

Conditional Structure versus Conditional Estimation in NLP Models

Dan Klein and Christopher D. Manning

Computer Science Department

Stanford University

Stanford, CA 94305-9040

{klein, manning}@cs.stanford.edu

Abstract

This paper separates conditional *parameter estimation*, which consistently raises test set accuracy on statistical NLP tasks, from conditional *model structures*, such as the conditional Markov model used for maximum-entropy tagging, which tend to lower accuracy. Error analysis on the POS tagging task shows that the actual tagging errors made by the conditionally structured model derive principally not from label bias, as has been claimed, but from other ways in which the independence assumptions of the conditional model structure are unsuited to linguistic sequences. The paper presents new word-sense disambiguation and POS tagging experiments, and integrates apparently conflicting reports from other recent work.

1 Introduction

The success and widespread adoption of probabilistic models in NLP has led to numerous variant methods for any given task, and it can be difficult to tell what aspects of a system have led to its relative successes or failures. As an example, maximum entropy taggers have achieved very good performance (Ratnaparkhi, 1998; Toutanova and Manning, 2000; Lafferty et al., 2001), but almost identical performance has also come from finely tuned HMM models (Brants, 2000; Thede and Harper, 1999). Are any performance gains due to the sequence model used, the maximum entropy approach to parameter estimation, or the features employed by the system?

Recent experiments have given conflicting recommendations about what is best. Johnson (2001) finds that a conditionally trained PCFG marginally outperforms a standard jointly trained PCFG, but that a conditional shift-reduce model performs worse than a joint formulation. Lafferty et al. (2001) suggest on abstract grounds that conditional models will suffer from a phenomenon called *label bias*

(Bottou, 1991), (see section 3), but is this a significant effect for real NLP problems?

We suggest that the results in the literature, along with the new results we present in this work, can be explained by the following generalizations:

- The ability to include better features in a well-founded fashion means better performance.
- For fixed features, assumptions implicit in the model structure have a large impact on errors.
- Maximizing the objective being evaluated has a reliably positive, but sometimes small, effect.

It is especially important to study these issues using NLP data sets: NLP tasks are marked by their complexity and sparsity, and, as we show, conclusions imported from the machine-learning literature do not always hold in these extreme contexts.

In previous work, the structure of a model and the method of parameter estimation are often both changed simultaneously (for reasons of naturalness or computational ease), but in this paper we seek to tease apart the separate effects of these two factors. Section 2 examines a fixed probabilistic model with varying objective functions, in the context of word-sense disambiguation (WSD) with Naive Bayes models. Section 3 examines the situation where the model structure is allowed to vary, in the context of part-of-speech (POS) tagging. Finally, we relate our discussions and experiments to recent findings by other researchers.

2 Objective Functions: Naive-Bayes

Throughout, we will assume we have a corpus $D = (O, S)$ of observed structure instances $o \in O$, each associated with an unobserved structure $s \in S$. A simple case is when we have a bag of observations making up each o and a single, atomic class for each s , for example in bag-of-words WSD. A particular model for this case is the familiar *Naive Bayes* (NB)

model (Gale et al., 1992), where we assume independence between each of the $o_i \in o$:

$$P(s, o) = P(s) \prod_i P(o_i | s)$$

The NB model gives a joint distribution over the $\{o_i\}$ and s variables.

A *loglinear model* is one where each (s, o) is given a non-negative score

$$\text{score}(s, o) = \prod_k \theta_k^{f_k(s, o)}$$

where each $f_k(s, o)$ is an indicator of the number of times a feature k occurs inside (s, o) , and θ_k is a weighting given to that feature. A product model in this form is a linear model in log space.

