

**BEYOND BAGS OF WORDS:
EFFECTIVELY MODELING DEPENDENCE AND
FEATURES IN INFORMATION RETRIEVAL**

A Dissertation Presented

by

DONALD A. METZLER JR.

Submitted to the Graduate School of the
University of Massachusetts Amherst in partial fulfillment
of the requirements for the degree of

DOCTOR OF PHILOSOPHY

September 2007

Computer Science

© Copyright by Donald A. Metzler Jr. 2007

All Rights Reserved

**BEYOND BAGS OF WORDS:
EFFECTIVELY MODELING DEPENDENCE AND
FEATURES IN INFORMATION RETRIEVAL**

A Dissertation Presented

by

DONALD A. METZLER JR.

Approved as to style and content by:

W. Bruce Croft, Chair

James Allan, Member

John Buonaccorsi, Member

Andrew McCallum, Member

Andrew Barto, Department Chair
Computer Science

To Shelley and my parents.

ACKNOWLEDGMENTS

This thesis would not have been possible without the immense support that I received from many people within the Computer Science Department at the University of Massachusetts Amherst (UMass). First and foremost, I would like to thank my advisor, W. Bruce Croft, for instilling in me all of the virtues necessary to conduct novel, meaningful research. Bruce taught me the importance of being cognizant of the past, asking the right questions, analyzing results critically, and striving to be the absolute best. I would also like to thank the other faculty of the Center for Intelligent Information Retrieval (CIIR), which includes James Allan, Andrew McCallum, and R. Manmatha, each of which I was fortunate enough to work with on some level or another throughout my studies. I must also thank my fellow CIIR graduate students and David Fisher for all of the guidance, help, and feedback they provided over the years. In particular, I would like to thank Trevor Strohman and Fernando Diaz for their nearly endless supply of thought-provoking discussions. I would also like to thank Kate Moruzzi, the CIIR secretary, for all of her help with the administrative aspects of my studies, and Andre Gauthier, the CIIR system administrator, for thanklessly supporting my resource intensive experiments.

I must also acknowledge all of the people that I was fortunate enough to collaborate with outside of UMass. I graciously thank Alistair Moffat, Justin Zobel, and Yaniv Bernstein for their hospitality and endless entertainment during my visit to Australia. Finally, I thank Susan Dumais and Chris Meek for mentoring me and giving me the opportunity to see what it was like to be a researcher in an industrial research lab during my summer at Microsoft Research.

ABSTRACT

BEYOND BAGS OF WORDS: EFFECTIVELY MODELING DEPENDENCE AND FEATURES IN INFORMATION RETRIEVAL

SEPTEMBER 2007

DONALD A. METZLER JR.

B.S., ROSE-HULMAN INSTITUTE OF TECHNOLOGY

M.S., UNIVERSITY OF MASSACHUSETTS AMHERST

Ph.D., UNIVERSITY OF MASSACHUSETTS AMHERST

Directed by: Professor W. Bruce Croft

Current state of the art information retrieval models treat documents and queries as bags of words. There have been many attempts to go beyond this simple representation. Unfortunately, few have shown consistent improvements in retrieval effectiveness across a wide range of tasks and data sets. Here, we propose a new statistical model for information retrieval based on Markov random fields. The proposed model goes beyond the bag of words assumption by allowing dependencies between terms to be incorporated into the model. This allows for a variety of textual and non-textual features to be easily combined under the umbrella of a single model. Within this framework, we explore the theoretical issues involved, parameter estimation, feature selection, and query expansion. We give experimental results from a number of information retrieval tasks, such as ad hoc retrieval and web search.

TABLE OF CONTENTS

	Page
ACKNOWLEDGMENTS	v
ABSTRACT	vi
LIST OF TABLES	xii
LIST OF FIGURES	xvii
 CHAPTER	
1. INTRODUCTION	1
1.1 From Archie to Google	1
1.2 The Academic and Industrial Perspectives	2
1.3 Paradigm Shifts	3
1.4 A Robust Retrieval Model	7
1.5 Outline	8
1.6 Contributions	10
2. RELATED WORK	12
2.1 Bag of Words Models	12
2.1.1 Binary Independence Retrieval Model	14
2.1.1.1 Estimation with Relevance Information	16
2.1.1.2 Estimation without Relevance Information	16
2.1.1.3 Tree Dependence Model	18
2.1.2 2-Poisson Model	19
2.1.3 BM25 Model	21
2.1.4 Unigram Language Models	22
2.1.5 Other Bag of Words Models	23
2.2 Models That Go Beyond the Bag of Words Assumption	24

2.2.1	<i>n</i> -Gram Language Models	25
2.2.2	Indri Inference Network Model	26
2.2.2.1	Document Representation	27
2.2.2.2	Language Models	28
2.2.2.3	Representation Nodes	29
2.2.2.4	Query Nodes	30
2.2.2.5	Explicit vs. Implicit Query Generation	31
2.2.3	Other Models That Go Beyond the Bag of Words Assumption	32
2.3	The Current State of the Art	33
3.	A MARKOV RANDOM FIELD MODEL FOR INFORMATION RETRIEVAL	34
3.1	Modeling Relevance	34
3.2	MRFs for IR	35
3.2.1	Graph Structure	37
3.2.2	Potential Functions	39
3.3	Building MRFs	42
3.3.1	Dependence Model Type	43
3.3.2	Clique Set Type	43
3.3.3	Weighting Function	47
3.3.3.1	Weighting Functions for T_{QD} , O_{QD} , and U_{QD}	47
3.3.3.2	Weighting Functions for T_Q , O_Q , and U_Q	50
3.3.3.3	Weighting Functions for D	50
3.3.4	Examples	50
3.4	Ranking	54
3.5	Discussion	55
4.	THE THEORY OF LINEAR FEATURE-BASED MODELS	57
4.1	Linear Feature-Based Models	58
4.2	Previous Uses of Linear Feature-Based Models in IR	60
4.3	Parameter Space	61
4.3.1	Reduction to Multinomial Manifold	62
4.3.2	Rank Equivalence	64
4.3.3	Distance Between Models	64

4.4	Parameter Estimation	66
4.4.1	Direct Search	67
4.4.1.1	Grid Search	68
4.4.1.2	Coordinate Ascent	68
4.4.1.3	Discussion	71
4.4.2	Optimization Using Surrogate Functions.....	72
4.4.2.1	Perceptron Learning	73
4.4.2.2	RankNet	73
4.4.2.3	Support Vector Machine Optimization	74
4.4.2.4	Discussion	75
4.5	Quantifying the Impact of Metric Divergence	75
5.	EVALUATION OF THE BASIC MRF MODEL.....	79
5.1	Ad Hoc Retrieval	80
5.1.1	MRF Models for Ad Hoc Retrieval	82
5.1.1.1	Full Independence	82
5.1.1.2	Sequential Dependence.....	84
5.1.1.3	Full Dependence	88
5.1.2	Evaluation	93
5.1.2.1	Smoothing	93
5.1.2.2	Collection Size	96
5.1.2.3	The Role of Features	98
5.1.2.4	Robustness	102
5.1.2.5	Long Queries	106
5.1.2.6	BM25 Weighting	109
5.1.2.7	Comparison to Bigram Model.....	111
5.1.2.8	Generalization	112
5.1.3	Summary of Results	117
5.2	Web Search	118
5.2.1	Previous Models for Web Search	120
5.2.2	Document Priors	123
5.2.2.1	Inlink Count	123
5.2.2.2	PageRank	124

5.2.3	MRF Models for Web Search	126
5.2.4	Results	129
6.	AUTOMATIC FEATURE SELECTION	135
6.1	Related Work	136
6.2	Automatic Feature Selection	137
6.2.1	Motivation	137
6.2.2	Algorithm	138
6.3	Evaluation	139
6.3.1	No Retraining vs. Retraining	141
6.3.2	Number of Features	142
6.3.3	Feature Analysis	142
6.3.4	Summary of Results	145
7.	LATENT CONCEPT EXPANSION	147
7.1	Related Work	148
7.2	Latent Concept Expansion	150
7.2.1	Query Expansion	154
7.2.2	Comparison to Relevance Models	155
7.3	Experimental Results	155
7.3.1	Ad Hoc Retrieval Results	156
7.3.1.1	Expansion with Single Term Concepts	156
7.3.1.2	Expansion with Multi-Term Concepts	158
7.3.2	Robustness	159
7.3.3	Multi-Term Concept Generation	161
7.4	Discussion	163
7.4.1	Relevance vs. Relevant Documents	163
7.4.2	The Role of Dependence	164
8.	CONCLUSIONS AND FUTURE WORK	166
8.1	Using the MRF Model	166
8.2	Summary of <i>Ad Hoc</i> Retrieval Results	167
8.3	Contributions	167
8.4	Future Work	171

APPENDICES

A. DATA SETS	173
A.1 Anatomy of a TREC Data Set	173
A.2 Summary of Data Sets	177
B. EVALUATION METRICS	178
 BIBLIOGRAPHY	 181

LIST OF TABLES

Table	Page
3.1 Example clique sets for the query $q1\ q2\ q3$ under full dependence model.	44
3.2 Summary of Dirichlet and BM25 weighting functions that can be used with cliques in the T_{QD} , O_{QD} , and U_{QD} clique sets. Here, M and N act as weighting function parameters that affect how matching is done, $tf_{e,D}$ is the number of times expression e matches in document D , $cf_{e,D}$ is the number of times expression e matches in the entire collection, df_e is the total number of documents that have at least one match for expression e , $ D $ is the length of document D , $ D _{avg}$ is the average document length, N is the number of documents in the collection, and $ C $ is the total length of the collection. Finally, μ^t , μ^w , k_1^t , k_1^w , b^t , and b^w are weighting function hyperparameters. The t and w superscripts indicate term and window hyperparameters, respectively.	49
3.3 Summary of ICF and IDF weighting functions that can be used with cliques in the T_Q , O_Q , and U_Q clique sets.	51
4.1 Features used in the bag of words experiments. $tf_{w,D}$ is the number of times term w occurs in document D , cf_w is the number of times term w occurs in the entire collection, df_w is the number of documents term w occurs in, $ D $ is the length (in terms) of document D , $ C $ is the length (in terms) of the collection, and N is the number of documents in the collection.	77
4.2 Training and test set mean average precision values for various <i>ad hoc</i> retrieval data sets and training methods. The † represents a statistically significant improvement over language modeling and ‡ denotes significant improvement over the balanced SVM model. Tests done using a one tailed paired t-test at the 95% confidence level.	77
5.1 Test set results for the MRF-FI model.	84

5.2	Mean average precision for various parameter settings for LM-U-N using the MRF-SD model.	86
5.3	Test set results for the MRF-SD model. A † indicates a statistically significant improvement over the MRF-FI model.	87
5.4	Mean average precision using the MRF-FD model over different combinations of term, ordered, and unordered features.	90
5.5	Test set results for the MRF-FD model. A † indicates a statistically significant improvement over the MRF-FI model and a ‡ indicates statistically significant improvement over the MRF-SD model.	91
5.6	Median values for various statistics computed across judged relevant (Rel) and non-relevant (Nonrel) documents.	99
5.7	The 10 most improved and 10 most hurt test set queries when using the MRF-SD model on the ROBUST04 data set. Effectiveness is measure in terms of average precision.	104
5.8	The 10 most improved and 10 most hurt test set queries when using the MRF-SD model on the GOV2 data set. Effectiveness is measure in terms of average precision.	105
5.9	Test set mean average precision for description-length queries using full and sequential dependence models. All improvements are statistically significant.	108
5.10	Test set results for the MRF-BM25 model. The †, ‡, and * indicate statistically significant improvements over the MRF-FI, BM25 and MRF-SD models, respectively. Recommended term and window hyperparameter values are also provided.	110
5.11	Test set results for the bigram language model. The † indicates a statistically significant improvement over the MRF-FI model and the ↓ indicates a statistically significant <i>decrease</i> in effectiveness compared to the MRF-SD model (i.e., MRF-SD > MRF-BM25). Recommended smoothing parameter values are also provided.	111
5.12	Intracollection generalization results for mean average precision. Values given are effectiveness ratios.	113
5.13	Inter-collection generalization results. Table includes mean average precision effectiveness ratios across all possible train/test splits using the F2EXP model.	114

5.14	Inter-collection generalization results. Table includes mean average precision effectiveness ratios across all possible train/test splits using the MRF-SD model.	115
5.15	Summary of test set mean average precision for the MRF-FI, MRF-SD, and MRF-FD models across all of the <i>ad hoc</i> retrieval data sets. Values in parenthesis denote percentage improvement over MRF-FI model. A † indicates a statistically significant improvement over the MRF-FI model, and a ‡ indicates a statistically significant improvement over the MRF-SD model. Recommended smoothing values are given for each collection, and recommended MRF model parameters are provided for each model.	117
5.16	URLs in the GOV2 collection with the largest raw PageRank scores. The number of inlinks for each URL is also shown.	126
5.17	Summary of named page weighting functions. The NP, NP-O-M, and NP-U-N weighting functions are based on mixtures of field language models, and INLINK and PAGERANK are based on document priors. The α_f values correspond to the mixing probabilities $P(f D)$. Term and window smoothing parameters are denoted by μ_f^t and μ_f^w , respectively.	128
5.18	Summary of named page finding results.	130
5.19	Results comparing the mean reciprocal rank of the LM-Mixture and MRF-NP models with and without document priors.	131
5.20	The 10 most improved and 10 most hurt queries on the TREC 2006 Terabyte Track named page finding data set. Effectiveness is measured in terms of reciprocal rank.	132
6.1	Training and test set mean average precision values for no retraining and retraining.	141
6.2	Comparison of test set mean average precision for language modeling (MRF-FI), BM25, MRF model using language modeling weighting (MRF-SD), MRF model using BM25 weighting (MRF-BM25), and MRF learned using our proposed feature selection algorithm (MRF-FS). A † indicates a statistically significant improvement over <i>both</i> the MRF-FI and BM25 models and a ‡ indicates a significant improvement over the MRF-BM25 model.	145

7.1	Test set mean average precision for MRF-FI, MRF-SD, relevance models (RM3), and latent concept expansion (LCE). The superscripts α , β , and γ indicate statistically significant improvements over MRF-FI, MRF-SD, and RM3, respectively.	156
7.2	Test set precision at 10 for MRF-FI, MRF-SD, relevance models (RM3), and latent concept expansion (LCE). The superscripts α , β , and γ indicate statistically significant improvements over MRF-FI, MRF-SD, and RM3, respectively.	157
7.3	Test set mean average precision values for multi-term concept LCE experiments.	158
7.4	One and two term expansion concepts for the query <i>price fixing</i> (ROBUST04 topic 622) and <i>tax evasion indicted</i> (ROBUST04 topic 650). Concepts are listed in descending order of $P(e Q)$ and $P(e_1, e_2 Q)$, respectively.	158
7.5	Fifteen most likely one, two, and three word concepts constructed using the top 25 documents retrieved for the query <i>hubble telescope achievements</i> on the ROBUST04 collection.	162
8.1	Test set mean average precision across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.	168
8.2	Test set geometric mean average precision for across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.	168
8.3	Test set precision at 10 across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.	169
8.4	Test set precision at R (R-prec) across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.	169
A.1	Overview of TREC collections and topics used in most of our experiments.	176

A.2	TREC data sets used in Chapter 4. The disk numbers refer to the TREC volumes used to construct the index.	177
B.1	Summary of common information retrieval evaluation metrics, where $R(1, k)$ is defined in Equation B.1 and $ R $ is the total number of judged relevant documents.	179
B.2	Overview of aggregate measures. For each aggregate measure we show how it is computed. Here, N refers to the total number of queries being aggregated and q_i is the i^{th} query in the set. Notice that GMAP is zero if any query has an average precision of zero. In order to correct this, $\text{AvgP}'(q_i)$ is defined to be $\max(\text{AvgP}(q_i), .00001)$	180

LIST OF FIGURES

Figure	Page
1.1 A summary of the three primary information retrieval paradigm shifts. They include the TF shift (aboutness), the IDF shift (informativeness), and the noise shift (noisiness).....	7
2.1 An example spanning tree for five terms, rooted at D.	19
2.2 Indri's inference network retrieval model	27
2.3 Example Indri document representation for the document A B C A B. The features correspond to the single terms A, B, C, and bigrams "A A", "A B", "A C", "B A", "B B", "B C", "C A", "C B", and "C C", respectively. The function f takes a document and set of features as input and outputs a document representation.	28
3.1 Examples of three ways to model the joint distribution $P(Q, D)$ using Markov random fields.....	36
3.2 Example Markov random field model for three query terms under various independence assumptions, including full independence (left), sequential dependence (middle), and full dependence (right).....	38
3.3 Illustration showing how the full independence model generalizes unigram language modeling and BM25 (top), and how the sequential dependence model generalizes bigram language modeling (bottom).	56
5.1 TREC topic number 744.	80
5.2 Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the AP data set.	94
5.3 Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the WSJ data set.....	95

5.4	Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the ROBUST04 data set.	95
5.5	Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the WT10G data set.	96
5.6	Relationship between the number of documents in a collection and the relative improvement in mean average precision of the MRF-FD model over unigram language modeling (MRF-FI). Note that the x -axis is log scaled.	97
5.7	Plot of average distance between sequential query terms for WT10G data set. Coll represents the entire collection, nonrel the set of judged non-relevant documents, and rel the set of judged relevant documents.	101
5.8	Robustness of MRF-SD and MRF-FD models for the AP, WSJ, ROBUST04, WT10G, and GOV2 test sets. The MRF-FI model is used as the baseline by which the improvements are computed. The evaluation metric used is average precision.	103
5.9	Mean average precision values plotted over MRF-SD parameter simplex for AP, WSJ, WT10g, and GOV2 collections.	116
5.10	Example TREC named page finding topics.	119
5.11	Inlink count prior.	124
5.12	PageRank prior.	127
5.13	Robustness of the MRF-NP models for the 2005 and 2006 Terabyte Track named page finding data sets. The LM-Mixture model is used as the baseline by which the improvements were computed. The evaluation metric used is reciprocal rank.	134
6.1	Mean average precision versus number of iterations for the training and test sets of the AP, WSJ, ROBUST04, WT10G and GOV2 data sets.	143
7.1	Graphical model representations of relevance modeling (top), latent concept expansion using single term concepts (middle), and latent concept expansion using two term concepts (bottom) for a three term query.	151

7.2	Histograms that demonstrate and compare the robustness of relevance models (RM3) and latent concept expansion (LCE) with respect to the MRF-FI model for the AP, WSJ, ROBUST04, and WT10G data sets.	160
A.1	Example TREC document.	174
A.2	Example TREC topic.....	175
A.3	Portion of a TREC relevance judgment file. The format of each line is: <code>query-id 0 doc-id judgment</code> . Judgments of 0, 1, and 2 refer to non-relevant, relevant, and highly relevant, respectively.	176

CHAPTER 1

INTRODUCTION

1.1 From Archie to Google

Information retrieval is a dynamic discipline with a rich history. However, for much of its history, it had little or no impact on people's everyday lives. Many of the earliest consumers of information retrieval technologies were government researchers, scientists, and librarians. That began to change after the invention of the World Wide Web (web) in the early 1990s.

Before the introduction of the web, a number of information sources were available online. Most of the online information was published and controlled by government organizations or academic institutions. It was uncommon for everyday citizens to use these online systems, let alone publish their own content. The web revolutionized the way that information was published. It allowed individuals and organizations to create content that was instantly available and easy to access. It also provided a way of linking content together, which was not possible with the older online systems. As computing costs decreased and online popularity increased, the amount of information available on the web exploded.

As more electronic documents started appearing online, a natural desire to search the content arose. Various search tools were developed to help users find relevant files and documents. The earliest Internet search tools, Archie, Gopher, Veronica, and Jughead allowed users to search FTP servers. However, the popularity of FTP waned after the introduction of the web. This ushered in a new era that gave rise to web search engines. Unlike their predecessors, which were used by small fractions of

the population, web search engines such as Google are used every day by millions of users across the globe. Therefore, what started as a small, relatively unknown field of study, has evolved into an integral part of modern society.

1.2 The Academic and Industrial Perspectives

Yahoo and Google were both grown out of academic research projects. They currently are the two most popular commercial web search engines in the United States. Clearly, the academic research community, in the early days of the web, was developing cutting edge search technologies. However, as the commercial search engines came of age, it became increasingly difficult for the academic researchers to keep up with the collection sizes and other critical research issues related to web search. This caused a divide to form between the information retrieval research being done within academia and industry.

There are several reasons for this divide. First, as commercial search engines mature, they are able to collect more data in the form of query logs, click-through patterns, and other types of user data which is invaluable to web search. The companies have little incentive to release this data to academia, especially amid growing privacy concerns. Second, commercial search engines have much more computing power than most academic research institutions. Therefore, they are able to crawl more web pages, build larger indexes, use real data streams, and experiment with much more costly computations. Finally, commercial search engines are very protective of their search algorithms and techniques and do not typically publish their findings in scholarly conferences and journals. This is not surprising, since revealing technical details of ranking functions may allow spammers and other malicious entities to adversely influence search results.

To put things into perspective, let us compare academic and industrial collection sizes. The Text REtrieval Conference (TREC), which was started in 1992, provides

a set of standard, reusable test collections (i.e., document collection, queries, and relevance judgments) that most academic information retrieval researchers use when evaluating retrieval systems. The largest TREC test collection, called GOV2, is a 2004 crawl of the .gov top level domain. It consists of approximately 25 million web pages (428GB of text). In comparison, it is believed that Google has upwards of 25 billion items in its index, which is 1,000 times larger than GOV2. In addition, many of the most widely used academic information retrieval models were initially developed for test collections that consist of fewer than 1 million documents.

One of the goals of this work is to reduce the divide in understanding that exists between academic and commercial information retrieval systems with respect to large data sets. Many of the techniques and ideas developed here have been inspired by large test collections, such as GOV2. While the GOV2 collection is admittedly not web-scale, it is a significant and sizeable improvement over the test collections that have been used to develop most of the current state of the art information retrieval models.

1.3 Paradigm Shifts

As we just alluded to, large collections, such as those handled by commercial search engines, provide a new set of challenges for information retrieval researchers. In this work, we develop highly effective information retrieval models for both smaller, classical data sets, and larger web collections. As we will show throughout this work, the current state of the art academic retrieval models are not robust enough to achieve consistently effective retrieval results on large collections.

Most of these models are based on the so-called “bag of words” assumption. Under this assumption, text (e.g., queries and documents) are represented as unordered sets of terms. This means that any notion of term ordering is lost. For example, under this representation, the texts *the bear ate the human* and *the hu-*

man ate the bear are identical. However, these pieces of text clearly have different meanings. While this is an overly simplistic representation, very few have been able to develop non-bag of words retrieval models that are consistently and significantly better than the state of the art bag of words models. Many researchers over the past few decades have tried in vain, but there has been very little success.

The ranking functions associated with bag of words retrieval models often consist of some combination of *term frequency* (TF) and *inverse document frequency* (IDF). The IDF component acts to discriminate between informative and non-informative query terms. Those terms that have a high IDF are considered more informative, because they rarely occur in the collection. On the other hand, terms that have a low IDF are considered uninformative, since they occur in many documents. As the number of documents in a collection increases, IDF becomes increasingly important in order to discriminate between those documents that only contain non-informative query terms and those that contain highly informative query terms.

On the other hand, the TF component, which is often normalized in some way with respect to the *document length*, is used to discriminate between documents that contain a query term several times and those that contain the term many times. This makes the assumption that documents that contain more mentions of a given query term are more “about” the given term and therefore are more likely to be relevant to the query. As we will discuss shortly, this is a bad assumption, especially as collection sizes increase and documents become noisier. The TF component becomes more important as documents get longer, since query terms are unlikely to occur more than one time in a very short document, and since long documents are more likely to contain more diverse term occurrence statistics.

Therefore, the TF and IDF components used within bag of words ranking functions, when combined together, discriminate along two dimensions – informativeness

(IDF) and aboutness (TF). However, when dealing with large web collections, a third dimension that we call noisiness enters the picture.

All collections, even small ones that consist entirely of news articles, contain some noise. However, large web collections are likely to contain abundant amounts of noise. The standard TF and IDF features are not enough to overcome this noise. In fact, these features may actually help amplify the noise in some cases. Let us consider the query *habitat for humanity* run against the GOV2 collection. Using a state of the art bag of words retrieval model, many of the top ranked results are relevant to the request. However, there are several results very high in the ranked list that do not contain a single occurrence of the term *humanity*. Instead, these documents contain hundreds of occurrences of the high IDF term *habitat*. These documents are ranked so highly because they contain many occurrences of a very high IDF term.

Documents that contain hundreds of occurrences of some high IDF term are going to result in poor, noisy matches for most bag of words models based on TF and IDF. Such documents may arise by coincidence, or a spammer who wishes to increase the ranking of a given web page may “stuff” the page with such terms. In either case, it is very undesirable for these documents to be ranked highly.

Another more subtle way that noise may be introduced into bag of words matches happens when two or more query terms match a document, but the matches are random or unrelated to the query. For example, in the *habitat for humanity* case, consider a document that contains a paragraph that discusses habitat changes caused by global warming and another paragraph that discusses the negative impacts of global warming on humanity. Both the terms *habitat* and *humanity* will match this document, but the matches are unrelated to the query. That is, the terms just happened to match by chance. This is another example of noisy matches that can arise in large collections. In fact, as collection size grows, so does the chance that any two query terms will randomly match within some document.

Hence, new ranking function components, above and beyond TF and IDF must be used in order to reduce the number of noisy matches. There are a few ways to address this issue. First, one of the simplest ideas is to cap the TF component and not allow it to grow unbounded. While this addresses some noise issues, it fails to address the problem of randomly matching query terms. Second, in order to address the so-called term stuffing problem, anti-spam techniques may be developed in order to automatically detect malicious or misleading content. However, like capping TF, this only addresses some of the noise issues. Finally, term proximity features may be used in order to ensure that matches are not random and that they are somehow related to the query. For example, these types of features could be used to promote documents that contain the exact phrase *habitat for humanity* as opposed to those that simply contain random occurrences of the terms *habitat* and *humanity* on their own. It is this third option that we heavily explore within this work in order to overcome the limitations imposed by TF and IDF alone. It is important to notice that by using term position information, we are abandoning the bag of words assumption and move to a richer, more realistic text representation.

Aboutness, informativeness, and noisiness reflect the three primary information retrieval paradigm shifts. Here, a paradigm shift is a new way of approaching a problem with a given set of characteristics. The paradigm shifts are summarized in Figure 1.1. The figure plots three data sets (CACM, TREC Disks 1 and 2, and GOV2) with respect to their average document length and the number of documents in the collection. As the figure shows, the TF paradigm shift moves along the average document length axis and the IDF shift moves along the number of documents axis. We also see that the noise shift moves along both axes, but is only present for large collections, such as GOV2.

We hypothesize that many of the previous attempts to go beyond the bag of words assumption have failed because of the small data sets used. In fact, most, if not all,

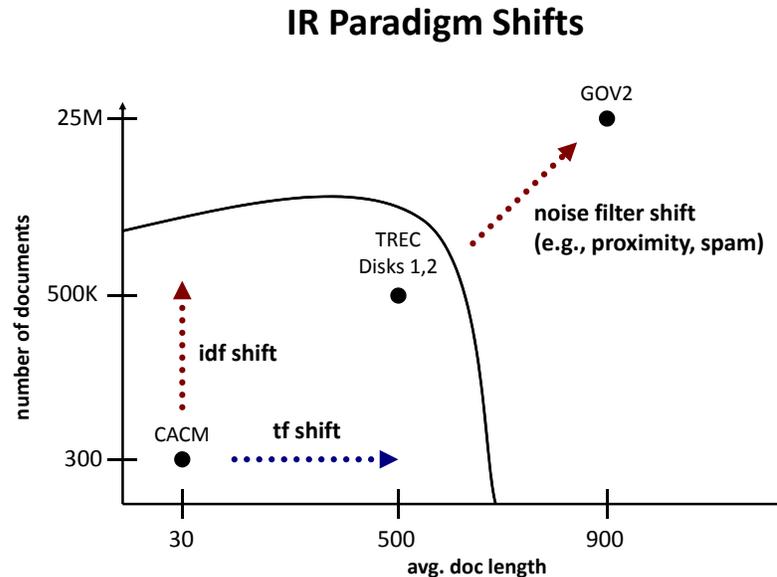


Figure 1.1. A summary of the three primary information retrieval paradigm shifts. They include the TF shift (aboutness), the IDF shift (informativeness), and the noise shift (noisiness).

of the previous research on non-bag of words model have been evaluated on test collections within the region shown in Figure 1.1. Poor, or inconclusive results, were achieved because the data sets did not exhibit the characteristics necessary to exploit the noise reducing features associated with non-bag of words models. Therefore, new models that go beyond the bag of words assumption should be tested on large, noisy data sets in order to properly evaluate their full potential.

1.4 A Robust Retrieval Model

In this work, we develop a robust statistical information retrieval model based on Markov random fields. In particular, the model is designed to support the following desiderata:

1. Support basic information retrieval tasks (e.g., ranking, query expansion, etc.).
2. Easily and intuitively model query term dependencies.

3. Handle arbitrary textual and non-textual features.
4. Consistently and significantly improve effectiveness over bag of words models across a wide range of tasks and data sets.

The model we develop goes beyond the bag of words assumption in two ways. First, the model can easily exploit various types of dependencies that exist between query terms. This eliminates the term independence assumption that often accompanies bag of words models. Second, arbitrary textual or non-textual features can be used within the model. Thus, it is possible to use simple features such as TF and IDF, or more complex features, such as those based on term proximity. Other possible features include PageRank, inlink count, readability, spam probability, among others. None of the current state of the art models allow arbitrary features to be incorporated as easily as our model.

As we will show, combining term dependencies and arbitrary features results in a very robust, powerful retrieval model. Within the model, we propose several extensions, such as an automatic feature selection algorithm and a query expansion framework. The resulting model and extensions provide a flexible framework for highly effective retrieval across a wide range of tasks and data sets. Our algorithm consistently and significantly outperforms the current state of the art models, especially when applied to very large web data sets.

1.5 Outline

The remainder of this work is laid out as follows:

- **Chapter 2 - Related Models** – We survey related information retrieval models. We include both bag of words models, such as language modeling and BM25, as well as models that have attempted to go beyond the bag of words assumption, like n -gram language models, the Indri retrieval model, and others.

- **Chapter 3 - A Markov Random Field Model for Information Retrieval**
 - We present the basics of our Markov random field (MRF) model for information retrieval, including the theoretical underpinnings of the model and how to use the model for ranking. In addition, we describe a novel way of compactly representing textual features that we use throughout the remainder of the work.
- **Chapter 4 - The Theory of Linear Feature-Based Models** – We present the theory of linear feature-based retrieval models, which is a broad class of retrieval models. We show that our MRF model, when used for ranking, is an instantiation of a linear-feature based model. Our theoretical analysis sheds novel insights into the parameter space of such models. We describe a novel parameter estimation technique for these methods and do a detailed literature survey of other parameter estimation techniques that have been recently proposed.
- **Chapter 5 - Evaluation of the Basic MRF Model** – We experimentally evaluate the MRF model on *ad hoc* and web search tasks. Our results show consistent and significant improvements in effectiveness over baseline bag of words models on both tasks. We also examine the generalization properties of the model. Our results show that estimated parameters generalize very well both within and across collections.
- **Chapter 6 - Automatic Feature Selection** – We present a fully supervised, automatic feature selection algorithm that can be used with the MRF model or any feature-based model. The algorithm adds great flexibility to the MRF model by removing the need to manually select features. We experimentally evaluate various aspects of the algorithm and comparing it against models with manually selected features. Our results show that the models automatically

learned using the algorithm are either indistinguishable from or significantly better than models with manually chosen features.

- **Chapter 7 - Latent Concept Expansion** – We develop a powerful query expansion framework for the MRF model. Our expansion technique is called latent concept expansion. It allows arbitrary features and term dependencies to be modeled both in the original query and the expansion concepts. This is the first query expansion model ever developed that provides this type of flexibility. Our results show that latent concept expansion significantly outperforms relevance modeling, which is currently the state of the art query expansion technique.
- **Chapter 8 - Conclusions and Future Work** – We summarize our contributions and provide a comprehensive table summarizing our primary results. In addition, we discuss potential areas of future work.

1.6 Contributions

The following is a summary of our primary contributions:

1. **Robust retrieval model.** We develop a new, formally motivated, statistical retrieval model based on Markov random fields that robustly and effectively handles term dependencies and the combination of arbitrary features.
2. **Better understanding of features for information retrieval.** By modeling dependencies between terms and encoding rich features, such as those based on phrases and term proximity, we are able to better understand how and when such features can improve retrieval effectiveness.
3. **Novel parameter estimation technique.** Our technique exploits the nature of rank-equivalence and works to directly maximize the underlying retrieval

metric, which leads to better performance than maximizing the data likelihood or margin. This avoids the problem of metric divergence.

4. **Automatic model learning.** We propose a supervised feature selection algorithm that can be used to automatically learn highly effective models. This eliminates the need for human experts to manually select model features on a per-task basis.
5. **Concept-based query expansion.** Our retrieval model provides an elegant mechanism for expanding queries using multi-term concepts in the context of relevance or pseudo-relevance feedback.
6. **State of the art retrieval effectiveness.** Our model shows consistent and significant improvements in retrieval effectiveness over current state of the art retrieval models on *ad hoc* retrieval and web search tasks.

CHAPTER 2

RELATED WORK

In this chapter we survey information retrieval models that are relevant to our work. There is no standard or formal way of describing retrieval models. However, most models can be uniquely described in terms of their *document representation*, *query representation*, and a *ranking function*. In the most common scenario, the ranking function takes a document and query representation as input, and outputs a score that is indicative of how relevant the document is to the query.

Throughout our discussion we will refer to documents, queries, and relevance. However, it should be noted that these terms are actually intended to be used quite generally. For example, in an automatic question answering system, a ‘document’ is an answer, a ‘query’ is a question, and ‘relevance’ is defined according to whether or not the answer, with regard to the question, is correct [112]. Other examples include image retrieval, where the ‘documents’ are images, and a query-query similarity task, where the ‘documents’ are queries [67, 69].

2.1 Bag of Words Models

We begin by looking at so-called bag of words retrieval models. As we will see, these models make use of many different types of document representations, query representations, and ranking functions. However, the one thing that all of these models have in common is the fact that term order is ignored when constructing the document and query representations. That is, given a document A and A^π , a permutation of A , the representations of A and A^π are identical.

This assumption is obviously overly simple. It conflates many texts that have very different semantic meanings into a single form. For example, the texts *the human ate the bear* and *the bear ate the human* have very different semantic meanings, but are represented the same way under the bag of words assumption. While it is not known how the human brain represents text, it is unlikely that term ordering is completely ignored. Indeed, term orderings play an important role in semantics. Therefore, it is difficult to expect a computer, using such a simple representation, to determine relevance as accurately as a human assessor can.

Despite the fact that the bag of words assumption is so simple, some believe that it is a reasonable first approximation [57, 97]. For example, in 1988, Salton and Buckley stated:

In reviewing the extensive literature accumulated during the past 25 years in the area of retrieval system evaluation, the overwhelming evidence is that the judicious use of single term identifiers is preferable to the incorporation of more complex entities...

Although this was written in 1988, there has been little, if any, conclusive evidence to discredit this claim.

Throughout this work, we argue, and aim to validate experimentally, that models based on the bag of words assumption are inferior to models that consider richer representations. We acknowledge that single term identifiers are likely to be the most powerful and discriminative types of identifiers. However, we propose a number of other identifiers (features) that can be used to significantly improve retrieval effectiveness well beyond the state of the art bag of words models discussed in this chapter.

Finally, one important thing to note is that the bag of words assumption does not force the underlying model to treat term occurrences as independent events. For example, it is possible for bag of words models to incorporate term co-occurrence

statistics into the ranking function, thereby modeling one very basic form of term dependence [7, 110, 116]. We return to the subject of term dependence models shortly.

2.1.1 Binary Independence Retrieval Model

We begin our discussion of bag of words models by describing the Binary Independence Retrieval (BIR) model, which is one of the earliest probabilistic models for information retrieval [92]. The model is also known by many other names including the Okapi model, the City model, the Robertson-Sparck Jones model, and the classical probabilistic model.

Under the BIR model, documents are represented as binary vectors d indexed over a fixed vocabulary \mathcal{V} (i.e., $d \in \{0, 1\}^{|\mathcal{V}|}$). The vocabulary typically consists of single terms, but may also include more general concepts, such as phrases. If a term, $w \in \mathcal{V}$, occurs one or more times in a document, then $d_w = 1$, otherwise $d_w = 0$. Documents are then ranked according to the *Probability Ranking Principle* (PRP) [87].

Probability Ranking Principle. If a reference retrieval system's response to each request is a ranking of the documents in the collection in order of decreasing probability of relevance to the user who submitted the request, where the probabilities are estimated as accurately as possible on the basis of whatever data have been made available to the system for this purpose, the overall effectiveness of the system to its user will be the best that is obtainable on the basis of those data.

The PRP states that documents should be ranked in decreasing order of probability of relevance given all of the available evidence. Thus, documents are ranked according to $P(R = 1|d)$, where R is a binary random variable that models relevance. Also, in order to make parameter estimation feasible, the model assumes that the term occurrence variables (d_w) are conditionally independent, given a relevance class. Using these assumptions, we now derive the BIR model's ranking function:

$$\begin{aligned}
P(R = 1|d) &\stackrel{\text{rank}}{=} \frac{P(R = 1|d)}{P(R = 0|d)} \\
&= \frac{P(d|R = 1)P(R = 1)}{P(d|R = 0)P(R = 0)} \\
&\stackrel{\text{rank}}{=} \frac{P(d|R = 1)}{P(d|R = 0)} \\
&= \prod_{w \in \mathcal{V}} \frac{P(d_w = 1|R = 1)^{\delta_w} P(d_w = 0|R = 1)^{1-\delta_w}}{P(d_w = 1|R = 0)^{\delta_w} P(d_w = 0|R = 0)^{1-\delta_w}} \\
&= \prod_{w : \delta_w=1} \frac{P(d_w = 1|R = 1)P(d_w = 0|R = 0)}{P(d_w = 0|R = 1)P(d_w = 1|R = 0)} \prod_{w \in \mathcal{V}} \frac{P(d_w = 0|R = 1)}{P(d_w = 0|R = 0)} \\
&\stackrel{\text{rank}}{=} \sum_{w : \delta_w=1} \log \frac{P(d_w = 1|R = 1)P(d_w = 0|R = 0)}{P(d_w = 0|R = 1)P(d_w = 1|R = 0)} \tag{2.1}
\end{aligned}$$

where $\stackrel{\text{rank}}{=}$ denotes rank equivalence¹ and δ_w is 1 if term w occurs in the document and 0 otherwise.

Although the model, at first glance, appears to make the relatively strong assumption of conditional independence, the model is actually much less restrictive [24]. Instead, as Cooper shows, the model allows dependencies to exist between terms, but requires the strength of the dependence to be equal in both the relevant ($R = 1$) and the non-relevant ($R = 0$) sets of documents, thereby linking the dependence. Therefore, instead of assuming conditional independence, the model actually assumes *linked dependence*. Even though this is a weaker assumption, it is still restrictive and unlikely to hold, in general.

One criticism of the model is that there is no explicit query involved. Rather, the model assumes that $P(R = 1|d)$ is implicitly conditioned on some underlying information need. This simplifies the model to a certain extent, but raises certain theoretical issues when researchers try to incorporate aspects of the query into the model.

¹Two functions are defined to be rank equivalent if they guaranteed to produce the same ordering of items when sorted according to function value. Monotonic transforms, scaling (by a constant), and translation (by a constant) are all examples of rank-preserving operations.

In order to use the model for ranking, $P(d_w = 1|R = 0)$ and $P(d_w = 1|R = 1)$ must be estimated for all w . This is a total of $2|\mathcal{V}|$ parameters, which for large vocabularies, is an immense number of parameters. We now describe how these parameters are typically estimated when we have relevance information and when we do not.

2.1.1.1 Estimation with Relevance Information

If we are given relevance information for a given information need then estimation is straightforward. Assume we are given a set of documents that are known to be relevant to the information need, as well as a set that are known to be non-relevant. In the simplest case, we can get this information from a user who has manually judged a set of documents. Using this information, we obtain the following estimates:

$$P(d_w = 1|R = 0) = \frac{nr_w + \alpha_{nr}}{NR + \alpha_{nr} + \beta_{nr}} \quad (2.2)$$

$$P(d_w = 1|R = 1) = \frac{r_w + \alpha_r}{R + \alpha_r + \beta_r} \quad (2.3)$$

where nr_w is the number of judged non-relevant documents that term w occurs in, NR is the total number of judged non-relevant documents, r_w is the number of judged relevant documents term w occurs in, R is the total number of documents judged relevant, and α and β are smoothing parameters that avoid zero probabilities and help overcome data sparsity. Historically, the smoothing parameters have been set to $\alpha_{nr} = \alpha_r = 0.5$ and $\beta_{nr} = \beta_r = 0$, although there is no reason to believe these are the optimal settings.

2.1.1.2 Estimation without Relevance Information

As we just showed, estimation with relevance information is rather straightforward. However, for most queries, a system will not have access to relevance information. Not surprisingly, parameter estimation under this scenario proves to be much

more challenging. In order to use the model when there is no relevance information, a number of assumptions must be made [28, 94]. These assumptions include:

1. $P(d_w = 1|R = 0) = P(d_w = 1|R = 1)$ for all terms that do not occur in the query.
2. $P(d_w = 1|R = 1) = 0.5$ for all terms that occur in the query.
3. $P(d_w = 1|R = 0) = \frac{df_w + \alpha_{nr}}{N + \alpha_{nr} + \beta_{nr}}$ for all terms that occur in the query, where df_w is the number of documents that w occurs in and N is the number of documents in the collection.

The first two assumptions are poor and very unlikely to hold true. However, they greatly simplify how the ranking function is computed, which is necessary, given the lack of relevance information.

The third assumption is the most reasonable. It assumes that we have observed no relevant documents and that every document in the collection is non-relevant. This information is treated as real relevance information and plugged into Equation 2.2 to derive the estimate of $P(d_w = 1|R = 0)$ shown above. Of course, this assumption is not completely accurate. However, it has been shown that the collection, as a whole, acts as a good proxy for modeling the set of non-relevant documents.

After invoking these assumptions, the BIR model ranking function simplifies to the following form:

$$P(R = 1|d) \stackrel{rank}{=} \sum_{w : \delta_w = 1 \wedge w \in Q} \log \frac{N - df_w + 0.5}{df_w + 0.5} \quad (2.4)$$

where $w \in Q$ indicates that term w occurs in the query and we set $\alpha_{nr} = \beta_{nr} = 0.5$, which is commonly used in this scenario. The term inside of the summation is the well-known Okapi IDF. This derivation provides one of the many theoretical justifications for the importance of IDF [88].

