

Feature Bagging: Preventing Weight Undertraining in Structured Discriminative Learning

Charles Sutton, Michael Sindelar, and Andrew McCallum

Department of Computer Science
University of Massachusetts Amherst
Amherst, MA 01003 USA

{casutton, mccallum}@cs.umass.edu, msindela@student.umass.edu

Abstract

Discriminatively-trained probabilistic models are very popular in NLP because of the latitude they afford in designing features. But training involves complex trade-offs among weights, which can be dangerous: a few highly-indicative features can swamp the contribution of many individually weaker features, causing their weights to be undertrained. Such a model is less robust, for the highly-indicative features may be noisy or missing in the test data. To ameliorate this *weight undertraining*, we propose a new training method, called *feature bagging*, in which separate models are trained on subsets of the original features, and combined using a mixture model or a product of experts. We evaluate feature bagging on linear-chain conditional random fields for two natural-language tasks. On both tasks, the feature-bagged CRF performs better than simply training a single CRF on all the features.

1 Introduction

Discriminative methods for training probabilistic models have enjoyed wide popularity in natural language processing, such as in part-of-speech tagging (Toutanova et al., 2003), chunking (Sha and Pereira, 2003), named-entity recognition (Florian et al., 2003; Chieu and Ng, 2003), and most recently parsing (Taskar et al., 2004). A discriminative probabilistic model is trained to maximize the conditional probability $p(y|x)$ of output labels y given input variables x , as opposed to modeling the joint probability $p(y, x)$, as in generative models such as the Naive Bayes classifier and hidden Markov models. The popularity of discriminative models stems from the great flexibility they allow in defining features: because the distribution over input features $p(x)$ is not modeled, it can contain rich, highly overlapping features without making the model intractable for training and inference.

In NLP, for example, useful features include word bigrams and trigrams, prefixes and suffixes, membership in domain-specific lexicons, and information from semantic databases such as WordNet. It is not uncommon to have hundreds of thousands or even millions of features.

But not all features, even ones that are carefully engineered, improve performance. Adding more features to a model can hurt its accuracy on unseen testing data. One well-known reason for this is overfitting: a model with more features has more capacity to fit chance regularities in the training data. In this paper, however, we focus on another, more subtle effect: adding new features can cause existing ones to be *underfit*. Training of discriminative models, such as regularized logistic regression, involves complex trade-offs among weights. A few highly-indicative features can swamp the contribution of many individually weaker features, even if the weaker features, taken together, are just as indicative of the output. Such a model is less robust, for the few strong features may be noisy or missing in the test data.

This effect was memorably observed by Dean Pomerleau (1995) when training neural networks to drive vehicles autonomously. Pomerleau reports one example when the system was learning to drive on a dirt road:

The network had no problem learning and then driving autonomously in one direction, but when driving the other way, the network was erratic, swerving from one side of the road to the other. ... It turned out that the network was basing most of its predictions on an easily-identifiable ditch, which was always on the right in the training set, but was on the left when the vehicle turned around. (Pomerleau, 1995)

The network had features to detect the sides of the road, and these features were active at training and test time, although weakly, because the dirt road was difficult to detect. But the ditch was so highly indicative that the network did not learn the dependence between the road edge and the desired steering direction.

In this paper, we examine a novel way to avoid feature undertraining in discriminative sequence models. We train separate models for groups of competing features—in the driving example, one model with the ditch features, and one with the side-of-the-road features—and then average them into a single model.

We test these methods on conditional random fields (CRFs) (Lafferty et al., 2001), which is the class of discriminatively-trained undirected models. On two natural-language tasks, we show that feature bagging performs significantly better than training a single CRF with all available features.

2 Conditional Random Fields

Conditional random fields (CRFs) (Lafferty et al., 2001) are undirected graphical models that are discriminatively trained. Let \mathbf{G} be an undirected graphical model over random vectors \mathbf{y} and \mathbf{x} . As a typical special case, $\mathbf{y} = \{y_t\}$ and $\mathbf{x} = \{x_t\}$ for $t = 1, \dots, T$, so that \mathbf{y} is a labeling of an observed sequence \mathbf{x} . For a given collection $C = \{\{\mathbf{y}_c, \mathbf{x}_c\}\}$ of cliques in \mathbf{G} , a CRF defines the conditional probability of an assignment to labels \mathbf{y} given the observed variables \mathbf{x} as:

$$p_{\Lambda}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in C} \Phi(\mathbf{y}_c, \mathbf{x}_c), \quad (1)$$

where Φ is a potential function and the partition function $Z(\mathbf{x}) = \sum_{\mathbf{y}} \prod_{c \in C} \Phi(\mathbf{y}_c, \mathbf{x}_c)$ is a normalization factor over all possible label assignments.