The scores of a loglinear model can always be viewed as conditional probabilities if normalized (via a “partition function”):

$$P(s|o) = \text{score}(s, o) / \sum_s \text{score}(s, o)$$

Each $P(s)$ or $P(o_i | s)$ in the NB model is a number $\theta \in \Theta$, and thus the Naive Bayes model is a log-linear classifier with a feature for each s and each (s, o_i) pair. In NLP applications of NB models, the o_i are typically multinomial distributions over words in a particular position in the document (cf. McCallum and Nigam 1998), and the typical way used to set the parameters is using (smoothed) *relative frequency estimators* (RFES):

$$\begin{aligned} \theta_s &= P(s) = \text{count}(s) / |D| \\ \theta_{o_i | s} &= P(o_i | s) = \text{count}(o_i, s) / \sum_{o'} \text{count}(o', s) \end{aligned}$$

These intuitive relative frequency estimators are the estimates for Θ which maximize the probability of the corpus D , according to the joint NB model. We will refer to this objective function as *joint likelihood* (JL). A NB model which has been trained by maximizing JL will be referred to as NB-JL.

We can set the parameters in other ways, without changing our model. We just won’t maximize JL. If we are doing classification, we may not care about JL, but about whatever kinds of errors we get charged for by the classification evaluation. If we either receive a 1 or 0 score on an instance, then we have what is called *0/1-loss* (Friedman, 1997). If, as in the SENSEVAL competition (Kilgarriff, 1998), we may guess a distribution $P(s)$ over answers and our score is $\sum P(s_{correct})$, then we have *linear loss*.¹ We refer to 0/1 loss simply as *accuracy* (Acc). Since

¹This is in fact a strange scoring criterion if one expects to use WSD systems as a step in a probabilistic process, rather

the score under linear loss is the sum, over the test instances o , of $P(s|o)$, we refer to that objective as *sum of conditional likelihoods* (SCL).

If we multiply, rather than add, our conditional estimates, then we have the likelihood of the corpus’ classes conditional on their observations. This is *conditional likelihood* (CL). To recap, we have the following objectives:

$$\begin{aligned} JL(\Theta, D) &= \prod_{(o, s) \in D} P(s, o) \\ CL(\Theta, D) &= \prod_{(o, s) \in D} P(s|o) \\ SCL(\Theta, D) &= \sum_{(o, s) \in D} P(s|o) \\ Acc(\Theta, D) &= \sum_{(o, s) \in D} \delta(s, \arg \max_{s'} P(s'|o)) \end{aligned}$$

Although in principle we can optimize any of these objectives, in practice some are harder to optimize than others. To optimize for a given objective function numerically, we need to decide what values of Θ we will allow. If we want to have a well-formed joint NB interpretation, we must have the inequalities $\sum_o \theta_{o|s} \leq 1$ and $\sum_s \theta_s \leq 1$. If we want to be guaranteed a non-deficient joint interpretation, we could insist on equality. However, if we relax the equality then we have a larger feasible space which may give better values of our objective.²

For the Naive-Bayes model, we performed the following experiments. We took as data the collection of SENSEVAL-2 English lexical sample WSD corpora.³ We set the parameters for the NB model in several ways. We optimized JL (by using the RFES). We also optimized SCL and (the log of) CL, using a conjugate gradient (CG) ascent method (Press et al., 1988).⁴ For CL and SCL, we optimized each objective both in an unconstrained fashion and with the

than in isolation. CL is a much better measure for filter processes, because it appropriately highly punishes wrongly assigning zero or near-zero probabilities to outcomes.

²With a NB model, we can always divide the vector Θ by any positive constant and leave our conditional scores unchanged. Thus, any parameter setting Θ has some Θ' which satisfies the linear inequalities and which makes the same conditional predictions. If we have a fixed number of observation variables per instance, then even strict equality constraints do not change the set of resulting conditional models.

³<http://www.sle.sharp.co.uk/senseval2/>

⁴All optimization was done using conjugate gradient ascent over log parameters $\lambda_i = \log \theta_i$, rather than the given parameters due to sensitivity near zero and improved quality of quadratic approximations during optimization. Linear constraints over θ are not linear in log space, and were enforced using a quadratic Lagrange penalty method (Bertsekas, 1995).