2.1.1.3 Tree Dependence Model

The tree dependence model is one of many extensions that have been proposed for the BIR model [110]. It attempts to model first order dependencies between terms by making use of the Chow expansion [19]. The Chow expansion is a method for approximating a joint distribution in terms of first order dependencies.

The tree dependence model constructs a weighted, undirected graph G for both the relevant and non-relevant sets, such that the vertices of G are the terms in the vocabulary, there exists an edge between every pair of terms, and the edge weights are the expected mutual information between the two terms, which is computed as:

$$I(t_1, t_2 | R) = \sum_{\delta_1 \in \{0,1\}} \sum_{\delta_2 \in \{0,1\}} P(d_{t_1} = \delta_1, d_{t_2} = \delta_2 | R) \log \frac{P(d_{t_1} = \delta_1, d_{t_2} = \delta_2 | R)}{P(d_{t_1} = \delta_1 | R)P(d_{t_2} = \delta_2 | R)} \quad (2.5)$$

where t_1 and t_2 are terms, and $P(d_t | R)$ and $P(d_{t_1}, d_{t_2} | R)$ are typically estimated using relevance information.

This graph is then used in the following way. First, a maximum spanning tree is constructed from the graph. Next, an arbitrary node is chosen as the root, which allows all of the edges to be directionalized. Finally, the directed graph is used to compute a first order approximation to the joint distribution over all of the terms by assuming that a term is dependent on its parent term².

Figure 2.1 show an example maximum spanning tree over a graph with five terms. As we see, the root node is term D . The first-order approximation of the joint distribution can then be written as:

$$P(A, B, C, D, E) = P(D)P(A|D)P(C|D)P(B|C)P(E|C) \quad (2.6)$$

²The root node, which has no parents, is assumed to not be dependent on any other terms.

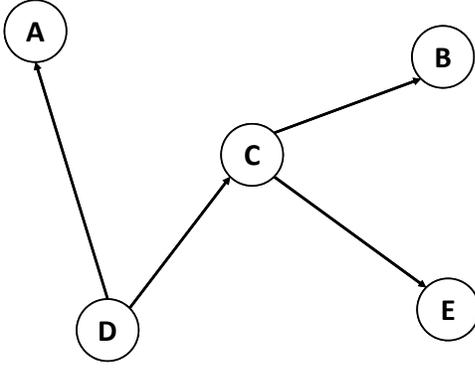


Figure 2.1. An example spanning tree for five terms, rooted at D.

In the context of the BIR model, this approximation can be used to compute $P(d|R = 0)$ and $P(d|R = 1)$, which results in the following ranking function:

$$\begin{aligned}
 P(R = 1|d) &\stackrel{rank}{=} \frac{P(d|R = 1)}{P(d|R = 0)} \\
 &= \prod_{w \in \mathcal{V}} \frac{P(d_w = 1|d_{\pi_w} = \delta_{\pi_w}, R = 1)^{\delta_w} P(d_w = 0|d_{\pi_w} = \delta_{\pi_w}, R = 1)^{1-\delta_w}}{P(d_w = 1|d_{\pi_w} = \delta_{\pi_w}, R = 0)^{\delta_w} P(d_w = 0|d_{\pi_w} = \delta_{\pi_w}, R = 0)^{1-\delta_w}} \\
 &\stackrel{rank}{=} \sum_{w : \delta_w=1} \log \frac{P(d_w = 1|d_{\pi_w} = \delta_{\pi_w}, R = 1)P(d_w = 0|d_{\pi_w} = \delta_{\pi_w}, R = 0)}{P(d_w = 0|d_{\pi_w} = \delta_{\pi_w}, R = 1)P(d_w = 1|d_{\pi_w} = \delta_{\pi_w}, R = 0)}
 \end{aligned} \tag{2.7}$$

where π_w is the parent term of w in the spanning tree.

Parameter estimation in the tree dependence model is even more difficult than in the BIR model because the data sparseness problem is only exacerbated by the large number of conditionals that must be estimated. Furthermore, the model has never shown consistent improvements in effectiveness over the BIR model [110]. Therefore, the model is interesting mostly from a theoretical and historical perspective.

2.1.2 2-Poisson Model

As we just explained, the BIR model represents documents as binary vectors. This representation can only capture whether or not a term occurs in a document. It

ignores the number of times the term occurs (term frequency), which, as we discussed in Chapter 1, is important, especially as document lengths increase.

The 2-Poisson model was proposed to overcome this limitation [48, 93]. Under this model, documents are represented as vectors of term frequencies. The vector is indexed over a fixed vocabulary \mathcal{V} , and thus $d \in \mathbb{N}^{|\mathcal{V}|}$ for every document d .

The 2-Poisson ranking function is derived analogously to the BIR ranking function, as follows:

$$\begin{aligned}
P(R = 1|d) &\stackrel{rank}{=} \frac{P(d|R = 1)}{P(d|R = 0)} \\
&= \prod_{w \in \mathcal{V}} \frac{P(d_w = tf_w|R = 1)}{P(d_w = tf_w|R = 0)} \\
&= \prod_{w : tf_w > 0} \frac{P(d_w = tf_w|R = 1)P(d_w = 0|R = 0)}{P(d_w = 0|R = 1)P(d_w = tf_w|R = 0)} \prod_{w \in \mathcal{V}} \frac{P(d_w = 0|R = 1)}{P(d_w = 0|R = 0)} \\
&\stackrel{rank}{=} \sum_{w : tf_w > 0} \log \frac{P(d_w = tf_w|R = 1)P(d_w = 0|R = 0)}{P(d_w = 0|R = 1)P(d_w = tf_w|R = 0)} \tag{2.8}
\end{aligned}$$

where term frequencies, as before, as assumed to be conditionally dependent³ and tf_w denotes the number of times that term w occurs in the document.

To use the model, $P(d_w = tf_w|R = 0)$ and $P(d_w = tf_w|R = 1)$ must be estimated. As its name implies, the 2-Poisson model assumes that term frequencies are distributed according to a mixture of two Poissons, as given by:

$$P(d_w = tf_w|R = 1) = P(E = 1|R = 1) \frac{e^{-\alpha_w} \alpha_w^{tf_w}}{tf_w!} + P(E = 0|R = 1) \frac{e^{-\beta_w} \beta_w^{tf_w}}{tf_w!} \tag{2.9}$$

$$P(d_w = tf_w|R = 0) = P(E = 1|R = 0) \frac{e^{-\alpha_w} \alpha_w^{tf_w}}{tf_w!} + P(E = 0|R = 0) \frac{e^{-\beta_w} \beta_w^{tf_w}}{tf_w!} \tag{2.10}$$

where α_w and β_w are parameters of the Poisson distributions, and E is a binary variable that represents *eliteness* [48, 89].

³As in the BIR model, the weaker linked dependence assumption holds.

Eliteness is a hidden, or latent, variable that reflects whether or not a given term is actually “about” a document or not. Therefore, elite terms are “about” a document, and non-elite terms are not. Given the subjective nature of its definition, it is very difficult to quantify or actually model eliteness in practice. For this reason, the model is interesting purely from a theoretical point of view.

2.1.3 BM25 Model

The BM25 model, proposed by Robertson and Walker, is an empirical, hand-crafted approximation of the 2-Poisson model [89]. It was first introduced at TREC in 1995, and has been widely used ever since [90]. The ranking function is given by:

$$P(R = 1|d) \approx \sum_{w \in Q \cap d} tf_{w,Q} \frac{(k_1 + 1)tf_{w,d}}{k_1 \left((1 - b) + b \frac{|d|}{|d|_{avg}} \right) + tf_{w,d}} \log \frac{N - df_w + 0.5}{df_w + 0.5} \quad (2.11)$$

where $tf_{w,Q}$ is the number of times term w occurs in the query, $tf_{w,d}$ is the number of times term w occurs in the document, $|d|$ is the number of terms in the document, $|d|_{avg}$ is the average document length, N is the number of documents in the collection, df_w is the number of documents term w occurs in the collection, and k_1 and b are model parameters.

The model is very simple to implement and, with carefully chosen model parameters, has been shown to consistently achieve state of the art effectiveness. Unfortunately, the model, despite being inspired by the 2-Poisson model, is heuristic and has no built-in mechanism for modeling term dependencies. It has been shown that it is possible to incorporate query independent features, such as PageRank into the model [26], as well as term proximity information [18], However, these improvements are heuristic and are unrelated to the assumed underlying 2-Poisson model. Furthermore, there is no convenient, formally motivated framework for easily adding other types of features to the model, which limits the usefulness of the model.

2.1.4 Unigram Language Models

Language modeling, a statistical technique first applied to speech recognition [96], has also been successfully applied to information retrieval [84]. Since its introduction, it has grown in popularity and has been proven to be a robust, highly effective retrieval model. In the context of information retrieval, language models are statistical models of text generation. In this section we will describe the simplest type of language model, the unigram model, which is based on the bag of words assumption. Later, we describe more complex language models.

The most common strategy for using language models for information retrieval is called the *query likelihood* approach. Given a query Q , documents are ranked according to the likelihood that the query was generated, given document D as evidence. Typically, for the sake of smoothing the document model, a Bayesian estimate is used. Using the approach, documents are ranked according to:

$$\begin{aligned}
 P(Q|D) &= \prod_{q \in Q} P(q|D) \\
 &= \prod_{q \in Q} \int_{\theta_D} P(q|\theta_D)P(\theta_D|D) \\
 &\propto \prod_{q \in Q} \int_{\theta_D} P(q|\theta_D)P(D|\theta_D)P(\theta_D)
 \end{aligned} \tag{2.12}$$

where the unigram language model (θ_D) is typically a multinomial distribution over a fixed vocabulary [101]. In addition, for computational simplicity, it is assumed that $P(\theta_D)$ is Dirichlet. This is typically called Bayesian or Dirichlet smoothing [123]. Under these assumptions, we derive the following estimate for $P(w|D)$:

$$P(w|D) = \frac{tf_{w,D} + \mu P(w|C)}{|D| + \mu} \tag{2.13}$$

where it is assumed that the Dirichlet parameters are $\alpha_w = \mu P(w|C)$, where μ is a model hyperparameter and $P(w|C) = \frac{cf_w}{|C|}$, with cf_w being the number of times term

w occurs in the collection and $|C|$ is the total number of terms in the collection. Using this estimate, documents are then ranked according to:

$$\begin{aligned}
P(Q|D) &\stackrel{rank}{=} \sum_{q \in Q} \log \frac{tf_{w,D} + \mu P(w|C)}{|D| + \mu} \\
&\stackrel{rank}{=} \sum_{q \in Q \cap D} \log \frac{tf_{w,D} + \mu P(w|C)}{\frac{\mu P(w|C)}{|D| + \mu}} + \sum_{q \in Q} \log \frac{\mu P(w|C)}{|D| + \mu} \\
&\stackrel{rank}{=} \sum_{q \in Q \cap D} \log \left[1 + \frac{tf_{w,D}}{\mu} \cdot \frac{|C|}{cf_w} \right] - |Q| \log(|D| + \mu) \quad (2.14)
\end{aligned}$$

which can be interpreted as another variant on the standard *tf.idf* formula with built-in document length normalization.

Although we described Bayesian smoothing here, it should be noted that many other types of smoothing are possible [123, 124, 125]. Also, distributions other than the multinomial have been proposed for modeling documents and queries, such as the multiple-Bernoulli distribution [68].

Even though language modeling is more formally motivated than BM25, the ranking functions are quite similar, with both relying on the standard *tf*, *idf*, and document length normalization components. Many of the problems with BM25 are also carried over to this model. For example, language models are models of *text* generation and therefore it is difficult to incorporate non-textual features into the model. Arbitrary query independent features are often encoded using a document prior [54], but this is not applicable to query dependent features. One possible solution is to use a general Naïve Bayes model, but such a model, by its very nature, is incapable of modeling term dependencies. As we will soon show, there have been a number of models proposed to address this problem using more complex language models.

2.1.5 Other Bag of Words Models

The axiomatic approach to retrieval [37] and the divergence from randomness model [2] are two recently proposed bag of words retrieval models. The ranking

functions of these two models are variants on the *tf.idf* theme. While the models shed interesting insights into retrieval modeling, they have not been shown to be significantly better or more flexible than language modeling or BM25.

The tree dependence model was one of the first bag of words models that attempted to capture the dependence that exists between terms. Several other bag of words models have been proposed to capture dependencies, as well. Early examples include the Bahadur Lazarsfeld expansion (BLE) [86] which is an exact, but computationally intensive, method of modeling high order dependencies, and the Generalized Dependence Model, which generalizes both the tree dependence model and the BLE expansion [119].

Other examples include latent semantic analysis [31, 49], term association models [7, 38, 102, 116], cluster-based language models [33, 55, 61], topic models [8, 9, 44, 115], and pseudo-relevance feedback [34, 58, 122].

Despite the increased complexity, many of these have failed to yield substantial improvements in effectiveness. Several of the models, including cluster-based language modeling and some of the topic models, have shown significant improvements in retrieval effectiveness. However, these models are often computationally intensive, making them impractical to apply to web-scale collections. The model that we propose in this work does not require any expensive computations and can easily be applied to very large collections. In addition, it allows other types of term dependence features, beyond simple co-occurrence statistics, to be used, thus making it more robust and practical for a wide range of tasks and data sets.

2.2 Models That Go Beyond the Bag of Words Assumption

We now describe a set of models that go beyond the bag of words assumption. These models are typically more complex, less efficient, and less effective. For these reasons they are not widely used within the information retrieval community. How-

ever, it is important to describe the breadth of work done in this area in order to provide a clear picture of the difficulty involved with developing highly effective models that go beyond the bag of words assumption.

2.2.1 *n*-Gram Language Models

In this section we describe *n*-gram language models (for $n > 1$). These models are simple generalizations of the unigram language model approach that take context into account. That is, *n*-gram language models generate terms by conditioning on the previous $n - 1$ terms encountered. In a unigram model, generating the term *Lincoln* is equally likely regardless of the previous term. In a bigram model ($n = 2$), *Lincoln* has a higher likelihood of being generated after *president* than after *brick*, for example. Therefore, *n*-gram models capture the sequential structure of language generation.

As with unigram language models, documents are ranked according to query likelihood, which is computed as:

$$P(Q|D) = \prod_{i=1}^{|Q|} P(q_i|q_{i-1}, \dots, q_{i-n+1}, D) \quad (2.15)$$

where $P(q_i|q_{i-1}, \dots, q_{i-n+1}, D)$ can be estimated in a number of ways [41, 101, 105], many of which include some form of backoff to a unigram model. One way of estimating bigram probabilities, from [41], is:

$$P(w_i|w_{i-1}, D) = (1 - \lambda_1) \left[(1 - \lambda_2) \frac{tf_{w_i, w_{i-1}, D}}{tf_{w_{i-1}, D}} + \lambda_2 \frac{tf_{w_i, D}}{|D|} \right] + \lambda_1 \left[(1 - \lambda_3) \frac{cf_{w_i, w_{i-1}}}{cf_{w_{i-1}}} + \lambda_3 \frac{cf_{w_i}}{|C|} \right] \quad (2.16)$$

where λ_1 , λ_2 , and λ_3 are free parameters that control the smoothing.

Gao et al. showed that this model consistently outperformed unigram language models across a number of data sets using description-length queries [41]. Unfortunately, the model, as described, performs poorly on title-length queries. The model

is a generalization of Jelinek-Mercer smoothing, which is known to work well on longer queries [125]. Therefore, the model must be adjusted to be more like Dirichlet smoothing in order to perform well on title queries. This can be achieved by setting $\lambda_1 = \frac{\mu_1}{\mu_1 + |D|}$. We note that this modification does not follow naturally or formally from some underlying model. Instead it is a heuristic modification that only works because it makes the ranking function more like the Dirichlet ranking function.

Although n -gram language models capture the relationship between terms better, they are still not adequately robust for our needs. One criticism of such models is that they rely on evidence from the previous $n - 1$ terms, when in fact the next $n - 1$ terms might provide just as strong evidence. Consider the text *white house rose garden*. In a reasonable bigram language model of general English, $P(\textit{house}|\textit{white})$ would be assigned a high probability, but $P(\textit{rose}|\textit{house})$ would not. Therefore, under the bigram model, the likelihood of this sequence may actually be underestimated. However, conditioning on both past and future words could overcome such a problem. It is also noted that n -grams are typically overly strict, in that they do not allow the modeling of longer-range, unordered dependencies, such as the fact that two terms tend to often appear, not necessarily in order, within close proximity to each other. As we will show, our model is capable of handling a wide range of dependencies, including those that n -gram language models are not capable of.

2.2.2 Indri Inference Network Model

The retrieval model implemented in the Indri search engine [107] is an enhanced version of the model described in [64], which combines the language modeling [101] and inference network [109] approaches to information retrieval. The resulting model allows rich, structured queries to be evaluated using language modeling estimates within the network. Figure 2.2 shows a graphical model representation of the network. Within the model, documents are ranked according to $P(I|D, \alpha, \beta)$, the belief

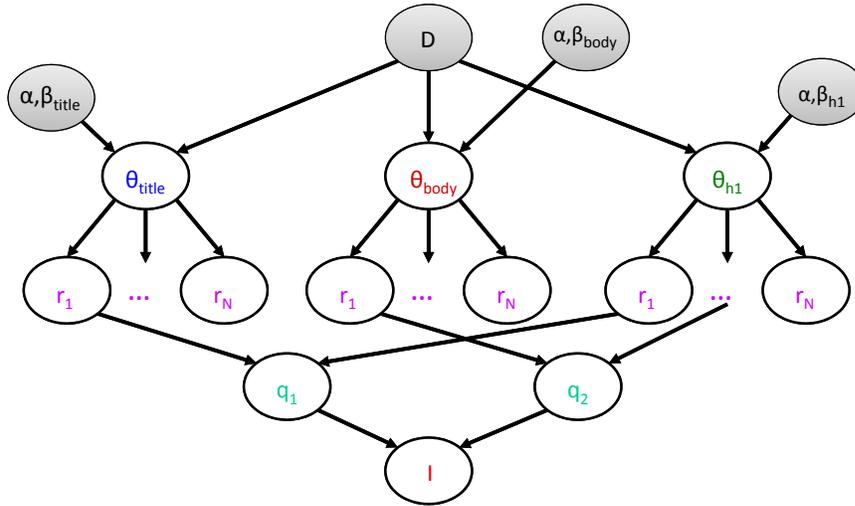


Figure 2.2. Indri's inference network retrieval model

the information need I is met given document D and hyperparameters α and β as evidence.

2.2.2.1 Document Representation

Typically, in the language modeling framework, documents are represented as a sequence of terms. Based on this sequence, a multinomial language model over the vocabulary is estimated. However, it is often the case that we wish to model more interesting text phenomenon, such as phrases or the *absence* of a term. Therefore, a different representation scheme is necessary. In the Indri model, documents are represented as multisets of binary feature vectors. Given a document, a feature vector is extracted for *every* position within the document. Figure 2.3 provides an example representation using this scheme. As the figure shows, the document is represented by a set of five vectors, one for each position within the document. This representation is very general and provides a way of modeling almost arbitrary textual features.

$$f(\text{A B C A B}, \begin{bmatrix} A \\ B \\ C \\ AA \\ AB \\ AC \\ BA \\ BB \\ BC \\ CA \\ CB \\ CC \end{bmatrix}) = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \right\}$$

Figure 2.3. Example Indri document representation for the document A B C A B. The features correspond to the single terms A, B, C, and bigrams “A A”, “A B”, “A C”, “B A”, “B B”, “B C”, “C A”, “C B”, and “C C”, respectively. The function f takes a document and set of features as input and outputs a document representation.

This representation moves away from modeling text towards modeling features of text. Throughout the remainder of this section we refer to such models as language models, although they really are better described as *language feature models*.

2.2.2.2 Language Models

Since the document event space is binary, it is no longer appropriate to estimate multinomial language models for each document. Instead, multiple-Bernoulli models are estimated, as in Model B of [68]. This overcomes the theoretical issues encountered in [64]. Note that the multiple-Bernoulli model imposes the assumption that the features (r_i 's) are independent, which of course may be a poor assumption depending on the feature set.

A Bayesian approach is taken and a multiple-Beta prior is imposed over the distribution of language models (θ). The Beta is chosen for simplicity, as it is the conjugate prior to the Bernoulli distribution. Thus, $P(D|\theta)$ is distributed according to *Multi-Bernoulli*(θ) and $P(\theta|\alpha, \beta)$ is distributed according to *Multi-Beta*(α, β).

Hence, our belief at node θ is computed as:

$$\begin{aligned}
 P(\theta_i|D, \alpha, \beta) &= \frac{P(D|\theta_i)P(\theta_i|\alpha_i, \beta_i)}{\int_{\theta_i} P(D|\theta_i)P(\theta_i|\alpha_i, \beta_i)} \\
 &= \text{Beta}(\#(r_i, D) + \alpha_i, |D| - \#(r_i, D) + \beta_i) \quad (2.17)
 \end{aligned}$$

for each i where $\#(r_i, D)$ is the number of times feature r_i is set to 1 in document D 's multiset of feature vectors.

Such a model is estimated using the entire document text. Additionally, if a document is structured, as HTML, SGML, and XML documents are, then language models can be estimated for each field. To do so, we treat all of the text in that appears within a given field as a pseudo-document. For example, a model can be estimated for all of the text that appears within the **h1** tags of a web page document.

2.2.2.3 Representation Nodes

The r_i nodes correspond to document features that can be represented in an Indri structured query. There is a one-to-one correspondence between r_i nodes and the features used to represent the document. Therefore, the r_i nodes represent binary events that feature i is observed.

Indri implements a number of textual features, including single terms, #N (ordered window N), and #uwN (unordered window N). Please refer to [64] for more details on these operators.

Using the framework developed thus far, the belief at a given representation node is computed as:

$$\begin{aligned}
 P(r_i = 1|D, \alpha, \beta) &= \int_{\theta_i} P(r_i = 1|\theta_i)P(\theta_i|D, \alpha_i, \beta_i) \\
 &= E[P(r_i|\theta_i)] \\
 &= \frac{\#(r_i, D) + \alpha_i}{|D| + \alpha_i + \beta_i} \quad (2.18)
 \end{aligned}$$

Furthermore, selecting $\alpha_i = \mu P(r_i = 1|C)$ and $\beta_i = \mu(1 - P(r_i = 1|C))$ results in the multiple-Bernoulli model equivalent of the multinomial model’s Dirichlet smoothing [123] estimate:

$$P(r_i|D, \alpha, \beta) = \frac{\#(r_i, D) + \mu P(r_i|C)}{|D| + \mu} \quad (2.19)$$

where μ acts as a tunable smoothing parameter.

2.2.2.4 Query Nodes

The query node operators are soft probabilistic operators that are used to combine evidence within the network. The operators are primarily used to combine evidence from representation nodes and other query nodes. The operators implemented in Indri are `#combine` (same as `#and`), `#weight` (same as `#wand`), `#or`, `#not`, `#sum`, `#wsum`, and `#max` [64].

The information need node, I , is also a query node that acts to combine all of the evidence of the query into a single belief. It is this belief (i.e., $P(I = 1|D, \alpha, \beta)$) that is used to rank documents. Therefore, the ranking function is defined in terms of query and representation nodes. For example, consider the following Indri query:

```
#weight( 1.0 #or( american #1(united states) ) 2.0 presidents )
```

For this query, the resulting ranking function would first compute beliefs for `american` and `#1(united states)`, then combine the beliefs using the probabilistic `#or` operator. The belief of this `#or` operator and the belief of `presidents` would be combined using the probabilistic `#weight` operator to produce the belief that the information need is satisfied. Within the inference network framework, queries are not explicitly defined. Instead, a structured query is used to construct a query network, which encodes a user’s information need.

Finally, we note that, since language modeling probabilities are used within the network, the `#wsum` operator no longer makes sense. Instead, the `#combine` (`#and`) and `#weight` (`#wand`) operators are more appropriate, since they produce an *idf*

effect [64]. It can be shown that the Indri query $\#combine(q_1 \dots q_N)$ using the estimates just described returns exactly the same ranked list as the query $q_1 \dots q_N$ using the traditional (multinomial with Dirichlet smoothing) query likelihood model.

2.2.2.5 Explicit vs. Implicit Query Generation

The Indri retrieval model can be used in two ways. First, given a simple keyword query, a system can be developed to convert the query into a structured Indri query. This process acts to transform the simple query into a richer representation. For example, phrases, synonyms, or task-specific operators may be automatically added to the query in order to improve effectiveness over the simple keyword query. The Indri retrieval model was successfully used in this capacity during the 2004-2006 TREC Terabyte Tracks [71, 72, 70], and the 2005 TREC Robust Track [66].

Alternatively, users can use the query language to manually construct complex queries. It has been shown that intelligently constructed manual queries can significantly outperform automatically generated queries [64]. However, the query language is too complex for novice users to use successfully. Only expert users, such as information analysts and librarians, are likely to benefit from such a query language. For this reason, algorithmic query construction is important.

Despite its success, the Indri retrieval model does not provide a formal mechanism for learning how to combine various types of evidence, making use of arbitrary evidence, or automatically converting a short keyword query into a rich structured query. The model that we present in this work is inspired by the Indri retrieval model, and attempts to overcome some of its limitations. As we will show, our proposed model does not require users to use a complex query language. It allows users to enter their information needs as short, keyword queries. Furthermore, the model can incorporate arbitrary textual and non-textual evidence, automatically learn the

best set of features to extract from the query, and provides a framework for learning how to combine the various types of evidence.

2.2.3 Other Models That Go Beyond the Bag of Words Assumption

There have been many models proposed to that go beyond the bag of words assumption [20, 29, 27, 35, 30]. We now briefly highlight several of these models.

Fagan examines how to identify and use non-syntactic (statistical) phrases [35]. Fagan identifies phrases using factors such as the number of times the phrase occurs in the collection and the proximity of the phrase terms. His results suggest no single method of phrase identification consistently yields improvements in retrieval effectiveness across a range of collections. For several collections, significant improvements in effectiveness are achieved when phrases are defined as any two terms within a query or document with unlimited proximity. That is, any two terms that co-occurred within a query or document were considered a phrase. However, for other collections, this definition proved to yield marginal or negative improvements.

Work done by Croft et al. shows similar results [29]. Their results showed phrases formed with a probabilistic AND operator slightly outperformed proximity phrases. The probabilistic AND operator boosts the scores of documents where the phrase terms co-occur. Therefore, little benefit was shown as the result of modeling term proximity.

In addition to the n -gram language models we described, several other language model variants have been proposed that attempt to model term dependencies [41, 77]. The dependence language model presented by Gao et al. in [41] showed consistent improvements over a baseline query likelihood system on a number of TREC collections. However, the model uses a link structure for each query which is not straightforward to construct. Our proposed model does not require a query link structure to be con-

structed. However, if such information is available, it can easily be incorporated into the model.

Recently, Mishne and de Rijke explored the use of proximity information to improve web retrieval [73]. Our model shares many closely related insights. Despite the high level similarity, the details of the models differ greatly, with our model allowing more general query dependencies and features to be considered in a more formally well-grounded framework.

Therefore, there have been many models proposed to go beyond the bag of words assumption, but none of them have allowed the use of arbitrary features, easy modeling of term dependencies, and yielded consistent, significant improvements in retrieval effectiveness. The model that we propose in the next chapter combines and generalizes the best aspects of these previous models within a robust, effective retrieval framework.

2.3 The Current State of the Art

Despite the large number of attempts to go beyond the bag of words assumption, there have been very few, if any, models that have been proven to be consistently better than the current best bag of words models (i.e., language modeling and BM25).

In fact, strong evidence that BM25 and language modeling are considered the state of the art retrieval models comes from looking at the models used by participants in recent years at TREC. Outside of several obscure, poor performing models, a majority of participant used either BM25 or language modeling, with some additional task-specific engineering added on top. Therefore, little progress has been made in advancing the state of the art of retrieval models since the advent of language modeling and BM25 a decade ago.

CHAPTER 3

A MARKOV RANDOM FIELD MODEL FOR INFORMATION RETRIEVAL

In this chapter we introduce our Markov random field retrieval model. Only the basics of the model are covered in this chapter. Subsequent chapters will go into more detail and describe various extensions of the basic model.

3.1 Modeling Relevance

We begin by describing what we seek to model. The four primary variables in most information retrieval systems are users (\mathcal{U}), queries (\mathcal{Q}), documents (\mathcal{D}), and relevance (\mathcal{R}). We define our event space to be $\mathcal{U} \times \mathcal{Q} \times \mathcal{D}$ and define relevance, $R \in \mathcal{R}$, to be a random variable over $\mathcal{U} \times \mathcal{Q} \times \mathcal{D}$. Thus, some relevance value is associated with every user, query, document tuple. Other factors, such as time and context are ignored.

These variables interact in real information systems in the following way. Users submit queries to the system and are presented a ranked list of documents. Some of the documents in the ranked are relevant, while others are non-relevant. Suppose that we were to collect a list of query/document pairs (Q, D) , such that some user found document D relevant to query Q . Imagine that such a list was collected across a large sample of users. The resulting list can be thought of as a sample from some underlying population of relevant query/document pairs that are aggregated across

users and conditioned on relevance. This, is then, a *relevance distribution*¹, which is similar in spirit to the one proposed by Lavrenko [57]. It is this distribution, $P(Q, D|R = 1)$, the joint distribution over query and document pairs, conditioned on relevance, that we focus on modeling. For notational convenience, we drop the explicit conditioning on relevance (i.e., $R = 1$) throughout the remainder of this work, unless otherwise noted.

3.2 MRFs for IR

There are many possible ways to model a joint distribution. In this work, we choose to use Markov random fields. Markov random fields, sometimes referred to as undirected graphical models, are commonly used in the statistical machine learning domain to model complex joint distributions. As we will show throughout the remainder of this section, there are many advantages and few, if any, disadvantages to using MRFs for information retrieval.

A Markov random field is constructed from a graph G . The nodes in the graph represent random variables, and the edges define the independence semantics between the random variables. The independence semantics are governed by the Markov property.

Markov Property. Let $G = (V, E)$ be the undirected graph associated with a Markov random field, then $P(v_i|v_{j \neq i}) = P(v_i|v_j : (v_i, v_j) \in E)$ for every random variable v_i associated with a node in V .

The Markov Property states that every random variable in the graph is independent of its non-neighbors given observed values for its neighbors. Therefore, different edge configurations impose different independence assumptions.

¹Note that we make the assumption that relevance is binary, which is commonly used for information retrieval tasks. If relevance is non-binary, then a different relevance distribution can be estimated for each relevance level.

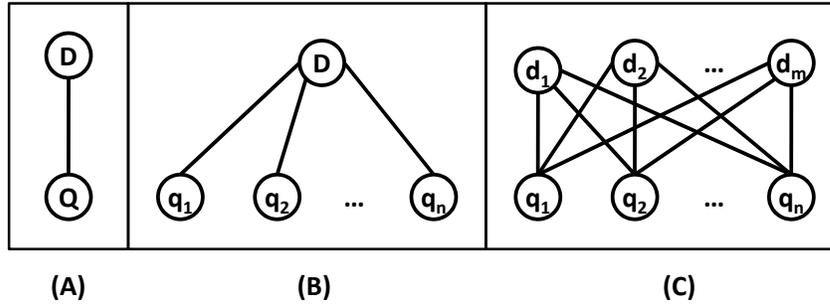


Figure 3.1. Examples of three ways to model the joint distribution $P(Q, D)$ using Markov random fields.

There are several ways to model the joint distribution $P(Q, D)$ using Markov random fields. Figure 3.1 summarizes the various options that are available. Option A constructs a graph with two nodes, a query node Q and a document node D . However, this model is too coarsely specified and does not provide any insight into the types of term dependencies that are being modeled since it models whole queries and documents. Option B breaks the query apart into individual terms and treats the document as a whole. Given a query of length n , this results in a graph with n query term nodes and a document node. This option provides more specific control over which query term dependencies are modeled. Finally, option C breaks apart both the document and the query into individual terms. Given a query of length n and a document of length m , our graph would contain n query term nodes and m document term nodes. This option provides the most flexibility for modeling both query and document term dependencies. However, the model is likely to be overly complex. Modeling dependencies between query terms is more feasible than modeling dependencies between document terms since queries are generally much shorter than documents and exhibit less complex dependencies between terms.

Option B satisfies our needs without being overly complex, and so it will be used throughout the remainder of this work. Thus, given a query of length n , our graph G

consists of n query term nodes and a single document node D . The random variables associated with the query term nodes are multinomials over the underlying vocabulary \mathcal{V} and the random variable associated with the document node is also a multinomial over the set of documents in the collection. We note that variations on this theme are possible. For example, it may be appropriate to include several document nodes or even other types of nodes, such as document structure nodes within the MRF.

The joint probability mass function over the random variables in G is defined by:

$$P_{G,\Lambda}(Q, D) = \frac{1}{Z_\Lambda} \prod_{c \in C(G)} \psi(c; \Lambda) \quad (3.1)$$

where $Q = q_1 \dots q_n$ are the set of query term nodes, D is the document node, $C(G)$ is the set of maximal cliques in G , each $\psi(\cdot; \Lambda)$ is a non-negative *potential function* over clique configurations parameterized by Λ and $Z_\Lambda = \sum_{Q,D} \prod_{c \in C(G)} \psi(c; \Lambda)$ normalizes the distribution. It is generally infeasible to compute Z_Λ due to the exponential number of terms in the summation.

Therefore, in order to compute the joint distribution we need a graph G , potential functions ψ , and the parameter vector Λ . Detailed descriptions of these components are given in the following sections.

3.2.1 Graph Structure

We have already described how the nodes of our MRF are chosen. We now must show how these nodes can be connected together. As explained before, the Markov Property dictates the dependence semantics of the MRF. Therefore, it is relatively straightforward to explore various independence assumptions by constructing MRFs with different graph structures.

We consider three generalized graph structures, each with different underlying independence assumptions. The three structures are *full independence* (FI), *sequential dependence* (SD), and *full dependence* (FD). Figure 3.2 shows graphical model

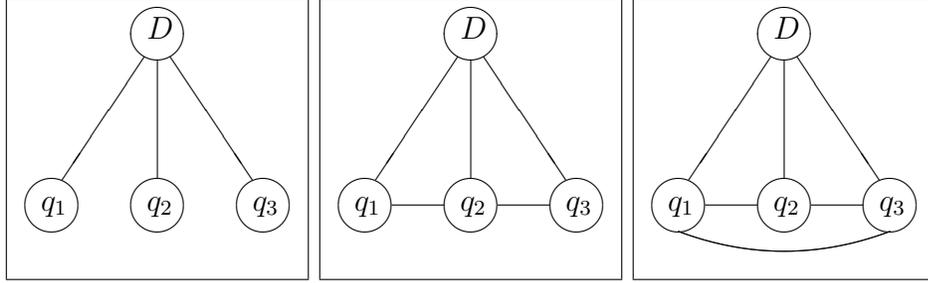


Figure 3.2. Example Markov random field model for three query terms under various independence assumptions, including full independence (left), sequential dependence (middle), and full dependence (right).

representations of each. These generalized structures are considered because of their significance to information retrieval. As we now show, each corresponds to a well-studied class of retrieval models.

The full independence structure makes the assumption that query terms q_i are independent given some document D . That is, the likelihood of observing query term q_i is not affected by the observation of any other query term, or more succinctly, $P(q_i|D, q_{j \neq i}) = P(q_i|D)$. This corresponds to the independence assumption made by many of the bag of words models that were described in Chapter 2.

As its name implies, the sequential dependence structure assumes a dependence between neighboring query terms. Formally, this assumption states that $P(q_i|D, q_{j \neq i}) = P(q_i|D, q_{i-1}, q_{i+1})$. Models of this form are similar in nature to bigram and biterm language models [101, 105].

The last structure we consider is the full dependence structure. In this structure, we assume all query terms are in some way dependent on each other. Graphically, a query of length n translates into the complete graph K_{n+1} , which includes edges from each query node to the document node D , as well. This model is an attempt to capture longer range dependencies than the sequential dependence structure. If such a model can accurately be estimated, it should be expected to perform at least as well as a model that ignores term dependence.

There are other reasonable ways of constructing G given a query, such as that proposed by Gao et al. [41], in which dependencies between terms are inferred using natural language processing techniques. The advantage of using one of the structures just described is that there is no need to rely on natural language processing techniques, which can often produce noisy output, especially on short segments of text. Of course, some of the dependencies imposed by our structure may be incorrect, but in general, they capture meaningful relationships between terms.

3.2.2 Potential Functions

In order to compute the MRF’s joint probability mass function (Equation 3.1), a set of potential functions must be defined over configurations of the maximal cliques in the underlying graph. These potential functions can be thought of as *compatibility functions*. That is, they are meant to reflect how compatible a given clique configuration is. How compatibility is defined and measured depends on the task and clique.

For example, in Figure 3.2, the nodes D and q_1 form a maximal clique in the full independence variant. The potential function defined over the clique should reflect how compatible the term q_1 is to D . Here, compatibility may be defined as “aboutness” and measured using some *tf.idf* score for the term q_1 in D .

Typically, potential functions are built top-down, starting with a maximal clique and defining a potential over it. However, within our model, we choose to build potential functions in a bottom-up fashion, which provides more fine grained control over the behavior of the functions. This is accomplished by first associating one or more real-valued *feature functions* with each (maximal or non-maximal) clique in the graph. Each feature function has a feature weight associated with it that is a free parameter in our model. Then, non-negative potential functions over the maximal

cliques are constructed from these feature functions and feature weights using an exponential form. We now formally describe the details of this process.

1. Assign one or more feature functions to each clique in G . This assignment can be encoded as a set of 3-tuples, $\mathcal{C} = \{(c, f(\cdot), \lambda)\}_{i=1}^n$, where c is a clique of G , $f(\cdot)$ is the feature function assigned to the clique, and λ is the weight (parameter) associated with the feature. Recall that the same clique may be associated with more than one feature function.
2. For every $(c, f(\cdot), \lambda) \in \mathcal{C}$, assign c to one of the maximal clique(s) in G that c is a subclique of. It is always possible to assign subcliques to maximal cliques, although this assignment is not guaranteed to be unique.
3. For every maximal clique in G , define its potential function as $\psi(\cdot) = \exp(\sum_c \lambda f(\cdot))$, where the sum goes over the cliques that were assigned to the maximal clique in Step 2.

We now provide an example to illustrate the process. Consider the full independence graph in Figure 3.2. Suppose that we make the following assignment of feature functions and parameters to the graph:

$$\begin{aligned}
 & (\{q_1, D\}, f_1(q_1, D), \lambda_1) \\
 & (\{q_1, D\}, f_4(q_1, D), \lambda_4) \\
 & (\{q_2, D\}, f_2(q_2, D), \lambda_2) \\
 & (\{q_2, D\}, f_4(q_2, D), \lambda_4) \\
 & (\{q_3, D\}, f_3(q_3, D), \lambda_3) \\
 & (\{q_3, D\}, f_4(q_3, D), \lambda_4) \\
 & (\{D\}, f_5(D), \lambda_5)
 \end{aligned}$$

where each f_i is some real-valued feature function defined over the configurations of the clique. The specific form of the feature functions is not important in this example.

After assigning each clique to a maximal clique, we construct the following potential functions:

$$\psi(q_1, D) = \exp [\lambda_1 f_1(q_1, D) + \lambda_4 f_4(q_1, D) + \lambda_5 f_5(D)] \quad (3.2)$$

$$\psi(q_2, D) = \exp [\lambda_2 f_2(q_2, D) + \lambda_4 f_4(q_2, D)] \quad (3.3)$$

$$\psi(q_3, D) = \exp [\lambda_3 f_3(q_3, D) + \lambda_4 f_4(q_3, D)] \quad (3.4)$$

This construction is not unique since the clique $\{D\}$ is a subclique of all three maximal cliques. Therefore, we can assign feature function $f_5(D)$ to any of the maximal cliques. In the previous set of potential functions, it was assigned to the maximal clique $\{q_1, D\}$. If it had been assigned to the maximal clique $\{q_3, D\}$ instead, the following potential functions would have been constructed:

$$\psi(q_1, D) = \exp [\lambda_1 f_1(q_1, D) + \lambda_4 f_4(q_1, D)] \quad (3.5)$$

$$\psi(q_2, D) = \exp [\lambda_2 f_2(q_2, D) + \lambda_4 f_4(q_2, D)] \quad (3.6)$$

$$\psi(q_3, D) = \exp [\lambda_3 f_3(q_3, D) + \lambda_4 f_4(q_3, D) + \lambda_5 f_5(D)] \quad (3.7)$$

It is critical to note that, even though the potential function definitions are not guaranteed to be unique using this formulation, the joint probability mass function will be unique. It is easy to see that, for this example, the joint, under all possible assignments is equal to:

$$\begin{aligned} P_{G,\Lambda}(Q, D) = Z_{\Lambda}^{-1} \exp[& \lambda_1 f_1(q_1, D) + \lambda_4 f_4(q_1, D) + \\ & \lambda_2 f_2(q_2, D) + \lambda_4 f_4(q_2, D) + \\ & \lambda_3 f_3(q_3, D) + \lambda_4 f_4(q_3, D) + \lambda_5 f_5(D)] \end{aligned} \quad (3.8)$$

This example also serves to illustrate that both functions and parameters can be shared across potential functions. Here, the feature function f_4 and parameter λ_4

were shared across three cliques. In order to share a feature function across cliques, we require that the input to the feature function be compatible with each clique. For example, a feature function that takes two term nodes and a document node as input can only be shared across cliques with two term nodes and a document node. We do not permit a feature function that only takes a document node as input to be shared with a clique that contains both a query term node and a document node. There are no restrictions on sharing parameters across cliques, however. By sharing parameters across cliques, we effectively tie parameters together, which reduces the number of free parameters and can help overcome data sparseness issues.

3.3 Building MRFs

As we just showed, potential functions are constructed by assigning feature functions and parameters to arbitrary cliques in the MRF. In this section, we describe how textual and non-textual features can be represented and assigned to cliques. Potentials can then be built from these features and be used to compute $P_{G,\Lambda}(Q, D)$.

In this section, we propose a novel method for representing MRFs for information retrieval. We represent MRFs using a *canonical form*. The canonical form is designed to be a compact, intuitive, and flexible method of representing MRFs. It can handle a wide variety of graph structures and features that are useful for information retrieval tasks. A canonical forms have the following structure:

$$\begin{aligned}
 &(\text{dependence model type, clique set type, weighting function})_1 : \lambda_1 \\
 &(\text{dependence model type, clique set type, weighting function})_2 : \lambda_2 \\
 &\dots \\
 &(\text{dependence model type, clique set type, weighting function})_n : \lambda_n
 \end{aligned}$$

Here, a 3-tuple represents how feature functions are assigned to cliques. Each 3-tuple assigns a feature function to one or more cliques within the graph. The variable after the colon represents the parameter associated with all of the feature functions

assigned by the 3-tuple. As we showed in the previous section, this ties the parameters of all of the feature functions associated with a 3-tuple together. The details of this assignment and tying process will become clearer later in this section when we work through several examples.