We assume the potentials factorize according to a set of features $\{f_k\}$, which are given and fixed, so that

$$\Phi(\mathbf{y}_c, \mathbf{x}_c) = \exp \left(\sum_k \lambda_k f_k(\mathbf{y}_c, \mathbf{x}_c) \right) \quad (2)$$

The model parameters are a set of real weights $\Lambda = \{\lambda_k\}$, one weight for each feature.

Many applications have used the *linear-chain CRF*, in which a first-order Markov assumption is made on the hidden variables. In this case, the cliques of the conditional model are the nodes and edges, so that there are feature functions $f_k(y_{t-1}, y_t, \mathbf{x}, t)$ for each label transition. (Here we write the feature functions as potentially depending on the entire input sequence.) Feature functions can be arbitrary. For example, a feature function $f_k(y_{t-1}, y_t, \mathbf{x}, t)$ could be a binary test that has value 1 if and only if y_{t-1} has the label “*adjective*”, y_t has the label “*proper noun*”, and x_t begins with a capital letter.

Linear-chain CRFs correspond to finite state machines, and can be roughly understood as conditionally-trained hidden Markov models (HMMs). This class of CRFs is also a globally-normalized extension to *Maximum Entropy Markov Models* (McCallum et al., 2000) that avoids the label bias problem (Lafferty et al., 2001).

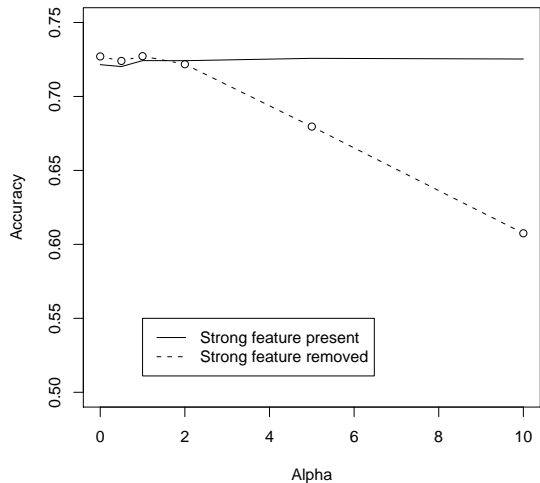


Figure 1: Effect of a single strong feature drowning out weaker features in logistic regression on synthetic data. The x -axis indicates the strength of the strong feature. In the top line, the strong feature is present at training and test time. In the bottom line, the strong feature is missing from the training data at test time.

Note that the number of state sequences is exponential in the input sequence length T . In linear-chain CRFs, the partition function $Z(\mathbf{x})$, the node marginals $p(y_i|\mathbf{x})$, and the Viterbi labeling can be calculated efficiently by variants of the dynamic programming algorithms for HMMs.

3 Weight Undertraining

In the section, we give a simple demonstration of weight undertraining. In a discriminative classifier, such as a neural network or logistic regression, a few strong features can drown out the effect of many individually weaker features, even if the weak features are just as indicative put together. To demonstrate this effect, we present an illustrative experiment using logistic regression, because of its strong relation to CRFs. (Conditional random fields are in fact the generalization of logistic regression to sequence data.)

Consider random variables $x_1 \dots x_n$, each distributed as independent standard normal variables. The output y is a binary variable whose probability depends on all the x_i ; specifically, we define its distribution as $y \sim \text{Bernoulli}(\text{logit}(\sum_i x_i))$. The correct decision boundary in this synthetic problem is the hyperplane tangent to the weight vector $(1, 1, \dots, 1)$. Thus, if n is large, each x_i contributes weakly to the output y . Finally, we include a highly indicative feature $x_S = \alpha \sum_i x_i + \mathcal{N}(\mu =$

$0, \sigma^2 = 0.04$). This variable alone is sufficient to determine the distribution of y . The variable α is a parameter of the problem that determines how strongly indicative x_S is; specifically, when $\alpha = 0$, the variable x_S is random noise.