TRAINING SET					
Optimization	Acc	MacroAcc	log JL	log CL	SCL
NB-JL	86.8	86.2	-22969684.7	-243184.1	4505.9
NB-CL*	98.5	96.2	-23366291.2	-973.0	5101.2
NB-CL	98.5	96.2	-23431010.0	-854.1	5115.1
NB-SCL*	94.2	93.7	-23054768.6	-226187.8	4884.4
NB-SCL	97.3	95.5	-23146735.3	-220145.0	5055.8

TEST SET					
Optimization	Acc	MacroAcc	log JL	log CL	SCL
NB-JL	73.6	55.0	-1816757.1	-55251.5	3695.4
NB-CL*	72.3	53.4	-1954977.1	-19854.1	3566.3
NB-CL	76.2	56.5	-1964498.5	-20498.7	3798.8
NB-SCL*	74.8	57.2	-1841305.0	-43027.8	3754.1
NB-SCL	77.5	59.7	-1872533.0	-33249.7	3890.8

Figure 1: Scores for the NB model trained according to various objectives. Scores are usually higher on both training and test sets for the objective maximized, and discriminative criteria lead to better test-set accuracy. The best scores are in bold.

linear equalities necessary to ensure that the resulting NB model would be non-deficient (giving CL^* and SCL^*). Acc was not directly optimized. Although it is the most common evaluation criterion in NLP, it is not continuous in Θ and is therefore unsuited to direct optimization (indeed, the problem of finding an optimum is NP-complete).

Unconstrained CL corresponds exactly to a conditional maximum entropy model (Berger et al., 1996; Lafferty et al., 2001). This particular case where there are multiple explanatory variables and a single categorical response variable is also precisely the well-studied statistical model of (multinomial) *logistic regression* (Agresti, 1990). Its optimization problem is concave (over log parameters) and therefore has a unique global maximum (in either parameter space). Adding the non-deficiency constraints to CL, or optimizing SCL, only guarantees us a local optimum in theory, but in practice we detected no cases of maxima which were not global over the feasible region.

Smoothing, of course, is vital in this task. So that the various models would receive as similar as possible smoothing, we smoothed by adding phantom data. We added one instance of each class occurring with the bag containing each vocabulary word once. This gave the same result as add-one smoothing on the RFES for NB-JL, and ensured that NB-CL and NB-CL* would not assign zero conditional probability to any unseen event. The phantom data did not, however, result in smoothed estimates for SCL and SCL^* , because any conditional probability will sum to one over the phantom instances.

As a result, we added a penalty term proportional to $\sum_{\theta} (\log \theta)^2$, which ensured that no conditional probabilities reached 0 or 1.

Figure 1 shows, for each objective maximized, the values of all objectives on both the training and test set. Optimizing for a given objective generally gave the best score for that objective for both the training set and the test set. The exception is NB-SCL and NB-SCL* which have lower SCL score than NB-CL and NB-CL*. This is due to the penalty used for smoothing the summed models.

Accuracy is higher when optimizing CL than JL (including for macro-averaging over senses). That NB-CL should beat NB-JL on accuracy is unsurprising, since Acc is a discretization of CL, not of JL. Both discriminative objectives (CL and SCL) give better accuracy and better CL and SCL. This supports the claim that maximizing conditional likelihood, or other discriminative objectives, improves test set accuracy for realistic NLP tasks. Interestingly, NB-SCL has even better test-set accuracy than NB-CL, though this seems to be due to the difference in smoothing methods used (Chen and Rosenfeld (1999) show that quadratic penalties are very effective in practice, while the phantom data method is quite crude). NB-CL* is somewhere between JL and CL for all objectives on the training data. Its behavior shows that the change from a standard NB approach (NB-JL) to a maximum entropy classifier (NB-CL) can be broken into two aspects: a change in objective and an abandonment of a non-deficiency constraint. It is worth noting that the JL score for NB-CL*, is not very much lower than for NB-JL, despite the dramatic change in CL.