Given a canonical form, it is easy to systematically build the corresponding MRF and derive both the joint probability mass function, as well as the ranking function. We represent all MRFs throughout the remainder of this work using these canonical forms. We now describe the meaning and details of each component in the 3-tuple.

3.3.1 Dependence Model Type

The first entry in the tuple is the *dependence model type*, which specifies the dependencies, if any, that are to be modeled between query terms. As we described before, dependencies are encoded by the edges in the MRF, with different edge configurations correspond to different types of dependence assumptions.

In this work, we only allow the dependence model type to be full independence (FI), sequential dependence (SD), or full dependence (FD), which are the three generalized graph structures described in Section 3.2.1 and illustrated in Figure 3.2.

For a given MRF, each feature function may have a different dependence model type. The dependence model type simply defines the graph structure that the current feature is applied to. The graph structure that the resulting MRF has depends on the dependence model types of all of its features combined.

3.3.2 Clique Set Type

The second entry in the tuple, the *clique set type*, describes the set of (maximal or non-maximal) cliques within the graph that the feature function is to be applied to. Thus, each feature function can be applied to one or more cliques within the graph, depending on the clique set.

Description	Notation	Example
Set of cliques containing the document node and exactly one query term	T_{QD}	$\{\{q_1, D\}, \{q_2, D\}, \{q_3, D\}\}$
Set of cliques containing the document node and two or more query terms that appear in sequential order within the query	O_{QD}	$\{\{q_1, q_2, D\}, \{q_2, q_3, D\}, \{q_1, q_2, q_3, D\}\}$
Set of cliques containing the document node and two or more query terms that appear unordered within the query	U_{QD}	$\{\{q_1, q_3, D\}\}$
Set of cliques containing exactly one query term	T_Q	$\{\{q_1\}, \{q_2\}, \{q_3\}\}$
Set of cliques containing two or more query terms that appear in sequential order within the query	O_Q	$\{\{q_1, q_2\}, \{q_2, q_3\}, \{q_1, q_2, q_3\}\}$
Set of cliques containing two or more query terms that appear unordered within the query	U_Q	$\{\{q_1, q_3\}\}$
Set containing only the singleton node D	D	$\{\{D\}\}$

Table 3.1. Example clique sets for the query $q_1 q_2 q_3$ under full dependence model.

We propose seven clique sets that can be used within our model. These sets are summarized in Table 3.1. In order to motivate these clique sets, we enumerate every possible type of clique that is of interest to us, beginning with cliques that contain the document node and one or more query term nodes.

First, the simplest type of clique that contains the document node and one or more query nodes is a 2-clique consisting of an edge between a query term q_i and the document D . A potential function over such a clique should measure how well, or how likely query term q_i describes the document.

Next, we consider cliques that contain two or more query terms. For such cliques there are two possible cases, either all of the query terms within the clique appear contiguously in the query or they do not. The fact that query terms appear contiguously within a query provides different (stronger) evidence about the information need than a set of non-contiguous query terms. For example, in the query *train station security measures*, if any of the subphrases, *train station*, *train station security*, *station security measures*, or *security measures* appear in a document then there is strong evidence in favor of relevance.

Although the occurrence of contiguous sets of query terms provide strong evidence of relevance, it is also the case that the occurrence of non-contiguous sets of query terms can provide valuable evidence. However, since the query terms are not contiguous we do not expect them to appear in order within relevant documents. Rather, we only expect the terms to appear ordered or unordered within a given proximity of each other. In the previous example, documents containing the terms *train* and *security* within some short proximity of one another also provide additional evidence towards relevance. This issue has been explored in the past by a number of researchers [29, 35].

Therefore, for cliques consisting of the document node and one or more query term nodes, we have the following clique sets:

- T_{QD} – set of cliques containing the document node and exactly one query term.
- O_{QD} – set of cliques containing the document node and two or more query terms that appear in sequential order within the query.
- U_{QD} – set of cliques containing the document node and two or more query terms that appear unordered within the query.

Note that the cliques that make up each set may change for different dependence model types. For example, O_{QD} and U_{QD} are empty under the full independence assumption since that would result in a graph where there are no cliques with two or more query term nodes. However, under the sequential dependence assumption, and with a query of length 2 or more, such cliques will exist and O_{QD} and U_{QD} will be non-empty.

Next, we consider cliques that only contain query term nodes. These clique sets are defined in an analogous way to those just defined, except the the cliques are only made up of query term nodes and do not contain the document node. Potential functions over these cliques should capture how compatible query terms are to one another. These clique potentials may take on the form of language models that impose well-formedness of the terms. Therefore, we define following query-dependent clique sets:

- T_Q – set of cliques containing exactly one query term.
- O_Q – set of cliques containing two or more query terms that appear in sequential order within the query.
- U_Q – set of cliques containing two or more query terms that appear unordered within the query.

Finally, there is the clique that only contains the document node. Potentials over this node can be used as a type of document prior, encoding document-centric properties. This trivial clique set is then:

- D – clique set containing only the singleton node D

We note that our clique sets form a partition over the cliques of G . This partition separates the cliques into sets that are meaningful from an information retrieval perspective. Thus, these clique sets make it easy to apply features in a very specific manner within the MRF.

Of course, the clique sets we defined here are not unique. It is possible to define many different types of clique sets. For example, another clique set may be defined as “the clique that contains the first query term and the document node”. Given enough training data, it may be possible to define such fine grained clique sets. However, given the limited amount of training data, we focus our attention on the coarse grained clique sets defined above.

3.3.3 Weighting Function

Finally, the third entry in the tuple is the *weighting function*, which defines the feature function that is applied to the cliques defined by the clique set. In this section we define weighting functions that can be used with the different clique sets we just defined. It is not our goal to provide a comprehensive list of possible feature functions. Instead, we simply seek to provide a few examples of the types of feature functions that are possible.

3.3.3.1 Weighting Functions for T_{QD} , O_{QD} , and U_{QD}

We first describe weighting functions that can be used with cliques in the T_{QD} , O_{QD} , and U_{QD} clique sets. These cliques consist of a set of query term nodes and a

document node. Therefore, the weighting functions applied to these cliques should measure how much the document is “about” the query terms.

The weighting functions we propose are based on language modeling estimates and the BM25 weighting model, which we described in Chapter 2. It is straightforward to use the standard forms for these weighting functions for the single term cliques (T_{QD}). However, we must define how to match the query terms within documents when applying these weighting functions to ordered term cliques (O_{QD}) and unordered term cliques (U_{QD}).

For ordered term cliques, we match terms in documents using the Indri ordered window operator ($\#M$), where the parameter M determines how many non-matching terms are allowed to appear between matched terms [64]. For clique $\{q_i, \dots, q_{i+k}, D\}$, we match documents according to $\#M(q_i \dots q_{i+k})$. This rewards documents for preserving the order that the query terms occur in.

In the unordered clique set case, we match terms using the Indri unordered window operator ($\#\text{uw}N$), where N defines the maximum size of the window that the terms may occur (ordered or unordered) in. For clique $\{q_i, \dots, q_j, D\}$ that contains k query terms, documents are matched according to $\#\text{uw}Nk(q_i \dots q_j)$. Notice that we multiply the number of terms in the clique set by N . If $N = 1$, then all k query terms must occur, ordered or unordered, within a window of k terms of each other within the document. As N increases, the matching becomes looser. If $N = \textit{unlimited}$, then any document that contains all k query terms is matched. By using this matching scheme, we reward documents in which subsets of query terms occur appear within close proximity of each other.

Table 3.2 summarizes these weighting functions. Of course, many different types of weighting functions can easily be used within the model. For example, if new, more effective term weighting functions are developed in the future, then they can be easily used instead of, or in addition to, the Dirichlet or BM25 weighting functions.

LM
$f_{LM,T}(q_i, D) = \log \left[\frac{tf_{q_i,D} + \mu^t \frac{cf_{q_i}}{ C }}{ D + \mu^t} \right]$
LM-O-M
$f_{LM,O,M}(q_1, \dots, q_k, D) = \log \left[\frac{tf_{\#M(q_1 \dots q_k),D} + \mu^w \frac{cf_{\#M(q_1 \dots q_k)}}{ C }}{ D + \mu^w} \right]$
LM-U-N
$f_{LM,U,N}(q_1, \dots, q_k, D) = \log \left[\frac{tf_{\#uwnk(q_1 \dots q_k),D} + \mu^w \frac{cf_{\#uwnk(q_1 \dots q_k)}}{ C }}{ D + \mu^w} \right]$
BM25
$f_{T,BM25}(q_i, D) = \frac{(k_1^t + 1)tf_{w,D}}{k_1^t \left((1-b^t) + b^t \frac{ D }{ D _{avg}} \right) + tf_{w,D}} \log \frac{N - df_w + 0.5}{df_w + 0.5}$
BM25-O-M
$f_{BM25,O,M}(q_1, \dots, q_k, D) = \frac{(k_1^w + 1)tf_{\#M(q_1 \dots q_k),D}}{k_1^w \left((1-b^w) + b^w \frac{ D }{ D _{avg}} \right) + tf_{\#M(q_1 \dots q_k),D}} \log \frac{N - df_{\#M(q_1 \dots q_k)} + 0.5}{df_{\#M(q_1 \dots q_k)} + 0.5}$
BM25-U-N
$f_{BM25,U,N}(q_1, \dots, q_k, D) = \frac{(k_1^w + 1)tf_{\#uwnk(q_1 \dots q_k),D}}{k_1^w \left((1-b^w) + b^w \frac{ D }{ D _{avg}} \right) + tf_{\#uwnk(q_1 \dots q_k),D}} \log \frac{N - df_{\#uwnk(q_1 \dots q_k)} + 0.5}{df_{\#uwnk(q_1 \dots q_k)} + 0.5}$

Table 3.2. Summary of Dirichlet and BM25 weighting functions that can be used with cliques in the T_{QD} , O_{QD} , and U_{QD} clique sets. Here, M and N act as weighting function parameters that affect how matching is done, $tf_{e,D}$ is the number of times expression e matches in document D , $cf_{e,D}$ is the number of times expression e matches in the entire collection, df_e is the total number of documents that have at least one match for expression e , $|D|$ is the length of document D , $|D|_{avg}$ is the average document length, N is the number of documents in the collection, and $|C|$ is the total length of the collection. Finally, μ^t , μ^w , k_1^t , k_1^w , b^t , and b^w are weighting function hyperparameters. The t and w superscripts indicate term and window hyperparameters, respectively.

3.3.3.2 Weighting Functions for T_Q , O_Q , and U_Q

Next, we consider weighting functions for the cliques in the T_Q , O_Q , and U_Q clique sets. These cliques consist of one or more query terms and no document nodes. Weighting functions defined over them should reflect their general importance or informativeness. Therefore, IDF-based measures are a natural set of feature functions to use for these types of cliques.

The two IDF measures that we use as feature functions are inverse collection frequency (ICF) and the Okapi IDF. Inverse collection frequency is very similar to IDF, except it considers the number of times an expression occurs, rather than the number of documents it occurs in. As with the weighting functions described in the previous section, it is straightforward to apply standard IDF features to the single term cliques (T_Q). We use the same matching semantics as described in the previous section for the ordered terms cliques (O_Q) and the unordered terms cliques (U_Q).

Our proposed feature functions are shown in Table 3.3. Other possible feature functions for these types of cliques include measures of how lexically cohesive the terms are and the average vocabulary level of the terms.

3.3.3.3 Weighting Functions for D

Depending on the task, there are a wide variety of weighting functions that can be applied to the document node clique. Some examples include document length [100], document quality [127], PageRank [10], URL depth [54], readability [99], sentiment [80], and opinionatedness [79].

Although we do not explore all of these query independent features in this work, we do make use of several of them for a web search task in Chapter 5.

3.3.4 Examples

Now that we have described each element that makes up the 3-tuple, we show how to construct MRFs from canonical forms. We do this by working through a number

ICF
$f_{ICF,T}(q_i, D) = -\log \frac{cf_{q_i}}{ C }$
ICF-O-M
$f_{ICF,U,M}(q_1, \dots, q_k, D) = -\log \frac{cf_{\#M(q_1 \dots q_k)}}{ C }$
ICF-U-N
$f_{ICF,O,N}(q_1, \dots, q_k, D) = -\log \frac{cf_{\#uNk(q_1 \dots q_k)}}{ C }$
IDF
$f_{IDF,BM25}(q_i, D) = \log \frac{N-df_w+0.5}{df_w+0.5}$
IDF-O-M
$f_{IDF,O,M}(q_1, \dots, q_k, D) = \log \frac{N-df_{\#M(q_1 \dots q_k)}+0.5}{df_{\#M(q_1 \dots q_k)}+0.5}$
IDF-U-N
$f_{IDF,U,N}(q_1, \dots, q_k, D) = \log \frac{N-df_{\#uNk(q_1 \dots q_k)}+0.5}{df_{\#uNk(q_1 \dots q_k)}+0.5}$

Table 3.3. Summary of ICF and IDF weighting functions that can be used with cliques in the T_Q , O_Q , and U_Q clique sets.

of examples. In all of the following examples examples, it is assumed that the query being evaluated is *new york city*.

Our first example is for the following canonical form:

$$(FI, T_{QD}, BM25) : \lambda$$

This canonical form includes a single feature function. The feature uses the full independence graph structure, is applied to the cliques in T_{QD} , and uses the BM25 weighting function. This expands to the following assignment of feature functions:

$$(\{\text{new}, D\}, f_{BM25,T}(\text{new}, D), \lambda)$$

$$(\{\text{york}, D\}, f_{BM25,T}(\text{york}, D), \lambda)$$

$$(\{\text{city}, D\}, f_{BM25,T}(\text{city}, D), \lambda)$$

Notice that all of the features share the same parameter.

This assignment can then be transformed into the following set of potential functions, using the process described in Section 3.2.2:

$$\psi(\text{new}, D) = \exp[\lambda f_{BM25,T}(\text{new}, D)] \quad (3.9)$$

$$\psi(\text{york}, D) = \exp[\lambda f_{BM25,T}(\text{york}, D)] \quad (3.10)$$

$$\psi(\text{city}, D) = \exp[\lambda f_{BM25,T}(\text{city}, D)] \quad (3.11)$$

where $f_{BM25,T}$ takes on the BM25 form as given in Table 3.2. The resulting probability mass function is then given by:

$$\begin{aligned} P_{G,\Lambda}(\text{new york city}, D) &= Z_{\Lambda}^{-1} \exp[\lambda f_{BM25,T}(\text{new}, D) + \\ &\quad \lambda f_{BM25,T}(\text{york}, D) + \\ &\quad \lambda f_{BM25,T}(\text{city}, D)] \end{aligned} \quad (3.12)$$

We see that this joint probability mass function is rank equivalent to the BM25 score of query for document D . Analogously, if $f_{BM25,T}$ is replaced with $f_{LM,T}$, the probability mass function is rank equivalent to query likelihood scoring in the language modeling framework.

Next, we consider the following canonical form:

$$(\text{SD}, O_{QD}, \text{LM-O-4}) : \lambda$$

which contains a single feature that uses the sequential dependence model, is applied to cliques in O_{QD} , and uses the Dirichlet weighting function. This expands into the following assignment of feature functions to cliques:

$$(\{ \text{new, york}, D \}, f_{LM,O,4}(\text{new, york}, D), \lambda)$$

$$(\{ \text{york, city}, D \}, f_{LM,O,4}(\text{york, city}, D), \lambda)$$

which is then transformed into the following set of potential functions:

$$\psi(\text{new, york}, D) = \exp[\lambda f_{LM,O,4}(\text{new, york}, D)] \quad (3.13)$$

$$\psi(\text{york, city}, D) = \exp[\lambda f_{LM,O,4}(\text{york, city}, D)] \quad (3.14)$$

where $f_{LM,O,4}$ takes on the Dirichlet form and M , the ordered window size, is set to 4.

Finally, we provide an example of a more complex canonical form. Consider the following canonical form:

$$\begin{aligned} (\text{FD}, O_{QD}, \text{LM-O-8}) &: \lambda_1 \\ (\text{FI}, T_Q, \text{IDF}) &: \lambda_2 \\ (\text{FI}, D, \text{PageRank}) &: \lambda_3 \end{aligned}$$

which then results in the following set of feature function assignments:

$$\begin{aligned} & (\{ \text{new, york}, D \}, f_{LM,O,8}(\text{new, york}, D), \lambda_1) \\ & (\{ \text{york, city}, D \}, f_{LM,O,8}(\text{york, city}, D), \lambda_1) \\ & (\{ \text{new, york, city}, D \}, f_{LM,O,8}(\text{new, york, city}, D), \lambda_1) \\ & (\{ \text{new} \}, f_{IDF,T}(\text{new}, D), \lambda_2) \\ & (\{ \text{york} \}, f_{IDF,T}(\text{york}, D), \lambda_2) \\ & (\{ \text{city} \}, f_{IDF,T}(\text{city}, D), \lambda_2) \\ & (\{ D \}, f_{PageRank}(D), \lambda_3) \end{aligned}$$

and the following potential function:

$$\begin{aligned} \psi(\text{new, york, city}, D) &= \exp[\lambda_1 f_{LM,O,8}(\text{new, york}, D) + \lambda_1 f_{LM,O,8}(\text{york, city}, D) + \\ & \lambda_1 f_{LM,O,8}(\text{new, york, city}, D) + \lambda_2 f_{IDF,T}(\text{new}, D) + \\ & \lambda_2 f_{IDF,T}(\text{york}, D) + \lambda_2 f_{IDF,T}(\text{city}, D) + \\ & \lambda_3 f_{PageRank}(D)] \end{aligned} \tag{3.15}$$

These examples illustrate that our proposed canonical form allows us to compactly define a large, rich set of MRFs for use with information retrieval tasks.

3.4 Ranking

Using our canonical feature representation, we derive the following simplified form of the joint distribution:

$$\begin{aligned} \log P_{G,\Lambda}(Q, D) = & \underbrace{\sum_{c \in T_{QD}} \lambda_c f_c(c) + \sum_{c \in O_{QD}} \lambda_c f_c(c) + \sum_{c \in U_{QD}} \lambda_c f_c(c)}_{\text{document + query dependent}} + \\ & \underbrace{\sum_{c \in T_Q} \lambda_c f_c(c) + \sum_{c \in O_Q} \lambda_c f_c(c) + \sum_{c \in U_Q} \lambda_c f_c(c)}_{\text{query dependent}} + \\ & \underbrace{\sum_{c \in D} \lambda_D f_c(c)}_{\text{document dependent}} - \underbrace{\log Z_\Lambda}_{\text{document + query independent}} \end{aligned} \quad (3.16)$$

$$(3.17)$$

where λ_c and f_c are the parameter and weighting (feature) function associated with clique c , respectively.

Given a query Q as evidence, we can use the model to rank documents in descending order according of the conditional $P_{G,\Lambda}(D|Q)$. Fortunately, properties of rankings allow us to significantly simplify the computation. That is,

$$\begin{aligned} P_{G,\Lambda}(D|Q) & \stackrel{\text{rank}}{=} \log P_{G,\Lambda}(D|Q) \\ & = \log \frac{P_{G,\Lambda}(Q, D)}{P_{G,\Lambda}(Q)} \\ & = \log P_{G,\Lambda}(Q, D) - \log P_{G,\Lambda}(Q) \\ & \stackrel{\text{rank}}{=} \log P_{G,\Lambda}(Q, D) \end{aligned} \quad (3.18)$$

After dropping document independent expressions from $\log P_{G,\Lambda}(Q, D)$, we derive the following ranking function:

$$P_{G,\Lambda}(D|Q) \stackrel{\text{rank}}{=} \sum_{c \in T_{QD}} \lambda_c f_c(c) + \sum_{c \in O_{QD}} \lambda_c f_c(c) + \sum_{c \in U_{QD}} \lambda_c f_c(c) + \sum_{c \in D} \lambda_c f_c(c) \quad (3.19)$$

which is a simple weighted linear combination of feature functions that can be computed efficiently for reasonable graphs since the partition function Z_Λ does not need to be computed. Later, in Chapter 7 we show how the reverse conditional, $P_{G,\Lambda}(Q|D)$, can be used for query expansion.

Furthermore, we note that, under the assumption of binary relevance, ranking documents according to $P_{G,\Lambda}(D|Q)$ adheres to the Probability Ranking Principle (PRP) [87]. This can be shown as follows:

$$\begin{aligned}
 P(R = 1|Q, D) &= \frac{P(D|Q, R = 1)P(Q|R = 1)P(R = 1)}{P(Q, D)} \\
 &\stackrel{rank}{=} P(D|Q, R = 1)
 \end{aligned}
 \tag{3.20}$$

The rank equivalence follows because $P(Q|R = 1)$ and $P(R = 1)$ are document independent, and $P(Q, D)$ is uniform across all query/document pairs [40]. Therefore, $P(R = 1|Q, D)$ is rank equivalent to $P(D|Q, R = 1)$, which is exactly $P_{G,\Lambda}(D|Q)$, as defined in our model². Hence, we have shown that ranking documents according to $P_{G,\Lambda}(D|Q)$ is equivalent to ranking documents in descending order of probability of relevance, which is the optimal strategy according to the PRP.

3.5 Discussion

In this chapter, we described the basics of our proposed Markov random field model for information retrieval. We explained our underlying model of relevance, basic MRF theory, and how MRF models can easily be constructed using a canonical form. The proposed model is very robust, as it can model a wide variety of dependencies between query terms, and can make use of arbitrary textual and non-textual

²Recall that we chose to implicitly condition on relevance (i.e., $R = 1$) in order to simplify the notation.

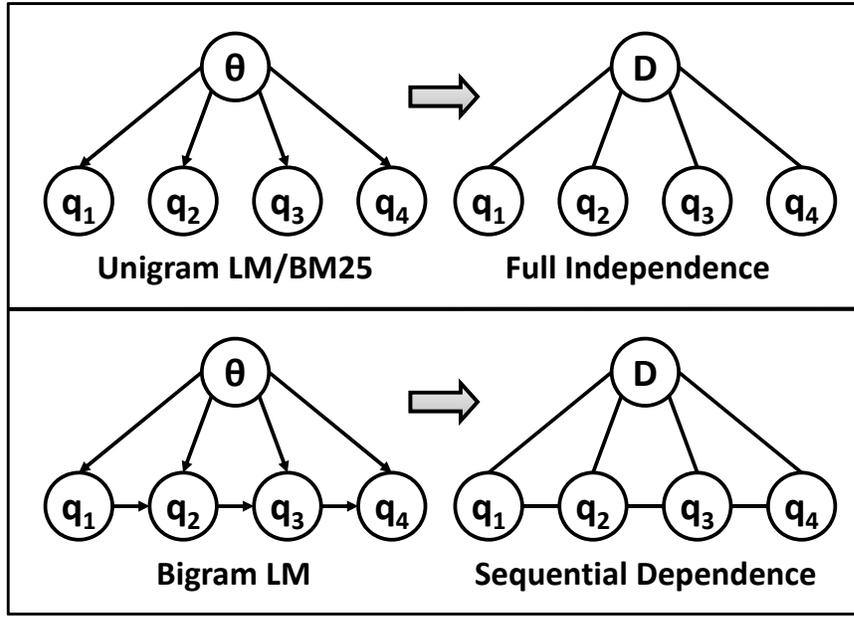


Figure 3.3. Illustration showing how the full independence model generalizes unigram language modeling and BM25 (top), and how the sequential dependence model generalizes bigram language modeling (bottom).

features, as well. This is the first model for information retrieval that has both of these important properties.

It is easy to show that our proposed model subsumes many previously proposed information retrieval models, which further proves the model’s flexibility. Figure 3.3 shows two simple examples of how the MRF model generalizes other models. Here, we see that the full independence model, with properly defined potential functions, gives rise to unigram language modeling and the BM25 model. Similarly, the sequential dependence model subsumes bigram and biterm language models. By studying previous retrieval models in the context of our MRF model, we gain fresh perspective and insight into the underlying principles of these models.

This chapter skirted the issue of parameter estimation (i.e., how to set Λ). Since this issue is critical to achieving good effectiveness, it gets a detailed treatment in the following chapter.

CHAPTER 4

THE THEORY OF LINEAR FEATURE-BASED MODELS

Features lie at the very heart of information retrieval models. The two features term frequency and inverse document frequency form the core of most modern retrieval models, including language modeling and the BM25 model. In most cases, the only difference between any two models is how these elementary features are combined. The resulting scoring functions are often non-linear and contain one or more free parameters that can be tuned in various ways.

In this chapter we describe the theory behind a class of models we call *linear feature-based models*. As its name implies, the scoring functions of these models are comprised of a linear, weighted combination of features. One of the main benefits of such models is their ability to combine many different kinds of features in a straightforward manner. There have been many models proposed that fall under this general framework, including our proposed MRF model. Other models include, but are not limited to, Gey’s logistic regression model [43], Nallapati’s discriminative model [76], and Gao et al.’s linear discriminate model [42]. We return to details of these models in Section 4.2.

Our goal is to synthesize the previous research done and present a general theory of these types of models. This includes investigating what existing scoring functions fit into the framework, characteristics of the underlying parameter space, and how model parameters can be effectively estimated.

This chapter, beyond providing a general theory of linear feature-based models, provides several interesting theoretical insights into our proposed MRF model, and

suggests ways to effectively estimate the model parameters, which we have not yet explained.

4.1 Linear Feature-Based Models

This section provides a general description of linear feature-based models and discusses characteristics of their underlying parameter space.

Suppose we are given a set of documents \mathcal{D} , queries $\mathcal{Q} = \{Q_i\}_{i=1}^N$, and training data \mathcal{T} . In addition, we are given a real-valued scoring function $S_\Lambda(D; Q)$ parameterized by Λ , a vector of parameters. Given a query Q_i , the scoring function $S_\Lambda(D; Q_i)$ is computed for each $D \in \mathcal{D}$ and documents are then ranked in descending order according to their score.

The scoring function induces a total ordering¹ (ranking) $R(\mathcal{D}, Q_i, S_\Lambda)$ on \mathcal{D} for each query Q_i . For simplicity, we rewrite $R(\mathcal{D}, Q_i, S_\Lambda)$ as $R_i(\Lambda)$ and let $\mathcal{R}_\Lambda = \{R_i(\Lambda)\}_{i=1}^N$ be the set of rankings induced over all of the queries.

Finally, in order to evaluate a parameter setting, we need an evaluation function $E(\mathcal{R}_\Lambda; \mathcal{T})$ that produces real valued output given a set of ranked lists and the training data. It should be noted that we require E to only consider the document *rankings* and not the document scores. The scores are only used to rank the documents and not used to evaluate the ranking. This is a standard characteristic among information retrieval evaluation metrics such as mean average precision, precision at 10, among others.

The general framework allows \mathcal{T} to be any type of data that can be used to compute the evaluation metric E over a set of ranked lists. For example, this data may come in the form of TREC relevance judgments or web click-through data [50, 52]. To estimate the parameters we only need to evaluate E , so models may even use

¹We assume ties are broken using some unique document identifier.

abstract concepts in place of \mathcal{T} , such as novelty [46] or aspect precision/recall [121], as long as E remains independent of the document scores.

Using this terminology, the generic information retrieval parameter estimation problem is to find the parameter setting Λ that maximizes the evaluation metric E over the parameter space. This can be formulated as the following optimization problem:

$$\begin{aligned} \hat{\Lambda} &= \arg \max_{\Lambda} E(\mathcal{R}_{\Lambda}; \mathcal{T}) \\ \text{s.t.} \quad &\mathcal{R}_{\Lambda} \sim S_{\Lambda}(D; Q) \\ &\Lambda \in M_{\Lambda} \end{aligned} \tag{4.1}$$

where $\mathcal{R}_{\Lambda} \sim S_{\Lambda}(D; Q)$ denotes that the orderings in \mathcal{R}_{Λ} are induced using scoring function S , and M_{Λ} is the parameter space over Λ .

Solving the general optimization problem proposed above is difficult. Therefore, we aim to solve a more constrained version. We restrict our focus to scoring functions from the following family:

$$\begin{aligned} \mathcal{S} = \{S_{\Lambda}(D; Q) : \exists l(\cdot) \text{ s.t. } l \text{ is strictly monotonically increasing and} \\ l(S_{\Lambda}(D; Q)) = \Lambda^T f(D, Q) + Z\} \end{aligned} \tag{4.2}$$

where $f(\cdot, \cdot)$ maps query/document pairs to real-valued feature vectors in \mathbb{R}^d , Z is a constant that does not depend on D (but may depend on Λ or Q). That is, we require there to exist some strictly monotonically increasing function l that, when applied to S , yields a function that is linear in our parameters Λ . The ranking functions in \mathcal{S} define the universe of linear feature-based models.

Examples of functions within this family include linear discriminants, such as those used with perceptrons, Support Vector Machines (SVMs) [16], and the so-called maximum entropy distribution [32]. In addition, many information retrieval

ranking functions proposed in the past, at their very core, also live within this family, including our proposed MRF model.

By definition, every $S \in \mathcal{S}$ can be reduced to a linear form via a strictly monotonically increasing function. Since such functions are rank preserving and subsequently evaluation metric preserving, we can always write the parameter estimation problem for any scoring function in \mathcal{S} as:

$$\begin{aligned}
 \hat{\Lambda} &= \arg \max_{\Lambda} E(\mathcal{R}_{\Lambda}; \mathcal{T}) \\
 s.t. \quad &\mathcal{R}_{\Lambda} \sim \Lambda^T f(D, Q) + Z \\
 &\Lambda \in M_{\Lambda}
 \end{aligned} \tag{4.3}$$

This general optimization problem fully describes how documents are ranked and how the parameters are estimated. Thus, linear feature-based models are instantiated by choosing an evaluation function E , training data \mathcal{T} , features f , and a parameter space M_{Λ} .

4.2 Previous Uses of Linear Feature-Based Models in IR

Many information retrieval models proposed in the past, at their very core, are linear feature-based models [42, 43, 65, 76]. The models typically differ in their formulation, features, or training. This section briefly summarizes several of these models.

In 1994, Gey proposed a logistic regression model for information retrieval [43]. In terms of our discussion above, the scoring function is in \mathcal{S} after application of the rank-preserving logit transformation and thus is a linear feature-based model. Six features were used in the model. The features were query absolute frequency, query relative frequency, document absolute frequency, document relative frequency, *idf*, and relative frequency in all documents. The maximum likelihood estimate was used

for the parameters. Results showed mixed improvements over a vector space baseline when trained on one collection and tested on another.

In [76], Nallapati argued for a discriminative model for information retrieval, focusing in particular on a SVM formulation. Like Gey, Nallapati also made use of six features. Table 4.1 shows the six features considered. In this case, the parameter vector is estimated by training a linear SVM, with relevant documents considered the “positive class” and non-relevant documents the “negative class”. Therefore, the ranking task is treated as a classification problem. Results were mixed when compared against a language modeling baseline.

Finally, Gao, et al. [42] described a linear discriminant model for information retrieval. The model combined a number of single term and phrase features. The ranking functions are linear feature-based models similar in spirit to those of Gey and Nallapati, but were found to significantly outperform baseline systems more consistently. One critical difference between this approach, and those proposed by Gey and Nallapati, is that the model was trained by directly maximizing a lower bound on mean average precision. As we will show shortly, careful, intelligent parameter estimation is critical to the success of any linear feature-based model.

4.3 Parameter Space

Thus far we have only talked abstractly about the parameter space M_Λ . There are many potential ways to choose a parameter space, each with its advantages and disadvantages. The most obvious choice, to not constrain the parameter space, occurs when $M_\Lambda = \mathbb{R}^d$. The advantage of this model is that the parameter weights can be either negative or positive. This allows the model to include features that convey both negative and positive evidence. On the downside, the search space is somewhat daunting, although not unmanageable.

Another option is to restrict the parameter values to be non-negative. Although this may seem too strict, there are several reasons why it is an acceptable assumption in most cases. In information retrieval, a majority of the features commonly used provide positive evidence. For example, large *tf* or *idf* values are evidence in support of relevance. If a feature is known *a priori* to provide negative evidence, then the feature can still be used, with its value negated. If we do not know if a feature provides positive or negative evidence and we ‘guess’ incorrectly, then the trained model will simply assign that feature a weight of 0. If this occurs, the feature value can be adjusted accordingly and the model retrained.

Therefore, the positivity constraint will typically have only a minor impact on the model. As we will now show, several nice results can be shown to hold under this assumption. Hence, for the remainder of this section, we assume that $M_\Lambda = \mathbb{R}_+^d \stackrel{def}{=} \{\Lambda \in \mathbb{R}^d : \lambda_i \geq 0\}$.

4.3.1 Reduction to Multinomial Manifold

Only considering positive parameter values allows us to map our problem onto a more intuitively appealing space with several nice characteristics. We will now show that the parameter estimation problem previously described, under the positivity constraint, is equivalent to the following constrained optimization problem:

$$\begin{aligned}
 \hat{\Lambda} &= \arg \max_{\Lambda} E(\mathcal{R}_\Lambda; \mathcal{T}) \\
 s.t. \quad &\mathcal{R}_\Lambda \sim \Lambda^T f(D, Q) + Z \\
 &\Lambda \in \mathbb{P}^{d-1}
 \end{aligned} \tag{4.4}$$

where \mathbb{P}^k is a multinomial manifold (also known as a k -simplex) described by:

$$\mathbb{P}^k = \left\{ \Theta \in \mathbb{R}^{k+1} : \forall j \theta_j \geq 0, \sum_{i=1}^{k+1} \theta_i = 1 \right\} \tag{4.5}$$

The multinomial manifold \mathbb{P}^k can be intuitively thought of as the space of all multinomial distributions over $k + 1$ potential outcomes. We now give proof of this equivalence.

Theorem. Any solution to the optimization problem over \mathbb{R}_+^d has a rank-equivalent solution to the optimization over \mathbb{P}^{d-1} .

Proof. Suppose that $\hat{\Lambda} \in \mathbb{R}_+^d$ is the solution to the original optimization problem. Now, consider the following transformation of $\hat{\Lambda}$ to $\hat{\Theta}$:

$$\hat{\theta}_i = \frac{\hat{\lambda}_i}{W} \tag{4.6}$$

where $W = \sum_i \hat{\lambda}_i$. If $W = 0$, then λ_i is mapped to $\theta_i = \frac{1}{d}$ for all i , which is the uniform distribution.

It is easy to see that the transformed parameter setting $\hat{\Theta}$ is in \mathbb{P}^{d-1} , and thus the transformation maps the original point onto the manifold. We must now show this transformation preserves rank-equivalence. Under $\hat{\Theta}$, the scoring function becomes:

$$\begin{aligned} S_{\hat{\Theta}}(D; Q) &= \sum_i \hat{\theta}_i f(D, Q)_i + Z \\ &= \sum_i \left(\frac{\hat{\lambda}_i}{W} \right) f(D, Q)_i + Z \\ &= \frac{1}{W} \sum_i \hat{\lambda}_i f(D, Q)_i + Z \\ &\stackrel{rank}{=} \sum_i \hat{\lambda}_i f(D, Q)_i \end{aligned} \tag{4.7}$$

where the last step follow from the fact that scaling (by $\frac{1}{W}$) and translating (by Z) all scores in the same way has no effect on ranking. Thus, the model under the transformed parameter $\hat{\Theta}$ ranks documents exactly the same as using $\hat{\Lambda}$, the original parameter value. We have therefore shown that any solution to the problem over \mathbb{R}_+^d can be transformed into a rank-equivalent solution over \mathbb{P}^d , thus completing the proof \square

It can be shown that this mapping from \mathbb{R}_+^d to \mathbb{P}^{d-1} is surjective (i.e., many-to-one and onto). This suggests that our original parameter space is inefficient, in that many (in fact, infinitely many) points within \mathbb{R}_+^d produce the same ranking. Within \mathbb{P}^{d-1} , all of these redundant points are conflated to a single point.

4.3.2 Rank Equivalence

We now show that the reduction to the multinomial manifold has an interesting connection with the notion of rank equivalence. Given two parameter settings of a linear, feature-based model, Λ_1 and Λ_2 in \mathbb{R}_+^d , with a fixed set of features, we define the binary relation “ \sim ” as: $\Lambda_1 \sim \Lambda_2$ if and only if Λ_1 and Λ_2 , under the model, are guaranteed to produce the same ranking for all queries. That is, “ \sim ” is the binary relation corresponding to rank equivalence between two parameter settings.

It is easy to see that “ \sim ” is an equivalence relation as it is reflexive ($\Lambda \sim \Lambda$), symmetric ($\Lambda_1 \sim \Lambda_2 \Rightarrow \Lambda_2 \sim \Lambda_1$), and transitive ($\Lambda_1 \sim \Lambda_2 \wedge \Lambda_2 \sim \Lambda_3 \Rightarrow \Lambda_1 \sim \Lambda_3$). It therefore induces equivalence classes over the original Euclidean parameter space. In fact, every parameter on the multinomial manifold corresponds to a *unique* equivalence class. Therefore, the set of parameters on the multinomial manifold can be thought of as *canonical* parameters that can be used to describe any possible parameter setting.

4.3.3 Distance Between Models

Our reduction provides another unique mechanism that is not available in most other retrieval models. The reduction allows us, for a fixed set of features, to quantitatively measure the *distance* between two models (i.e., two parameter settings). For the BM25 retrieval model, there is no straightforward way of measuring the distance between two parameter settings. In language modeling, there exists the notion of distance in terms of KL-divergence between two models [56], but it is not the same

as the distance we can compute here. Instead, we can compute the distance between the actual scoring functions themselves.

The naïve way to measure the distance between two parameter vectors which live in \mathbb{R}^d is to use the Euclidean distance. However, this leads to unappealing results. For example, consider the following two parameter vectors in \mathbb{R}^2 :

$$\Lambda_1 = \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \Lambda_2 = \begin{bmatrix} 2 \\ 6 \end{bmatrix}$$

The Euclidean distance between these two vectors is $\sqrt{10}$, despite the fact that the two parameter settings produce precisely the same ranking. Therefore, the intuitive distance between these two parameter settings is 0. Now, suppose we apply the mapping to the multinomial manifold (\mathbb{P}^1), that was defined above, to these two points. The mapped points can be shown to be:

$$\Theta_1 = \begin{bmatrix} 0.25 \\ 0.75 \end{bmatrix}, \Theta_2 = \begin{bmatrix} 0.25 \\ 0.75 \end{bmatrix}$$

Both of the original parameters are mapped to the same point in \mathbb{P}^1 , and thus they have no distance between them. Now that it is clear the manifold better captures the *intrinsic* properties of the parameters, we still have to answer the question of how to properly measure the distance between arbitrary points on the manifold. Results from information geometry tell us that the multinomial manifold follows a non-Euclidean geometry, and therefore the Euclidean distance does not hold. Thus, we must use a more appropriate distance metric, known as the geodesic distance, which is a generalization of the Euclidean distance to non-Euclidean geometries. The geodesic distance between two points, $\Theta, \Theta' \in \mathbb{P}^d$, is computed as:

$$d(\Theta, \Theta') = 2 \arccos \left(\sum_{i=1}^{d+1} \sqrt{\theta_i \theta'_i} \right) \quad (4.8)$$

where $d(\cdot, \cdot)$ ranges from 0 to π . Although we do not explore specific applications of this distance in this work, a great deal of work considering various properties of the multinomial manifold exists [60, 126]. Future applications may find uses for this unique property.

We also note that the cosine distance is an equally valid measure of the distance between two linear feature-based parameter settings. However, it is less theoretically motivated in terms of the underlying intrinsic geometry of the parameter space.

Finally, it is important to state again that the properties discussed in this section are only valid if the parameters are constrained to be non-negative. We feel the theoretical benefits gained by imposing the constraint outweigh any potential disadvantages.

4.4 Parameter Estimation

In this section we describe how to estimate the parameters of linear feature-based models, which corresponds to solving the optimization proposed earlier in this chapter. This is often called the *learning to rank* problem. The topic has recently gained popularity in the machine learning and information retrieval fields, generating a large number of tools and techniques. While many of these techniques take a purely machine learning approach to the problem, we attempt to focus more on the important information retrieval aspects of the problem. We believe this is important, because critical information retrieval issues are often ignored during the development of new learning to rank techniques. This is somewhat unfortunate, since valuable insights can often be gleaned by attacking these types of problems from an information retrieval perspective (e.g., via feature engineering, failure analysis, etc.) instead of a machine learning perspective.

Many of the linear feature-based models that have been used for information retrieval use parameters that are estimated using maximum likelihood, maximum

a posteriori, or maximum margin techniques. These techniques, however, do not maximize the correct metric. Classification accuracy, likelihood, and margin size are generally of little concern when ranking documents. It may be argued that estimating parameters by maximizing the likelihood of some training data or minimizing classification error is optimizing a function that is correlated with the underlying retrieval metric, such as mean average precision. However, this has been shown to be experimentally invalid, and it can also be shown to be theoretically invalid, as well. This phenomenon, where the metric being optimized diverges from the actual metric of interest, is known as *metric divergence* [74]. Hence, the appropriate way to estimate the parameters of linear feature-based models is to directly solve the optimization problem given in Equation 4.3.

The remainder of this section describes various approaches for solving this optimization problem, including two novel algorithms that perform a direct search over the multinomial manifold. These algorithms can easily handle large training sets, such as those that typically arise in information retrieval. They can also handle the highly unbalanced nature of the training data. Although these techniques have nice properties, they also have several pitfalls, as we will show. We will also briefly describe other learning to rank approaches that have been proposed recently.

4.4.1 Direct Search

We now describe two novel direct search techniques that operate on the multinomial manifold. Here, direct search refers to the fact that the optimization problem is solved in the original metric space. It is important to remember that these metric spaces are typically non-smooth with respect to the parameter settings, which prevents us from applying standard, gradient-based optimization techniques, such as gradient ascent. Our proposed techniques are very simple, yet have been shown to be very effective.

4.4.1.1 Grid Search

The most naïve approach to solving the optimization problem is to perform an exhaustive grid search over the parameter space. That is, we place a grid over the parameter space and evaluate $E(\mathcal{R}_\Lambda; \mathcal{T})$ at every grid intersection, outputting the parameter setting that yields the maximum at the end.

A grid search over \mathbb{R}^d is unbounded and ill-defined. For this reason, we restrict our discussion to the case where our parameter space is the multinomial manifold. For this case, the grid search is bounded and can be easily implemented.

Given a parameter $\epsilon = \frac{1}{K}$ for $K \in \mathbb{Z}^+$ that controls how fine grained our grid is, we define \mathcal{G} , the set of grid points in \mathbb{P}^{d+1} that we search over as:

$$\begin{aligned} \mathcal{G} &= \left\{ \Lambda = (k_1\epsilon \dots k_d\epsilon) : \sum_i k_i\epsilon = 1, k_i \in \mathbb{N} \right\} \\ &= \left\{ \Lambda = (k_1\epsilon \dots k_d\epsilon) : \sum_i k_i = K, k_i \in \mathbb{N} \right\} \end{aligned} \quad (4.9)$$

It is clear from this definition that, $|\mathcal{G}|$, the number of parameter values we must evaluate E at, depends both on d (the number of parameters) and K (how fine grained our grid is). In fact, $|\mathcal{G}| \in \Theta(K^{d-1})$. Therefore, a grid search is feasible only if both d and K are relatively small. For larger values, other training methods that do not require as many sorting operations and metric evaluations must be used. Although the grid search algorithm is relatively costly, it is guaranteed to find a global maximum as K gets large. Algorithm 1 provides a simple implementation of the grid search algorithm.