We choose this synthetic model by analogy to Pomerleau’s observations. The x_i correspond to the side of the road in Pomerleau’s case—the weak features present at both testing and training—and x_S corresponds to the ditch—the strongly indicative feature that is corrupted at test time.

We examine how badly the learned classifier is degraded when x_S feature is present at training time but missing at test time. For several values of the weight parameter α , we train a regularized logistic regression classifier on 1000 instances with $n = 10$ weak variables. In Figure 1, we show how the amount of error caused by ablating x_S at test time varies according to the strength of x_S . Each point in Figure 1 is averaged over 100 randomly-generated data sets. When x_S is weakly indicative, it does not affect the predictions of the model at all, and the classifier’s performance is the same whether it appears at test time or not. When x_S becomes strongly indicative, however, the classifier learns to depend on it, and performs much more poorly when x_S is ablated, even though exactly the same information is available in the weak features.

4 Feature Bagging

In this section, we describe the feature bagging method. We divide the set of features $F = \{f_k\}$ into a collection of possibly overlapping subsets $\mathcal{F} = \{F_1, \dots, F_M\}$, which we call *feature bags*. We train individual CRFs on each of the feature bags using standard MAP training, yielding individual CRFs $\{p_1, \dots, p_M\}$.

We average the individual CRFs into a single combined model. This averaging can be performed in several ways: we can average probabilities of entire sequences, or of individual transitions; and we can average using the arithmetic mean, or the geometric mean. This yields four combination methods:

1. *Per-sequence mixture*. The distribution over label sequences \mathbf{y} given inputs \mathbf{x} is modeled as a mixture of the individual CRFs. Given nonnegative weights $\{\alpha_1, \dots, \alpha_m\}$ that sum to 1, the combined model is given by

$$p_{\text{SM}}(\mathbf{y}|\mathbf{x}) = \sum_{i=1}^M \alpha_i p_i(\mathbf{y}|\mathbf{x}). \quad (3)$$

It is easily seen that if the sequence model is defined as in Equation 3, then the pairwise marginals

are mixtures as well:

$$p_{\text{SM}}(y_t, y_{t-1}|\mathbf{x}) = \sum_{i=1}^M \alpha_i p_i(y_t, y_{t-1}|\mathbf{x}). \quad (4)$$

The probabilities $p_i(y_t, y_{t-1}|\mathbf{x})$ are pairwise marginal probabilities in the individual models, which can be efficiently computed by the forward-backward algorithm.

It is because of Equation 4 that we can efficiently compute the maximum-probability state sequence in the mixture model. Because the sequence mixture is still a linear-chain, it can be written in terms of transition probabilities as $p_{\text{SM}}(\mathbf{y}|\mathbf{x}) = \prod_t p_{\text{SM}}(y_t|y_{t-1}, \mathbf{x})$. And each of these transition probabilities can be easily computed by dividing Equation 3 by the single-node marginal $p_{\text{SM}}(y_{t-1}|\mathbf{x})$.

So we can efficiently compute the maximum-probability state sequence in the combined model by first running forward-backward on each of the individual models, combining the pairwise marginals as in Equation 4, computing the transition probabilities for the mixture, and then using the standard Viterbi algorithm.

The mixture weights can be selected in a variety of ways, including equal voting, as in traditional bagging, or EM.

2. *Per-sequence product of experts*. The distribution over label sequences \mathbf{y} given inputs \mathbf{x} is modeled as a product of experts (Hinton, 2000). In a product of experts, instead of summing the probabilities from the individual models, we multiply them together. Essentially we take a geometric mean instead of an arithmetic mean. Given nonnegative weights $\{\alpha_1, \dots, \alpha_m\}$ that sum to 1, the product model is

$$p(\mathbf{y}|\mathbf{x}) \propto \prod_{i=1}^M (p_i(\mathbf{y}|\mathbf{x}))^{\alpha_i}. \quad (5)$$

The combined model can also be viewed as a conditional random field whose features are the log probabilities from the original models:

$$p(\mathbf{y}|\mathbf{x}) \propto \exp \left\{ \sum_{i=1}^M \alpha_i \log p_i(\mathbf{y}|\mathbf{x}) \right\} \quad (6)$$

By substituting in the CRF definition, it can be seen that the model in Equation 6 is simply a single CRF whose parameters are a weighted average of the original parameters. So feature bagging using the product method does not increase the family of models that are considered: standard training of a single

CRF on all available features could potentially pick the same parameters as the bagged model.