It would be too strong to state without qualification that maximizing CL (in particular) and discriminative objectives (in general) is better than maximizing JL, as far as test-set accuracy. CL strictly beat JL in accuracy for only 15 of the 24 words. Figure 2 shows a plot of the relative accuracy for CL: $(Acc_{CL} - Acc_{JL}) / Acc_{JL}$. The x -axis is the average number of training instances per sense, weighted by the frequency of that sense in the test data. There is a clear trend that larger training sets saw a larger benefit from using NB-CL. The scatter of this trend is partially due to the wide range in data set conditions. These data sets exhibit an unusual amount of drift between training and test distributions. For example, the test data for *amaze* consists of 70 instances of the less frequent of its two

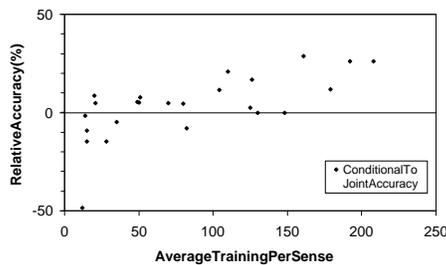


Figure 2: Conditional NB has higher accuracy than Joint NB for WSD on most SENSEVAL word sets. The relative improvement gained by switching to conditional estimation is positively correlated to training set size.

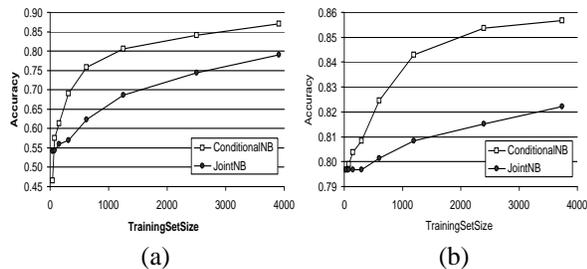


Figure 3: Conditional NB is better than Joint NB for WSD given all but possibly the smallest training sets, and the advantage increases with training set size. (a) “line” (b) “hard”

training senses, and represents the highest point on this graph, with NB-CL having a relative accuracy increase of 28%. This situation in general favors conditional estimates. On the other hand, many of these data sets are very small, individually, and 6 of the 7 sets where NB-JL wins are among the 8 smallest, 4 of them in fact being the 4 smallest. Ng and Jordan (2002) show that theoretically, for generative NB-JL vs. discriminative NB-CL, the discriminative model generally has a lower asymptotic error, but the generative model performs better in low-data situations. They in fact argue that unless one has a relatively large data set, one is likely to be better off with the generative estimate. Their claim seems too strong here; even smaller data sets often show benefit to accuracy from CL estimation, although all would qualify as small by their criteria.

Since the number of senses and skew towards common senses is so varied between SENSEVAL words, we turned to larger data sets to test the effective “break-even” size for WSD data, using the *hard* and *line* data from (Leacock et al., 1998). Figure 3 shows the accuracy of NB-CL and NB-JL as the amount of training data increases. Conditional beats joint for all but the smallest training sizes, and

the improvement is greater with larger training sets. Only for the *line* data does the conditional model ever drop below the joint model.

In short, for this task, NB-CL is performing better than it would be expected to. This appears to be due to two ways in which CL estimation is suited to linguistic data. First, the (Ng and Jordan, 2002) results do not involve smoothed data – their data sets do not require it like linguistic data does – and smoothing largely prevents the low-data overfitting that can plague conditional models.