4.4.1.2 Coordinate Ascent

Coordinate ascent is another technique that can be used to solve non-smooth optimization problems. One of the benefits of the algorithm is that it reduces multivariate search problems into a set of single variable problems, which are often easier to tackle

Algorithm 1 Grid Search

```
1:  $(k_1, k_2, \dots, k_d) \leftarrow (0, 0, \dots, 0)$ 
2:  $\Lambda^* \leftarrow \{\}$ 
3:  $E^* \leftarrow -\infty$ 
4: while  $k_d \leq K$  do
5:    $i \leftarrow 1$ 
6:    $k_i \leftarrow k_i + 1$ 
7:   while  $i < d$  and  $k_i = K$  do
8:      $k_i \leftarrow 0$ 
9:      $i \leftarrow i + 1$ 
10:     $k_i \leftarrow k_i + 1$ 
11:   end while
12:   if  $\sum_{i'} k_{i'} = K$  then
13:      $\Lambda \leftarrow (\frac{k_1}{K}, \frac{k_2}{K}, \dots, \frac{k_d}{K})$ 
14:      $E \leftarrow E(\mathcal{R}_\Lambda; \mathcal{T})$ 
15:     if  $E > E^*$  then
16:        $E^* \leftarrow E$ 
17:        $\Lambda^* \leftarrow \Lambda$ 
18:     end if
19:   end if
20: end while
21: Return  $\Lambda^*$ 
```

in non-smooth spaces. The algorithm first chooses one parameter to be free. It then holds all other parameters as fixed and optimizes the objective function over the single free parameter. This produces an uphill step along one coordinate dimension. This process is repeated for all parameters over a number of iterations. The technique is known to converge slowly on objective functions with long ridges. Variations of the method, including Powell's method, have been proposed to overcome this issue [85].

Coordinate ascent can be applied to the optimization problem under consideration regardless of whether we choose to optimize in the original Euclidean parameter space (\mathbb{R}^d) or the mapped multinomial parameter space (\mathbb{P}^{d-1}). Optimizing over the manifold may be beneficial due to the reduction in the number of repeated local extrema.

If coordinate ascent is performed over the multinomial manifold, then only a minor modification to the original algorithm is necessary. All one dimensional searches done

by the algorithm will be performed as if they were being done in \mathbb{R}^d . However, this does not ensure that the updated parameter estimate will be a point on the manifold. Therefore, after a step is taken in \mathbb{R}^d , we project the point back onto the manifold, which we showed is always possible. Note that this projection preserves the function value since the unnormalized and projected parameter estimates lead to equivalent rankings. Therefore, the optimization is implicitly being done in a space that we know how to optimize over (\mathbb{R}^d), but is continually being projected back onto to the manifold.

Algorithm 2 Coordinate Ascent

- 1: Initialize $\Lambda^0 \leftarrow (\lambda_1^0, \lambda_1^0, \dots, \lambda_d^0)$
 - 2: $t \leftarrow 1$
 - 3: **repeat**
 - 4: $\Lambda^t \leftarrow \Lambda^{t-1}$
 - 5: **for** i from 1 to d **do**
 - 6: $\lambda_i^t \leftarrow \arg \max_{\lambda_i} E(\mathcal{R}_{\Lambda^t}; \mathcal{T})$
 - 7: $\lambda_j^t \leftarrow \frac{\lambda_j^t}{\sum_i \lambda_i^t} \forall j$ (optional)
 - 8: **end for**
 - 9: **until** $|E(\mathcal{R}_{\Lambda^t}; \mathcal{T}) - E(\mathcal{R}_{\Lambda^{t-1}}; \mathcal{T})| > \epsilon$
 - 10: Return Λ^t
-

The coordinate ascent algorithm is given in Algorithm 2. In the algorithm, Λ^t denotes the parameter setting during iteration t . The algorithm first initializes Λ^0 , which can be done uniformly (i.e., $\lambda_i^0 = \frac{1}{d}$), randomly, or in a more informed manner using prior knowledge about the importance of each feature. Then, each λ_i^t is updated according to $\arg \max_{\lambda_i} E(\mathcal{R}_{\Lambda^t}; \mathcal{T})$, which holds all parameter values fixed except λ_i , and finds the setting of λ_i that results in the maximum evaluation metric score. This single dimensional search problem can be solved using a line search. The line search can be done exhaustively or finite difference derivatives can be estimated and used in lieu of exact derivatives. Of course, if E is partially differentiable with respect to each parameter then it may be possible to solve the single dimensional search problem

analytically. After λ_i^t is updated, the updated parameter vector Λ^t is optionally projected back onto the multinomial manifold.

The algorithm continues until the evaluation metric does not change more than ϵ between subsequent iterations. If ϵ is set to 0, then the algorithm stops only when no more uphill moves are possible. Setting $\epsilon > 0$ may cause the algorithm to return a solution that is not maximal. However, it also results in faster training and can help minimize overfitting.

Our implementation of the algorithm uses a line search to solve the single dimensional search problem, projects updated parameters onto the multinomial manifold, and uses $\epsilon = 0.0001$. After extensive evaluation, it was found that these settings resulted in the most well-behaved configuration of the algorithm.

4.4.1.3 Discussion

Finding the maximum of an arbitrary evaluation function E using direct search can be very difficult and error-prone, especially in high-dimensional space. Only a grid search method, with a suitably chosen granularity, is guaranteed to find a global maxima. Coordinate ascent is a local search technique that only finds a global maxima if the evaluation function E is concave. Our experiments using this approach, show that, for a certain set of term and phrase features, mean average precision is approximately concave over a wide range of collections. This may be the case for many related applications and feature sets, but is not true in general, as was pointed out in [42]. For functions with many local maxima, a multiple random restart strategy can be used to increase the chances of finding a global solution.

Another potential disadvantage of the direct search techniques presented here is the fact that they are fully supervised. If little or no training data exists, then unsupervised and active learning [98] techniques from machine learning can potentially be employed. However, such methods are out of the scope of the current work.

Despite these disadvantages, our approach has the advantage that it can make use of all of the training data and does not suffer in the face of highly unbalanced training data. When training using maximum likelihood or SVMs, it is often important to have balanced training data. However, in information retrieval it is very often the case that there are many more relevant documents compared to non-relevant documents for a given query. For this reason, the training data is very unbalanced. Nallapati found that the data needed to be balanced in order to achieve good generalization performance [76]. Balancing was done by undersampling the majority (non-relevant) class. Although this led to improved performance over the unbalanced case, it had the negative effect of throwing away valuable training data. Other solutions to the unbalanced data problem for SVMs exist that do not require training data to be compromised, such as allowing separate costs for training errors in the positive and negative classes [75]. Since our coordinate ascent approach does not make any implicit or explicit assumptions about the underlying distribution, we can use the training data in its entirety.

4.4.2 Optimization Using Surrogate Functions

A number of other estimation techniques have been proposed for training linear feature-based models for information retrieval. The common thread among these approaches is that they do not attempt to directly search in the evaluation metric space. Instead, they typically optimize some *smooth* surrogate function that is related to the evaluation metric space. By using a smooth, typically convex, surrogate function, it is possible to apply standard optimization machinery to the problem. However, it is not always the case that the optimum of the surrogate is equal to the optimum of the actual metric. For example, the surrogate function may be a lower or upper bound on the metric. Although these surrogates may yield reasonable estimates, many exhibit metric divergence. In certain cases, however, the optimum of the surrogate is equal

to the optimum of the evaluation metric, thereby eliminating any metric divergence. In the remainder of this section we summarize the various techniques proposed and describe the properties of their surrogate functions.

4.4.2.1 Perceptron Learning

Gao et al. proposed using a perceptron-based algorithm to optimize mean average precision [42]. The technique uses pairwise preferences as training data [117]. That is, training data, denoted by \mathcal{R} , comes in the form of tuples of the form $(d_i, d_j)_q$, which indicates that document i should be ranked higher than document j for query q . There are various ways to derive such preferences from manual relevance judgments or click-through data [50]. The perceptron learning algorithm is demonstrated in Algorithm 3. The perceptron approach, like our proposed coordinate ascent algo-

Algorithm 3 Perceptron Learning

```

1: Initialize  $\Lambda^0 \leftarrow (1, 0, \dots, 0)$ 
2: for  $t$  from 1 to  $MAX\_ITERATIONS$  do
3:   for each  $(d_i, d_j)_q \in \mathcal{R}$  do
4:     if  $\Lambda^T f(d_j, q) > \Lambda^T f(d_i, q)$  then
5:        $\lambda_i \leftarrow \lambda_i + \eta(f(d_i, q) - f(d_j, q))$ 
6:     end if
7:   end for
8: end for
9: Return  $\Lambda$ 

```

rithm, is not guaranteed to find a global maxima. Instead, the perceptron learning algorithm optimizes a lower bound on mean average precision. Therefore, metric divergence may be a problem. In addition, this technique can not easily be applied to other evaluation metrics. Despite this, the algorithm has been shown to produce reasonable effectiveness over a wide range of data sets.

4.4.2.2 RankNet

Another approach, based on neural networks, called RankNet, has recently been proposed [15]. A RankNet is trained using gradient descent over the following differ-

entiable cost function:

$$C(\mathcal{Q}, \mathcal{R}) = \sum_{q \in \mathcal{Q}} \sum_{(i,j) \in \mathcal{R}} -\hat{P}_{q,i,j} \log P_{q,i,j} - (1 - \hat{P}_{q,i,j}) \log(1 - P_{q,i,j}) \quad (4.10)$$

where \mathcal{Q} is the set of queries, \mathcal{R} is the set of pairwise preferences used for training, $P_{q,i,j} = \frac{1}{1 + \exp[s(q,d_j) - s(q,d_i)]}$ is the probability that document i ranks higher than document j for query q under the current neural network ($s(q, d_i)$ is the score of document d for query q as computed by the neural network), and $\hat{P}_{q,i,j}$ is the target probability of document i being ranked higher than j (obtained from training data).

The gradients of this smooth cost function can be computed analytically, which makes it easy to apply gradient descent to solve the optimization problem. However, the model suffers from standard neural network training issues, such as local minima. In addition, the cost function, which is simply the cross entropy between the target distribution and the distribution modeled using the neural network, does not minimize (or maximize) any specific retrieval metric. Therefore, it is vulnerable to metric divergence. Several recent studies have attempted to address this problem [17, 62].

We also note that RankNet is a linear feature-based model when the underlying neural network has no hidden layers. The underlying ranking function, in the presence of hidden layers, may be highly non-linear.

4.4.2.3 Support Vector Machine Optimization

Finally, Joachims proposed a large margin training technique for multivariate performance measures [50, 51]. The technique uses a surrogate objective function based on SVM using structured outputs that can be solved using quadratic programming. The approach can maximize a variety of information retrieval metrics, such as precision at k , precision-recall breakeven, and area under the ROC curve. In fact, any metric that can be computed based solely on a contingency table can be maximized efficiently. One particularly nice property of these metrics is that the maximum found

is the actual maximum, and therefore, there is no metric divergence. However, this is not the case for any arbitrary metric. Recently, approaches based on this work have been proposed to optimize lower bounds of average precision [120] and nDCG [59]. Furthermore, since the method is based on SVMs, it is easy to employ the “kernel trick” and implicitly project inputs into a higher, possibly infinite, dimensional space.

4.4.2.4 Discussion

There are two downsides to these types of approaches. First, they specifically work for one metric or a family of metrics. Second, many of them suffer from metric divergence, even though the resulting optimization problems are easier and more efficient to solve using well-established optimization techniques. The grid search and coordinate ascent algorithms, however, do not suffer from either of these problem.

There currently is no well developed understanding of best practices for estimating the parameters of linear feature-based information retrieval models. Most studies have looked at traditional machine learning problems, which typically differ from information retrieval tasks. Therefore, an interesting, and necessary, direction of future work is to undertake a comprehensive evaluation of these techniques, in terms of how effective they are across a wide range of retrieval data sets and metrics, how well they generalize, and how efficient they are. Another critical question that must be answered is whether or not non-linear models, such as those that arise when using neural networks or SVMs with non-linear kernels, are more effective than simple linear models.

4.5 Quantifying the Impact of Metric Divergence

In this section, we run several simple experiments against a number of *ad hoc* retrieval data sets. These experiments are chosen to empirically show that training by directly maximizing mean average precision is superior to training using a naïve

SVM approach. This allows us to determine, quantitatively, how significant the metric divergence problem really is.

The goal of *ad hoc* retrieval is to retrieve as many relevant documents as high in the ranked list as possible. Relevance is binary (relevant/non-relevant) and is assessed according to whether the document is topically related to the query or not. More details of *ad hoc* retrieval will be provided in the following chapter. For the experiments in this chapter, we use mean average precision as our primary evaluation metric. Please refer to Appendix B for further details on mean average precision.

We compare the effectiveness of three types of models. The first is a linear feature-based model trained by directly maximizing mean average precision using the coordinate ascent approach proposed earlier in this chapter. The second is Nallapati’s SVM model [76], which treats relevant documents as the positive class and non-relevant documents as the negative class and learns a (linear) SVM classifier. The model was discussed in detail earlier in this chapter. Both an unbalanced and balanced version of the model are investigated. The last model is a language modeling system [84], which provides a baseline and sanity check for the effectiveness of the learned models. Three standard newswire TREC collections are used. For each collection, 50 topics are used for training and 50 for testing. Only the title portion of the TREC topics are used. A summary of the collections used and the training and test topics are given in Table A.2.

For these experiments, E = mean average precision, \mathcal{T} = TREC relevance judgments, f = set of 6 features shown in Table 4.1, and $M_\Lambda = \mathbb{R}^d$. Each model is trained using only the data in the relevance judgments. That is, when the model is being trained it only “knows” about the documents contained in the relevance judgments and not about any of the unjudged documents in the collection. However, when being tested, *all* documents in the collection are ranked. In the case of the balanced SVM model, the non-relevant judgments from the relevance file were undersampled. The

Number	Feature
1	$\sum_{w \in Q \cap D} \log(tf_{w,D})$
2	$\sum_{w \in Q \cap D} \log(1 + \frac{tf_{w,D}}{ D })$
3	$\sum_{w \in Q \cap D} \log(\frac{N}{df_w})$
4	$\sum_{w \in Q \cap D} \log(\frac{ C }{cf_w})$
5	$\sum_{w \in Q \cap D} \log(1 + \frac{tf_{w,D}}{ D } \frac{N}{df_w})$
6	$\sum_{w \in Q \cap D} \log(1 + \frac{tf_{w,D}}{ D } \frac{ C }{cf_w})$

Table 4.1. Features used in the bag of words experiments. $tf_{w,D}$ is the number of times term w occurs in document D , cf_w is the number of times term w occurs in the entire collection, df_w is the number of documents term w occurs in, $|D|$ is the length (in terms) of document D , $|C|$ is the length (in terms) of the collection, and N is the number of documents in the collection.

	Disks 1,2		Disk 3		Disks 4,5	
	Train	Test	Train	Test	Train	Test
Unbalanced SVM	0.0955	0.1091	0.1501	0.1336	0.1421	0.1434
Balanced SVM	0.1577	0.1849	0.1615	0.1361	0.1671	0.1897
Coordinate Ascent	0.1955†‡	0.2327†‡	0.2080†‡	0.1773‡	0.2238†‡	0.2328†‡
Language modeling	0.1883‡	0.2155‡	0.1875‡	0.1642‡	0.1819	0.1995

Table 4.2. Training and test set mean average precision values for various *ad hoc* retrieval data sets and training methods. The † represents a statistically significant improvement over language modeling and ‡ denotes significant improvement over the balanced SVM model. Tests done using a one tailed paired t-test at the 95% confidence level.

feature-based model is trained using the same features as the SVM, and therefore has no additional power. The language modeling system ranks documents via query likelihood, with document models estimated using Bayesian (Dirichlet) smoothing [123] and is trained by finding the smoothing parameter that maximizes the mean average precision on the training data.

The results of the experiments are given in Table 4.2. As we see from the results, our parameter estimation technique based on coordinate ascent results in consistently and statistically significant improvements over both the unbalanced and balanced SVM estimates. Furthermore, it significantly outperforms language modeling on 4

out of 6 runs. Language modeling, on the other hand, significantly outperforms the SVM model on 4 out of the 6 runs.

The results indicate that language modeling, despite its simplicity, stands up very well compared to sophisticated feature-based machine learning techniques. The results also provide empirical proof that SVM parameter estimation is simply not the correct paradigm here, because it optimizes the wrong evaluation metric function. Our estimation technique, however, is directly maximizing the evaluation metric under consideration and results in stable, effective parameter estimates across the collections. Therefore, metric divergence can significantly hinder the effectiveness of a model.

CHAPTER 5

EVALUATION OF THE BASIC MRF MODEL

Our discussion, up until this point, has focused entirely on the theoretical issues surrounding the Markov random field model. We now shift our focus to more practical matters. In this chapter, we empirically evaluate the retrieval effectiveness of our proposed MRF model. This requires us to choose one or more tasks to evaluate the model against. There are a large number of important information retrieval tasks, such as web search [10], enterprise search [25], question answering (QA) [112], blog search [79], legal search [5], desktop search [83], and image search. Rather than evaluating our model on all of these tasks, we restrict our focus to *ad hoc* retrieval and web search. As we will describe in more detail shortly, these two tasks are the the most common and widely used in information retrieval.

This chapter has three primary contributions. First, we define the *basic MRF models for IR*. These models are our first attempts at using the MRF framework for retrieval. These models are purposefully simple and designed to show the robustness and functionality of the MRF framework. Later on, in Chapters 6 and 7, we describe extensions of these basic models. Second, we compare the effectiveness of these basic MRF models to the state of the art bag of words models. Finally, we do a detailed analysis of various aspects of the basic models in order to develop a richer, more complete understanding of the strengths and weaknesses of the framework.

```
<top>
<num> Number: 744

<title>
Counterfeit ID punishments

<desc> Description:
What punishments or sentences have been given in the U.S. for
making or selling counterfeit IDs?

<narr> Narrative:
Relevant documents will describe punishments for manufacturing or
selling counterfeit identification, such as drivers licenses,
passports, social security cards, etc. Fake professional
certifications and fake credit cards are relevant. Counterfeit goods
or auto serial numbers not relevant. Counterfeit checks are not
relevant. "Counterfeiting" with no indication of type is relevant.

</top>
```

Figure 5.1. TREC topic number 744.

5.1 Ad Hoc Retrieval

Ad hoc retrieval is one of the most important information retrieval tasks. In the task, a user submits a query, and the system returns a ranked list of documents that are *topically* relevant to the query. Therefore, the goal of the task is to find topically relevant documents in response to a query. It is critical to develop highly effective *ad hoc* retrieval models since such models often play important roles in other retrieval tasks. For example, most QA systems use an *ad hoc* retrieval system to procure documents that are topically relevant to some question. The QA systems then employ various techniques to extract answers from the document retrieved [112]. Thus, by improving on the current state of the art *ad hoc* retrieval models, it is possible to positively impact the effectiveness of a wide range of tasks.

In the remainder of this section we describe experiments using three different basic MRF models. Our aim is to analyze and compare the retrieval effectiveness of each

model across collections of varying size and type. We make use of the AP, WSJ, and ROBUST04 data sets, which are smaller collections that consist of news articles that are mostly homogeneous, and two web data sets, WT10g and GOV2, which are considerably larger and less homogeneous. Further details about the data sets are provided in Appendix A.

Each of these are TREC data sets. A TREC data set consists of a collection of documents, a set of topics, and human relevance assessments. An example *ad hoc* topic is shown in Figure 5.1. A TREC topic typically consists of a title, description, and narrative. It is important to note that a topic is not the same thing as a query, although the two terms are often used interchangeably. A query must be distilled from a topic. This is typically done by using the text contained in one or more of the topic fields as the query. For all of the experiments in this section, except where noted otherwise, we follow the common TREC procedure of using only the title portion of the topic as our query. Therefore, the query that we distill for the topic given in Figure 5.1 is *counterfeit id punishments*.

TREC relevance judgments are done by human assessors. When determining relevance, the entire TREC topic is considered. The assessors judge documents using a binary¹ scale, where rating 0 indicates not relevant and rating 1 indicates relevant.

All of the evaluation metrics that we consider in this section are based on binary judgments. For all of our experiments, we return a ranked list of no more than 1000 documents per query, as is standardly done during TREC evaluations. Furthermore, the primary evaluation metric that we use to evaluate *ad hoc* retrieval is mean average precision. Further details about the retrieval metrics we use can be found in Appendix B.

¹Although some TREC collections actually do have ternary (i.e., not relevant, relevant, and highly relevant) judgments, they have never been used during official evaluations. When ternary judgments do exist, all relevant (rating 1) and highly relevant (rating 2) documents are considered relevant, which thereby binarizes the judgments.

Finally, for all of our experiments, documents were stemmed using the Porter stemmer and a standard list of 418 stopwords was applied. All model parameters were estimated by maximizing mean average precision using our proposed coordinate ascent algorithm (see Algorithm 2).

Throughout all of our experiments, statistical significance is always tested using a one-tailed, paired t -test at significance level $p < 0.05$.

5.1.1 MRF Models for Ad Hoc Retrieval

We now define the three basic MRF models for the *ad hoc* retrieval task. The models correspond to the three dependence model types shown in Figure 3.2. Each model represents a different set of underlying dependence assumptions and makes use of different features. We define each model in terms of its canonical form, provide its joint probability mass function, and show its ranking function.

5.1.1.1 Full Independence

The first basic model that we consider makes use of the full independence model shown in Figure 3.2 (left). Recall that, under this model, query term nodes are independent of each other given a document as evidence. This model, therefore, shares many properties with standard bag of words retrieval models.

We now introduce the first basic MRF model, which we call the *Full Independence MRF Model* (MRF-FI model). It is constructed using the following canonical form:

$$(\text{FI}, T_{QD}, \text{LM}) : \lambda_{T_D}$$

$$(\text{FI}, T_Q, \text{ICF}) : \lambda_{T_Q}$$

The model defines two features. One feature is defined over the T_{QD} clique set and the other is defined over the T_Q clique set. Both features use the full independence assumption and language modeling features.

Notice that no feature is defined over D , the document node clique set. By not defining a feature over the document node clique, we are enforcing the constraint that documents, in isolation, provide no useful information for the *ad hoc* retrieval task. While this may seem like an extreme assumption, it is actually quite valid. No single document prior has ever been shown to significantly improve effectiveness across a wide range of data sets. Therefore, in order to keep the model as simple as possible, we simply do not define a feature over this clique. However, we note that for specific tasks it may be beneficial to define such a feature.

The model results in the following joint probability mass function:

$$P_{G,\Lambda}(Q, D) = Z^{-1} \exp \left[\lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \lambda_{T_Q} \sum_{q_i \in T_Q} \log \frac{|C|}{cf_{q_i}} \right] \quad (5.1)$$

Furthermore, it is easy to see that we obtain the following linear feature-based model when ranking according to $P(D|Q)$:

$$P_{G,\Lambda}(D|Q) \stackrel{rank}{=} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} \quad (5.2)$$

which shows that the MRF-FI model reduces exactly to the unigram query likelihood language modeling approach with Dirichlet smoothing (see Equation 2.14).

Although the MRF-FI model is not technically a bag of words model, we consider it as our bag of words baseline. This is appropriate, since, as we just showed, the model is rank equivalent to the unigram language modeling approach, which is a bag of words model. Therefore, we use the MRF-FI model as a baseline by which we compare other MRF models that actually go beyond the bag of words assumption.

Table 5.1 shows the test set results for the MRF-FI model across data sets. In the table, MAP refers to mean average precision, GMAP is geometric mean average

	AP	WSJ	ROBUST04	WT10G	GOV2
MAP	0.2077	0.3258	0.2920	0.1861	0.2984
GMAP	0.1219	0.2267	0.1970	0.1176	0.1891
P@10	0.3460	0.4860	0.4293	0.3204	0.5180
R-Prec	0.2448	0.3558	0.3291	0.2199	0.3515
μ^t	1750	2000	1000	1000	1500

Table 5.1. Test set results for the MRF-FI model.

precision, P@10 is precision at 10 ranked documents, R-Precision is precision at R (number of judged relevant documents), and μ^t denotes the smoothing parameter learned on the training set. All models were trained to maximize mean average precision. These numbers serve as our baselines, which we attempt to improve upon by employing more complex models.

5.1.1.2 Sequential Dependence

The second of the basic MRF models corresponds to the sequential dependence model shown in Figure 3.2 (center). It is the *Sequential Dependence MRF Model* (MRF-SD model), which is constructed according to the following canonical form:

$$(\text{FI}, T_{QD}, \text{LM}) : \lambda_{T_D}$$

$$(\text{FI}, T_Q, \text{ICF}) : \lambda_{T_Q}$$

$$(\text{SD}, O_{QD}, \text{LM-O-1}) : \lambda_{O_D}$$

$$(\text{SD}, O_Q, \text{ICF-O-1}) : \lambda_{O_Q}$$

$$(\text{SD}, O_{QD}, \text{LM-U-4}) : \lambda_{U_D}$$

$$(\text{SD}, O_Q, \text{ICF-U-4}) : \lambda_{U_Q}$$

which defines features over single term (i.e., T_{QD} and T_Q) clique sets, as well as ordered term clique sets (i.e., O_{QD} and O_Q). Unlike the MRF-SI model, this model makes use of some of the MRF model’s strengths. As we see, the model defines ordered and unordered window features over the ordered cliques in the graph. By doing so, we

go beyond the bag of words assumption. The joint probability mass function for the model is:

$$\begin{aligned}
P_{G,\Lambda}(Q, D) = Z^{-1} \exp \left[\right. & \lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \\
& \lambda_{T_Q} \sum_{q_i \in T_Q} \log \frac{|C|}{cf_{q_i}} + \\
& \lambda_{O_D} \sum_{(q_1, q_2, D) \in O_{QD}} \log \frac{tf_{\#1(q_1 q_2), D} + \mu^w \frac{cf_{\#1(q_1 q_2)}}{|C|}}{|D| + \mu^w} + \\
& \lambda_{O_Q} \sum_{(q_1, q_2) \in O_Q} \log \frac{|C|}{cf_{\#1(q_1 q_2)}} + \\
& \lambda_{U_D} \sum_{(q_1, q_2, D) \in U_{QD}} \log \frac{tf_{\#uw8(q_1 q_2), D} + \mu^w \frac{cf_{\#uw8(q_1 q_2)}}{|C|}}{|D| + \mu^w} + \\
& \left. \lambda_{U_Q} \sum_{(q_1, q_2) \in U_Q} \log \frac{|C|}{cf_{\#uw8(q_1 q_2)}} \right] \tag{5.3}
\end{aligned}$$

and the ranking function simplifies to the following linear feature-based model:

$$\begin{aligned}
P_{G,\Lambda}(D|Q) \stackrel{rank}{=} & \lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \\
& \lambda_{O_D} \sum_{(q_1, q_2, D) \in O_{QD}} \log \frac{tf_{\#1(q_1 q_2), D} + \mu^w \frac{cf_{\#1(q_1 q_2)}}{|C|}}{|D| + \mu^w} + \\
& \lambda_{U_D} \sum_{(q_1, q_2, D) \in U_{QD}} \log \frac{tf_{\#uw8(q_1 q_2), D} + \mu^w \frac{cf_{\#uw8(q_1 q_2)}}{|C|}}{|D| + \mu^w} \tag{5.4}
\end{aligned}$$

Recall that both the ordered (LM-O- M) and unordered (LM-U- N) features have free parameters that allow the size of the unordered window (scope of proximity) to vary. We now motivate our choices for selecting these specific values.

First, for M , the parameter that controls the ordered window matching, we decided to use 1, because it results in an “exact phrase” feature that does not allow any room in the ordered matching of query terms. This choice is motivated by the

N	AP	WSJ	WT10g	GOV2
2	0.1860	0.2776	0.2148	0.2697
8	0.1867	0.2763	0.2167	0.2832
50	0.1858	0.2766	0.2154	0.2817
Unlimited	0.1857	0.2759	0.2138	0.2714

Table 5.2. Mean average precision for various parameter settings for LM-U- N using the MRF-SD model.

fact that exact phrases are commonly used in many different applications. Furthermore, there has been little previous research that looked at relaxing such phrases. Therefore, choosing 1 is the most reasonable choice.

The other value, N , which controls the window width for unordered matching, was chosen after careful consideration of previous research. Fagan shows that the best choice of N varies across collections [35]. Optimal values found included setting N to either 2, the length of a sentence, or “unlimited” (matches any co-occurrences of the terms within a document). Croft et al. showed improvements could be achieved with passage-sized windows of 50 terms [29]. Therefore, since there were no strong conclusions, we experimented with window sizes of 2, 50, sentence, and “unlimited” to see what impact each had on effectiveness. Instead of segmenting sentences at index time, we observe that the average length of an English sentence is 8-15 terms, and choose a window size of 8 terms to model sentence-level proximity.

The results, which were evaluated on the entire data set, are given in Table 5.2. The results show very little difference across the various window sizes. However, for the AP, WT10g, and GOV2 collection the sentence-sized windows performed the best. For the WSJ collection, $N = 1$ performed the best. The only collection where mean average precision varies noticeably is the GOV2 collection. These results suggest that a limited scope of proximity (2-50 terms) performs reasonably, but can be approximated rather well by an “unlimited” scope, which reaffirms past research into dependence models based on co-occurrences. However, it appears as though

	AP	WSJ	ROBUST04	WT10G	GOV2
MAP	0.2147†	0.3425	0.3096†	0.2053†	0.3325†
GMAP	0.1265	0.2399†	0.2196†	0.1286†	0.2449†
P@10	0.3340	0.5080	0.4566†	0.3245	0.5680†
R-Prec	0.2580†	0.3633	0.3363	0.2374†	0.3716†

Table 5.3. Test set results for the MRF-SD model. A † indicates a statistically significant improvement over the MRF-FI model.

smaller scopes of proximity may provide better performance for larger collections, as evidenced by the GOV2 results. Therefore, given this experimental evidence, we decide to set $N = 4$ for use with our basic MRF-SD model.

Now that we have describe why the rationale behind the manual construction of our model, we must see how well it performs compared to the simple MRF-FI model. The results are given in Table 5.3. Results that are statistically significantly better than the MRF-FI model are indicated by a †.

The results show that the MRF-SD model is significantly better than the MRF-FI model on every data set except WSJ for mean average precision, which is our primary evaluation metric. The improvements in mean average precision are 3.4% for AP, 5.1% for WSJ, 6.0% for ROBUST04, 10.3% for WT10G, and 11.4% for GOV2. These results indicate very strong, consistent improvements over the bag of words baseline.

Similar results are exhibited for geometric mean average precision. GMAP heavily penalizes queries with a low average precision. Therefore, GMAP is often used to measure robustness [114]. As our results show, the MRF-SD model is quite robust and significantly improves GMAP for every data set except AP. We do a deeper analysis of the robustness of the MRF model later in this chapter.

We see that the precision at 10 is improved across most data sets, but is only significant on two of them (ROBUST04 and GOV2). Therefore, it appears as though most of the boost in mean average precision that is achieved from using the MRF-SD

model does not come from the very top of the ranked list. Instead, the improvement is likely coming from lower in the ranked list, where the ordered and unordered window features are bringing in more relevant documents and filtering out many of the low ranked, poorly matching documents.

Finally, it is important to recall that training is done to maximize mean average precision. It is likely that more significant improvements could be achieved for the other metrics if the model were trained to optimize them.

5.1.1.3 Full Dependence

The third basic MRF model is derived from the full dependence model. The model, which is shown in Figure 3.2 (right), is called the *Full Dependence MRF Model* (MRF-FD model). The model attempts to incorporate dependencies between every subset of query terms and is the most general of the basic models. Here, the number of cliques is exponential in the number of query terms, which restricts the application of this variant to shorter queries. This is not a problem for the MRF-FI and MRF-SD models, which have a linear number of cliques. The model is constructed according to the following canonical form:

$$\begin{aligned}
 (\text{FI}, T_{QD}, \text{LM}) &: \lambda_{T_D} \\
 (\text{FI}, T_Q, \text{ICF}) &: \lambda_{T_Q} \\
 (\text{FD}, O_{QD}, \text{LM-O-1}) &: \lambda_{O_D} \\
 (\text{FD}, O_Q, \text{ICF-O-1}) &: \lambda_{O_Q} \\
 (\text{FD}, O_{QD}, \text{LM-U-4}) &: \lambda_{U_D} \\
 (\text{FD}, O_Q, \text{ICF-U-4}) &: \lambda_{U_Q} \\
 (\text{FD}, U_{QD}, \text{LM-U-4}) &: \lambda_{U_D} \\
 (\text{FD}, U_Q, \text{ICF-U-4}) &: \lambda_{U_Q}
 \end{aligned}$$

which is similar to the MRF-SD model. However, the models differ in several key ways. First, the MRF-FD model uses the full dependence model type. Second, the

MRF-FD model defines two new features for the unordered clique sets (U_{QD} and U_Q). These clique sets are empty in the MRF-SD model. Furthermore, the parameters for all the unordered features are tied together. While this is not required, it simplifies the model.

The resulting joint probability mass function for the model is then given by:

$$\begin{aligned}
P_{G,\Lambda}(Q, D) = Z^{-1} \exp \left[\right. & \lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \\
& \lambda_{T_Q} \sum_{q_i \in T_Q} \log \frac{|C|}{cf_{q_i}} + \\
& \lambda_{O_D} \sum_{(q_1, \dots, q_k, D) \in O_{QD}} \log \frac{tf_{\#1(q_1 \dots q_k), D} + \mu^w \frac{cf_{\#1(q_1 \dots q_k)}}{|C|}}{|D| + \mu^w} + \\
& \lambda_{O_Q} \sum_{(q_1, \dots, q_k) \in O_Q} \log \frac{|C|}{cf_{\#1(q_1 \dots q_k)}} + \\
& \lambda_{U_D} \sum_{(q_1, \dots, q_k, D) \in O_{QD} \cup U_{QD}} \log \frac{tf_{\#uw8(q_1 \dots q_k), D} + \mu^w \frac{cf_{\#uw8(q_1 \dots q_k)}}{|C|}}{|D| + \mu^w} + \\
& \left. \lambda_{U_Q} \sum_{(q_1, \dots, q_k) \in O_Q \cup U_Q} \log \frac{|C|}{cf_{\#uw8(q_1 \dots q_k)}} \right] \quad (5.5)
\end{aligned}$$

and the resulting ranking function is then:

$$\begin{aligned}
P_{G,\Lambda}(D|Q) \stackrel{rank}{=} & \lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \\
& \lambda_{O_D} \sum_{(q_1, \dots, q_k, D) \in O_{QD}} \log \frac{tf_{\#1(q_1 \dots q_k), D} + \mu^w \frac{cf_{\#1(q_1 \dots q_k)}}{|C|}}{|D| + \mu^w} + \\
& \lambda_{U_D} \sum_{(q_1, \dots, q_k, D) \in O_{QD} \cup U_{QD}} \log \frac{tf_{\#uw4k(q_1 \dots q_k), D} + \mu^w \frac{cf_{\#uw4k(q_1 \dots q_k)}}{|C|}}{|D| + \mu^w} \quad (5.6)
\end{aligned}$$

Now that we have defined the MRF-FD model, we would like to understand what effect the ordered and unordered features have on the model's effectiveness and how well the models learned on one collection generalize to another. In order to measure

	Term + Ordered				Term + Unordered			
Train \ Test	AP	WSJ	WT10g	GOV2	AP	WSJ	WT10g	GOV2
AP	0.1847	0.2718	0.2180	0.2669	0.1840	0.2674	0.2179	0.2754
WSJ	0.1842	0.2733	0.2167	0.2613	0.1840	0.2674	0.2179	0.2754
WT10G	0.1847	0.2718	0.2180	0.2669	0.1838	0.2674	0.2189	0.2783
GOV2	0.1840	0.2705	0.2150	0.2675	0.1838	0.2674	0.2189	0.2783

	Term + Ordered + Unordered			
Train \ Test	AP	WSJ	WT10g	GOV2
AP	0.1866	0.2716	0.2226	0.2839
WSJ	0.1841	0.2738	0.2195	0.2694
WT10G	0.1865	0.2719	0.2231	0.2783
GOV2	0.1852	0.2709	0.2201	0.2844

Table 5.4. Mean average precision using the MRF-FD model over different combinations of term, ordered, and unordered features.

this, we train on one data set and then use the parameter values found to test on the other data sets. Results for models trained using terms and ordered features, terms and unordered features, and terms, ordered, and unordered features are given in Table 5.4.

For the AP collection, there is very little difference between using ordered and unordered features. However, there is a marginal increase when both ordered and unordered features are used together. The results for the WSJ collection are different. For that collection, the ordered features produce a clear improvement over the unordered features, but there is very little difference between using ordered features and the combination of ordered and unordered. The results for the two web collections, WT10g and GOV2, are similar. In both, unordered features perform better than ordered features, but the combination of both ordered and unordered features led to noticeable improvements in mean average precision.

From these results we can conclude that strict matching via ordered window features is more important for the smaller newswire collections. This may be due to the homogeneous, clean nature of the documents, where an ordered window match is likely to be a high quality match instead of noise. For the web collections, the

	AP	WSJ	ROBUST04	WT10G	GOV2
MAP	0.2128	0.3429†	0.3092†	0.2140†‡	0.3360†
GMAP	0.1257	0.2404†	0.2196†	0.1361†‡	0.2421†
P@10	0.3540	0.5080†	0.45605†	0.3469†‡	0.5720†
R-Prec	0.2543†	0.3694†	0.3394†	0.2417†‡	0.3763†

Table 5.5. Test set results for the MRF-FD model. A † indicates a statistically significant improvement over the MRF-FI model and a ‡ indicates statistically significant improvement over the MRF-SD model.

opposite is true. Here, the fuzzy unordered window matches provide better evidence. In these less homogeneous, noisy collections, an ordered window match is less likely to be a high quality match and more likely to be a noisy match. Instead, fuzzy matches are appropriate because they deal better with the noise inherent in web documents.

These results also suggest that parameters trained on any of the data sets generalize well to other data sets. This result is somewhat surprising; we expected parameters trained on newswire (web) data would generalize better to newswire (web) test data. However, this is not the case. It appears as though the parameters trained on any reasonable data set will generalize well, which allows one to use a single setting of the parameters across multiple data sets. This may imply that the features used here only capture general aspects of the text and that more domain-specific features may yield further improvements. We return to the issue of parameter generalization later in this chapter.

We conclude our discussion of the MRF-FD model by reporting test set effectiveness results. The results are given in Table 5.5. The improvements over the MRF-FI model are highly consistent, even more so than the improvements we saw with the MRF-SD model. Consistent and significant improvements in mean average precision and geometric mean average precision are observed on every data set except AP. Furthermore, both precision at 10 and R-prec are consistently improved across nearly all of the data sets, as well.

These results indicate that the MRF-FD model is better at improving precision at the top of the ranked list. This suggests that modeling dependencies between non-adjacent query terms, via the use of ordered and unordered features, enhances precision more so than modeling dependencies between adjacent query terms. By using the full dependence model, we impose a more global (i.e., across all query terms) type of proximity constraint on the query terms, whereas the sequential dependence model imposes more of a local (i.e., only between adjacent query terms) proximity constraint. Hence, the MRF-FD model promotes documents where all of the query terms occur within a close proximity to each other, and the MRF-SD model only promotes documents based on the proximity of pairs of adjacent query terms. The MRF-SD model, therefore, may result in lower quality matches that do not satisfy the global proximity constraints imposed by the MRF-FD model, which may lead to fewer relevant documents returned at the top of the ranked list.

Despite the fact that the MRF-SD model only enforces local proximity constraints, it is only significantly worse than the MRF-FD model on the WT10G data set. The two models are statistically indistinguishable for all other metrics and data sets.

This is an interesting result with practical ramifications. If a system builder had to decide whether to implement the MRF-SD model or the MRF-FD model, they would need to analyze this efficiency/effectiveness tradeoff closely. Our results show that, statistically, there is often no difference between the two models. However, as we showed, the MRF-FD model does tend to produce better results across all data sets and metrics. In terms of efficiency, the MRF-SD model requires less computation in order to rank documents, since there are only a linear number of cliques. The MRF-FD model, on the other hand, has an exponential number of cliques. Therefore, the key practical factors to consider are average query length, importance of excellent effectiveness, and computational resources.

Recent advances in inverted indexing technology and query evaluation may be able to significantly improve the efficiency by which both MRF-SD and MRF-FD models can be evaluated. These new techniques, based on impact ordered indexes, pre-compute complicated features and store them directly in the index [3]. Then, rather than computing an exponential number of feature functions per query, the aggregated feature value can be read directly from the index. Of course, applying such an indexing strategy requires a large amount of disk space to store the “feature lists”, but could result in very fast query evaluation, especially using recently developed query optimization techniques [4, 106].

5.1.2 Evaluation

In this section, we delve deeper into a number of issues related to the three basic MRF models. By analyzing these issues, we develop a better understanding of the MRF model. Many of the insights described here can be widely applied to other information retrieval models.

5.1.2.1 Smoothing

All three of the basic models have one or more model hyperparameters that control smoothing. These parameters live outside of the MRF model and must be tuned separately. Previous research has shown that language modeling effectiveness is often sensitive to the setting of the smoothing parameters [123]. Therefore, it is important to consider how sensitive the effectiveness of the basic models are to the setting of the hyperparameters.

There are two different hyperparameters associated with the basic models. They are μ^t , which controls single term smoothing, and μ^w , which controls both ordered and unordered window smoothing. We choose to smooth single terms different from windows because terms tend to behave differently than windows and have different occurrence statistics. Although not explored here, it is also possible to smooth the

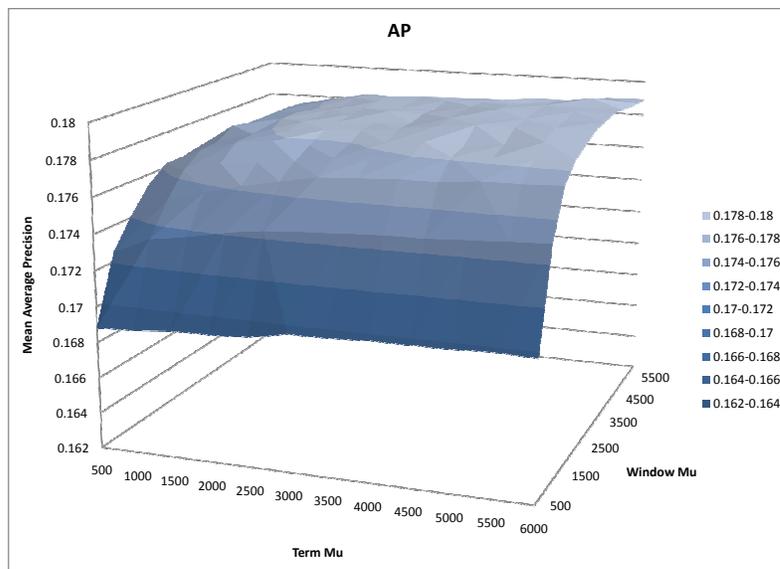


Figure 5.2. Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the AP data set.

ordered window features differently than the unordered windows. However, we feel that the two window types are similar enough that they can be smoothed in the same way.