Nevertheless, in Section 5, we show that this feature bagging method performs better than standard CRF training.

The previous two combination methods combine the individual models by averaging probabilities of entire sequences. Alternatively, in a sequence model we can average probabilities of individual transitions $p_i(y_t|y_{t-1}, \mathbf{x})$. Computing these transition probabilities requires performing probabilistic inference in each of the original CRFs, because $p_i(y_t|y_{t-1}, \mathbf{x}) = \sum_{\mathbf{y} \setminus y_t, y_{t+1}} p(\mathbf{y}|y_{t-1}, \mathbf{x})$.

This yields two other combination methods:

3. *Per-transition mixture.* The transition probabilities are modeled as

$$p_{\text{TM}}(y_t|y_{t-1}, \mathbf{x}) = \sum_{i=1}^M \alpha_i p_i(y_t|y_{t-1}, \mathbf{x}) \quad (7)$$

Intuitively, the difference between per-sequence and per-transition mixtures can be understood generatively. In order to generate a label sequence \mathbf{y} given an input \mathbf{x} , the per-sequence model selects a mixture component, and then generates \mathbf{y} using only that component. The per-transition model, on the other hand, selects a component, generates y_1 from that component, selects another component, generates y_2 from the second component given y_1 , and so on.

4. *Per-transition product of experts.* Finally, we can combine the transition distributions using a product model

$$p_{\text{SP}}(y_t|y_{t-1}, \mathbf{x}) \propto \prod_{i=1}^M p(y_t|y_{t-1}, \mathbf{x})^{\alpha_i} \quad (8)$$

Each transition distribution is thus—similarly to the per-sequence case—an exponential-family distribution whose features are the log transition probabilities from the individual models. Unlike the per-sequence product, there is no weight-averaging trick here, because the probabilities $p(y_t|y_{t-1}, \mathbf{x})$ are marginal probabilities.

Considered as a sequence distribution $p(\mathbf{y}|\mathbf{x})$, the per-transition product is a local-normalized maximum-entropy Markov model (McCallum et al., 2000). It would not be expected to suffer from label bias, however, because each of the features take the future into account; they are marginal probabilities from CRFs.

Those are the four combination methods we propose. Although for concreteness we describe them in terms of sequence models, they may be generalized to arbitrary graphical structures.

5 Results

We evaluate feature bagging on two natural language tasks, named entity recognition and noun-phrase chunking. The purpose of the named entity recognition task is to identify all of the named entities corresponding to people, locations, organizations and other miscellaneous entities. We use the standard CoNLL 2003 English data set, which is taken from Reuters newswire and consists of a training set of 14987 sentences, a development set of 3466 sentences, and a testing set of 3684 sentences. For noun-phrase chunking the goal is to extract base noun phrases from sentences. We use the standard CoNLL 2000 data set, which consists of 8936 sentences for training and 2012 sentences for testing, taken from Wall Street Journal articles annotated by the Penn Treebank project.

As is standard, we compute precision and recall for both tasks based upon the chunks (or named entities for CoNLL 2003) as

$$P = \frac{\# \text{ correctly labeled chunks}}{\# \text{ labeled chunks}}$$

$$R = \frac{\# \text{ correctly labeled chunks}}{\# \text{ actual chunks}}$$

We report the harmonic mean of precision and recall as $F_1 = (2PR)/(P + R)$.

For both tasks, we use per-sequence product-of-experts feature bagging with two feature bags which we manually choose based on prior experience with the data set. For each experiment, we report two baseline CRFs, one trained on union of the two feature sets, and one trained only on the features that were present in both bags, such as lexical identity and regular expressions.

For the named entity task, we use two feature bags based upon character ngrams and lexicons. Both bags contain a set of baseline features, such as word identity and regular expressions (Table 4). The ngram CRF includes binary features for character ngrams of length 2, 3, and 4 and word prefixes and suffixes of length 2, 3, and 4. The lexicon CRF includes membership features for a variety of lexicons containing people names, places, and company names. The combined model had 2,342,543 features. The mixture weight α is selected using the development set.