There is another, more interesting reason why we might expect conditional estimation for this model to work better for an NLP task like WSD than for a general machine learning task. One signature difficulty in NLP is that the data contains a great many rare observations. In the case of WSD, the issue is in telling the kinds of rare events apart. Consider a word w which occurs only once, with a sense s . In the joint model, smoothing ensures that w does not signal s too strongly. However, every w which occurs only once with s will receive the same θ . Ideally, we would want to be able to tell the accidental singletons from true indicator words. The conditional model implicitly does this. If w occurs with s in an example where other good indicator words are present, then those other words’ large weights will explain the occurrence of w and, without w having to have a large weight, its expected count with s in that instance will become close to 1. On the other hand, if no trigger words occur in that instance, there will be no other explanation for w ’s presence, and w ’s expected count with s in that example will not be near one unless w ’s own weight helps cause it to be.

As a concrete illustration, we isolated two senses of “line” into a two-sense data set. Sense 1 was “a queue” and sense 2 was “a phone line.” In this data, the words *transatlantic* and *flowers* both occur only once, and only with the “phone line” sense (plus once with each sense in the phantom data). However, *transatlantic* occurs in the instance *thanks, anyway, the transatlantic line_2 died.*, while *flowers* occurs in the longer instance *... phones with more than one line_2, plush robes, exotic flowers, and complimentary wine.* In the first case, the only non-singleton word which could indicate sense 2 is *died* which occurs once with sense 1 and twice with sense 2. However, in the other case, *phone* occurs 191 times with sense 2 and only 4 times with

sense 1. Additionally, there are more words in the second instance with which *flowers* can share the responsibility of increasing its expectation. Experimentally, the parameters for (*transatlantic*, 2) and (*flowers*, 2) were, of course, equal by joint estimation, and gave a logscore of 0.63 more to sense 2 than to sense 1 per occurrence, indicating it equally strongly. However, with conditional estimation, (*transatlantic*, 2) gave a relative logscore advantage of 1.32 to sense 2, while (*flowers*, 2) gave only 0.72. Given that the conditional estimation is implicitly differentially weighting rare events in a fairly reasonable way, it is perhaps unsurprising that a task like WSD would see the benefits on smaller corpus sizes than would be expected on standard machine-learning data sets.⁵

In any case, we would like to emphasize that these trends are reliable, but sometimes small. In practice, one must decide if a 5% error reduction is worth the added work: CG optimization, especially with constraints, is considerably harder to implement than simple RFE estimates for JL. It is also considerably slower: the total training time for the entire SENSEVAL corpus was less than 3 seconds for NB-JL, but two hours for NB-CL.

3 Model Structure: HMMs and CMMs

We now consider sequence data, with POS tagging as the concrete NLP example. In the previous section, we had a single model, but several ways of estimating parameters. In this section, we have two different model structures.

First is the classic hidden Markov model (HMM), shown in figure 5a. For an instance ($s = \{s_i\}, o = \{o_i\}$) we write the following joint model:

$$P(s, o) = P(s)P(o|s) = \prod_i P(s_i|s_{i-1})P(o_i|s_i)$$

where we make a special start state s_0 to simplify the notation. Forgetting that the parameters have probabilistic interpretations, this is again a loglinear model with features for the transitions and emissions in (s, o) :

$$\text{score}(s, o) = \prod_i \theta_{s_i, s_{i-1}} \theta_{o_i, s_i}$$

There are various options within this model space. If you insist that $\sum_{s'} \theta_{s', s} \leq 1$ and $\sum_o \theta_{o, s} \leq 1$,

⁵Interestingly, the common approach of discarding low-count events (for both training speed and overfitting reasons) when estimating the conditional models used in maxent taggers robs the system of the opportunity to exploit this effect of conditional estimation.

Objective	Model		
	HMM	MEMM	MEMM*
JL	91.23	89.22	90.44
CL*	91.41	89.22	90.44
CL	91.44	89.22	90.44

Figure 4: Tagging accuracy: For a fixed model, conditional estimation is slightly advantageous. For a fixed objective, the MEMM is inferior, though it can be improved by *unobserving* unambiguous words.