Since nobody has ever applied smoothing to ordered and unordered windows in this way, it is important to analyze how sensitive effectiveness is with regard to the window smoothing parameters. In Figures 5.2, 5.3, 5.4, and 5.5 we plot mean average precision surfaces over a wide range of settings for μ^t and μ^w using the MRF-SD model. Results are given for the AP, WSJ, ROBUST04, and WT10G data set.

The surfaces show that, in general, effectiveness is more sensitive to the setting of the window smoothing parameter (μ^w) than the term smoothing parameter (μ^t). The results suggest that it is important to tune the smoothing parameters, especially the window smoothing parameter. These surfaces also support our decision to smooth windows and terms differently, as it is apparent that setting $\mu^t = \mu^w$ is often far from the optimal setting.

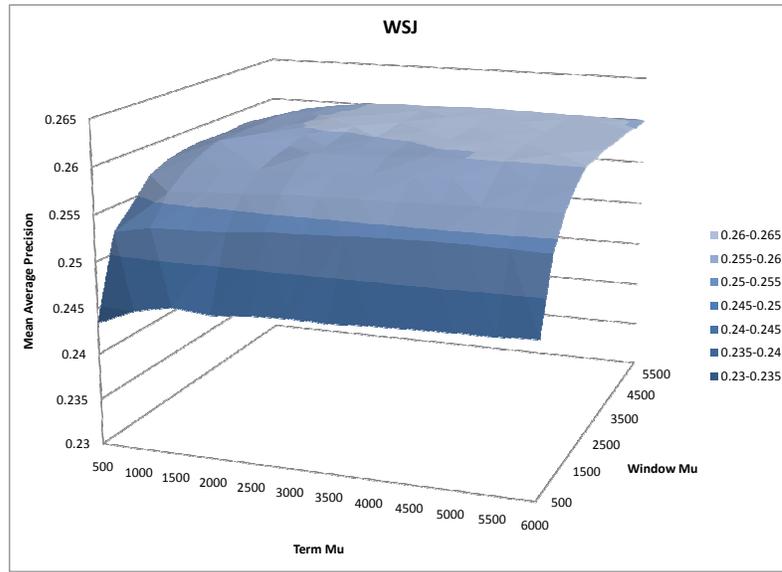


Figure 5.3. Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the WSJ data set.

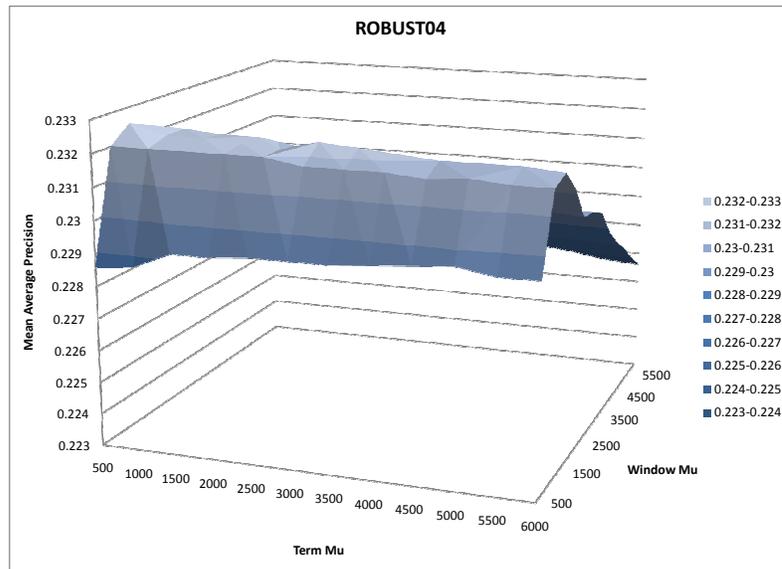


Figure 5.4. Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the ROBUST04 data set.

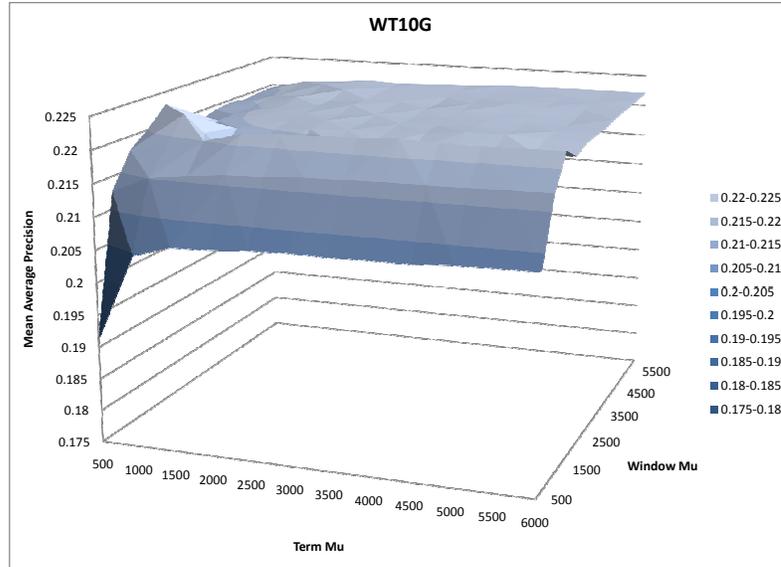


Figure 5.5. Training set mean average precision as a function of term and window smoothing parameters using the sequential dependence model on the WT10G data set.

Finally, we note that the AP, WSJ, and WT10G curves are shaped similarly. However, the ROBUST04 surface has a very distinct shape to it. The difference appears to be that it is difficult to “saturate” the window smoothing parameter on the AP, WSJ, and WT10G data sets, but that the window smoothing parameter quickly saturates on the ROBUST04 collection. This result may have to do with the fact that the ROBUST04 query set was specifically chosen to be difficult (for retrieval systems built before 2004). It may be that these queries were “hard” because the models that were applied to them did not take term proximity into account [12, 113]. Therefore, when we apply our model to these queries, it may be possible to over smooth the window features, which reduces the effect of term proximity on the ranking function and decreases effectiveness.

5.1.2.2 Collection Size

In Chapter 1, we described the various paradigm shifts that have occurred in information retrieval. We argued that as collection sizes grew and average document

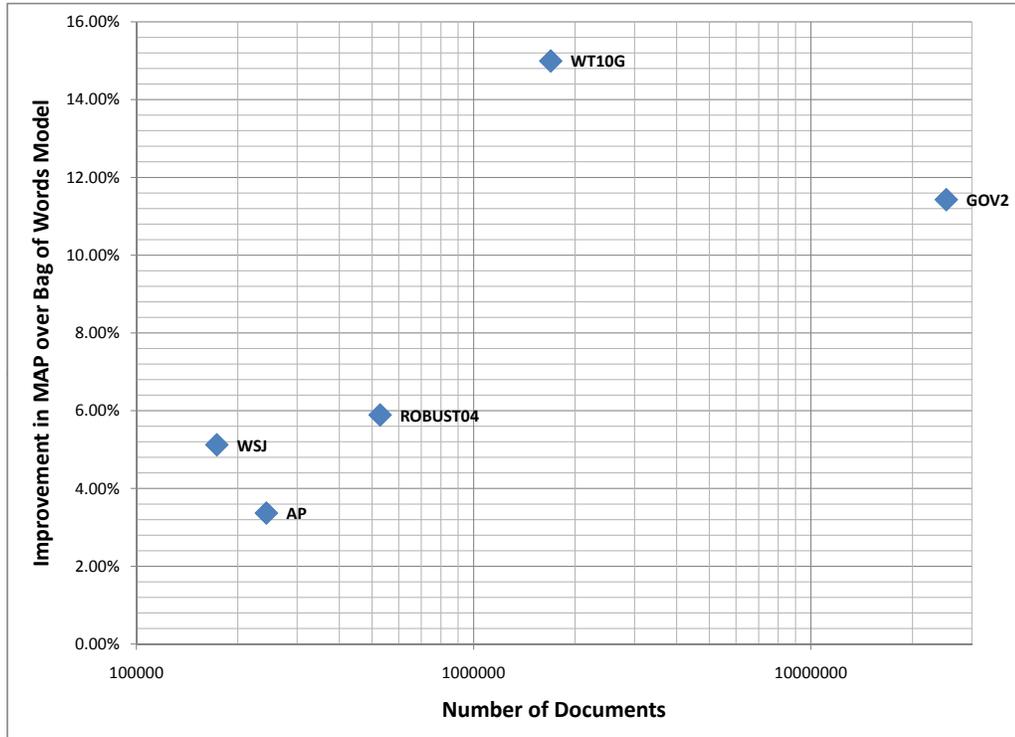


Figure 5.6. Relationship between the number of documents in a collection and the relative improvement in mean average precision of the MRF-FD model over unigram language modeling (MRF-FI). Note that the x -axis is log scaled.

lengths increased, new types of features beyond term frequency and inverse document frequency would become important. We argued that these new types of features that go beyond bags of words would act to filter out all of the noisy matches that were made by chance.

We test this hypothesis by analyzing how much better the MRF-FD model is compared to the MRF-FI model across a range of data set sizes. This data is plotted in Figure 5.6. In the Figure, the x -axis represents the number of documents in the collection (log scale) and the y -axis represents the relative improvement in mean average precision of MRF-FD over MRF-FI.

Although there are only 5 data points, there is clearly an increasing trend in the data, with bigger improvements seen for the larger collections. The trend, of course, is not perfect, but it does help validate our hypothesis that bag of words features

begin to fail as collection sizes increase. It will be interesting to see whether this trend continues as larger data sets are made available. Clearly, bag of words models, such as language modeling or BM25, are not well suited for *ad hoc* retrieval against web-scale collections.

5.1.2.3 The Role of Features

We have just shown that going beyond the bag of words assumption and making use of term proximity features is highly effective, especially on very large collections. In this section, we take this analysis one step further and investigate the role of various types of features across the data sets. This analysis provides insights into why certain features are more effective than others on a given data set and helps us understand which other types of features may be important in the future.

To aid our analysis, we compute statistics for various information retrieval features. The statistics are computed in the following manner. For every query, we compute the feature of interest for every document in \mathcal{D} , the document set of interest. The average feature value is then computed across all of the documents in \mathcal{D} . We then use the median of the average values as our primary statistic. We use the median, rather than the average, because many of the feature distributions are highly skewed. This procedure is carried out for the following features:

- **Overlap** – $\frac{|Q \cap D|}{|Q|}$, which is the fraction of query terms that occur in the document. If this value is 1, then every query term occurs in the document. We note that this is similar in spirit to Buckley et al.’s *titlestat* measure [13].
- **Average TF** – $\frac{\sum_{w \in Q} tf_{w,D}}{|Q|}$, which is the average term frequency of the query terms in the document.
- **Average Distance (Sequential)** – The average distance (with respect to term positions) between every pair of query terms that are adjacent to each other. Single term queries are ignored when computing this feature.

	Overlap		Avg. TF		Avg. Dist. (Seq)]		Avg. Dist. (Tot)	
	Rel	Nonrel	Rel	Nonrel	Rel	Nonrel	Rel	Nonrel
AP	0.63	0.52	2.6	1.84	212	207	215	210
WSJ	0.67	0.53	2.7	1.93	370	460	387	471
ROBUST04	0.67	0.53	2.9	2.42	430	4933	452	5998
WT10G	0.80	0.66	6.5	5.16	1389	10357	1414	10357
GOV2	0.94	0.76	18.6	18.79	7090	7826	7164	8017

Table 5.6. Median values for various statistics computed across judged relevant (Rel) and non-relevant (Nonrel) documents.

- **Average Distance (Total)** – The average distance between *every* pair of query terms. Again, single term queries are ignored when computing this feature.

These features are meant to capture various bag of words features (overlap and average TF), as well as notions of term proximity (average distances).

The medians, as computed using the procedure described above, are given in Table 5.6. Results are given for the AP, WSJ, ROBUST04, WT10G, and GOV2 data sets. The statistics are computed for both the set of judged relevant documents and the set of judged non-relevant documents. By comparing the median values of these features in both sets, we are able to better understand which features discriminate well between relevant and non-relevant documents.

Of course, the judged relevant and judged non-relevant documents are heavily biased because of the pooling procedure used at TREC. However, these statistics still provide valuable insights into the fine line between relevant and non-relevant documents and what types of features are important for data sets with varying characteristics.

We first analyze the overlap feature. As the results show, the overlap is higher in the relevant set than in the non-relevant set. This is to be expected, as relevant documents typically contain most of the query terms. However, there is a noticeable increasing trend in the value as the collection size increases. This suggests that as collections get larger, relevant documents that appear high in the ranked list (i.e.,

those that would get pooled and judged) will contain most, if not all, of the query terms. This suggests that it might be useful to run the query as a simple conjunctive Boolean query first, and then apply a more complex ranking function to the filtered set of documents.

A similar trend exists for the average term frequency feature, with larger average TF values for larger collections that contain longer documents. This, again, should not be surprising, since collections that contain longer documents will naturally contain more term occurrences. Furthermore, since many of the judgment pools include a large number of runs that use bag of words models based on *tf.idf* scoring, it is only natural for the results to be biased towards high term frequencies. The gap in TF from the relevant to the non-relevant set is not very large, and therefore average TF alone cannot be used as a very good discriminator. Indeed, it is very interesting to observe that the median average TF of non-relevant documents for GOV2 is larger than the median average TF of relevant documents. This suggests that many of the bag of words models return documents that contain many chance occurrences of the query terms. Since these are not meaningful occurrences of the query terms, the documents were actually non-relevant. This is where the term proximity and other types of features begin to become important, as we will now show.

The two term proximity features show the greatest discriminative potential of any of the features we looked at, especially as collection sizes grow. For the AP and WSJ collections, there is little difference between the term proximity features in the relevant and non-relevant sets. However, for the ROBUST04, WT10G, and GOV2 data sets, there is a noticeable divide between the term proximity characteristics of the relevant and non-relevant document sets. The biggest divide occurs for the WT10G data set, which, not surprisingly, showed the biggest boost in effectiveness when the MRF-SD and MRF-FD models were used. These statistics validate our arguments as to the importance of term proximity features, especially on larger collections.

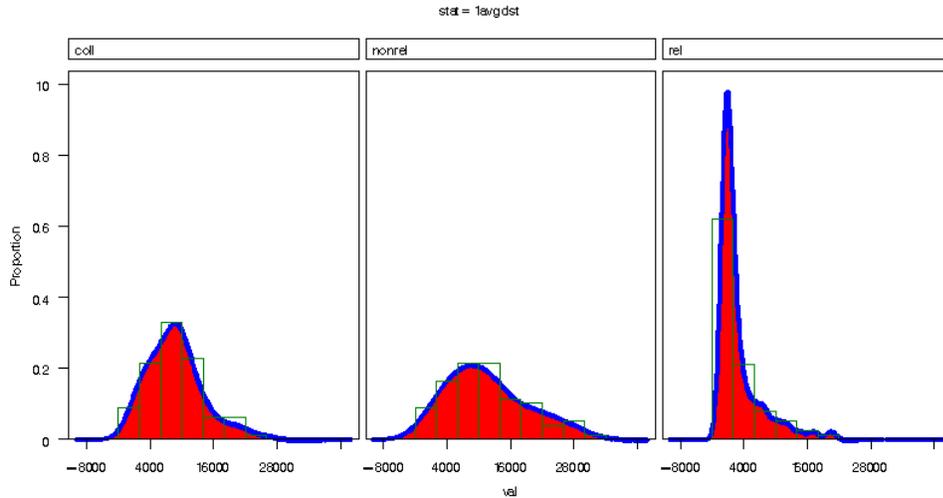


Figure 5.7. Plot of average distance between sequential query terms for WT10G data set. Coll represents the entire collection, nonrel the set of judged non-relevant documents, and rel the set of judged relevant documents.

They show that both local proximity, as modeled by the MRF-SD model, as well as global proximity, as modeled by the MRF-FD model, actually model discriminative characteristics of query terms that discriminate between relevant and non-relevant documents.

Finally, in order to give an idea of the distribution of the average distance (sequential) feature, we plot a histogram and smoothed density estimate for the WT10g data set in Figure 5.7. In addition to the statistics about the judged relevant and non-relevant documents, the same statistics were also computed for the entire set of documents. The entire set of documents is a much more realistic, less biased model of “not relevant” than the judged non-relevant documents. As the figure shows, the relevant distribution is highly skewed towards small values. The non-relevant distribution is skewed, but not as much as the relevant distribution. The collection distribution does not appear to be as skewed as the other distributions, but it is clear that the average distances compute across the collection are generally much larger than both the judged relevant and non-relevant documents.

5.1.2.4 Robustness

Next, we investigate the robustness of the MRF-SD and MRF-FD models. Here, we define robustness as the number queries whose effectiveness is improved/hurt (and by how much) as the result of applying these methods. A highly robust model will significantly improve many queries over the baseline and only minimally hurt a few.

In order to evaluate the robustness of the MRF-SD and MRF-FD models, we plot histograms that show how many queries were helped or hurt by a given amount. These plots are given in Figure 5.8. The bin labels indicate the relative change in mean average precision with regard to the baseline MRF-FI model. We see from the results that the distributions are skewed in the direction of positive improvements. In fact, for most of the data sets, there are very few queries that are hurt by more than 50%. Similarly, many queries are often improved by over 50% on every data set.

These histograms are only useful for aggregate data analysis. However, we would like to know which queries were the most helped and most hurt by these more complex models. In Tables 5.7 and 5.8 we provide the 10 most improved and 10 most hurt queries when using the MRF-SD model on the ROBUST04 and GOV2 data sets, respectively.

The first observation we make about these results is that the most improved queries are often those that have poor MRF-FI average precision (e.g., below 0.1). Of course, since these queries are so poor, it is very easy to achieve large relative improvements. However, some queries, such as *price fixing* (ROBUST04 topic 622), *big dig pork* (GOV2 topic 835), and *spanish civil war support* (GOV2 topic 829) are significantly improved and have “acceptable” MRF-SD average precision values. There are very few cases of queries with large MRF-FI average precisions being significantly hurt when the MRF-SD model is applied to them.

The second observation is that queries consisting of meaningful, common two word phrases, such as *gasoline tax* (ROBUST04 topic 700), *price fixing* (ROBUST04 topic

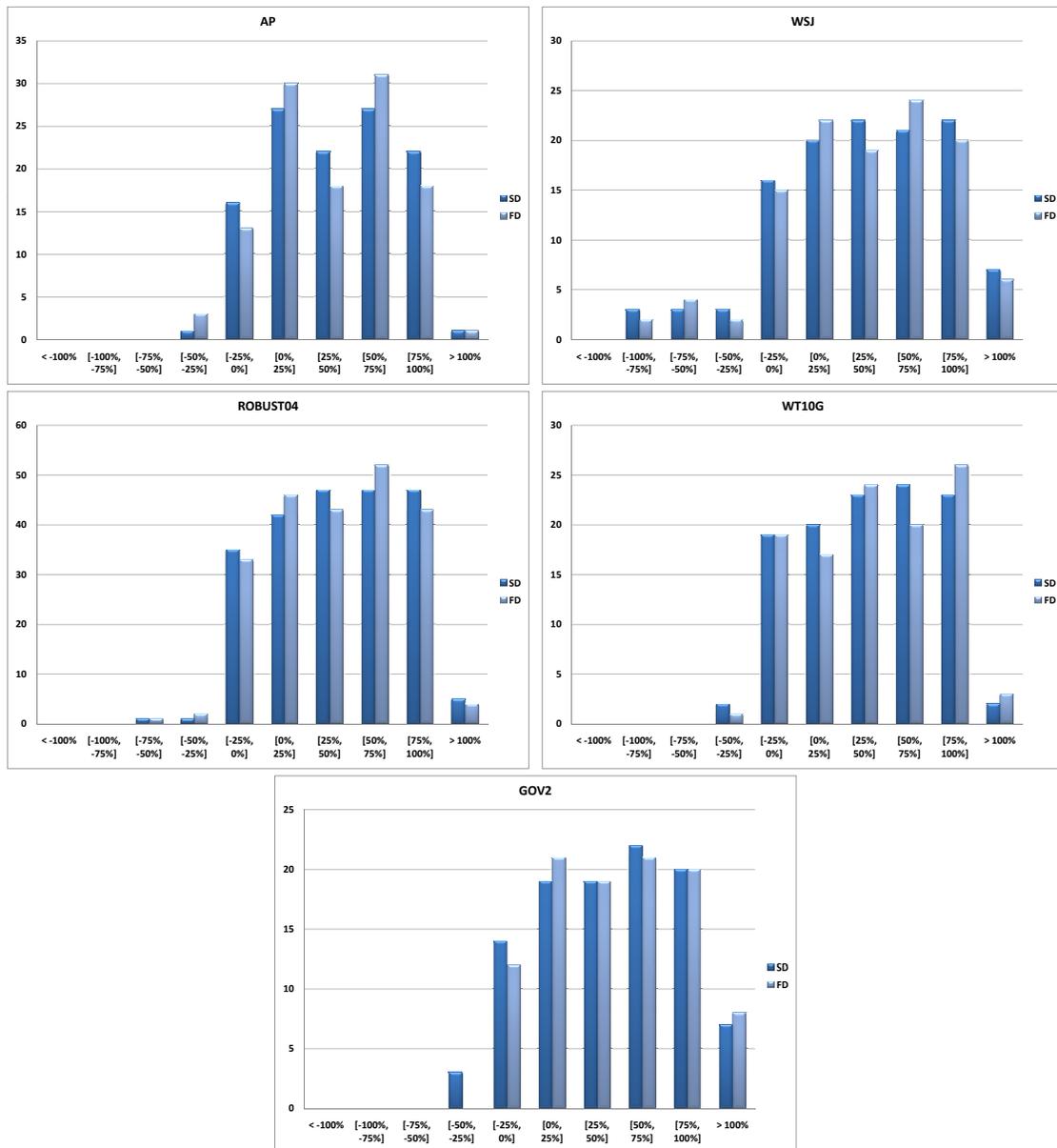


Figure 5.8. Robustness of MRF-SD and MRF-FD models for the AP, WSJ, ROBUST04, WT10G, and GOV2 test sets. The MRF-FI model is used as the baseline by which the improvements are computed. The evaluation metric used is average precision.

Topic	Query	MRF-FI	MRF-SD	% Change
615	timber exports asia	0.1477	0.0697	-52.81%
685	oscar winner selection	0.2689	0.1740	-35.29%
623	toxic chemical weapon	0.2982	0.2238	-24.95%
651	ethnic population	0.0381	0.0305	-19.95%
659	cruise health safety	0.3216	0.2589	-19.50%
644	exotic animals import	0.1719	0.1392	-19.02%
698	literacy rates africa	0.5278	0.4314	-18.26%
612	tibet protesters	0.4528	0.3733	-17.56%
695	white collar crime sentence	0.2746	0.2316	-15.66%
693	newspapers electronic media	0.3108	0.2680	-13.77%
...				
639	consumer line shopping	0.1353	0.2094	54.77%
629	abortion clinic attack	0.1568	0.2527	61.16%
700	gasoline tax	0.2811	0.4579	62.90%
684	part time benefits	0.0881	0.1482	68.22%
627	russian food crisis	0.0102	0.0175	71.57%
638	wrongful convictions	0.0231	0.0468	102.60%
690	college education advantage	0.0027	0.0062	129.63%
689	family planning aid	0.0224	0.0583	160.27%
666	thatcher resignation impact	0.0078	0.0333	326.92%
622	price fixing	0.0287	0.1354	371.78%

Table 5.7. The 10 most improved and 10 most hurt test set queries when using the MRF-SD model on the ROBUST04 data set. Effectiveness is measure in terms of average precision.

Topic	Query	MRF-FI	MRF-SD	% Change
822	custers stand	0.0815	0.0562	-31.04%
833	iceland government	0.5150	0.3669	-28.76%
847	portugal world war ii	0.2513	0.1859	-26.02%
837	eskimo history	0.0622	0.0486	-21.86%
838	urban suburban coyotes	0.2778	0.2402	-13.53%
846	heredity obesity	0.1734	0.1517	-12.51%
816	usaid assistance galapagos	0.7751	0.7056	-8.97%
850	mississippi river flood	0.1799	0.1663	-7.56%
808	north korean counterfeiting	0.7212	0.6717	-6.86%
845	new jersey tomato	0.3892	0.3635	-6.60%
...				
825	national guard involvement iraq	0.0755	0.1183	56.69%
843	pol pot	0.3158	0.5012	58.71%
849	scalable vector graphics	0.2080	0.4029	93.70%
842	david mccullough	0.0806	0.1799	123.20%
805	identity theft passport	0.0412	0.0961	133.25%
844	segmental duplications	0.0543	0.1486	173.66%
830	model railroads	0.0327	0.0992	203.36%
829	spanish civil war support	0.0637	0.2183	242.70%
835	big dig pork	0.0593	0.2141	261.05%
806	doctors borders	0.0061	0.0750	1129.51%

Table 5.8. The 10 most improved and 10 most hurt test set queries when using the MRF-SD model on the GOV2 data set. Effectiveness is measure in terms of average precision.

622), *pol pot* (GOV2 topic 843), and *model railroads* (GOV2 topic 830), are more likely to be improved than two word phrases that are not as common or meaningful, such as *ethnic population* (ROBUST04 topic 651), *tibet protesters* (ROBUST04 topic 612), *eskimo history* (GOV2 topic 837), and *heredity obesity* (GOV2 topic 846). This may be the result of how ordered and unordered feature weights are computed in the MRF-SD model. It may be useful in the future to include notions of lexical cohesiveness in the computation of ordered and unordered phrase features in order to rectify this issue [111].

Finally, we note that parsing/stopword removal errors may have contributed to the reduction in effectiveness observed for some of the queries. One clear example of such an error is the query *custers stand* (GOV2 topic 822). The original title is *custers's last stand*. However, our query distillation process removes a large set of stopwords, including the term *last*. When the MRF-SD model is applied to the query *custers stand*, the exact phrase feature becomes a very poor feature, since the actual exact phrase features that we want are *custers last* and *last stand*, instead of *custers stand*. Interestingly, the query *doctors borders* (GOV2 topic 806), which is originally *doctors without borders* is the most improved query for the GOV2 data set. A better understanding of stopword removal from within phrases is needed in order to deal with these cases in a more consistent manner.

5.1.2.5 Long Queries

Up until this point, we have only considered queries that were constructed from the title portion of the TREC topics. These queries tend to be very short, high quality queries that contain few, if any, function words. In this section, we examine whether or not the MRF model maintains its effectiveness on long queries. In particular, we are interested in evaluating the model on queries constructed from the description portion of the TREC topic. The description field contains a longer natural language

description of the underlying information need and often contains more useful terms. However, it also includes many function words. For this reason, we developed a special stopword list that contains common function words that often occur in TREC description fields. This stopword list is then applied to queries in order to remove most of the noisy query terms, while keeping the important content terms.

In order to evaluate long queries, we define two new basic MRF models designed specifically for long queries. The first, *MRF-FI-L*, is a variant of the MRF-FI model. The *LM* weighting function is replaced with another language modeling weighting function based on Jelinek-Mercer smoothing. This is done because previous research has suggested that Jelinek-Mercer smoothing is more effective on longer queries than Dirichlet smoothing [123]. This small change results in the following ranking function:

$$P_{G,\Lambda}(D|Q) \stackrel{rank}{=} \sum_{(q_i,D) \in T_{QD}} \log \left[(1 - \delta^t) \frac{tf_{q_i,D}}{|D|} + \delta^t \frac{cf_{q_i}}{|C|} \right] \quad (5.7)$$

where δ^t is the term smoothing parameter. In a similar fashion, we modify the MRF-SD model to use Jelinek-Mercer smoothing instead of Dirichlet smoothing. This model has this ranking function:

$$\begin{aligned} P_{G,\Lambda}(D|Q) \stackrel{rank}{=} & \lambda_{T_D} \sum_{(q_i,D) \in T_{QD}} \log \left[(1 - \delta^t) \frac{tf_{q_i,D}}{|D|} + \delta^t \frac{cf_{q_i}}{|C|} \right] + \\ & \lambda_{O_D} \sum_{(q_1,q_2,D) \in O_{QD}} \log \left[(1 - \delta^w) \frac{tf_{\#1(q_1q_2),D}}{|D|} + \delta^w \frac{cf_{\#1(q_1q_2)}}{|C|} \right] + \\ & \lambda_{U_D} \sum_{(q_1,q_2,D) \in U_{QD}} \log \left[(1 - \delta^w) \frac{tf_{\#uw8(q_1q_2),D}}{|D|} + \delta^w \frac{cf_{\#uw8(q_1q_2)}}{|C|} \right] \end{aligned} \quad (5.8)$$

where δ^w is the ordered and unordered window smoothing parameter. Both δ^t and δ^w must be in the range $[0, 1]$.

The long query results are given in Table 5.9. These results show that the MRF-SD-L model significantly outperforms the MRF-FI-L model on all data sets. In fact,

	MRF-FI-L	MRF-SD-L
AP	0.1778	0.1956 (+10.0%)
WSJ	0.2395	0.2544 (+6.2%)
ROBUST04	0.2910	0.3120 (+7.2%)
WT10G	0.1288	0.1639 (+27.3%)
GOV2	0.2003	0.2412 (+20.4%)

Table 5.9. Test set mean average precision for description-length queries using full and sequential dependence models. All improvements are statistically significant.

the improvements here are much larger than those observed for the short queries, which signifies that modeling dependencies between adjacent query terms is even more important for longer queries. One potential reason for this behavior is the fact that longer queries often contain many more spurious terms that, when matched in a bag of words setting, will return many poor documents. Instead, when the adjacency constraint is enforced, the number of these poor matches is reduced.

These results also indicate that the longer queries are much less effective than the shorter version of the queries (see Table 5.3). This poor effectiveness is likely caused by the increased noise in the query. Although many function words are removed from the queries during pre-processing, some make it into the query. Similar results have been observed at TREC, as well [113, 114]. This brings up the question of whether or not more information, in the form of longer natural language queries, should be expected to return better results than short keyword queries. If the search engine were replaced by a librarian, then it is obvious that the more information that you were to provide, the better the results would ultimately be. However, natural language processing techniques have failed to have a positive effect on retrieval effectiveness, especially for longer queries. Perhaps once natural language techniques are improved, it may be reasonable to expect better effectiveness from longer queries. However, it remains to be seen whether or not users of information retrieval systems will be willing to enter long descriptive queries. Instead, the most likely answer in some technology

that lies between short keyword queries and fully descriptive queries. For example, a user interface that allows users to enter a keyword query and then provides a set of simple options to focus their search results. Such technologies are starting to show up in web search engines, but whether or not they will enter the mainstream remains yet to be seen.

5.1.2.6 BM25 Weighting

All of the basic MRFs proposed so far have used language modeling weighting functions. As described in Chapter 2, the BM25 weighting function has been shown to have effectiveness comparable to language modeling. For this reason, we are interested in examining the effectiveness of MRF models built using BM25 weighting functions. It is straightforward to modify one of the previously proposed MRF models to use BM25 weighting. In order to keep things relatively simple and provide for an easy comparison, we choose to modify the MRF-SD model. The resulting MRF-BM25 model is then given by the following canonical form:

$$\begin{aligned}
 (\text{FI}, T_{QD}, \text{BM25}) &: \lambda_{T_D} \\
 (\text{FI}, T_Q, \text{IDF}) &: \lambda_{T_Q} \\
 (\text{SD}, O_{QD}, \text{BM25-O-1}) &: \lambda_{O_D} \\
 (\text{SD}, O_Q, \text{IDF-O-1}) &: \lambda_{O_Q} \\
 (\text{SD}, O_{QD}, \text{BM25-U-4}) &: \lambda_{U_D} \\
 (\text{SD}, O_Q, \text{IDF-U-4}) &: \lambda_{U_Q}
 \end{aligned}$$

where the weighting functions are defined in Section 3.3.3. The MRF-BM25 model has the same form as the MRF-SD model, but replaces the LM and ICF weighting functions with analogous BM25 and IDF ones. This results in the following ranking function:

	AP	WSJ	ROBUST04	WT10G	GOV2
MAP	0.2210†	0.3512†‡	0.3101†‡	0.2129†‡	0.3476†‡
GMAP	0.1366†*	0.2471†‡	0.2199†‡	0.1181‡	0.2817† ‡ *
P@10	0.3140	0.5140†	0.4525†	0.3388	0.6100† ‡ *
R-Prec	0.2666†	0.3698†	0.3366‡	0.2508†‡	0.3834†
(k_1^t, b^t)	(1.75, 0.3)	(1.5, 0.3)	(0.5, 0.3)	(0.5, 0.2)	(1.0, 0.4)
(k_1^w, b^w)	(0.25, 0.1)	(0.25, 0.1)	(0.25, 0.0)	(0.25, 0.0)	(0.25, 0.0)

Table 5.10. Test set results for the MRF-BM25 model. The †, ‡, and * indicate statistically significant improvements over the MRF-FI, BM25 and MRF-SD models, respectively. Recommended term and window hyperparameter values are also provided.

$$\begin{aligned}
P_{G,\Lambda}(D|Q) \stackrel{\text{rank}}{=} & \\
\lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} & \frac{(k_1^t + 1)tf_{w,D}}{k_1^t \left((1 - b^t) + b^t \frac{|D|}{|D|_{\text{avg}}} \right) + tf_{w,D}} \log \frac{N - df_w + 0.5}{df_w + 0.5} + \\
\lambda_{O_D} \sum_{(q_1, q_2, D) \in O_{QD}} & \frac{(k_1^w + 1)tf_{\#1(q_1 q_2), D}}{k_1^w \left((1 - b^w) + b^w \frac{|D|}{|D|_{\text{avg}}} \right) + tf_{\#1(q_1 q_2), D}} \log \frac{N - df_{\#1(q_1 q_2)} + 0.5}{df_{\#1(q_1 q_2)} + 0.5} + \\
\lambda_{U_D} \sum_{(q_1, q_2, D) \in U_{QD}} & \frac{(k_1^w + 1)tf_{\#\text{uw}8(q_1 q_2), D}}{k_1^w \left((1 - b^w) + b^w \frac{|D|}{|D|_{\text{avg}}} \right) + tf_{\#\text{uw}8(q_1 q_2), D}} \log \frac{N - df_{\#\text{uw}8(q_1 q_2)} + 0.5}{df_{\#\text{uw}8(q_1 q_2)} + 0.5}
\end{aligned} \tag{5.9}$$

which has four hyperparameters, as opposed to the two hyperparameters in the MRF-SD model. While the two extra parameters make the model more flexible, to a certain extent, it also makes it more difficult to properly tune.

The results of our experiments using the MRF-BM25 model are shown in Table 5.10. Test set results are given and significance tests are done comparing the retrieval effectiveness against MRF-FI, BM25 (see Equation 2.11) and MRF-SD. According to our primary evaluation metric, mean average precision, we see that the MRF-BM25 is always significantly better than the MRF-FI model and significantly better than BM25 on all data sets, except AP. Furthermore, the MRF-BM25 model is statistically indistinguishable from the MRF-SD model. These results indicate that the improvement in effectiveness we observed when using the MRF-SD model was

	AP	WSJ	ROBUST04	WT10G	GOV2
MAP	0.2116†	0.3319† ↓	0.3012† ↓	0.2165†	0.3371†
GMAP	0.1229	0.2313↓	0.2076† ↓	0.1347†	0.2137† ↓
P@10	0.3420	0.4880↓	0.4343↓	0.3061	0.5500
R-Prec	0.2537†	0.3561↓	0.3339	0.2499†	0.3744†
μ_1	2750	2250	1000	2750	1250
λ_2	0.995	0.998	0.970	0.950	0.990
λ_3	1.00	1.00	1.00	0.99	1.00

Table 5.11. Test set results for the bigram language model. The † indicates a statistically significant improvement over the MRF-FI model and the ↓ indicates a statistically significant *decrease* in effectiveness compared to the MRF-SD model (i.e., MRF-SD > MRF-BM25). Recommended smoothing parameter values are also provided.

not specific to the language modeling weights used. Indeed, as we just showed, similar improvements can be obtained using BM25 weights. Therefore, this general form of model can be used in a “plug ’n play” manner, using any reasonable weighting function.

5.1.2.7 Comparison to Bigram Model

We now compare our model against another non-bag of words model. The model that we choose to compare against is the bigram language model (see Equation 2.16), ranked using query likelihood. This model has recently been shown to be one of the most consistently effective non-bag of words models to date [41]. We compare the effectiveness of the model against the MRF-FI and MRF-SD models. We choose the MRF-SD model since it is a direct generalization of the bigram model and models no additional dependencies, thus making it the most similar model to compare against.

In our experiments, we train the bigram model’s smoothing parameters to maximize mean average precision. The test set results are shown in Table 5.11. For each data set, we provide a set of recommended parameter settings. Furthermore, we indicate statistically significant improvements over the MRF-FI model and statistically significant *decreases* in effectiveness versus the MRF-SD model.

The results show that the bigram model is significantly better than the MRF-FI model across all data sets. This result is consistent with previous results [41]. However, the model is significantly worse than the MRF-SD model on the WSJ and ROBUST04 data sets. We note that the bigram model is never significantly better than the MRF-SD model for any metric. This result indicates that while the bigram model can be highly effective, the MRF-SD model is still a better choice, based purely on effectiveness.

Furthermore, we argue that the MRF framework, in general, is always a better choice than the bigram model. Since the MRF model clearly generalizes and supercedes the bigram model, it will always be more flexible and provide more modeling options. Furthermore, the bigram model uses a very rigid set of unigram and bigram features that cannot be changed across tasks. However, the MRF model provides an easy mechanism for including a wide range of arbitrary features. Therefore, there is little reason to choose the bigram model over the MRF model.

One particularly interesting result of the bigram experiments is that the improvement over the MRF-FI model increases as the collection size grows in a similar manner to the MRF-SD model. This result further supports our claim that term dependence and term proximity features are of the utmost importance when collection sizes grow, document lengths increase, and collections become noisier.

5.1.2.8 Generalization

Finally, we investigate several aspects of how well the MRF model parameters generalize. An underlying goal of parameter selection strategies is to produce a model that generalizes well. A model is said to generalize well if, when trained on one set of data, remains effective on an unseen test set. A model that is capable of achieving excellent effectiveness on a training set but performs poorly on a test set is of minimal value. Therefore, if some parameter selection method results in effectiveness \hat{m} , and

	AP	WSJ	ROBUST04	WT10G	GOV2	Avg.
LM	99.33	99.57	100.0	94.28	99.42	98.52
BM25	99.35	99.67	98.67	98.34	99.79	99.17
F2EXP	100.0	99.97	97.95	95.50	99.66	98.62
MRF-SD	97.93	97.39	99.23	100.0	100.0	98.91

Table 5.12. Intracollection generalization results for mean average precision. Values given are effectiveness ratios.

the optimal effectiveness is m^* , we then compute the following:

$$G = \frac{\hat{m}}{m^*} \quad (5.10)$$

which we define as the *effectiveness ratio*. An ideal model, that generalizes perfectly, would achieve an effectiveness ratio of 1 for every unseen data set. In information retrieval, even a 2-5% change in some measures, such as mean average precision, can be statistically significant, and therefore effectiveness ratios below 0.90 indicate a model’s inability to generalize can severely hinder its effectiveness. Most reasonable retrieval models will have an effectiveness ratio greater than 0.95.

We are particularly interested in *intracollection* and *intercollection* generalization, which are two different ways of measuring the generalization properties of a model, which we now describe.

Intracollection generalization deals with how well a model trained on a set of topics from some collection generalizes to another set of topics on that same collection. This is a common setting in TREC evaluations, where collections are often reused from year to year, and systems are typically trained on the topics from the previous year(s).

We ran a number of experiments to test the intracollection properties of various retrieval models. The retrieval models considered are language modeling (LM), BM25, F2EXP (an axiomatic retrieval model that was designed to be less sensitive to parameter estimation [37]), and the MRF-SD model. In these experiments, parameters are estimated by maximizing mean average precision on the training set. Models

Train\Test	AP	WSJ	ROBUST04	WT10G	GOV2	Avg.
AP	-	99.2	94.3	93.0	93.6	95.0
WSJ	99.1	-	97.4	96.4	96.2	97.3
ROBUST04	95.0	97.7	-	97.5	99.6	97.4
WT10G	91.8	92.7	96.5	-	93.4	93.4
GOV2	95.6	98.2	99.3	97.2	-	97.6
Avg.	95.4	96.9	96.9	96.0	95.8	96.2

Table 5.13. Inter-collection generalization results. Table includes mean average precision effectiveness ratios across all possible train/test splits using the F2EXP model.

are evaluated according to the effectiveness ratio on the test set. The metric used to compute the effectiveness ratio is mean average precision.

The results are given in Table 5.12. The table lists the effectiveness ratios of each model across each data set, as well as the average effectiveness ratio across all data sets. A model with perfect intracollection generalization would have an effectiveness ratio of 100. The results indicate that all of the models do a relatively good job of generalizing, with average effectiveness ratios well above 98%. We note that the F2EXP model tends to generalize better within newswire collections, while the dependence model generalizes better for web collections. The BM25 model, however, has the best average effectiveness ratio, which indicates its parameters do a particularly good job of capturing collection-dependent characteristics, rather than topic set-specific ones.

The other type of generalization we consider is inter-collection generalization. This type of generalization measures how well a model trained on a topic set from one collection generalizes to a different topic set on a different collection. This is a practical scenario for ‘off the shelf’ retrieval systems that may be used across a wide range of different collections. It is unlikely that the end users of these systems will be willing or able to provide training data to the system, and therefore the system must be shipped with a very solid set of pre-tuned, highly generalizable parameters.

Train\Test	AP	WSJ	ROBUST04	WT10G	GOV2	Avg.
AP	-	100	99.7	98.4	98.9	99.3
WSJ	100	-	99.7	98.4	98.9	99.3
robust04	99.6	99.6	-	99.7	99.3	99.5
WT10G	98.1	98.9	99.8	-	97.0	98.5
GOV2	99.6	99.4	99.7	98.0	-	99.2
Avg.	99.3	99.5	99.7	98.7	98.5	99.1

Table 5.14. Inter-collection generalization results. Table includes mean average precision effectiveness ratios across all possible train/test splits using the MRF-SD model.

In order to measure the inter-collection generalization, we compute the effectiveness ratio for every possible combination of training/test splits. The results for the F2EXP and MRF-SD models are shown in Tables 5.13 and 5.14, respectively.