For the chunking task, the two feature sets were selected based upon part of speech and lexicons. Again, a set of baseline features was used, similar to the regular expressions and word identity features used on the named

Model	F1
Per-sequence Product of Experts	86.61
Per-transition Product of Experts	86.58
Per-sequence Mixture	86.46
Per-transition Mixture	86.42

Table 1: Results of various bagging methods on the CoNLL 2003 Named Entity Task.

entity task (Table 4). The first bag also includes part-of-speech tags generated by the Brill tagger and the conjunctions of those tags used by Sha and Pereira (2003). The second bag uses lexicon membership features for lexicons containing names of people, places, and organizations. In addition, we use part-of-speech lexicons generated from the entire Treebank, such as a list of all words that appear as nouns. These lists are also used by the Brill tagger (Brill, 1994). There were 536,203 features used in the combined model. The mixture weight α is selected using 2-fold cross validation.

In both data sets, the bagged model performs better than the single CRF trained with all of the features. For the named entity task, bagging improves performance from 85.45% to 86.61%, with a substantial error reduction of 8.32%. This is lower than the best reported results for this data set, which is 88.76% (Florian et al., 2003), but our model can still be improved by the inclusion of more extensively engineered features, such as the ones described in (Chieu and Ng, 2003). For the chunking task, bagging improved the performance from 94.34% to 94.77%, with an error reduction of 7.60%. In both data sets, the improvement is statistically significant (McNemar’s test; $p < 0.01$).

On the chunking task, the bagged model also outperforms the models of Kudo and Matsumoto (2001) and Sha and Pereira (2003), the best previously reported results on this data set. Although we used lexicons that were not included in the previous models, the additional features actually did not help the original CRF. It was only with feature bagging that these lexicons improved performance.

Finally, we compare the four bagging methods of Section 4: pre-transition mixture, pre-transition product of experts, and per-sequence mixture. On the named entity data, all four models performed in a statistical tie, with statistically significant difference between their results (Table 1).

6 Previous Work

Within classification, there is literature on combining models trained on feature subsets. Ho (1995) creates an ensemble of decision trees by randomly choosing a feature subset on which to grow each tree using standard decision tree learners. Other work along these lines

Model	F1
Single CRF(Base Feat.)	81.52
Single CRF(All Feat.)	85.45
Combined CRF	86.61

Table 2: Results for the CoNLL 2003 Named Entity Task. The bagged CRF performs significantly better than a single CRF with all available features (McNemar’s test; $p < 0.01$).

Model	F1
Single CRF(Base Feat.)	89.60
Single CRF(All Feat.)	94.34
(Sha and Pereira, 2003)	94.38
(Kudo and Matsumoto, 2001)	94.39
Combined CRF	94.77

Table 3: Results for the CoNLL 2000 Chunking Task. The bagged CRF performs significantly better than a single CRF (McNemar’s test; $p < 0.01$), and also better than previously published results.

include Bay’s (1998) using nearest-neighbor classifiers, and more recently Bryll et al (2003). Also, in Breiman’s work on random forests (2001), ensembles of random decision trees are constructed by choosing a random feature at each node. This literature mostly has the goal of improving accuracy by reducing the classifier’s variance.

In contrast, O’Sullivan et al. (2000) specifically focus on increasing robustness by training classifiers to use all of the available features. Their algorithm FeatureBoost is analogous to AdaBoost, except that the meta-learning algorithm maintains weights on features instead of on instances. Feature subsets are automatically sampled based on which features, if corrupted, would most affect the ensemble’s prediction. They show that FeatureBoost is more robust than AdaBoost on synthetically corrupted UCI data sets. Their method does not easily extend to sequence models, especially natural-language models with hundreds of thousands of features.

We are not aware of previous work on feature subspaces for ensembles of sequence models. Altun, Hofmann, and Johnson (2003) describe a boosting algorithm for sequence models, but they boost instances, not features. In fact, the main advantage of their technique is increased model sparseness, whereas in this work we aim to fully use *more* features to increase accuracy and robustness.

7 Conclusion

Discriminatively-trained probabilistic models have had much success in applications because of their flexibility in defining features, but sometimes even highly-

$w_t = w$
w_t begins with a capital letter
w_t contains only capital letters
w_t is a single capital letter
w_t contains some capital letters and some lowercase
w_t contains a numeric character
w_t contains only numeric characters
w_t appears to be a number
w_t is a string of at least two periods
w_t ends with a period
w_t contains a dash
w_t appears to be an acronym
w_t appears to be an initial
w_t is a single letter
w_t contains punctuation
w_t contains quotation marks
$P_t = P$
All features for time $t + \delta$ for all $\delta \in [-2, 2]$

Table 4: Baseline features used in all bags. In the above w_t is the word at position t , P_t is the POS tag at position t , w ranges over all words in the training data, and P ranges over all chunk tags supplied in the training data. The “appears to be” features are based on hand-designed regular expressions.

indicative features can fail to increase performance. We have shown that this can be due to feature undertraining, where highly-indicative features prevent training of many weaker features. One solution to this is feature bagging: repeatedly selecting feature subsets, training separate models on each subset, and averaging the individual models.