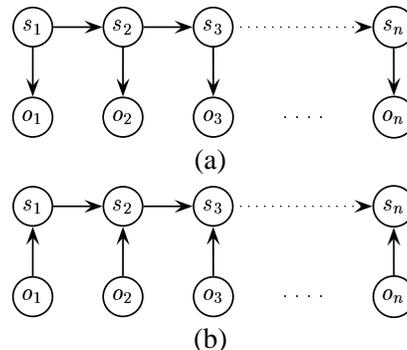


Figure 5: Graphical models: (a) the downward HMM, and (b) the upward conditional Markov model (CMM).

you can think about the terms as conditional probabilities for $P(s_i|s_{i-1})$ and $P(o_i|s_i)$, and you get a generative model: a (possibly deficient) HMM. If you require equality, you get a non-deficient vanilla HMM. Either way, the ML estimates are the standard RFEs. If you maximize CL, without requiring equality, you get (possibly deficient) HMMs which are instances of the Conditional Random Fields of Lafferty et al. (2001).⁶

Figure 4 shows the tagging accuracy of an HMM trained to maximize each objective. JL is the standard HMM. CL duplicates the simple CRFs in (Lafferty et al., 2001). CL* is again an intermediate, where we optimized conditional likelihood but required the HMM to be non-deficient. This separates out changes brought by a different objective function from changes brought by adopting a deficient model.⁷ In accordance with our observations in the last section, and consistent with the results of (Lafferty et al., 2001), the CL accuracy is slightly higher than JL for this fixed model.

Another model often used for sequence data is the upward Conditional Markov Model (CMM), shown

⁶The general class of CRFs is more expressive and reduces to deficient HMMs only when they have just these features.

⁷If one is only interested in classification, there is no reason to disprefer conditional models.

as a graphical model in figure 5b. This is the model used in maximum entropy tagging. The graphical model shown gives a *joint* distribution over (s, o) , just like an HMMs. It is a conditional model, in the sense that that distribution can be written as $P(s, o) = P(s|o)P(o)$. Since tagging only uses $P(s|o)$, we can ignore what the model says about $P(o)$. In particular, the model drawn assumes that each observation is independent, but we could add any arrows we please among the o_i without changing the conditional predictions. Therefore, it is common to think about this model as if the joint interpretation were absent, and not to model the observations at all. For models which are conditional in the sense of the factorization above, the JL and CL estimates for $P(s|o)$ will always be the same. It is therefore tempting to believe that since one can find closed-form CL estimates (the RFES) for these models, one can gain the benefit of conditional estimation. We will show that this is not true, at least not in this case.

Adopting a maximum entropy model has three kinds of assumptions. First, one has decided on the upward graphical model. This has effects in and of itself. The ML estimate for this model is the RFE for $P(s_i|s_{i-1}, o_i)$. For tagging, sparsity makes this impossible to reliably estimate directly, but assume we could do so. Still, we would have a graphical model with several defects. Every graphical model embodies conditional independence assumptions. The NB model assumes that observations are independent given the class. The HMM assumes the Markov property that future observations are independent from past ones given the intermediate state. These are both obviously false in the data, but the models do well enough for the tasks we ask of them. However, the assumptions in this upward model are far worse in how they impact performance. It is a conditional model, in that the model can be factored as $P(o)P(s|o)$. As a result, it makes no useful statement about the distribution of the data, making it useless, for example, for generation or language modeling. But more subtly note that states are independent of future observations. As a result, future cues are unable to influence past decisions in certain cases. For example, imagine observing an entire sentence where the first word is unknown. We might expect to be able to query our model for a distribution over tags for that word. With this model, that query will return the marginal distribution over

(sentence-initial) unknown words' tags, regardless of what the following words are.