As we see from the table, the cross-collection effectiveness ratios for the MRF-SD model are higher for every training/test set pair, with very few exceptions. In fact, on average, the MRF-SD comes within 1% of the optimal setting regardless of which collection is used for training, whereas the F2EXP model only comes within 4% of the optimal on average. The Dirichlet and BM25 models (not shown) have average effectiveness ratios of 98.9% and 96.9%, respectively. Therefore, the MRF-SD model and Dirichlet models are more robust when it comes to cross-collection generalization and make them good candidates for “out of the box” implementations that require a single parameter setting to work well across a wide range of collections.

As further evidence of the model’s generalization properties, Figure 5.9 illustrates the well-behaved, nearly concave surfaces that arise when mean average precision is plotted over the multinomial parameter simplex of the MRF-SD ranking function for various data sets. Each of the mean average precision surfaces has the same general form, which indicates that the features capture an inherent property that persists across different types of collections. Although there is no guarantee that such a nicely concave surface will exist for all features and all evaluation metrics, it provides

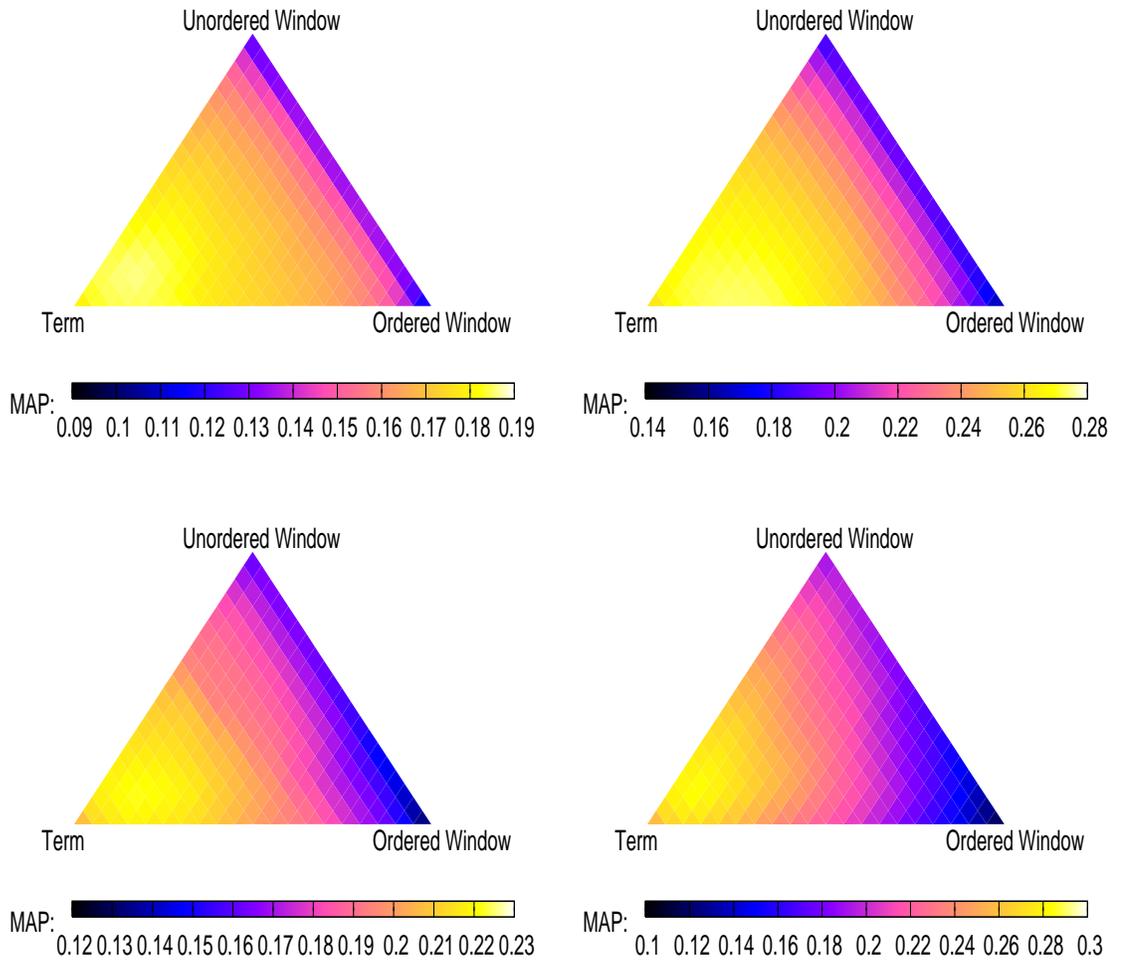


Figure 5.9. Mean average precision values plotted over MRF-SD parameter simplex for AP, WSJ, WT10g, and GOV2 collections.

	MRF-FI	MRF-SD	MRF-FD	(μ^t, μ^w)
AP	0.2077	0.2147 (+3.4%) [†]	0.2128 (+2.5%)	(1750, 5000)
WSJ	0.3258	0.3425 (+4.8%)	0.3429 (+5.2%) [†]	(2000, 1000)
ROBUST04	0.2920	0.3096 (+6.0%) [†]	0.3092 (+5.9%) [†]	(1000, 750)
WT10g	0.1861	0.2053 (+10.3%) [†]	0.2140 (+15.0%) ^{†‡}	(1000, 6000)
GOV2	0.2984	0.3325 (+11.4%) [†]	0.3360 (+12.6%) [†]	(1500, 4500)
$(\lambda_{T_D}, \lambda_{O_D}, \lambda_{U_D})$	N/A	(0.85, 0.10, 0.05)	(0.80, 0.10, 0.10)	

Table 5.15. Summary of test set mean average precision for the MRF-FI, MRF-SD, and MRF-FD models across all of the *ad hoc* retrieval data sets. Values in parenthesis denote percentage improvement over MRF-FI model. A [†] indicates a statistically significant improvement over the MRF-FI model, and a [‡] indicates a statistically significant improvement over the MRF-SD model. Recommended smoothing values are given for each collection, and recommended MRF model parameters are provided for each model.

some evidence that the functions we are maximizing over the simplex are not too difficult to optimize using simple algorithms, such as the coordinate ascent method we proposed in Chapter 4.

5.1.3 Summary of Results

For completeness, we summarize the main results of our *ad hoc* retrieval experiments in Table 5.15. The table shows test set mean average precision across all data sets for the MRF-FI, MRF-SD, and MRF-FD basic models. It also includes recommended smoothing parameters for each collection, as well as recommended MRF model parameters for each model.

These results show that the hand constructed non-bag of words MRF models (MRF-SD and MRF-FD) consistently outperform the bag of words model (MRF-FI). Although we only show results here that make use of language modeling weighting functions, we note that the results achieved using the MRF-BM25 model, which uses BM25 weighting functions, did consistently and significantly outperform the MRF-FI model in terms of mean average precision across all data sets, thereby satisfying one of the goals of our work. As we will show in Chapter 6, moving away from hand built

MRFs towards automatically constructed ones yields even more significant increases in effectiveness.

5.2 Web Search

Web search is one of the most popular and widely used information retrieval applications. The goal of a web search system is to return a set of web pages that are relevant to a user's query. There are several important differences between *ad hoc* retrieval and web search. First, not all web search queries are content-based, or *informational*, searches. For example, a user who enters the query *mcdonalds locations* is not looking for documents that are *about* McDonalds locations. Instead, they are likely searching for the web page on the McDonald's web page that lists where their restaurants are located. This type of query, where a user seeks out a specific page that they either know exists or think is very likely to exist, is known as a *navigational* query. Additionally, a user who enters the query *cheap digital cameras* is neither seeking information *about* digital cameras nor seeking a specific page. Instead, the user is likely interested in purchasing a digital camera. Such queries, which are intended to lead to an online transaction, are known as *transactional* queries. A study done by Broder in 2002 reports that approximately 50% of web queries are informational, 20% navigational, and 30% transactional [11].

Since these three query types are so different, they are often evaluated differently. Content-based (informational) web retrieval was evaluated in the previous section with our experiments on the WT10G and GOV2 data sets. Evaluation of transactional queries is difficult, since it often requires product databases and query clickthrough logs, which are not currently publicly available for privacy and intellectual property reasons. In this section, we focus on navigational queries, for which there are publicly available TREC data sets.

```
<num> Number: NP1048
<title> us embassy vietnam
</top>

<top>
<num> Number: NP1049
<title> phil english biography
</top>

<top>
<num> Number: NP1050
<title> tenet 9/11 testimony
</top>

<top>
<num> Number: NP1051
<title> hubble timeline
</top>

<top>
<num> Number: NP1052
<title> cdc west nile 2003 statistics by state
</top>
```

Figure 5.10. Example TREC named page finding topics.

These data sets were used during the TREC Web Track (1999-2004) and the TREC Terabyte Track (2004-2006). During these tracks, there were several navigational search-related subtasks. Of particular interest to us is the *named page finding* task that was run during the TREC Terabyte Track in 2005 and 2006. We are interested in this task because of the large data (GOV2) set and the large number of queries available to experiment with (252 from 2005, 181 from 2006).

The named page finding task requires systems to find “bookmarkable” pages that users either know exist or presume are likely to exist [21]. One of the main differences between named page finding and *ad hoc* retrieval is that there is typically only one relevant document for every named page finding query. Another key difference is that the primary evaluation metric is mean reciprocal rank, instead of mean average precision². Several example named page finding topics are shown in Figure 5.10.

In the remainder of this section, we first review previously proposed approaches to web search (named page finding). We then propose our basic MRF model for web search. Finally, we evaluate our model on the TREC Terabyte Track named page finding topics.

5.2.1 Previous Models for Web Search

Web data sets are very different than standard *ad hoc* retrieval data sets. They are typically larger and tend to be noisier, because users may publish their own content. Furthermore, web pages contain HTML markup and can link to other pages. These additional pieces of information play a particularly important role for navigational queries. For this reason, navigational search models often focus on *link structure* and *document structure*.

²Average precision is equal to reciprocal rank for queries with only one relevant document, so the two measures will only differ for those topics that have more than one relevant document

Link structure refers to the hyperlink structure of the web. Graph-based algorithms are typically applied to the link structure of the web in order to find hubs and authoritative pages. Examples of these algorithms include PageRank and HITS [10, 53]. These algorithms, along with other web-specific features such as inlink count, and URL depth have been shown to be useful for improvement the effectiveness of navigational queries [54].

Methods that make use of document structure often treat certain HTML fields in a special way. For example, a common technique is to weight the importance of text occurring in various fields differently, such as the `title` field or the anchor text pointing to the document [78, 91].

Recently, other aspects of web search, such as user behavior, have been found to be useful, but is beyond the scope of this this work [1].

In this section, we use a bag of words language modeling approach as our baseline. The model, which we refer to as LM-Mixture, makes use of both link structure and document structure, has been shown to be highly effective in the past [78]. Given a query, documents are ranked under the model according to:

$$\begin{aligned}
 P(D|Q) &= \frac{P(D) \prod_{w \in Q} \sum_f P(f|D)P(w|D, f)}{\sum_D P(D) \prod_{w \in Q} \sum_f P(f|D)P(w|D, f)} \\
 &\stackrel{rank}{=} P(D) \prod_{w \in Q} \sum_f P(f|D)P(w|D, f)
 \end{aligned} \tag{5.11}$$

where $P(w|D, f)$ is the probability of generating term w from field f in document D (i.e., this is a language model built from field f in document D), $P(f|D)$ is a mixing probability, and $P(D)$ is the document’s prior probability. It is easy to see that this model is closely related to the standard query likelihood ranking function, except the monolithic document model is replaced with a mixture of field models and a document prior is introduced.

There are other models that attempt to combine evidence from multiple fields, including The BM25F model, which is a field weighted variant of BM25 [91]. The model is similar in nature to the mixture of language models approach described here, but field weighting is done differently. Rather than weighting terms after document length normalization is done, as is done in Equation 5.11, term weights are incorporated before document length normalization. How to properly combine evidence and handle document length normalization in the presence of multiple fields is still an open question [103].

In order to fully specify the model we must describe how to estimate $P(w|D, f)$, $P(f|D)$, and $P(D)$. We begin with the field language model, $P(w|D, f)$. This probability is estimated in a straightforward manner by treating all of the text that appears in field f of document D as a pseudo-document. By doing so, we induce the pseudo-document D_f for field i and the pseudo-collection C_f , which is made up of all of the pseudo-documents constructed from field f . Now, standard language modeling estimation can be applied. We choose to model each field language model as a Dirichlet smoothed language model, which results in the following estimate:

$$P(w|D, f) = \frac{tf_{w,D_f} + \mu_f^t \frac{c_{f,w,f}}{|C_f|}}{|D_f| + \mu_f^t} \quad (5.12)$$

where all of the f subscripts refer to statistics computed in the pseudo-document/pseudo-collection and μ_f^t specifies the smoothing parameter for the field. Since there is a single smoothing parameter per field, accurate estimation may be difficult. Hence, smoothing parameter values are typically chosen to be two times the average length of pseudo-documents of type f . We follow this general rule of thumb in our experiments here.

The mixing probabilities, $P(f|D)$ can either be set to $\frac{|D_f|}{|D|}$ or uniformly. Alternatively, they can be hand or automatically tuned in order to maximize mean reciprocal

rank. In this work, we use a set of hand tuned values that have been found to be effective in previous experiments.

5.2.2 Document Priors

There are many different document priors that can be estimated for navigational web search tasks. In this section, we describe how to estimate priors based on inlink count and PageRank. The priors are estimated using TREC relevance judgments, although they may also be estimated in a completely unsupervised or semi-supervised setting, given other resources, such as query click logs.

When computing $P(D)$, we really are computing the prior probability that document D is relevant given some external piece of evidence about D , such as the number of links pointing to D or the PageRank of D . Therefore, instead of estimating $P(D)$ directly, we estimate $P(R = 1|evidence)$, where evidence is some random variable that only depends on the document itself.

5.2.2.1 Inlink Count

The inlink count of document D is the number of web pages that point to D . Inlink count is often a good feature to use for navigational queries because the pages that are “bookmarkable” often have a high inlink count associated with them. Therefore, we expect documents with larger numbers of inlinks to have a higher prior probability of relevance.

Following the framework we described above, we compute the probability of relevance, given the log of the inlink count (simply denoted l). We choose to apply the log function in order to compress the range of values to a more reasonable set. The resulting probability estimate, using Bayes’ rule, is:

$$P(R = 1|l = X) = \frac{P(l = X|R = 1)P(R = 1)}{P(l = X|R = 0)P(R = 0) + P(l = X|R = 1)P(R = 1)} \quad (5.13)$$

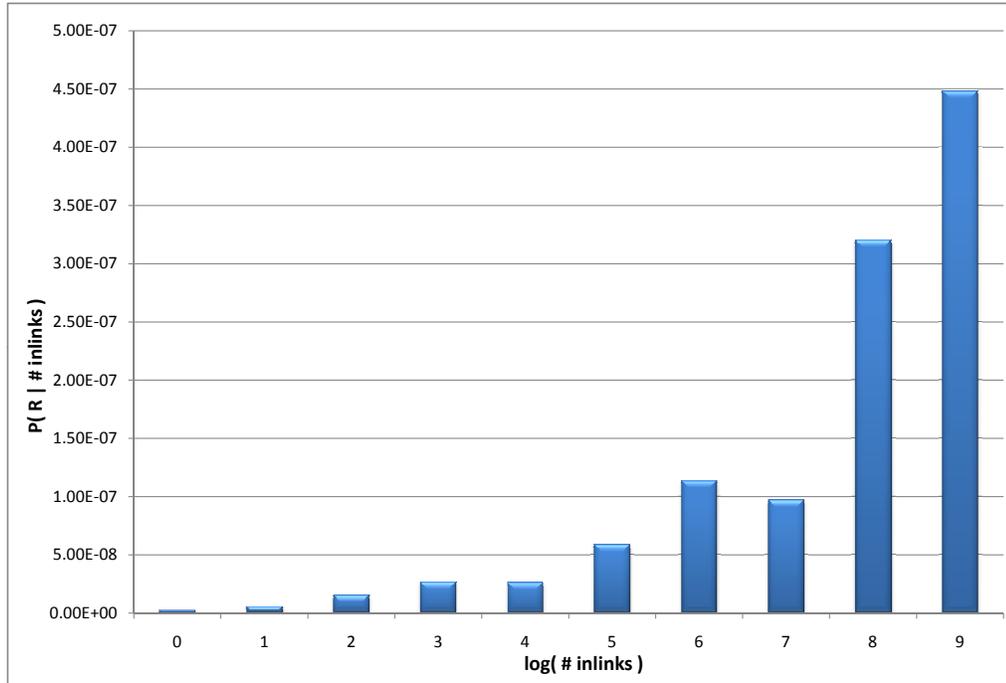


Figure 5.11. Inlink count prior.

where $P(l = X|R)$ and $P(R)$ are estimated empirically from TREC relevance judgments. The relevance judgments used are from the TREC 2001 and 2002 Web Track data set. Later, we will apply these priors to the larger TREC 2005-2006 Terabyte Track data set (GOV2). It is unknown how well these probabilities will generalize to this newer, larger data set, but we feel that the estimates should be fairly reasonable.

Figure 5.11 shows the estimated priors across a range of log inlink counts. We see that the prior probability of relevance increases as the number of inlinks increases, as expected.

5.2.2.2 PageRank

The problem with using inlink count alone is that there is no notion of authority involved. It is very easy for a spammer to create thousands of fake web pages and have them all point to each other (this is a so-called link farm). This results in a large number of inlinks, but none of the pages are actually authoritative at all.

The PageRank algorithm is based on the notion of spreading authority throughout a graph. The basic idea is that if a page has many inlinks from highly authoritative pages, then that page is likely to be authoritative, as well.

The PageRank of document D , denoted $r_{D,t}$, can be computed iteratively. Let t denote the current iteration. Then, the PageRank is computed as follows:

$$r_{D,t+1} = \alpha + (1 - \alpha) \sum_{p:p \rightarrow D} \frac{r_{p,t}}{\text{deg}(p)} \quad (5.14)$$

where $\alpha \in (0, 1]$ affects the amount of “random surfing” done, $p \rightarrow D$ indicates that page p links to document D , and $\text{deg}(p)$ is the number of links out of page p . The values are iteratively updated and renormalized until they converge to the raw PageRank value. There are other ways of computing PageRank, for example, by solving an eigenvector problem, but we use the iterative approach in our work for simplicity.

The raw PageRank is a value between 0 and 1. We assume that the true authoritativeness, or importance, of web pages is Zipfian in nature. That is, there are a small number of highly authoritative pages, a larger number of less authoritative pages, all the way down to a very large number of non-authoritative pages. Therefore, in order to impose such a distribution, we sort the documents by their raw PageRank scores and then geometrically bin the documents into 11 bins. This idea was inspired by Anh and Moffat’s work on document-centric impact weighting [3]. This results in each document being assigned a binned PageRank value between 0 and 10. We use these binned PageRank values in order to estimate the document prior.

The PageRank prior is computed in a similar fashion to the inlink count prior, as follows:

$$P(R = 1 | PR = X) = \frac{P(PR = X | R = 1)P(R = 1)}{P(PR = X | R = 0)P(R = 0) + P(PR = X | R = 1)P(R = 1)} \quad (5.15)$$

URL	# Inlinks
http://www.usgs.gov	1,329,036
http://www.ca.gov	471,819
http://www.nih.gov	1,502,324
http://www.epa.gov	1,386,329
http://es.epa.gov/cgi-bin/ncercqamail.pl	1,349,131
http://es.epa.gov/ncer/rfa	1,344,630
http://www.hhs.gov	778,652
http://www.ornl.gov	639,570
http://www.doi.gov	761,456
http://www.medicare.gov	186,948

Table 5.16. URLs in the GOV2 collection with the largest raw PageRank scores. The number of inlinks for each URL is also shown.

where PR denotes the binned PageRank value. Table 5.16 shows a list of the web pages in the GOV2 collection with the highest raw PageRank values and the number of inlinks those pages have. Although many of the pages with the highest PageRank have very many inlinks, this is not always the case. The best example of this is the fact that the Medicare web page has such a high PageRank, but only has 186,948 inlinks.

Lastly, Figure 5.12 shows the estimated document priors for each binned PageRank value. The plot has an interesting shape to it. Documents with very low PageRank are given a very low prior probability of relevance. The prior probability dramatically increases as the PageRank reaches 8 and 9. Interestingly, a higher prior is assigned to documents with PageRank 8 than 9 and that a zero probability is assigned to documents with PageRank 10. This is very likely the result of data sparseness and the fact that there are so few documents with PageRank 9 or 10 in our relevance judgments.

5.2.3 MRF Models for Web Search

We now describe one of the many possible ways to use the MRF model for web search. The mixture of language modeling approach described earlier (see Equ-

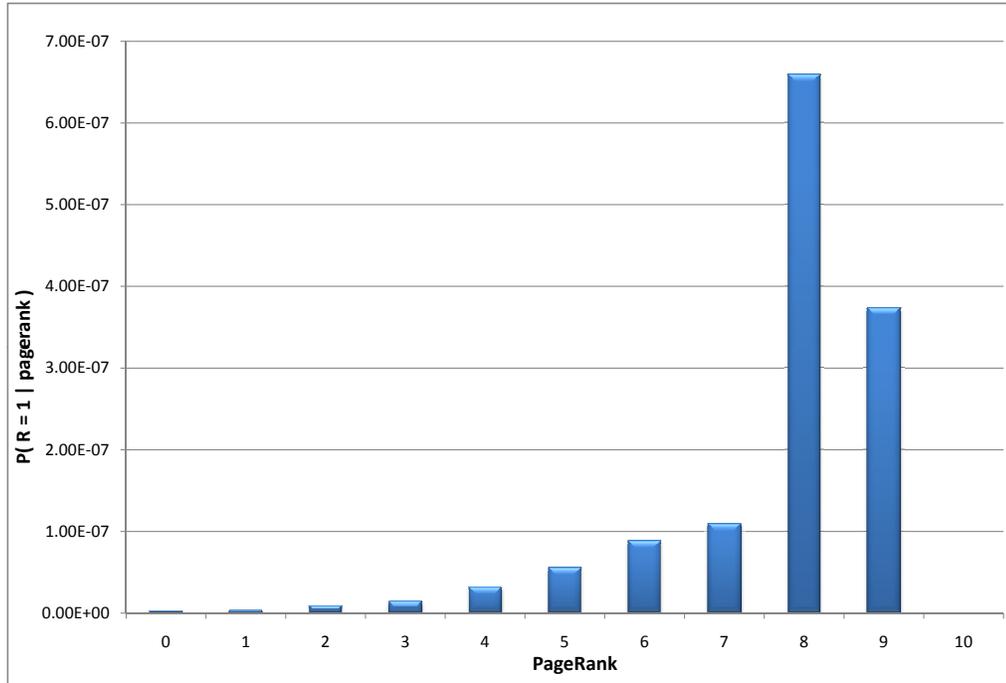


Figure 5.12. PageRank prior.

tion 5.11) has been shown to be highly effective. Therefore, we wish to use the MRF framework to generalize this model. By doing so, we will be able to model dependencies between query terms, which allows us to make use of phrase and term proximity weighting functions, which were shown to be very valuable for the *ad hoc* retrieval task.

To achieve this, we modify the MRF-SD model by replacing the standard language modeling feature weights (i.e., LM, LM-O-1, and LM-U-4) with analogous mixture of language modeling feature weights. In addition, we add the inlink count and PageRank document priors as feature weights defined over the document clique, as well. These new feature weights, named NP, NP-O-*M*, NP-U-*N*, INLINK, and PAGERANK are defined in Table 5.17. This gives rise to the basic MRF for web search, which we call the *MRF-NP* model, defined by the following canonical form:

$$(FI, T_{QD}, NP) : \lambda_{T_D}$$

<p>NP</p> $f_{NP,T}(q_i, D) = \log \left[\sum_f \alpha_f \frac{t f_{q_i, D_f} + \mu_f^t \frac{c f_{q_i, f}}{ C_f }}{ D_f + \mu_f^t} \right]$
<p>NP-O-M</p> $f_{NP,O,M}(q_1, \dots, q_k, D) = \log \left[\sum_f \alpha_f \frac{t f_{\#M(q_1 \dots q_k), D_f} + \mu_f^w \frac{c f_{\#M(q_1 \dots q_k), f}}{ C_f }}{ D_f + \mu_f^w} \right]$
<p>NP-U-N</p> $f_{NP,U,N}(q_1, \dots, q_k, D) = \log \left[\sum_f \alpha_f \frac{t f_{\#uNk(q_1 \dots q_k), D_f} + \mu_f^w \frac{c f_{\#uNk(q_1 \dots q_k), f}}{ C_f }}{ D_f + \mu_f^w} \right]$
<p>INLINK</p> $f_{INLINK}(D) = \log P(R = 1 l = l(D))$
<p>PAGERANK</p> $f_{PR}(D) = \log P(R = 1 PR = PR(D))$

Table 5.17. Summary of named page weighting functions. The NP, NP-O-M, and NP-U-N weighting functions are based on mixtures of field language models, and INLINK and PAGERANK are based on document priors. The α_f values correspond to the mixing probabilities $P(f|D)$. Term and window smoothing parameters are denoted by μ_f^t and μ_f^w , respectively.

$$(FI, T_Q, ICF) : \lambda_{T_Q}$$

$$(SD, O_{QD}, NP-O-1) : \lambda_{O_D}$$

$$(SD, O_Q, ICF-O-1) : \lambda_{O_Q}$$

$$(SD, O_{QD}, NP-U-4) : \lambda_{U_D}$$

$$(SD, O_Q, ICF-U-4) : \lambda_{U_Q}$$

$$(FI, D, INLINK) : \lambda_{IN}$$

$$(FI, D, PAGERANK) : \lambda_{PR}$$

Using this canonical form, the MRF-NP ranking function is given by:

$$\begin{aligned}
P_{G,\Lambda}(D|Q) \stackrel{rank}{=} & \lambda_{T_D} \sum_{(q_i,D) \in T_{QD}} \log \left[\sum_f \alpha_f \frac{t f_{q_i,D_f} + \mu_f^t \frac{c f_{q_i,f}}{|C_f|}}{|D_f| + \mu_f^t} \right] + \\
& \lambda_{O_D} \sum_{(q_1,q_2,D) \in O_{QD}} \log \left[\sum_f \alpha_f \frac{t f_{\#1(q_1q_2),D_f} + \mu_f^w \frac{c f_{\#1(q_1q_2),f}}{|C_f|}}{|D_f| + \mu_f^w} \right] + \\
& \lambda_{U_D} \sum_{(q_1,q_2,D) \in U_{QD}} \log \left[\sum_f \alpha_f \frac{t f_{\#uw8(q_1q_2),D_f} + \mu_f^w \frac{c f_{\#uw8(q_1q_2),f}}{|C_f|}}{|D_f| + \mu_f^w} \right] + \\
& \lambda_{IN} \log P(R = 1 | l = l(D)) + \\
& \lambda_{PR} \log P(R = 1 | PR = PR(D))
\end{aligned} \tag{5.16}$$

The model provides a mechanism for weighting fields (α_f), single term matches (λ_{T_D}), ordered phrases (λ_{O_D}), unordered phrases (λ_{U_D}), inlink prior (λ_{IN}), and PageRank prior (λ_{PR}). Thus, the model encompasses many of the features that have been shown to be effective for navigational web search.

There are many alternative ways that this model may have been constructed. For example, rather than mixing the field language models within the NP weighting functions, each field language model could be its own feature weight. This would promote the α_f hyperparameters to full-fledged MRF model parameters, which would allow them to be estimated using the machinery described in Chapter 4. Although we do not attempt to empirically analyze the differences in the formulations, we believe that there would be little, if any, difference in the effectiveness of the two models.

5.2.4 Results

We now empirically evaluate the retrieval effectiveness of the MRF-NP model. For our experiments, we use the TREC 2005 and 2006 Terabyte Track named page finding queries. These queries are evaluated against the GOV2 collection. Please refer to Appendix A for further information on this data set.

In our experiments, we consider four fields in the mixture models. These fields are **body** (full text of the web page), **title** (text within HTML elements) **heading** (text

	LM-Mixture			MRF-NP		
	MRR	S@10	Not Found	MRR	S@10	Not Found
TREC 2005	0.414	0.563	0.175	0.441	0.583	0.171
TREC 2006	0.472	0.657	0.133	0.512	0.696	0.138

Table 5.18. Summary of named page finding results.

within **h1**, **h2**, **h3**, and **h4** elements), and **anchor** (all of the anchor text that points to a page). The collection was stemmed using the Porter stemmer and a standard list of 418 stopwords was applied. The MRF model parameters were estimated by maximizing mean reciprocal rank. No more than 1000 results were returned for each query.

Little research has been done on term dependence and non-bag of words models for named page finding. Hence, for experimental purposes, we use the mixture of language models approach as our baseline. In order to ensure fairness, the LM-Mixture model uses the same fields, mixing parameters (α_f), and smoothing parameters (μ_f^t) as the MRF-NP model. In addition, we also use the inlink count and PageRank document priors with the LM-Mixture model. The two priors are combined into a single prior as described in [54].

The results of our experiments are given in Table 5.18. We report results for mean reciprocal rank (MRR), success at rank 10 (S@10), and not found. See Appendix B for definitions of these metrics.

The results show that the MRF-NP model outperforms the baseline language modeling mixture model in terms of MRR on both data sets. There is 6.5% improvement on the 2005 queries and a 8.5% improvement on the 2006 queries. Both of these results are statistically significant. Furthermore, the S@10 metric is improved on both data sets, as well, indicating that the MRF-NP model pulls relevant documents into the top 10. The last metric, not found, does not change significantly for the two models, which indicates that they are both relatively stable in terms of completely

LM-Mixture		MRF-NP	
No Prior	Prior	No Prior	Prior
0.463	0.472	0.498	0.512

Table 5.19. Results comparing the mean reciprocal rank of the LM-Mixture and MRF-NP models with and without document priors.

failing to find any relevant documents. Although these numbers are relatively good, there is still considerable room for improvement.

One thing that these results do not reveal, however, is how much of the increase in effectiveness comes from the term proximity features and how much comes from the document priors. In Table 5.19, we report the results of an ablation test that attempts to quantify the importance of each type of feature. The results show that document priors improve effectiveness 1.9% on for the LM-Mixture model and 2.8% for the MRF-NP model. This indicates that the MRF model does a better job at combining the evidence from the document priors than the LM-Mixture model. As for the term proximity features, there is an improvement of 7.6% when no priors are used, and an improvement of 8.5% when priors are used. These results show that the term proximity features account for most of the improvement in effectiveness and that, when used in conjunction with the document priors, there is a small additive effect. Therefore, despite what Google would like you to think about PageRank being the heart of their ranking function³, it appears as though, in reality, PageRank is far less important than fundamental information retrieval features, such as term proximity.

As with our previous analysis, we are interested in developing a better understanding of the types of queries the MRF-based model excels at, and those it fails at. Table 5.20 lists the 10 most helped and 10 most hurt queries from the 2006 data set. These examples do not show any clear trends as to which types of queries are likely to be helped or hurt by the model. In the future, it may be valuable to do an

³<http://www.google.com/technology/>

Topic	Query	LM-Mixture	MRF-NP	% Change
980	medline search	0.00	0.00	-100.0%
962	us iraq rebuilding accomplish- ments	0.01	<0.00	-80.0%
1038	USPTO Guide and manuals	0.33	0.09	-72.7%
922	marine mammal gray whale	0.33	0.11	-66.6%
906	kickstart deviceprobe	0.33	0.17	-50.00%
943	american folklife center home- page	1.00	0.50	-50.0%
1033	Coastal & Marine Geology In- foBank	0.06	0.04	-33.3%
916	Texas Department of Banking, Agency Philosophy	0.02	0.01	-31.5%
1005	bay trail map	0.20	0.14	-28.6%
1006	olympic games salt lake city new jobs	0.01	0.01	-26.6%
...				
1041	CDC homepage	<0.00	0.01	204.2%
936	patent DRAM cell constructions	0.14	0.50	250.0%
939	Sun Earth student section	0.14	0.50	250.0%
951	us embassy vienna	0.14	0.50	250.0%
984	informal personal caregiver em- ployment	0.13	0.50	300.0%
1030	Space Shuttle Mission #75	0.03	0.14	442.9%
912	Tips for Mobile Homes Residents in Wisconsin	0.03	0.25	650.0%
903	reasons to reduce waste	0.13	1.00	700.0%
1008	land use bill december 2003	0.02	0.17	816.7%
1027	1997 Surface Flows to Nevada from Canadian Province of Ori- gin, by Truck	0.10	1.00	900.0%

Table 5.20. The 10 most improved and 10 most hurt queries on the TREC 2006 Terabyte Track named page finding data set. Effectiveness is measured in terms of reciprocal rank.

analysis to determine if the most helped queries can be automatically detected using more sophisticated techniques, such as the so-called notion of concept density [34].

Finally, we examine the robustness of the MRF-NP method. The results are given in Figure 5.13. These results are similar to the *ad hoc* retrieval results, with many queries experiencing a large increase in reciprocal rank, and a small number experiencing less significant decreases. The reason why there is a large peak in the [0%,25%] bin is because the effectiveness of a large number of queries do not change at all. The large peak around [50%,75%] is likely the result of how the reciprocal rank is computed. If a relevant document moves up in the ranked list by a very small number of positions, then, depending on where in the ranked list it original appeared, the increase in reciprocal rank is likely to fall into this range.

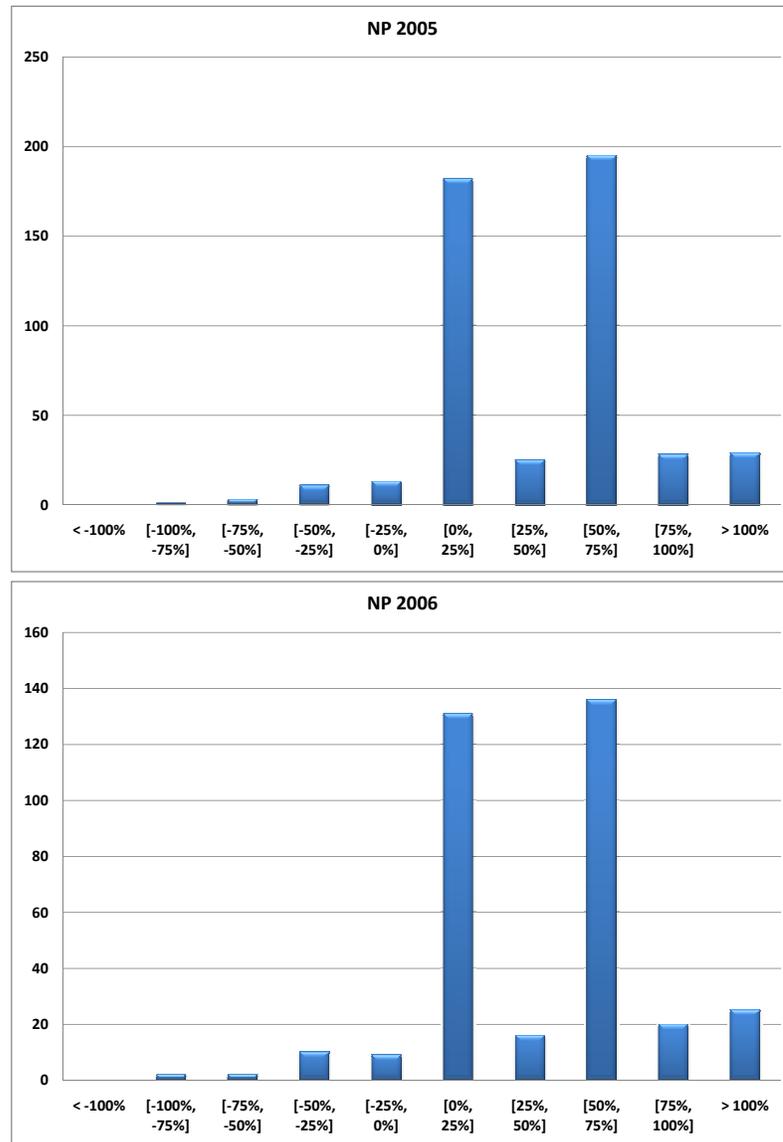


Figure 5.13. Robustness of the MRF-NP models for the 2005 and 2006 Terabyte Track named page finding data sets. The LM-Mixture model is used as the baseline by which the improvements were computed. The evaluation metric used is reciprocal rank.

CHAPTER 6

AUTOMATIC FEATURE SELECTION

As we showed in Chapter 2, many different types of retrieval models have been proposed throughout the years. These include Boolean, vector space, logic-based, probabilistic, and feature-based. One critical factor that must be considered when developing information retrieval models is the type of features to be used or modeled. Term frequency, inverse document frequency, document length, and term proximity are the fundamental features that are used in most of the modern information retrieval models including BM25 [89], language modeling [101], divergence from randomness (DFR) [2], axiomatic approach to IR [37], and our proposed MRF model.

However, most of these models make use of hand selected, probabilistically-inspired, or implicit features. Therefore, it is often difficult to adapt these types of models to new tasks, especially when the task has new, completely different types of features associated with it. Applying these models to new tasks typically requires an information retrieval expert to modify the underlying model in some way in order to properly account for the new types of features. This is a common theme in information retrieval modeling. Examples include incorporating PageRank as a prior into the BM25 model [26], allowing term proximity information as evidence in BM25 [18], modeling document structure in both language modeling and BM25 [78, 91], including term dependence in the DFR model [82], and allowing term associations in the axiomatic model [38]. These examples illustrate that incorporating new types of evidence and features into existing retrieval models is often non-trivial and can require significant amounts of human involvement.

Therefore, it is desirable for models to be flexible and robust enough to easily handle a wide range of features and provide a mechanism for automatically selecting relevant features. Then, given a large pool of candidate features, it would be possible to automatically learn the best model. Under this model learning paradigm, there would no longer be a need to manually tune or modify some existing retrieval model whenever a new task or data set is encountered. Instead, attention could be paid to developing a rich pool of features that are widely applicable.

We argue that our proposed MRF model, and similar types of models, such as Gao et al.’s linear discriminant model [42], are the correct types of models to use when model flexibility and robustness are important. In this chapter, we propose an automatic, supervised feature selection algorithm that can be used in conjunction with these types of models. Our proposed algorithm is general and can be applied to a wide range of feature sets, evaluation metrics, and methods for learning to rank. Besides being more robust and flexible, we also show that models constructed using our algorithm are often significantly more effective than the hand built basic MRF models we evaluated in Chapter 5.

6.1 Related Work

A number of feature selection techniques for random field models have been proposed in the machine learning literature [32, 63]. Our proposed algorithm is an adaptation of the feature induction technique proposed by Pietra et al. in [32]. Pietra et al. propose a greedy approach for adding induced features to the underlying model. During each iteration, the information gain for each induced feature is computed. The feature with the highest information gain is then added to the model and the entire model is retrained. Although we do not actually induce new features in our present work, we use a similar algorithm for selecting from a large pool of features. Another difference is that our algorithm scores each feature according to any infor-

mation retrieval metric of interest. The feature that improves the metric the most is the one that is added to the model.

There has also been some information retrieval research into automatically learning ranking function using genetic programming [36]. These algorithms attempt to find a locally optimal ranking function by iteratively “evolving” a population of ranking functions using mutations and crossovers. Ranking functions are represented as arithmetic trees that consist of arithmetic operators and standard bag of words information retrieval features (e.g., term frequency, document length, etc.). The learned ranking functions have been shown to be significantly more effective than baseline ranking algorithms for several data sets [36].

Finally, result fusion techniques are another way of combining evidence from multiple types of features [6, 39]. If each individual feature is used as a ranking function, then data fusion techniques can be used to determine the best way to combine the rankings. However, using these techniques in this way does not directly address the feature selection problem, which is our primary focus.

6.2 Automatic Feature Selection

As we described before, feature selection techniques are commonly used in the machine learning community. In this section, we propose a feature selection algorithm that can be used with our MRF model. The algorithm is specifically designed to be used for information retrieval tasks.

6.2.1 Motivation

Feature selection is important for a number of reasons. First, it provides a general, robust way of building models when there is little *a priori* knowledge about the types of features that may be important for a given task or data set. By using a feature selection algorithm, the model designer can focus less on building the best model and

can instead focus on designing good features. Second, feature selection can reduce the number of noisy or redundant features in a large feature set. Such features may reduce training efficiency and may result in a model that contains a number of non-identifiable parameters. Non-identifiable parameters are those that cannot be reasonably estimated given the training data. This often results from having redundant or highly correlated parameters. Feature selection helps overcome the problems associated with non-identifiable parameters. Finally, feature selection can provide insights into the important features for a given task or data set. By inspecting the order in which features are selected, we can often learn what characteristics of a given task are the most important or the most exploitable. This knowledge can then be used by the feature engineer to construct better features.

6.2.2 Algorithm

We now describe our automatic feature selection algorithm. While our discussion will focus on how the algorithm can be applied to the MRF model for IR, it should be noted that it can also be applied to a variety of other models. In particular, it can be easily applied to linear feature-based models (see Chapter 4).

Let M_t denote the model learned after iteration t . Features are denoted by f and the weight (parameter) associated with feature f is denoted by λ_f . The candidate set of features is denoted by \mathcal{F} . The entire set of feature weights for a model is denoted by Λ . A model, then, is represented as set of feature/weight pairs. Finally, we assume that $SCORE(M)$ returns the utility or ‘goodness’ of model M with respect to some training data. The utility function and the form of the training data largely depends on the underlying task. For example, for *ad hoc* retrieval, it is likely that $SCORE(\cdot)$ would return the mean average precision of using model M against some set of training data, such as TREC topics and relevance judgments. For a homepage finding task, $SCORE(\cdot)$ might be another metric, such as mean reciprocal rank. The

important thing to note here is that any utility function, regardless of whether or not it is differentiable with respect to the model parameters, can be used. The ultimate goal of our feature selection algorithm is to select features and set feature weights in such a manner as to maximize the metric imposed by $SCORE(\cdot)$.

The algorithm begins with an empty model (i.e., $M_0 = \{\}$). Then, we temporarily add a feature f to the model. We then hold all weights except λ_f fixed and find the setting for λ_f that maximizes the utility of the augmented model. This step can be done using any number of learning to rank techniques or parameter estimation techniques, including the ones described in Chapter 4. The utility of feature f ($SCORE_f$) is defined to be the maximum utility obtained during training. The feature's utility measures how good the current model would be if the feature were added to it. This process is repeated for every $f \in \mathcal{F}$, resulting in a utility being computed for every feature in the candidate pool. The feature with the maximum utility is then added to the model and removed from \mathcal{F} . After the new feature is added, we can, optionally, retrain the entire set of weights. The entire process is then repeated until either some fixed number of features have been added to the model or until the change in utility between consecutive iterations drops below some threshold. Algorithm 4 provides pseudo-code for this algorithm.

Note that our algorithm is not guaranteed to find the global maximum for $SCORE(M)$. Instead, we are only guaranteed to find a local maxima. Many factors, including properties of $SCORE(M)$, the number of features used, and the properties of the feature used, will affect the quality of the learned model.