On large, real-world natural-language processing tasks, feature bagging significantly improves performance, even with only two feature subsets. Areas for future work include automatically determining the feature subsets.

Acknowledgments

We thank Andrew Ng, Hanna Wallach, and Jerod Weinman for helpful conversations. This work was supported in part by the Center for Intelligent Information Retrieval, in part by The Central Intelligence Agency, the National Security Agency and National Science Foundation under NSF grant numbers IIS-0326249 and IIS-0427594. Any opinions, findings and conclusions or recommendations expressed in this material are the author(s) and do not necessarily reflect those of the sponsor.

References

Yasemin Altun, Thomas Hofmann, and Mark Johnson. 2003. Discriminative learning for label sequences via boosting. In *Advances in Neural Information Processing Systems (NIPS*15)*.

- Stephen D. Bay. 1998. Combining nearest neighbor classifiers through multiple feature subsets. In *ICML '98: Proceedings of the Fifteenth International Conference on Machine Learning*, pages 37–45. Morgan Kaufmann Publishers Inc.
- Leo Breiman. 2001. Random forests. *Machine Learning*, 45(1):5–32, October.
- Eric Brill. 1994. Some advances in transformation-based part of speech tagging. In *AAAI '94: Proceedings of the twelfth national conference on Artificial intelligence (vol. 1)*, pages 722–727. American Association for Artificial Intelligence.
- Robert Bryll, Ricardo Gutierrez-Osuna, and Francis Quek. 2003. Attribute bagging: improving accuracy of classifier ensembles by using random feature subsets. *Pattern Recognition*, 36:1291–1302.
- Hai Leong Chieu and Hwee Tou Ng. 2003. Named entity recognition with a maximum entropy approach. In Walter Daelemans and Miles Osborne, editors, *Proceedings of CoNLL-2003*, pages 160–163. Edmonton, Canada.
- Radu Florian, Abe Ittycheriah, Hongyan Jing, and Tong Zhang. 2003. Named entity recognition through classifier combination. In *Proceedings of CoNLL-2003*.
- G.E. Hinton. 2000. Training products of experts by minimizing contrastive divergence. Technical Report 2000-004, Gatsby Computational Neuroscience Unit.
- T. K. Ho. 1995. Random decision forests. In *Proc. of the 3rd Int'l Conference on Document Analysis and Recognition*, pages 278–282, Montreal, Canada, August.
- T. Kudo and Y. Matsumoto. 2001. Chunking with support vector machines. In *Proceedings of NAACL-2001*.
- J. Lafferty, A. McCallum, and F. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. *Proc. 18th International Conf. on Machine Learning*.
- Andrew McCallum, Dayne Freitag, and Fernando Pereira. 2000. Maximum entropy Markov models for information extraction and segmentation. In *Proc. 17th International Conf. on Machine Learning*, pages 591–598. Morgan Kaufmann, San Francisco, CA.
- Joseph O’Sullivan, John Langford, Rich Caruana, and Avrim Blum. 2000. Featureboost: A meta learning algorithm that improves model robustness. In *International Conference on Machine Learning*.
- Dean Pomerleau. 1995. Neural network vision for robot driving. In M. Arbib, editor, *The Handbook of Brain Theory and Neural Networks*.

Fei Sha and Fernando Pereira. 2003. Shallow parsing with conditional random fields. In *Proceedings of HLT-NAACL 2003*. Association for Computational Linguistics.

Ben Taskar, Dan Klein, Michael Collins, Daphne Koller, and Chris Manning. 2004. Max-margin parsing. In *Empirical Methods in Natural Language Processing (EMNLP04)*.

Kristina Toutanova, Dan Klein, Christopher D. Manning, and Yoram Singer. 2003. Feature-rich part-of-speech tagging with a cyclic dependency network. In *HLT-NAACL 2003*.