6.3 Evaluation

In this section we experimentally evaluate various aspects of our proposed feature selection algorithm.

Algorithm 4 Feature selection algorithm.

```
1:  $t \leftarrow 0$ 
2:  $M_t \leftarrow \{\}$ 
3: while  $SCORE(M_t) - SCORE(M_{t-1}) > \epsilon$  do
4:   for  $f \in \mathcal{F}$  do
5:      $\hat{\lambda}_f \leftarrow \arg \max_{\lambda_f} SCORE(M \cup \{(f, \lambda_f)\})$ 
6:      $SCORE_f \leftarrow SCORE(M \cup \{(f, \hat{\lambda}_f)\})$ 
7:   end for
8:    $f^* \leftarrow \arg \max_f SCORE_f$ 
9:    $M \leftarrow M \cup \{(f^*, \hat{\lambda}_{f^*})\}$ 
10:   $\Lambda \leftarrow \arg \max_{\Lambda} SCORE(M)$  (optional)
11:   $\mathcal{F} \leftarrow \mathcal{F} - \{f^*\}$ 
12:   $t \leftarrow t + 1$ 
13: end while
```

In order to investigate the strengths and weaknesses of the algorithm, we evaluate its effectiveness on a wide range of *ad hoc* retrieval data sets. The TREC data sets used in our experiments are summarized in Appendix A.

All collections were stopped using a standard list of 418 common terms and stemmed using a Porter stemmer. Only the title portion of the TREC topics are used to construct queries. Our primary evaluation metric is mean average precision. Statistical significance is determined using a one-tailed paired *t* – *test* evaluated at the $p < 0.05$ level.

We now describe our feature candidate pool in terms of the feature representation scheme we proposed in Section 3.3. For dependence model type, the features may be either *FI* (full independence), *SD* (sequential dependence), or *FD* (full dependence). The clique set type may either be T_{QD} , O_{QD} , or U_{QD} . The weighting functions include LM, BM25, [LM, BM25]-O-[1, 2, 4, 8, 16, or 32], and [LM, BM25]-U-[1, 2, 4, 8, 16, 32, or unlimited]. As we see, the pool is very robust and covers many different types of important features. It includes features that span all three type of dependence, use all three types of clique sets, allow both Dirichlet or BM25 weighting, and vary the window sizes for the ordered and unordered window matchings across a wide range

	No Retrain		Retrain	
	Train	Test	Train	Test
AP	0.1863	0.2266	0.1865	0.2246
WSJ	0.2700	0.3553	0.2703	0.3543
ROBUST04	0.2387	0.3079	0.2391	0.3065
WT10G	0.2344	0.2129	0.2357	0.2140

Table 6.1. Training and test set mean average precision values for no retraining and retraining.

of values. In total, after removing trivial and duplicate features, our candidate pool consists of 48 features.

For all of our experiments, features are added until there is no change in training set mean average precision between iterations (i.e., $\epsilon = 0$) or until we have added 5 features. Preliminary experiments showed that adding more than 5 features never resulted in significantly different training or test set results.

6.3.1 No Retraining vs. Retraining

We wish to analyze what effect, if any, retraining (see Algorithm 4, line 10) has on training and generalization properties of the model. Table 6.1 summarizes the mean average precision obtained on the training and test set when retraining is used and when it is not.

We first investigate whether or not the models learned with retraining vary significantly from those learned without retraining. As Table 6.1 shows, the training set mean average precision values for no retraining and retraining are nearly equivalent for every data set. In fact, the differences are statistically indistinguishable. In addition, we discovered that the same set of features were added regardless of whether or not retraining was done or not. Therefore, it appears as though retraining has little effect on the learned model, both in terms of the features selected and the training set mean average precision.

Next, we study the effect of retraining on the generalization properties of the model. As the test set results in Table 6.1 show, there is very little difference in mean average precision for no retraining versus retraining. The results, again, are statistically indistinguishable for every data set. Hence, retraining does not significantly affect how well the model generalizes to unseen data.

Therefore, given that retraining requires more computational power, has no effect on either the learned model or the generalization properties of the model, we conclude that there is no need to retrain the model each iteration.

6.3.2 Number of Features

We now analyze how sensitive the models are to the number of parameters, both in terms of potential overfitting, and in terms of test set effectiveness.

Figure 6.1 plots the training and test set mean average precision versus the number of features that have been added to the model. As the figure indicates, there appears to be little, if any overfitting happening. The test set mean average precision never significantly drops as more features are added to the model.

6.3.3 Feature Analysis

The greedy nature of our feature selection algorithm provides us with a mechanism for analyzing the importance of different types of features across data sets. By looking at the order in which features are selected, and the weight assigned to each, we can develop deeper insights into the role that features play for a given task and/or data set.

For example, for the WT10G data set, with no retraining, the features are selected in the following order:

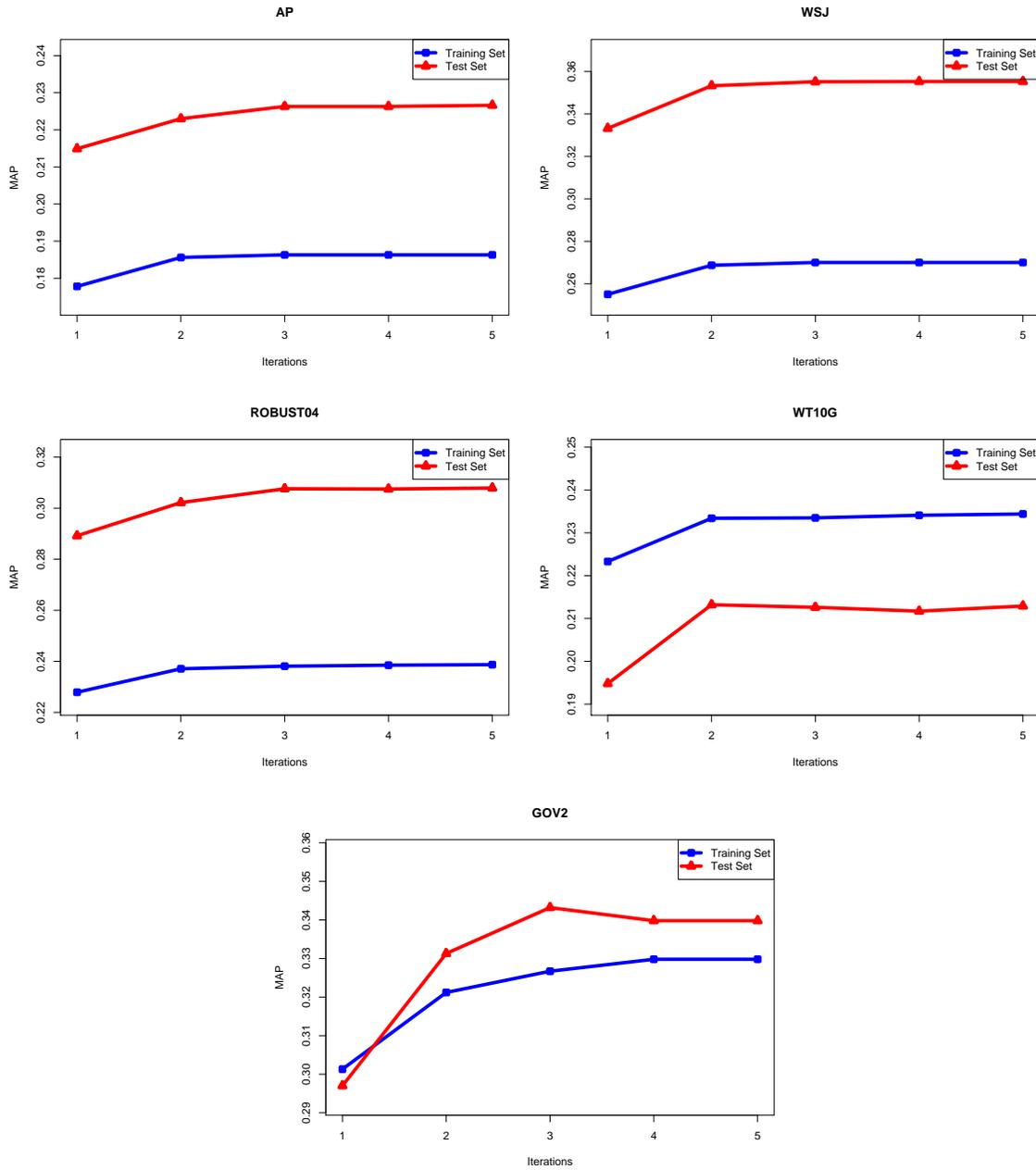


Figure 6.1. Mean average precision versus number of iterations for the training and test sets of the AP, WSJ, ROBUST04, WT10G and GOV2 data sets.

(FI, T_{QD} , BM25) :	0.8138
(FD, U_{QD} , LM-U-8) :	0.0001
(SD, U_{QD} , BM25-U-unlimited) :	0.0090
(FD, U_{QD} , BM25-U-8) :	0.1575
(SD, O_{QD} , BM25-O-8) :	0.0196

where the numbers after the colons are the weights assigned to each feature in the final model.

As Figure 6.1 shows, there is a large increase in both training and test set mean average precision after the second feature, LM-U-8, is added to the model. This large increase, which is also exhibited for the GOV2 data set, reiterates the importance of term proximity models for large web collections [65]. We see that simply adding a single proximity feature increases mean average precision substantially. However, there is a much smaller effect observed after further term proximity / dependence model features are added to the model.

To provide a different example, we consider the order in which features are selected for the WSJ collection. The features selected, in order, are:

(FI, T_{QD} , BM25) :	0.5864
(SD, U_{QD} , BM25-U-1) :	0.3746
(FD, U_{QD} , BM25-U-32) :	0.0193
(FI, T_{QD} , LM) :	0.0196
(FD, U_{QD} , BM25-U-unlimited) :	0.0001

As with the WT10G model, the first feature selected is the full independence, single term, BM25 feature. In fact, this feature was the first selected for every data

	MRF-FI	BM25	MRF-SD	MRF-BM25	MRF-FS
AP	0.2077	0.2149	0.2128	0.2210	0.2266†‡
WSJ	0.3258	0.3332	0.3429†	0.3512†	0.3553†‡
ROBUST04	0.2920	0.2892	0.3092†	0.3101†	0.3079†
WT10G	0.1861	0.1948	0.2140†	0.2129†	0.2129†
GOV2	0.2984	0.2971	0.3360†	0.3476†	0.3398†

Table 6.2. Comparison of test set mean average precision for language modeling (MRF-FI), BM25, MRF model using language modeling weighting (MRF-SD), MRF model using BM25 weighting (MRF-BM25), and MRF learned using our proposed feature selection algorithm (MRF-FS). A † indicates a statistically significant improvement over *both* the MRF-FI and BM25 models and a ‡ indicates a significant improvement over the MRF-BM25 model.

set. This is not surprising, however, since the overwhelming importance of single term features has long been understood.

However, no other strong regularities were observed across data sets. This indicates that the each data set has unique characteristics that make certain features more discriminative than others. Such characteristics may include things like query length, noise, document length distribution, and properties of the underlying vocabulary. This suggests that no single model, with a fixed feature set and fixed feature weights, can be applied to every possible task and data set. Instead, adaptive models and techniques, such as the one presented here, can provide a means for automatically and robustly learning the best set of features to use on a task-by-task basis.

6.3.4 Summary of Results

Finally, we compare the retrieval effectiveness of the models automatically learned using our feature selection algorithm (MRF-FS) with several other retrieval models, including language modeling (MRF-FI model), BM25, and two MRF models with hand selected features (MRF-SD and MRF-BM25, as defined in Chapter 5) that we have been shown to be highly effective. For each model, parameters are tuned on the training set to maximize mean average precision. Therefore, every model is properly

trained in accordance with the same evaluation metric. This allows us to compare the effectiveness of automatically learned models with models that use manually chosen features and have been proven to be highly effective.

Our results, which are summarized in Table 6.2, support previous observations that show that using MRF models with hand chosen features are generally more effective than bag of words models for *ad hoc* retrieval. However, we are interested in how effective the automatically learned models are. For both the AP and WSJ data sets, the mean average precision of the automatically learned model is statistically significantly better than all of the other models, including the MRF models with manually chosen features. The improvement in mean average precision over BM25 for the AP data set is 5.4% and 6.6% on the WSJ data set.

On the ROBUST04, WT10G, and GOV2 data sets, the automatically learned models are statistically significantly better than language modeling and BM25, but statistically indistinguishable from the two models with hand selected features. Despite the lack of a statistically significant improvement over the models with hand selected features, the results still provide evidence that the learned model is highly effective. Indeed, compared to BM25, the automatically learned model is 6.5% better for ROBUST04, 9.3% better for WT10G, and 14.4% better for GOV2.

Therefore, our results show that our proposed feature selection algorithm produces very effective models that are competitive with, and often significantly better than models with hand selected features. The most important result, however, is that *every* automatically learned model was significantly better than the two state of the art bag of words models. This result is very powerful, as it shows that using this framework, models can be automatically learned from a rich set of features and are very likely to be significantly better than the best hand crafted bag of words models.

CHAPTER 7

LATENT CONCEPT EXPANSION

Users of information retrieval systems are required to express complex information needs in terms of Boolean expressions, a short list of keywords, a sentence, a question, or possibly a longer narrative. A great deal of information is lost during the process of translating from the information need to the actual query. For this reason, there has been a strong interest in query expansion techniques. Such techniques are used to augment the original query to produce a representation that better reflects the underlying information need.

Query expansion techniques have been well studied for various models in the past and have shown to significantly improve effectiveness in both the relevance feedback and pseudo-relevance feedback setting [58, 95, 118, 122].

Until now, our MRF-based models have been solely used for ranking documents in response to a given query. In this chapter, we show how these models can be extended and used for query expansion using a technique that we call *latent concept expansion* (LCE). There are three primary contributions of our work.

First, LCE provides a mechanism for combining term dependence with query expansion. Previous query expansion techniques are based on bag of words models. Therefore, by performing query expansion using the MRF model, we are able to study the dynamics between term dependence and query expansion.

Next, query expansion techniques in the past have implicitly only made use of term occurrence features. By using more powerful feature sets, such as those we have

proposed earlier, it is possible to produce better expansion terms that discriminate between relevant and non-relevant documents better.

Finally, our proposed approach seamlessly provides a mechanism for generating both single and multi-term concepts. Most previous techniques, by default, generate terms independently. There have been several approaches that make use of generalized concepts, however such approaches were somewhat heuristic and done outside of the model [81, 118]. Our approach is both formally motivated and a natural extension of the underlying model.

Before describing the details of LCE and formally evaluating it, we review related work in the area of query expansion.

7.1 Related Work

One of the classic and most widely used approaches to query expansion is the Rocchio algorithm [95]. Rocchio’s approach, which was developed within the vector space model, reweights the original query vector by moving the weights towards the set of relevant or pseudo-relevant documents and away from the non-relevant documents. Unfortunately, it is not possible to formally apply Rocchio’s approach to a statistical retrieval model, such as language modeling for information retrieval.

A number of formalized query expansion techniques have been developed for the language modeling framework, including Zhai and Lafferty’s model-based feedback and Lavrenko and Croft’s relevance models [58, 122]. Both approaches attempt to use pseudo-relevant or relevant documents to estimate a better query model.

Model-based feedback finds the model that best describes the relevant documents while taking a background (noise) model into consideration. This separates the content model from the background model. The content model is then interpolated with the original query model to form the expanded query.

The other technique, relevance models, is more closely related to our work. Therefore, we go into the details of the model. Much like model-based feedback, relevance models estimate an improved query model. The only difference between the two approaches is that relevance models do not explicitly model the relevant or pseudo-relevant documents. Instead, they model a more generalized notion of relevance, as we now show.

Given a query Q , a relevance model is a multinomial distribution over the vocabulary, $P(\cdot|Q)$, that encodes the likelihood of each term given the query as evidence. It is computed as:

$$\begin{aligned}
 P(w|Q) &= \int_D P(w|D)P(D|Q) \\
 &\approx \frac{\sum_{D \in \mathcal{R}_Q} P(w|D)P(Q|D)P(D)}{\sum_w \sum_{D \in \mathcal{R}_Q} P(w|D)P(Q|D)P(D)} \tag{7.1}
 \end{aligned}$$

where \mathcal{R}_Q is the set of documents that are relevant or pseudo-relevant to query Q . In the pseudo-relevant case, these are the top ranked documents for query Q . Furthermore, it is assumed that $P(D)$ is uniform over this set. These mild assumptions make computing the Bayesian posterior more practical.

After the model is estimated, documents are ranked by clipping the relevance model by choosing the k most likely terms from $P(\cdot|Q)$. This clipped distribution is then interpolated with the original, maximum likelihood query model [71]. This can be thought of as expanding the original query by k weighted terms. Throughout the remainder of this work, we refer to this instantiation of relevance models as RM3.

There has been relatively little work done in the area of query expansion in the context of dependence models [47]. However, there have been several attempts to expand using multi-term concepts. Xu and Croft’s local context analysis (LCA) method combined passage-level retrieval with concept expansion, where concepts were single terms and phrases [118]. Expansion concepts were chosen and weighted using

a metric based on co-occurrence statistics. However, it is not clear based on the analysis done how much the phrases helped over the single terms alone.

Papka and Allan investigate using relevance feedback to perform multi-term concept expansion for document routing [81]. The concepts used in their work are more general than those used in LCA, and include Indri query language structures, such as #UW50(white house), which corresponds to the concept “the terms white and house occur, in any order, within 50 terms of each other”. Results showed that combining single term and large window multi-term concepts significantly improved effectiveness. However, it is unclear whether the same approach is also effective for *ad hoc* retrieval, due to the differences in the tasks.

7.2 Latent Concept Expansion

We now describe how our MRF framework can be used in a novel way to generate single and multi-term concepts that are topically related to some original query. As we will show, the concepts generated using our technique can be used for query expansion or other tasks, such as suggesting alternative query formulations.

We assume that when a user formulates their original query, they have some set of concepts in mind, but are only able to express a small number of them in the form of a query. We treat the concepts that the user has in mind, but did not explicitly express in the query, as latent concepts. These latent concepts can consist of a single term, multiple terms, or some combination of the two. It is our goal to recover these latent concepts given some original query.

This can be accomplished within our framework by first expanding the original MRF graph G to include the type of concept we are interested in generating. We call this expanded graph H . In Figure 7.1, the middle graph provides an example of how to construct an expanded graph that can generate single term concepts. Similarly, the graph on the bottom illustrates an expanded graph that generates two term concepts.

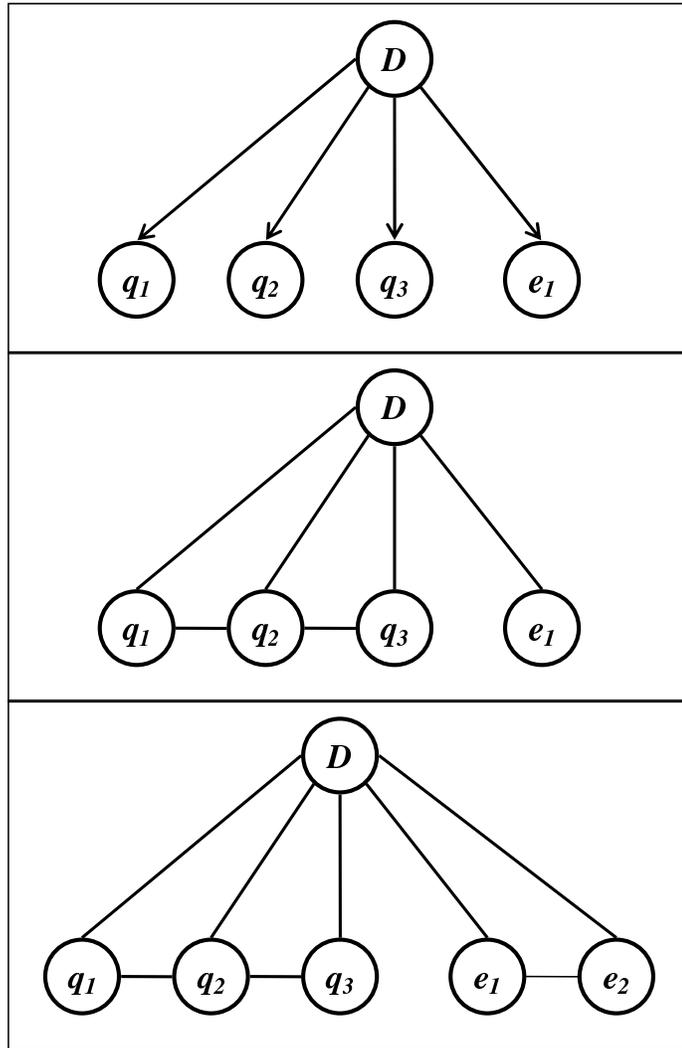


Figure 7.1. Graphical model representations of relevance modeling (top), latent concept expansion using single term concepts (middle), and latent concept expansion using two term concepts (bottom) for a three term query.

Although these two examples make use of the sequential dependence model (i.e., dependencies between adjacent query terms), it is important to note that both the original query and the expansion concepts can use any dependence structure.

After H is constructed, we compute $P_{H,\Lambda}(E|Q)$, a probability distribution over latent concepts, according to:

$$P_{H,\Lambda}(E|Q) = \frac{\sum_{D \in \mathcal{R}} P_{H,\Lambda}(Q, E, D)}{\sum_{D \in \mathcal{R}} \sum_E P_{H,\Lambda}(Q, E, D)} \quad (7.2)$$

where \mathcal{R} is the universe of all possible documents and E is some latent concept that may consist of one or more terms. Since it is not practical to compute this summation, we must approximate it. We notice that $P_{H,\Lambda}(Q, E, D)$ is likely to be peaked around those documents D that are highly ranked according to query Q . Therefore, we approximate $P_{H,\Lambda}(E|Q)$ by only summing over a small subset of relevant or pseudo-relevant documents for query Q . This is computed as follows:

$$P_{H,\Lambda}(E|Q) \approx \frac{\sum_{D \in \mathcal{R}_Q} P_{H,\Lambda}(Q, E, D)}{\sum_{D \in \mathcal{R}_Q} \sum_E P_{H,\Lambda}(Q, E, D)} \quad (7.3)$$

where \mathcal{R}_Q is a set of relevant or pseudo-relevant documents for query Q and all clique sets are constructed using H .

As an example, suppose that we were to perform LCE on a MRF-SD model and expand by single term concepts. Then, the graph H would be like the middle one in Figure 7.1. Under this setting, expansion concepts would be generated in the following way:

$$\begin{aligned}
P_{H,\Lambda}(e|Q) \propto \sum_{D \in \mathcal{R}_Q} \exp \left[\right. & \lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \\
& \lambda_{O_D} \sum_{(q_1, q_2, D) \in O_{QD}} \log \frac{tf_{\#1(q_1 q_2), D} + \mu^w \frac{cf_{\#1(q_1 q_2)}}{|C|}}{|D| + \mu^w} + \\
& \lambda_{U_D} \sum_{(q_1, q_2, D) \in U_{QD}} \log \frac{tf_{\#uw8(q_1 q_2), D} + \mu^w \frac{cf_{\#uw8(q_1 q_2)}}{|C|}}{|D| + \mu^w} + \\
& \lambda'_{T_D} \log \frac{tf_{e, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \\
& \left. \lambda'_{T_Q} \log \frac{|C|}{cf_e} \right] \tag{7.4}
\end{aligned}$$

As we see, the likelihood contribution for each document in \mathcal{R}_Q is a combination of the original query's score for the document (first 3 components, refer to Equation 5.4), concept E 's score for the document (fourth component), and E 's document-independent score (fifth component). Therefore, this equation can be interpreted as measuring how well Q and E account for the top ranked documents and the “goodness” of E , independent of the documents. For maximum flexibility, we introduce a new set of parameters for the expansion concept (i.e., λ'_{T_D} and λ'_{T_Q}), which allows us to weight the expansion features differently than those in the original query.

For the sake of completeness, we show a more complex example, where the original query uses the MRF-FD model and the expansion concept, which consists of two terms, uses the MRF-SD model. Under this model, concept probabilities are computed according to:

$$\begin{aligned}
P_{H,\Lambda}(e_1, e_2|Q) \propto \sum_{D \in \mathcal{R}_Q} \exp & \left[\lambda_{T_D} \sum_{(q_i, D) \in T_{QD}} \log \frac{tf_{q_i, D} + \mu^t \frac{cf_{q_i}}{|C|}}{|D| + \mu^t} + \right. \\
& \lambda_{O_D} \sum_{(q_1, \dots, q_k, D) \in O_{QD}} \log \frac{tf_{\#1(q_1 \dots q_k), D} + \mu^w \frac{cf_{\#1(q_1 \dots q_k)}}{|C|}}{|D| + \mu^w} + \\
& \lambda_{U_D} \sum_{(q_1, \dots, q_k, D) \in O_{QD} \cup U_{QD}} \log \frac{tf_{\#\text{uw}8(q_1 \dots q_k), D} + \mu^w \frac{cf_{\#\text{uw}8(q_1 \dots q_k)}}{|C|}}{|D| + \mu^w} + \\
& \lambda'_{T_D} \sum_{i=1}^2 \log \frac{tf_{e_i, D} + \mu^t \frac{cf_{e_i}}{|C|}}{|D| + \mu^t} + \\
& \lambda'_{O_D} \log \frac{tf_{\#1(e_1 e_2), D} + \mu^w \frac{cf_{\#1(e_1 e_2)}}{|C|}}{|D| + \mu^w} + \\
& \lambda'_{U_D} \log \frac{tf_{\#\text{uw}8(e_1 e_2), D} + \mu^w \frac{cf_{\#\text{uw}8(e_1 e_2)}}{|C|}}{|D| + \mu^w} + \\
& \lambda'_{T_Q} \sum_{i=1}^2 \log \frac{|C|}{cf_{e_i}} + \\
& \lambda'_{O_Q} \log \frac{|C|}{cf_{\#1(e_1 e_2)}} + \\
& \left. \lambda'_{U_Q} \log \frac{|C|}{cf_{\#\text{uw}8(e_1 e_2)}} \right] \tag{7.5}
\end{aligned}$$

where the first three components compute the MRF-FD score of the original query, the next three components compute the MRF-SD score of the expansion concept, and the last three components compute an IDF-based “goodness” score of the concept.

7.2.1 Query Expansion

To use this framework for query expansion, we first choose an expansion graph H that encodes the latent concept structure we are interested in expanding the query using. We then select the k latent concepts with the highest likelihood given by Equation 7.3. A new graph G' is constructed by augmenting the original graph G with the k expansion concepts e_1, \dots, e_k . Finally, documents are ranked according to $P_{G',\Lambda}(D|Q, e_1, \dots, e_k)$.

7.2.2 Comparison to Relevance Models

Inspecting Equations 7.1 and 7.3 reveals the close connection that exists between LCE and relevance models. Both models essentially compute the likelihood of a term (or concept) in the same manner. It is easy to see that just as the MRF model can be viewed as a generalization of language modeling, so too can LCE be viewed as a generalization of relevance models. In fact, by setting $\lambda_{T_D} = \lambda'_{T_D} = 1$ and all other parameters to 0 in Equation 7.4, we derive the relevance model formula.

There are important differences between MRFs/LCE and unigram language models/relevance models. See Figure 7.1 for graphical model representations of both models. Unigram language models and relevance models are based on the multinomial distribution. This distributional assumption locks the model into the bag of words representation and the implicit use of term occurrence features. However, the distribution underlying the MRF model allows us to move beyond both of these assumptions, by modeling both dependencies between query terms and allowing arbitrary features to be explicitly used.

Moving beyond the simplistic bag of words assumption in this way results in a general, robust model and, as we show in the next section, translates into significant improvements in retrieval effectiveness.

7.3 Experimental Results

In order to better understand the strengths and weaknesses of our technique, we evaluate it on a wide range of data sets. Appendix A provides a summary of the TREC data sets considered. For each data set, we split the available topics into a training and test set, where the training set is used solely for parameter estimation and the test set is used for evaluation purposes.

In all cases, only the title portion of the TREC topics are used to construct queries. We construct G using the sequential dependence assumption for all data sets [65].

	MRF-FI	MRF-SD	RM3	LCE
AP	.2077	.2147 ^{α}	.2518 ^{$\alpha\beta$}	.2692 ^{$\alpha\beta\gamma$}
WSJ	.3258	.3425 ^{α}	.3493 ^{α}	.3943 ^{$\alpha\beta\gamma$}
ROBUST04	.2920	.3096 ^{α}	.3382 ^{$\alpha\beta$}	.3601 ^{$\alpha\beta\gamma$}
WT10g	.1861	.2053 ^{α}	.1944 ^{α}	.2269 ^{$\alpha\beta\gamma$}
GOV2	.3234	.3520 ^{α}	.3656 ^{α}	.3924 ^{$\alpha\beta\gamma$}

Table 7.1. Test set mean average precision for MRF-FI, MRF-SD, relevance models (RM3), and latent concept expansion (LCE). The superscripts α , β , and γ indicate statistically significant improvements over MRF-FI, MRF-SD, and RM3, respectively.

7.3.1 Ad Hoc Retrieval Results

We now investigate how well our model performs in practice in a pseudo-relevance feedback setting. We compare the MRF-FI model, the MRF-SD model (without expansion), relevance models, and LCE to better understand how each model performs across the various data sets.

For the MRF models, we train the model parameters (i.e., Λ) and model hyperparameters (i.e., μ^t , μ^w). For RM3 and LCE, we also train the number of pseudo-relevant feedback documents used and the number of expansion terms.

7.3.1.1 Expansion with Single Term Concepts

We begin by evaluating how well our model performs when expanding using only single terms. The expansion term likelihoods are computed according to Equation 7.4. The equation clearly shows how LCE differs from relevance models. As we stated before, when we set $\lambda_{T_D} = \lambda'_{T_D} = 1$ and all other parameters to 0, we obtain the exact formula that is used to compute term likelihoods in the relevance modeling framework. Therefore, LCE adds two very important factors to the equation. First, it adds the ordered and unordered window features that are applied to the original query. Second, it applies an intuitive *tf.idf*-like form to the candidate expansion term w . The *idf* factor, which is not present in relevance models, plays an important role in expansion term selection.

	MRF-FI	MRF-SD	RM3	LCE
AP	.3460	.3340	.3640 ^{β}	.3720 ^{β}
WSJ	.4860	.5080	.4980	.5400 ^{$\alpha\beta\gamma$}
ROBUST04	.4293	.4566 ^{α}	.4576 ^{α}	.4798 ^{$\alpha\gamma$}
WT10g	.3204	.3245	.3265	.3633 ^{$\alpha\beta\gamma$}
GOV2	.5180	.5680 ^{α}	.6160 ^{α}	.6060 ^{α}

Table 7.2. Test set precision at 10 for MRF-FI, MRF-SD, relevance models (RM3), and latent concept expansion (LCE). The superscripts α , β , and γ indicate statistically significant improvements over MRF-FI, MRF-SD, and RM3, respectively.

The results, evaluated using mean average precision, are given in Table 7.1. As we see, the MRF-SD model, relevance models, and LCE always significantly outperform the bag of words MRF-FI model. In addition, LCE shows significant improvements over relevance models across all data sets. The relative improvements over relevance models is 6.9% for AP, 12.9% for WSJ, 6.5% for ROBUST04, 16.7% for WT10G, and 7.3% for GOV2.

We also note the interesting result that the MRF-SD model is statistically equivalent to relevance models on the two web data sets. In fact, the MRF-SD model outperforms relevance models on the WT10g data set. This reiterates the importance of non-unigram, proximity-based features for content-based web search observed previously [65, 72].

We also evaluate the methods according to precision at 10 in Table 7.2. Not surprisingly, the precision at 10 improvements achieved using LCE were not as significant as those observed when using mean average precision. It is well understood that expansion is recall enhancing, rather than precision enhancing. The MRF model significantly improves upon the unigram language model on every data set, but only in a few cases do relevance models or LCE significantly improve upon the MRF model.

It should also be noted that, despite the fact that our model has more free parameters than relevance models, there is surprisingly little overfitting. Instead, it exhibits good generalization properties.

	One Term	One + Two Term
AP	0.2692	0.2648
WSJ	0.3943	0.3945
ROBUST04	0.3601	0.3464
WT10G	0.2269	0.2187
GOV2	0.3924	0.3900

Table 7.3. Test set mean average precision values for multi-term concept LCE experiments.

<i>price fixing</i>		<i>tax evasion indicted</i>	
price	fixed pricing	evasion	tax evasion
fixed	price fixing	tax	income tax
market	control control	indicted	tax charges
report	market price	federal	los angeles
vote	united states	charges	san diego

Table 7.4. One and two term expansion concepts for the query *price fixing* (ROBUST04 topic 622) and *tax evasion indicted* (ROBUST04 topic 650). Concepts are listed in descending order of $P(e|Q)$ and $P(e_1, e_2|Q)$, respectively.

7.3.1.2 Expansion with Multi-Term Concepts

We also investigated expanding using both single and two word concepts. For each query, we expanded using a set of single term concepts, as well as a set of two term concepts. The sets were chosen independently. The results of the experiments are shown in Table 7.3. As the results show, there is little or no improvement when including two term concepts. In fact, the results are statistically indistinguishable for all data sets.

This unexpected result may be due to the fact that strong correlations exist between the single term expansion concepts. We found that the two word concepts chosen often consisted of two highly correlated terms that are also chosen as single term concepts. For example, for the query *price fixing*, there was a great deal of redundancy. Table 7.4 shows the single and two term concepts that were selected for expansion. We see “fixed pricing” and “price fixing”, “price”, and “fixed” all occur in the list. While these are excellent expansion terms, the fact that they appear so many

times likely results in terms “price” and “fixing” being overemphasized. Additionally, this query adds “control control” and “united states”, which are poor expansion concepts. These two concepts are given high probability because of the features used within the model. A similar analysis holds for the query *tax evasion indicted*. Therefore, other feature sets may ultimately yield different results, especially if they reduce the correlation among the expansion concepts. While the IDF-like feature may be appropriate for choosing single terms, it may be better to apply some other type of feature to the ordered and unordered cliques within the graph. A better understanding of phrases and named entities (e.g., San Diego, Los Angeles, United States in Table 7.4), and their effect on query expansion is still needed.

Furthermore, due to the computational complexity involved in computing multi-term expansion, we had to limit the number of multi-term expansion concepts that we used. We were unable to perform a full sweep over all such values. However, we made sure to focus on the most promising values, and, based on previous experience, it is unlikely that the results would have been significantly better than the one term expansion concepts, but perhaps would have been slightly improved over their current value.

Therefore, our experiments yield no conclusive results with regard to expansion using multi-term concepts. Instead, the results introduce interesting open questions and directions for future exploration.

7.3.2 Robustness

As we have shown, relevance models and latent concept expansion can significantly improve retrieval effectiveness over the baseline MRF-FI model. In this section we analyze the robustness of these two methods.

Figure 7.2 provides an analysis of the robustness of relevance modeling and latent concept expansion for the AP, WSJ, ROBUST04, and WT10G data sets. The

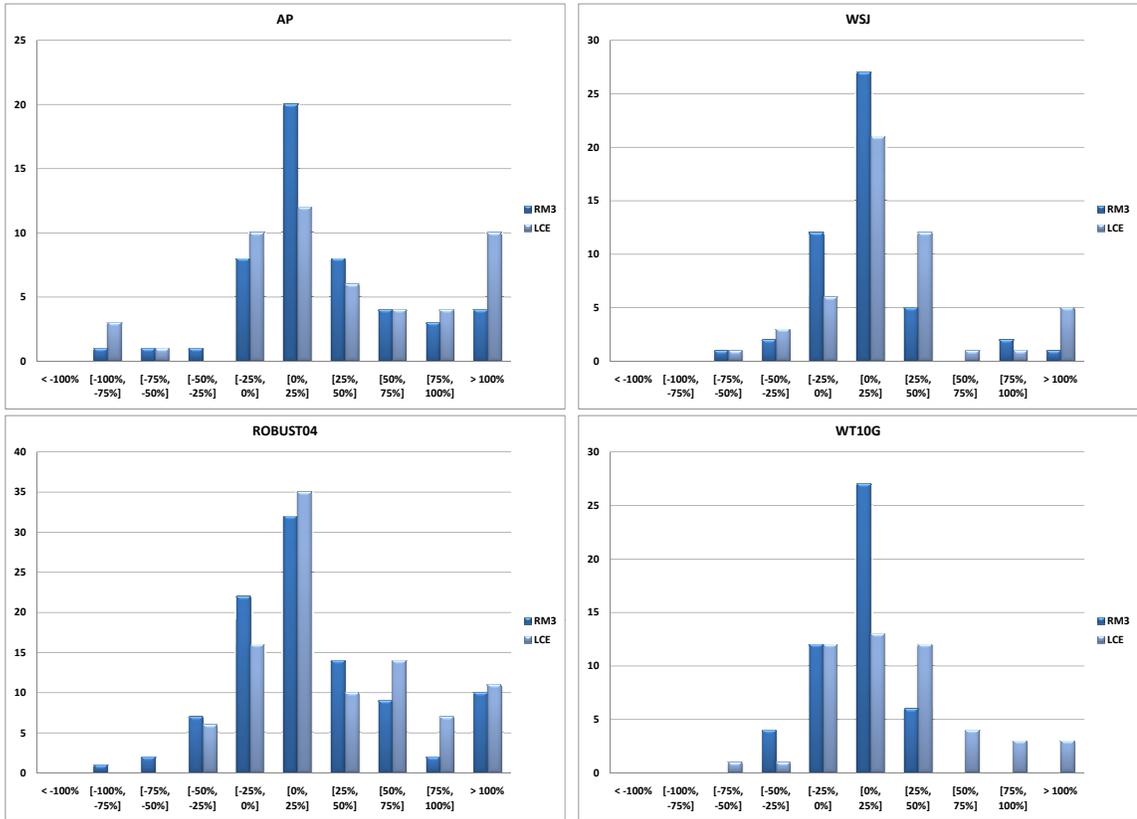


Figure 7.2. Histograms that demonstrate and compare the robustness of relevance models (RM3) and latent concept expansion (LCE) with respect to the MRF-FI model for the AP, WSJ, ROBUST04, and WT10G data sets.

analysis for GOV2 is similar. The histograms provide, for various ranges of relative decreases/increases in mean average precision, the number of queries that were hurt/improved with respect to the MRF-FI baseline.

As the results show, LCE exhibits strong robustness for each data set. For AP, relevance models improve 38 queries and hurt 11, whereas LCE improves 35 and hurts 14. Although relevance models improve the effectiveness of 3 more queries than LCE, the relative improvement exhibited by LCE is significantly larger. For the ROBUST04 data set, relevance models improve 67 queries and hurt 32, and LCE improves 77 and hurts 22. Finally, for the WT10G collection, relevance models improve 32 queries and hurt 16, and LCE improves 35 and hurts 14. As with AP, the amount of improvement exhibited by the LCE versus relevance models is significantly larger for both the ROBUST04 and WT10G data sets. In addition, when LCE does hurt performance, it is less likely to hurt as much as relevance modeling, which is a desirable property.

Overall, LCE improves effectiveness for 65%-80% of queries, depending on the data set. When used in combination with a highly accurate query performance prediction system, it may be possible to selectively expand queries and minimize the loss associated with sub-baseline performance.

7.3.3 Multi-Term Concept Generation

Although we found that expansion using multi-term concepts failed to produce conclusive improvements in effectiveness, there are other potential tasks that these concepts may be useful for, such as query suggestion/reformulation, summarization, and concept mining. For example, for a query suggestion task, the original query can be used to generate a set of latent concepts which correspond to alternative query formulations.

1 word concepts	2 word concepts	3 word concepts
telescope	hubble telescope	hubble space telescope
hubble	space telescope	hubble telescope space
space	hubble space	space telescope hubble
mirror	telescope mirror	space telescope NASA
NASA	telescope hubble	hubble telescope astronomy
launch	mirror telescope	NASA hubble space
astronomy	telescope NASA	space telescope mirror
shuttle	telescope space	telescope space NASA
test	hubble mirror	hubble telescope mission
new	NASA hubble	mirror mirror mirror
discovery	telescope astronomy	space telescope launch
time	telescope optical	space telescope discovery
universe	hubble optical	shuttle space telescope
optical	telescope discovery	hubble telescope flaw
light	telescope shuttle	two hubble space

Table 7.5. Fifteen most likely one, two, and three word concepts constructed using the top 25 documents retrieved for the query *hubble telescope achievements* on the ROBUST04 collection.

Although evaluating our model on these tasks is beyond the scope of this work, we wish to show an illustrative example of the types of concepts generated using our model. In Table 7.5, we present the most likely one, two, and three term concepts generated using LCE for the query *hubble telescope achievements* using the top 25 ranked documents from the ROBUST04 collection.

It is well known that generating multi-term concepts using a unigram-based model produces unsatisfactory results, since it fails to consider term dependencies. This is not the case when generating multi-term concepts using our model. Instead, a majority of the concepts generated are well-formed and meaningful. There are several cases where the concepts are less coherent, such as *mirror mirror mirror*. In this case, the likelihood of the term *mirror* appearing in a pseudo-relevant document outweighs the ICF features, which causes this non-coherent concept to have a high likelihood. Such examples are in the minority, however.

Not only are the concepts generated well-formed and meaningful, but they are also topically relevant to the original query. As we see, all of the concepts generated are on topic and in some way related to the Hubble telescope. It is interesting to see that the concept *hubble telescope flaw* is one of the most likely three term concepts, given that it is somewhat contradictory to the original query. Despite this contradiction, documents that discuss the telescope flaws are also likely to describe the successes, as well, and therefore this is likely to be a meaningful concept.

One important thing to note is that the concepts LCE generates are of a different nature than those that would be generated using a bigram relevance model. For example, a bigram model would be unlikely to generate the concept *telescope space NASA*, since none of the bigrams that make up the concept have high likelihood. However, since our model is based on a number of different features over various types of cliques, it is more general and robust than a bigram model.

Although we only provided the concepts generated for a single query, we note that the same analysis and conclusions generalize across other data sets, with coherent, topically related concepts being consistently generated using LCE.

7.4 Discussion

We conclude this chapter by discussing several theoretical issues involved with regard to query expansion.

7.4.1 Relevance vs. Relevant Documents

There are other reasonable ways of using the MRF model for query expansion beyond using Equation 7.3. For example, given a set of relevant or pseudo-relevant documents D_1, \dots, D_n , we could generate expansion concepts by computing $P(E|D_1, \dots, D_n)$.

The concepts generated in this way tend to show excellent effectiveness on the training set. However, they fail to yield significant improvements in retrieval ef-

fectiveness on the test set. In fact, using this expansion model shows little-to-no improvements over relevance models. It is somewhat surprising that this formulation performs so much worse on the test set. Buckley and Salton noted a similar observation when optimizing term weights using a Rocchio-based expansion technique [14]. They argued that there is a subtle, yet important difference between modeling relevance and modeling relevant documents.

Under this alternative formulation, we are directly modeling relevant documents. By doing so, we are overfitting the model. This explains the poor generalization behavior exhibited using this alternative formulation. This, however, is not an issue when using the original formulation. Using that formulation, $P_{H,\Lambda}(E|Q)$ can be thought of as a model of relevance, since it generates terms that are highly likely given the query, which is our best evidence of relevance. This avoids the problem of overfitting the model to the relevant documents by modeling a more generalized notion of relevance.

Therefore, the results obtained by Salton and Buckley using the Rocchio-based expansion technique apply to statistical query expansion techniques, as well. Thus, future query expansion models, to avoid overfitting, should focus on modeling more generalized notions of relevance, rather than the relevant documents themselves.

7.4.2 The Role of Dependence

Our latent concept expansion technique captures two semi-orthogonal types of dependence. In information retrieval, there has been a long-term interest in understanding the role of term dependence. Out of this research, two broad types of dependencies have been identified.

The first type of dependence is *syntactic dependence*. This type of dependence covers phrases, term proximity, and term co-occurrence [22, 27, 29, 35, 110]. These

methods capture the fact that queries implicitly or explicitly impose a certain set of positional dependencies.

The second type is *semantic dependence*. Examples of semantic dependence are relevance feedback, pseudo-relevance feedback, synonyms, and to some extent stemming [23]. These techniques have been explored on both the query and document side. On the query side, this is typically done using some form of query expansion, such as relevance models or LCE. On the document side, this is done as document expansion or document smoothing [55, 61, 108].

Although there may be some overlap between syntactic and semantic dependencies, they are mostly orthogonal. Our model uses both types of dependencies. The use of phrase and proximity features within the model captures syntactic dependencies, whereas LCE captures query-side semantic dependence. This explains why the initial improvement in effectiveness achieved by using the MRF model is not lost after query expansion. If the same types of dependencies were captured by both syntactic and semantic dependencies, LCE would be expected to perform about equally as well as relevance models. Therefore, by modeling both types of dependencies we see an additive effect, rather than an absorbing effect.

An interesting area of future work is to determine whether or not modeling document-side semantic dependencies can add anything to the model. Previous results that have combined query- and document-side semantic dependencies have shown mixed results [61, 115].

CHAPTER 8

CONCLUSIONS AND FUTURE WORK

This chapter provides a broad summary of our work. We begin by summarizing how to use the MRF model. We then present a comprehensive summary of results and reiterate our primary contributions. Finally, we conclude by proposing several potential directions of future work.

8.1 Using the MRF Model

To utilize the Markov random field model in a general retrieval setting, the following steps should be carried out:

1. **Construct graph** The graph G should be constructed in a manner that captures dependencies between query terms. We proposed the full independence, sequential dependence, and full dependence as three generalized dependency constructs that are meaningful to information retrieval.
2. **Select features** A set of potential functions over the cliques of G must be defined. We proposed using clique sets to reduce the number of potential functions and parameters within the model. The clique sets defined are based on cliques containing single query terms, ordered query terms, and unordered sets of query terms. We also proposed a set of primitive and high-level features to be used with these clique sets. These features can either be manually chosen or automatically chosen using the method we proposed in Section 6

3. **Parameter estimation** We must estimate the model parameters in some way. This can be done completely unsupervised, semi-supervised, or completely supervised. In Section 4.4, we proposed a coordinate-ascent based method that directly maximizes the metric under consideration and described other related methods for learning to rank.
4. For each query Q ,
 - **Rank documents** Rank documents in descending order of $P_{G,\Lambda}(D|Q)$, which can be computed efficiently using Equation 3.19.
 - **Expand query** (optional) Use the method proposed in Section 7 to expand the original query using either true or blind feedback based on the top ranked documents.

8.2 Summary of *Ad Hoc* Retrieval Results

Ad hoc retrieval was the primary task that we used to evaluate the effectiveness of various retrieval models. Hence, we now summarize all of our key *ad hoc* retrieval results. The results, which are in Tables 8.1, 8.2, 8.3, and 8.4, include MAP, GMAP, P@10, and R-Prec results for language modeling, BM25, bigram language modeling, MRF-SD, MRF-FD, MRF-BM25, MRF-FS, RM3, and LCE. The bold value in each column represents the best effectiveness achieved using a model that does not use pseudo-relevance feedback. In every case, the best effectiveness is achieved by a non-bag of words model. Furthermore, in every case but one, the model with the best effectiveness is a model based on our MRF framework. This nicely summarizes the superior effectiveness possible when using these types of models.

8.3 Contributions

We now summarize the major contributions developed throughout this work:

	AP	WSJ	ROBUST04	WT10G	GOV2
LM	0.2077	0.3258	0.2920	0.1861	0.2984
BM25	0.2149	0.3332	0.2892	0.1948	0.2971
BIGRAM	0.2116	0.3319	0.3012	0.2165	0.3371
MRF-SD	0.2147	0.3425	0.3096	0.2053	0.3325
MRF-FD	0.2128	0.3429	0.3092	0.2140	0.3360
MRF-BM25	0.2210	0.3512	0.3101	0.2129	0.3476
MRF-FS	0.2264	0.3595	0.3079	0.2077	0.3398
RM3	0.2518	0.3493	0.3382	0.1944	0.3656
LCE	0.2692	0.3943	0.3601	0.2269	0.3924

Table 8.1. Test set mean average precision across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.

	AP	WSJ	ROBUST04	WT10G	GOV2
LM	0.1219	0.2267	0.1970	0.1176	0.1891
BM25	0.1342	0.2344	0.1937	0.1055	0.1926
BIGRAM	0.1229	0.2313	0.2076	0.1347	0.2137
MRF-SD	0.1268	0.2399	0.2196	0.1284	0.2448
MRF-FD	0.1255	0.2405	0.2196	0.1361	0.2424
MRF-BM25	0.1366	0.2471	0.2199	0.1181	0.2817
MRF-FS	0.1411	0.2565	0.2108	0.1132	0.2526
RM3	0.1436	0.2392	0.2235	0.1204	0.2612
LCE	0.1420	0.2728	0.2516	0.1450	0.2842

Table 8.2. Test set geometric mean average precision for across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.

	AP	WSJ	ROBUST04	WT10G	GOV2
LM	0.3460	0.4860	0.4293	0.3204	0.5180
BM25	0.3200	0.5020	0.4354	0.3224	0.5440
BIGRAM	0.3420	0.4880	0.4343	0.3061	0.5500
MRF-SD	0.3340	0.5080	0.4566	0.3245	0.5680
MRF-FD	0.3540	0.5080	0.4505	0.3469	0.5720
MRF-BM25	0.3140	0.5140	0.4525	0.3388	0.6100
MRF-FS	0.3340	0.5200	0.4455	0.3347	0.6140
RM3	0.3640	0.4980	0.4576	0.3265	0.6160
LCE	0.3720	0.5400	0.4798	0.3633	0.6060

Table 8.3. Test set precision at 10 across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.

	AP	WSJ	ROBUST04	WT10G	GOV2
LM	0.2448	0.3558	0.3291	0.2199	0.3515
BM25	0.2632	0.3640	0.3226	0.2324	0.3437
BIGRAM	0.2537	0.3561	0.3339	0.2499	0.3744
MRF-SD	0.2580	0.3633	0.3362	0.2374	0.3716
MRF-FD	0.2543	0.3694	0.3394	0.2417	0.3763
MRF-BM25	0.2666	0.3698	0.3366	0.2508	0.3834
MRF-FS	0.2753	0.3825	0.3324	0.2448	0.3740
RM3	0.2846	0.3730	0.3589	0.2291	0.3753
LCE	0.2969	0.4066	0.3769	0.2548	0.4010

Table 8.4. Test set precision at R (R-prec) across a range of retrieval models. The model parameters were trained to maximize mean average precision. Bold value indicates the best technique that does not make use of pseudo-relevance feedback.

1. **Robust retrieval model.** We develop a new, formally motivated, statistical retrieval model based on Markov random fields that robustly and effectively handles term dependencies and the combination of arbitrary features.
2. **Better understanding of features for information retrieval.** By modeling dependencies between terms and encoding rich features, such as those based on phrases and term proximity, we are able to better understand how and when such features can improve retrieval effectiveness.
3. **Novel parameter estimation technique.** Our technique exploits the nature of rank-equivalence and works to directly maximize the underlying retrieval metric, which leads to better performance than maximizing the data likelihood or margin. This avoids the problem of metric divergence.
4. **Automatic model learning.** We propose a supervised feature selection algorithm that can be used to automatically learn highly effective models. This eliminates the need for human experts to manually select model features on a per-task basis.
5. **Concept-based query expansion.** Our retrieval model provides an elegant mechanism for expanding queries using multi-term concepts in the context of relevance or pseudo-relevance feedback.
6. **State of the art retrieval effectiveness.** Our model shows consistent and significant improvements in retrieval effectiveness over current state of the art retrieval models on *ad hoc* retrieval and web search tasks.

8.4 Future Work

In this work, we laid out the basic foundations of the MRF model for information retrieval. Although we tackled many critical issues, several challenges remain. We now briefly describe some of these challenges.

- **Feature engineering** - We presented a flexible framework for constructing MRFs for information retrieval tasks. However, for the *ad hoc* task, we focused on a very specific set of features/weighting functions based on term proximity features using both language modeling and BM25 weighting schemes. For the web search task, we used similar features plus priors based on inlink count and PageRank-based priors. Since the effectiveness of the MRF model depends so heavily on the features used within the model, it is important to engineer better features that can be applied across various tasks. Such feature engineering should, of course, not be done in a heuristic, uninformed manner. Instead, we advocate that a formal framework be developed for discovering new features and analyzing their usefulness.
- **Better understanding of multi-term concept expansion** - In Chapter 7, we presented LCE, our general query expansion framework for the MRF model. The framework is unique in that it allows term dependence and arbitrary features to be used. Another appealing characteristic of the framework is that it allows queries to be expanded with multi-term concepts. However, our preliminary results showed there was no benefit of including such concepts and that single term concepts remained the best option. We believe that it should be possible to use multi-term concepts in order to improve effectiveness. However, a deeper understanding of multi-term concepts and their role in expansion must first be developed. There are a number of ways to tackle such a problem, but one of the most straightforward approaches is to develop a better set of query-dependent features to be used with LCE. We believe that such features

should somehow take IDF into account, some notion of lexical cohesion, and some measure of redundancy.

- **Unified learning to rank framework** - The importance of careful tuning and parameter estimation cannot be stressed enough. If a good model is poorly trained, there is little hope of achieving good effectiveness. As models begin to incorporate more features, the need to develop robust learning to rank techniques grows. As we described in Chapter 4, there have been a number of approaches proposed to solve the learning to rank problem. The coordinate ascent approach that we proposed is able to handle arbitrary evaluation metrics, but is not guaranteed to find a global maxima and may not scale well when there are hundreds or thousands of features. Many of the machine learning approaches maximize some surrogate function that may or may not be related to some evaluation metric. While these approaches often are guaranteed to find a global maxima and scale well, they typically can not be applied to different evaluation metrics easily. Therefore, there is a need to combine the pros and the cons from both types of approaches into a powerful, flexible, efficient framework for learning to rank.
- **Efficient indexing architectures** - In order to support efficient query evaluation of MRF model queries, there is a need to develop efficient indexing architectures. Such indexing architectures must be able to index arbitrary features in a compact, efficient manner. Furthermore, query evaluation strategies must be developed to be used with such indexing architectures. We believe this is a critical component to the wide spread adoption of such models.

APPENDIX A

DATA SETS

In order to properly analyze the various dimensions of retrieval effectiveness, it is important to evaluate new models and techniques against a diverse set of data sets. We make use of TREC collections in all of our experiments, Therefore, we begin this Appendix by describing the anatomy of a TREC data set. We then provide the specific details of the data sets used in the experiments throughout this work.

A.1 Anatomy of a TREC Data Set

Every TREC data set consists of a set of *documents*, *topics*, and *relevance judgments*. We now provide an overview of each.

TREC documents, which are in SGML format, contain a number of metadata fields, as well as content. See Figure A.1 for an example TREC document. Examples of metadata fields include `DOCNO`, which is a unique document identifier and `HEAD`, which is the headline of the news article. The document’s textual content is contained within the `TEXT` field. In our experiments, we throw away all of the metadata, except for the document identifier, and index the content contained in the `TEXT` field¹.

A TREC topic typically consists of a `title`, `description`, and a `narrative`, although some topics may contain other fields. See Figure A.2 for an example TREC topic. It is important to note that a topic represents an information need and is not a query itself. Instead, researchers must distill a query from a topic. There are a

¹TREC web documents do not have a `TEXT` field. For these document, we index all of the text, except that contained in the `DOCHDR` field.

```

<DOC>
<DOCNO> AP890101-0010 </DOCNO>
<FILEID>AP-NR-01-01-89 1123EST</FILEID>
<FIRST>r i AM-Thatcher-Women 01-01 0287</FIRST>
<SECOND>AM-Thatcher-Women,0295</SECOND>
<HEAD>Thatcher Says Male Prime Ministers May Eventually Be Fashionable
Again</HEAD>
<DATELINE>LONDON (AP) </DATELINE>
<TEXT>
Prime Minister Margaret Thatcher, who made history
in 1979 when she became Europe's first woman prime minister, noted
Sunday she is not alone and joked that male politicians may one day
come back into fashion.
"We're getting more women prime ministers," she said in a
television interview, referring to the recent election of Benazir
Bhutto as prime minister of Pakistan.
"And don't forget ... Mrs. Gandhi was a very able, charming,
formidable prime minister of India."
Mrs. Thatcher now is the longest serving leader in the West.
Before she came to power, women had governed in Sri Lanka and
Israel. Part of the British leader's tenure in office coincided with
that of Mrs. Gandhi, who was assassinated in 1984.
"I think male prime ministers one day will come back into
fashion," she joked with interviewer David Frost on Britain's
commercial TV-am channel.
Asked about combining her job and her domestic life with her
husband, Denis, a retired oil executive, she said "women have run
both the home and work for a very long time."
"I mean, every working wife knows that if you decide to have
steak and kidney pie for supper, it's no better if you took 20
minutes thinking about it than if you took 20 seconds."
Frost recalled Mrs. Thatcher's comment "I never did understand
men," during an acrimonious meeting last year with fellow leaders
of the European Economic Community. Asked if she understood men
better now, Mrs. Thatcher replied:
"It may not be understanding of the deepest kind, but I do know
what they're likely to do and say. So, one has a certain
predictability about it."
</TEXT>
</DOC>

```

Figure A.1. Example TREC document.

```
<top>
<num> Number: 763

<title> Hunting deaths

<desc> Description:
Give information on human deaths associated with hunting for game.

<narr> Narrative:
Accidental deaths, murders, and suicides are relevant. Deaths can be
from any cause. Fatalities of people not in the hunting party are
relevant, but the deaths must be connected with hunting. Relevant
hunting must be for live prey. Deaths related to submarine hunting
are not relevant.

</top>
```

Figure A.2. Example TREC topic.

number of ways of doing this, but the most common way is to use one of the fields as the query. In most of the experiments done throughout this work, we distill a query by using only the text that appears in the title field.

Finally, relevance judgments are provided. For each topic, a set of documents are manually judged for relevance. Different definitions and scales of relevance are used for different tasks. Please refer to Sections 5.1 and 5.2 for examples. Furthermore, due to limited resources, not all documents are judged for every topic. Instead, TREC uses a number of techniques, including document pooling, in order to reduce the number of documents judged, while ensuring the test collection is reusable [45, 104]. Figure A.3 shows an excerpt from a TREC relevance judgment file. These judgments are used as the “ground truth” for the purposes of evaluation, which we cover in Appendix B.

763 0 GX018-79-11508357 2
763 0 GX018-79-15198972 0
763 0 GX019-18-8997173 0
763 0 GX019-40-2241667 0
763 0 GX019-56-11902532 1
763 0 GX019-97-10746106 0
763 0 GX020-45-4344240 0
763 0 GX020-99-3364231 0
763 0 GX021-19-6737921 0
763 0 GX021-41-4253461 0

Figure A.3. Portion of a TREC relevance judgment file. The format of each line is: query-id 0 doc-id judgment. Judgments of 0, 1, and 2 refer to non-relevant, relevant, and highly relevant, respectively.

Name	Description	# Docs	Train Topics	Test Topics
WSJ	Wall St. Journal '87-'92	173,252	51–150	151–200
AP	Associated Press '88-'90	242,918	51–150	151–200
ROBUST04	Newswire articles	528,155	301–450	601–700
WT10g	Small web crawl	1,692,096	451–500	501–550
GOV2	2004 crawl of .gov domain	25,205,179	701–750	751–800

Table A.1. Overview of TREC collections and topics used in most of our experiments.

	Disks 1,2	Disk 3	Disks 4,5
Num. Docs	741,856	336,310	556,077
Training topics	101-150	51-100	301-350
Test topics	151-200	101-150	401-450

Table A.2. TREC data sets used in Chapter 4. The disk numbers refer to the TREC volumes used to construct the index.

A.2 Summary of Data Sets

Now that we have explained the composition of a TREC data set, we summarize the details of the data sets considered in this work. First, Table A.1 gives an overview of five primary data sets used in our experiments. The table includes the data set name, a short description, the number of documents in the data set, the topics used for training and the topics used for testing. As we see, these data sets vary in size, homogeneity, and noisiness. The AP and WSJ data sets are small collections of newswire articles from a single source. Therefore, they are homogeneous and contain very little noise. The ROBUST04 data set is medium sized and consists of news stories from a number of sources, thus making it less homogeneous, and not very noisy. Finally, the WT10G and GOV2 data sets are (very) large web crawls, thus making them both heterogeneous and noisy. Hence, the data sets are diverse across a number of characteristics, making this a suitable evaluation testbed. Nearly all of the experiments done throughout this work use the data sets listed in this table.

Another set of data sets is also considered, but only used for one experiment. The data sets shown in Table A.2 are used in the experiments described in Chapter 4. All of these data sets consist primarily of news articles, are somewhat heterogeneous, and contain very little noise. The collections are also relatively small, as well. The table also lists the TREC topics used for training and testing purposes.

APPENDIX B

EVALUATION METRICS

This appendix provides a brief summary of metrics that are commonly used to evaluate the effectiveness of information retrieval systems. Each of these evaluation measures are rank-based. That is, they depend exclusively on the ranking of documents produced by a system. The scores of the documents do not influence the metrics in any way. In addition, a binary model of relevance is assumed, which states that an item is either *relevant* to a request or it is *not relevant*.

Given a ranked list of documents in response to a query, we define the function R such that $R(i) = 1$ if the document at rank i is relevant and $R(i) = 0$ if the document at rank i is not relevant. Furthermore, as a matter of convenience, we define the function $R(1,k)$ that counts the number of relevant documents between rank 1 and k (inclusive) as

$$R(1,k) = \sum_{i=1}^k R(i) \tag{B.1}$$

In Table B.1 the functional forms for precision at rank k , R-Prec, success at rank k , average precision, and reciprocal rank are given.

Precision at rank k measures the proportion of relevant documents that are ranked in the top k , whereas success at rank k is 1 if *any* relevant documents appear in the top k and 0 otherwise. Typical values of k range from 5 to 100. Precision at rank 5 or 10, for example, is typically used to evaluate a system's ability to return relevant documents at the top of the ranked list, which is important for many applications, including web search.

Metric Name	Functional Form
Precision at rank k	$P@k = \frac{R(1,k)}{k}$
Success at rank k	$S@k = \delta(R(1,k) \geq 1)$
R-Prec	$R\text{-Prec} = \frac{R(1, R)}{ R }$
Average Precision	$\text{AvgP} = \frac{1}{ R } \sum_{i:R(i)=1} \frac{R(1,i)}{i}$
Reciprocal Rank	$RR = \max \left\{ \frac{1}{i} : R(i) = 1 \right\}$

Table B.1. Summary of common information retrieval evaluation metrics, where $R(1, k)$ is defined in Equation B.1 and $|R|$ is the total number of judged relevant documents.

R-prec is a special case of precision at rank k . R-Prec computes the precision at rank $|R|$, where $|R|$ is the total number of judged relevant documents. This measure is more adaptive, as it computes precision to a different depth for every query.

Average precision can be thought of as a weighted precision measure that gives higher weight to relevant documents that appear near the top of the ranked list. The measure is computed by averaging the precision at k for every k such that the document at rank k is relevant. Average precision is standardly computed to a depth of 1000 documents¹. Any relevant documents that do not appear in the ranked list are assumed to be at rank “infinity”, and therefore contribute nothing to the average. Therefore, the measure rewards systems that rank relevant documents high in the ranked list and those that return more relevant documents. In this way, the measure implicitly accounts for both precision and recall (coverage of relevant documents). Average precision is typically used to evaluate *ad hoc* retrieval tasks and other tasks where both precision and recall are important factors.

The last measure, reciprocal rank, is computed according to the reciprocal of the highest ranked relevant document. Note that the measure quickly decays as the rank of the first relevant document increases (e.g., 1, 0.5, 0.33, 0.25, 0.2, etc.). We note

¹Throughout this work, all evaluation metrics are computed to a depth of 1000.

Name	Value
Precision at rank k	$P@k = \frac{1}{N} \sum_{i=1}^N P@k(q_i)$
Success at rank k	$S@k = \frac{1}{N} \sum_{i=1}^N S@k(q_i)$
R-Prec	$R\text{-Prec} = \frac{1}{N} \sum_{i=1}^N R\text{-Prec}(q_i)$
Mean Average Precision	$MAP = \frac{1}{N} \sum_{i=1}^N \text{AvgP}(q_i)$
Geometric Mean Average Precision	$GMAP = \prod_{i=1}^N \text{AvgP}'(q_i)^{1/N}$
Mean Reciprocal Rank	$MRR = \frac{1}{N} \sum_{i=1}^N RR(q_i)$

Table B.2. Overview of aggregate measures. For each aggregate measure we show how it is computed. Here, N refers to the total number of queries being aggregated and q_i is the i^{th} query in the set. Notice that GMAP is zero if any query has an average precision of zero. In order to correct this, $\text{AvgP}'(q_i)$ is defined to be $\max(\text{AvgP}(q_i), .00001)$.

that reciprocal rank is a special case of average precision when there is only a single relevant document for a given request. This measure is commonly used to evaluate known-item finding queries, where it is critical to return the relevant web page very high in the ranked list.

Each of these measures we just described are computed on a query-by-query basis. However, for the sake of comparison, we often need to compute a single measure from a set of query-by-query measures. Given a set of queries, the *arithmetic average* or the *geometric average* are often used to combine, or aggregate, the measures computed on a query-by-query basis. Table B.2 summarizes the most common aggregate measures used in information retrieval. In order to avoid confusion, the most commonly used names of the evaluation metrics were used. This results in aggregate measures having the same name as their non-aggregate counterparts (e.g., R-Prec is the name of both the aggregate and non-aggregate measure). We note, however, that throughout this work, it should be clear from the context of the discussion whether the measure in question has been computed for a single query or if it was aggregated over a set of queries.

BIBLIOGRAPHY

- [1] Agichtein, E., Brill, E., and Dumais, S. Improving web search ranking by incorporating user behavior information. In *Proc. 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 19–26.
- [2] Amati, G., and van Rijsbergen, C. J. Probabilistic models of information retrieval based on measuring the divergence from randomness. *ACM Transactions on Information Systems* 20, 4 (2002), 357–389.
- [3] Anh, V. N., and Moffat, A. Simplified similarity scoring using term ranks. In *Proc 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 226–233.
- [4] Anh, V. N., and Moffat, A. Pruned query evaluation using pre-computed impacts. In *Proc 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 372–379.
- [5] Baron, J.R., Lewis, D.D., and Oard, D. TREC 2006 legal track overview. In *Proc. 15th Text REtrieval Conference* (2006).
- [6] Bartell, B. T., Cottrell, G. W., and Belew, R. K. Automatic combination of multiple ranked retrieval systems. In *Proc. 17th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1994), pp. 173–181.
- [7] Berger, A., and Lafferty, J. Information retrieval as statistical translation. In *Proc. 22nd Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1999), pp. 222–229.
- [8] Blei, D., Griffiths, T., Jordan, M., and Tenenbaum, J. Hierarchical topic models and the nested Chinese restaurant process. In *Proc. 16th Proc. of Advances in Neural Information Processing Systems* (2003).
- [9] Blei, D., Ng, A., and Jordan, M. Latent dirichlet allocation. *Journal of Machine Learning Research* 3 (2003), 993–1022.
- [10] Brin, S., and Page, L. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems* 30, 1–7 (1998), 107–117.
- [11] Broder, A. A taxonomy of web search. *SIGIR Forum* 36, 2 (2002), 3–10.
- [12] Buckley, C. Why current IR engines fail. In *Proc. 27th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2004), pp. 584–585.

- [13] Buckley, C., Dimmick, D., Soboroff, I., and Voorhees, E. Bias and the limits of pooling. In *Proc. 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 619–620.
- [14] Buckley, C., and Salton, G. Optimization of relevance feedback weights. In *Proc. 18th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1995), pp. 351–357.
- [15] Burges, C., Shaked, T., Renshaw, E., Lazier, A., Deeds, M., Hamilton, N., and Hullender, G. Learning to rank using gradient descent. In *Proc. 22nd Proc. Intl. Conference on Machine Learning* (2005), pp. 89–96.
- [16] Burges, C. J. C. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery* 2, 2 (1998), 121–167.
- [17] Burges, C. J.C., Ragno, R., and Le, Q. V. Learning to rank with nonsmooth cost functions. In *Proc. 19th Proc. of Advances in Neural Information Processing Systems* (2007), pp. 193–200.
- [18] Büttcher, S., Clarke, C. L. A., and Lushman, B. Term proximity scoring for ad-hoc retrieval on very large text collections. In *Proc. 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 621–622.
- [19] Chow, C., and Liu, C. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory* 14, 3 (1968), 462–467.
- [20] Clarke, C., Cormack, G., and Burkowski, F. Shortest substring ranking (Multi-Text experiments for TREC-4). In *Proc. 4th Text REtrieval Conference* (1995).
- [21] Clarke, C., Scholar, F., and Soboroff, I. Overview of the TREC 2005 terabyte track. In *Proc. 14th Text REtrieval Conference* (2006).
- [22] Clarke, C. L. A., and Cormack, G. V. Shortest-substring retrieval and ranking. *ACM Trans. Inf. Syst.* 18, 1 (2000), 44–78.
- [23] Collins-Thompson, K., and Callan, J. Query expansion using random walk models. In *Proc. 14th Intl. Conf. on Information and Knowledge Management* (2005), pp. 704–711.
- [24] Cooper, W. S. Some inconsistencies and misnomers in probabilistic information retrieval. In *Proc. 14th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1991), pp. 57–61.
- [25] Craswell, N., de Vries, A.P., and Soboroff, I. Overview of the TREC 2005 enterprise track. In *Proc. 14th Text REtrieval Conference* (2005).

- [26] Craswell, N., Robertson, S., Zaragoza, H., and Taylor, M. Relevance weighting for query independent evidence. In *Proc. 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 416–423.
- [27] Croft, W. B. Boolean queries and term dependencies in probabilistic retrieval models. *Journal of the American Society for Information Science* 37, 4 (1986), 71–77.
- [28] Croft, W. B., and Harper, D. Using probabilistic models of information retrieval without relevance information. *Journal of Documentation* 35, 4 (1979), 285–295.
- [29] Croft, W. B., Turtle, H., and Lewis, D. The use of phrases and structured queries in information retrieval. In *Proc. 14th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1991), pp. 32–45.
- [30] de Kretser, O., and Moffat, A. Effective document presentation with a locality-based similarity heuristic. In *Proc. 22nd Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1999), pp. 113–120.
- [31] Deerwester, S., Dumais, S., Furnas, G. W., Landauer, T. K., and Harshman, R. Indexing by latent semantic analysis. *Journal of the Society for Information Science* 41, 6 (1990), 391–407.
- [32] Della Pietra, S., Della Pietra, V., and Lafferty, J. Inducing features of random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19, 4 (1997), 380–393.
- [33] Diaz, F. Regularizing ad hoc retrieval scores. In *Proc. 14th Intl. Conf. on Information and Knowledge Management* (2005), pp. 672–679.
- [34] Diaz, F., and Metzler, D. Improving the estimation of relevance models using large external corpora. In *Proc. 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 154–161.
- [35] Fagan, J. Automatic phrase indexing for document retrieval: An examination of syntactic and non-syntactic methods. In *Proc. tenth Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1987), pp. 91–101.
- [36] Fan, W., Gordon, M. D., and Pathak, P. A generic ranking function discovery framework by genetic programming for information retrieval. *Inf. Process. Manage.* 40, 4 (2004), 587–602.
- [37] Fang, H., and Zhai, C. An exploration of axiomatic approaches to information retrieval. In *Proc. 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 480–487.

- [38] Fang, H., and Zhai, C. Semantic term matching in axiomatic approaches to information retrieval. In *Proc. 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 115–122.
- [39] Fox, E. A., and Shaw, J. A. Combination of multiple searches. In *Proc. 2nd Text REtrieval Conference* (1993), pp. 243–252.
- [40] Fuhr, N. Probabilistic models in information retrieval. *The Computer Journal* 35, 3 (1992), 243–255.
- [41] Gao, J., Nie, J., Wu, G., and Cao, G. Dependence language model for information retrieval. In *Proc. 27th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2004), pp. 170–177.
- [42] Gao, J., Qi, H., Xia, X., and Nie, J. Linear discriminant model for information retrieval. In *Proc. 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 290–297.
- [43] Gey, F. Inferring probability of relevance using the method of logistic regression. In *Proc. 17th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1994).
- [44] Griffiths, T. L., Steyvers, M., Blei, D. M., and Tenenbaum, J. B. Integrating topics and syntax. In *Proc. 17th Proc. of Advances in Neural Information Processing Systems* (2005), pp. 537–544.
- [45] Harman, D. Overview of the first TREC conference. In *Proc. 16th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1993), pp. 36–47.
- [46] Harman, D. Overview of the TREC 2002 novelty track. In *Proc. 11th Text REtrieval Conference* (2002).
- [47] Harper, D., and van Rijsbergen, C. J. An evaluation of feedback in document retrieval using co-occurrence data. *Journal of Documentation* 34, 3 (1978), 189–216.
- [48] Harter, S. P. A probabilistic approach to automatic keyword indexing. *Journal of the American Society for Information Science* 26, 5 (1975), 197–206 and 280–289.
- [49] Hofmann, T. Probabilistic latent semantic indexing. In *Proc. 22nd Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1999), pp. 50–57.
- [50] Joachims, T. Optimizing search engines using clickthrough data. In *Proc. 8th Ann. Intl. ACM SIGKDD Conf. on Knowledge Discovery and Data Mining* (2002), pp. 133–142.

- [51] Joachims, T. A support vector method for multivariate performance measures. In *Proc. 22nd Proc. Intl. Conference on Machine Learning* (2005), pp. 377–384.
- [52] Joachims, T., Granka, L., Pan, B., Hembrooke, H., and Gay, G. Accurately interpreting clickthrough data as implicit feedback. In *Proc. 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 154–161.
- [53] Kleinberg, J. Authoritative sources in a hyperlinked environment. *Journal of the ACM* 46, 5 (1999), 604–632.
- [54] Kraaij, W., Westerveld, T., and Hiemstra, D. The importance of prior probabilities for entry page search. In *Proc. 25th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2002), pp. 27–34.
- [55] Kurland, O., and Lee, L. Corpus structure, language models, and ad hoc information retrieval. In *Proc. 27th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2004), pp. 194–201.
- [56] Lafferty, J., and Zhai, C. Document language models, query models, and risk minimization for information retrieval. In *Proc. 24th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2001), pp. 111–119.
- [57] Lavrenko, V. *A Generative Theory of Relevance*. PhD thesis, University of Massachusetts, Amherst, MA, 2004.
- [58] Lavrenko, V., and Croft, W. B. Relevance-based language models. In *Proc. 24th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2001), pp. 120–127.
- [59] Le, Q., and Smola, A. Direct optimization of ranking measures. <http://www.citebase.org/abstract?id=oai:arXiv.org:0704.3359>, 2007.
- [60] Lebanon, G., and Lafferty, J. Hyperplane margin classifiers on the multinomial manifold. In *Proc. 21st Proc. Intl. Conference on Machine Learning* (2004), pp. 66–71.
- [61] Liu, X., and Croft, W. B. Cluster-based retrieval using language models. In *Proc. 27th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2004), pp. 186–193.
- [62] Matveeva, I., Burges, C., Burkard, T., Laucius, A., and Wong, L. High accuracy retrieval with multiple nested ranker. In *Proc. 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 437–444.
- [63] McCallum, A. Efficiently inducing features of conditional random fields. In *Proc. 19th Conf. on Uncertainty in Artificial Intelligence* (2003).

- [64] Metzler, D., and Croft, W. B. Combining the language model and inference network approaches to retrieval. *Information Processing and Management* 40, 5 (2004), 735–750.
- [65] Metzler, D., and Croft, W. B. A markov random field model for term dependencies. In *Proc. 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 472–479.
- [66] Metzler, D., Diaz, F., Strohman, T., and Croft, W. B. UMass robust 2005: Using mixtures of relevance models for query expansion. In *Proc. 14th Text REtrieval Conference* (2005).
- [67] Metzler, D., Dumais, S., and Meek, C. Similarity measures for short segments of text. In *Proc. 29th European Conf. on Information Retrieval* (2007), pp. 16–27.
- [68] Metzler, D., Lavrenko, V., and Croft, W. B. Formal multiple bernoulli models for language modeling. In *Proc. 27th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2004), pp. 540–541.
- [69] Metzler, D., and Manmatha, R. An inference network approach to image retrieval. In *Proc. 3rd Intl. Conf. on Image and Video Retrieval* (2004), pp. 42–50.
- [70] Metzler, D., Strohman, T., and Croft, W. B. Lessons learned from three terabyte tracks. In *Proc. 15th Text REtrieval Conference* (2006).
- [71] Metzler, D., Strohman, T., Turtle, H., and Croft, W. B. Indri at TREC 2004: Terabyte track. In *Proc. 13th Text REtrieval Conference* (2004).
- [72] Metzler, D., Strohman, T., Zhou, Y., and Croft, W. B. Indri at TREC 2005: Terabyte track. In *Proc. 14th Text REtrieval Conference* (2005).
- [73] Mishne, G., and de Rijke, M. Boosting web retrieval through query operations. In *Proc. 27th European Conf. on Information Retrieval* (2005), pp. 502–516.
- [74] Morgan, W., Greiff, W., and Henderson, J. Direct maximization of average precision by hill-climbing with a comparison to a maximum entropy approach. Tech. rep., MITRE, 2004.
- [75] Morik, K., Brockhausen, P., and Joachims, T. Combining statistical learning with a knowledge-based approach - a case study in intensive care monitoring. In *Proc. 16th Proc. Intl. Conference on Machine Learning* (1999).
- [76] Nallapati, R. Discriminative models for information retrieval. In *Proc. 27th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2004), pp. 64–71.
- [77] Nallapati, R., and Allan, J. Capturing term dependencies using a language model based on sentence trees. In *Proc. 11th Intl. Conf. on Information and Knowledge Management* (2002), pp. 383–390.

- [78] Ogilvie, P., and Callan, J. Combining document representations for known-item search. In *Proc. 26th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2003), pp. 143–150.
- [79] Ounis, I., de Rijke, M., MacDonald, C., Mishne, G., and Soboroff, I. Overview of the TREC-2006 blog track. In *Proc. 15th Text REtrieval Conference* (2006).
- [80] Pang, B., Lee, L., and Vaithyanathan, S. Thumbs up? Sentiment classification using machine learning techniques. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (2002), pp. 79–86.
- [81] Papka, R., and Allan, J. Why bigger windows are better than smaller ones. Tech. rep., University of Massachusetts, Amherst, 1997.
- [82] Peng, J., Macdonald, C., He, B., Plachouras, V., and Ounis, I. Incorporating term dependency in the dfr framework. In *Proc 30th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2007), p. To appear.
- [83] Peng, Y., and He, D. Direct comparison of commercial and academic retrieval system: an initial study. In *Proc. 15th Intl. Conf. on Information and Knowledge Management* (2006), pp. 806–807.
- [84] Ponte, J., and Croft, W. Bruce. A language modeling approach to information retrieval. In *Proc. 21st Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1998), pp. 275–281.
- [85] Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P. *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge University Press, 1992.
- [86] R. Losee, Jr. Term dependence: Truncating the Bahadur Lazarsfeld expansion. *Information Processing and Management* 30, 2 (1994), 293–303.
- [87] Robertson, S. The probability ranking principle in IR. *Journal of Documentation* 33, 4 (1977), 294–303.
- [88] Robertson, S. Understanding inverse document frequency: on theoretical arguments for IDF. *Journal of Documentation* 60, 5 (2004), 503–520.
- [89] Robertson, S., and Walker, S. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *Proc. 17th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1994), pp. 232–241.
- [90] Robertson, S., Walker, S., Jones, S., Hancock-Beaulieu, M. M., and Gatford, M. Okapi at TREC-3. In *Proc. 3rd Text REtrieval Conference* (1994), pp. 109–126.

- [91] Robertson, S., Zaragoza, H., and Taylor, M. Simple bm25 extension to multiple weighted fields. In *Proc. 13th Intl. Conf. on Information and Knowledge Management* (2004), pp. 42–49.
- [92] Robertson, S. E., and Spärck Jones, K. Relevance weighting of search terms. *Journal of the American Society for Information Science* 27, 3 (1976), 129–146.
- [93] Robertson, S. E., van Rijsbergen, C. J., and Porter, M. F. Probabilistic models of indexing and searching. In *Proc. 3rd Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1980), pp. 35–56.
- [94] Robertson, S. E., and Walker, S. On relevance weights with little relevance information. In *Proc. 20th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1997), pp. 16–24.
- [95] Rocchio, J. J. *Relevance Feedback in Information Retrieval*. Prentice-Hall, 1971, pp. 313–323.
- [96] Rosenfeld, R. Two decades of statistical language modeling: Where do we go from here? *Proceedings of IEEE* 88, 8 (2000), 1270–1278.
- [97] Salton, G., and Buckley, C. Term-weighting approaches in automatic text retrieval. *Information Processing and Management* 24, 5 (1988), 513–523.
- [98] Shen, X., and Zhai, C. Active feedback in ad hoc information retrieval. In *Proc. 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 59–66.
- [99] Si, L., and Callan, J. A statistical model for scientific readability. In *Proc. 10th Intl. Conf. on Information and Knowledge Management* (2001), pp. 574–576.
- [100] Singhal, A., Buckley, C., and Mitra, M. Pivoted document length normalization. In *Proc. 19th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1996), pp. 21–29.
- [101] Song, F., and Croft, W. B. A general language model for information retrieval. In *Proc. 8th Intl. Conf. on Information and Knowledge Management* (1999), pp. 316–321.
- [102] Spärck Jones, K. *Automatic keyword classification for information retrieval*. Butterworths, 1971.
- [103] Spärck Jones, K. Wearing proper combinations. Tech. rep., University of Cambridge, 2005.
- [104] Spärck Jones, K., and van Rijsbergen, C. J. Information retrieval test collections. *Journal of Documentation* 32, 1 (1976), 59–72.

- [105] Srikanth, M., and Srihari, R. Biterm language models for document retrieval. In *Proc. 25th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2002), pp. 425–426.
- [106] Strohman, T., and Croft, W. B. Efficient document retrieval in main memory. In *Proc 30th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2007), p. To appear.
- [107] Strohman, T., Metzler, D., Turtle, H., and Croft, W. B. Indri: A language model-based search engine for complex queries. In *Proceedings of the International Conference on Intelligence Analysis* (2004).
- [108] Tao, T., Wang, X., Mei, Q., and Zhai, C. Language model information retrieval with document expansion. In *Proc. of HLT/NAACL* (2006), pp. 407–414.
- [109] Turtle, H., and Croft, W. B. Evaluation of an inference network-based retrieval model. *ACM Transactions on Information Systems* 9, 3 (1991), 187–222.
- [110] van Rijsbergen, C. J. A theoretical basis for the use of cooccurrence data in information retrieval. *Journal of Documentation* 33, 2 (1977), 106–119.
- [111] Vechtomova, O., Karamuftuoglu, M., and Robertson, S. E. On document relevance and lexical cohesion between query terms. *Information Processing and Management* 42, 5 (2006), 1230–1247.
- [112] Voorhees, E. The TREC-8 question answering track report. In *Proc. 8th Text REtrieval Conference* (1999), pp. 77–82.
- [113] Voorhees, E. Overview of the TREC 2004 robust retrieval track. In *Proc. 13th Text REtrieval Conference* (2004).
- [114] Voorhees, E. Overview of the TREC 2005 robust retrieval track. In *Proc. 14th Text REtrieval Conference* (2005).
- [115] Wei, X., and Croft, W. B. LDA-based document models for ad-hoc retrieval. In *Proc. 29th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2006), pp. 178–185.
- [116] Wei, X., and Croft, W. B. Modeling term associations for ad-hoc retrieval performance within language modeling framework. In *Proc. 29th European Conf. on Information Retrieval* (2007), pp. 52–63.
- [117] Wong, S. K.M., and Yao, Y. Y. Linear structure in information retrieval. In *Proc. 11th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (1988), pp. 219–232.
- [118] Xu, J., and Croft, W. B. Improving the effectiveness of information retrieval with local context analysis. *ACM Trans. Inf. Syst.* 18, 1 (2000), 79–112.

- [119] Yu, C. T., Buckley, C., Lam, K., and Salton, G. A generalized term dependence model in information retrieval. Tech. rep., Cornell University, 1983.
- [120] Yue, Y., Finley, T., Radlinski, F., and Joachims, T. A support vector method for optimizing average precision. In *Proc. 30th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2007), p. To appear.
- [121] Zhai, C. *Risk Minimization and Language Modeling in Information Retrieval*. PhD thesis, Carnegie Mellon University, 2002.
- [122] Zhai, C., and Lafferty, J. Model-based feedback in the language modeling approach to information retrieval. In *Proc. 10th Intl. Conf. on Information and Knowledge Management* (2001), pp. 403–410.
- [123] Zhai, C., and Lafferty, J. A study of smoothing methods for language models applied to ad hoc information retrieval. In *Proc. 24th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2001), pp. 334–342.
- [124] Zhai, C., and Lafferty, J. Two-stage language models for information retrieval. In *Proc. 25th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2002), pp. 49–56.
- [125] Zhai, C., and Lafferty, J. A study of smoothing methods for language models applied to information retrieval. *ACM Trans. Inf. Syst.* 22, 2 (2004), 179–214.
- [126] Zhang, D., Chen, X., and Lee, W. S. Text classification with kernels on the multinomial manifold. In *Proc. 28th Ann. Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval* (2005), pp. 266–273.
- [127] Zhou, Y., and Croft, W. B. Document quality models for web ad hoc retrieval. In *Proc. 14th Intl. Conf. on Information and Knowledge Management* (2005), pp. 331–332.