We constructed two taggers. One was an HMM, as in figure 5b. It was trained for JL, CL*, and CL. The second was a CMM, as in figure 5c. We used a maximum entropy model over the (word,tag) and (previous-tag,tag) features to approximate the $P(s_i|o_i, s_{i-1})$ conditional probabilities. This CMM is referred to as an MEMM. The Penn treebank was used as the data corpus. To smooth these models as equally as possible and give as unified a treatment of unseen words as possible, we took all words which occurred only once in training and mapped them to an unknown token. All new words in the test data were also mapped to this token.⁸

Using these taggers, we examined what kinds of errors actually occurred. One kind of error tendency in CMMS which has been hypothesized in the literature is called *label bias* (Bottou, 1991; Lafferty et al., 2001). Label bias is a type of explaining away phenomenon (Pearl, 1988) which can be attributed to the local conditional modeling of each state. The idea is that states whose following-state distributions have low entropy will be preferred. Whatever mass arrives at a state must be pushed to successor states; it cannot be dumped on alternate observations as in an HMM. In theory, this means that the model can get into a dysfunctional behavior where a trajectory has no relation to the observations but will still stumble onward with high conditional probability. The sense in which this is an explaining-away phenomenon is that the previous state explains the current state so well that the observation at the current state is effectively ignored. What we found in the case of POS tagging was quite the opposite. The state-state distributions are on average nowhere near as sharply distributed as the state-observation distributions. This gives rise to the reverse explaining-away effect. The observations explain the states above them so well that the previous states are effectively ignored. We call this *observation bias*.

As an example, consider what happens when a word has only a single tag. Modulo smoothing effects, which only make this effect soft rather than categorical, the conditional distribution for the tag above that word will always assign conditional probability one to that single tag, regardless of the

⁸Doing so lowered our accuracy by approximately 2% for all models, but gave better-controlled experiments.

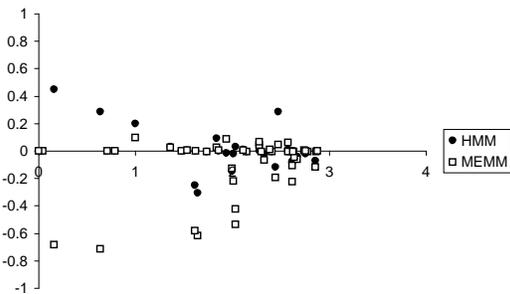


Figure 6: State transition entropy does not appear to be positively correlated with the relative over-proposal frequency of the tags for the MEMM model, though it is slightly so with the HMM model.

					HMM	MEMM	MEMM*
Correct States	PDT	DT	NNS	VBZ	-0.0	-1.3	-0.0
Incorrect States	DT	DT	NNS	VBZ	-5.4	-0.3	-5.7
Observations	All	the	indexes	dove			

Figure 7: The MEMM exhibits observation bias: knowing that *the* is a DT makes the quality of the DT-DT transition irrelevant, and *All* receives its most common tag (DT).

previous tag. Figure 7 shows the sentence *All the indexes dove .*, in which *All* should be tagged as a predeterminer (PDT).⁹ Most occurrences of *All*, however, are as a determiner (DT, 106/135 vs 26/135), and it is much more common for a sentence to begin with a determiner than a predeterminer. The other words occur with only one tag in the treebank.¹⁰ The HMM tags this sentence correctly, because two determiners in a row is rarer than *All* being a predeterminer (and a predeterminer beginning a sentence). However, the MEMM shows exactly the effect described above, choosing the most common tag (DT) for *All*, since the choice does not effect the deterministic conditional tagging distribution for *the*. The MEMM parameters do assign a lower score to DT DT than to PDT DT, but the *the* perfectly explains the second DT, regardless.

Exploiting the joint interpretation of the CMM, what we can do is to *unobserve* words. For example, we can retain our knowledge that the state above *the* is DT, but “forget” that we know that the word at that position is *the*. If we do inference in this example with *the* unobserved, then the conditional distribution over tag sequences changes as shown

⁹The treebank predeterminer tag is meant for when words like *All* are followed by a determiner, as in this case.

¹⁰For the sake of clarity, this example has been slightly doctored by the removal of several non-DT occurrences of *the* in the treebank – all incorrect.

under MEMM*: the correct tagging has once again become most probable. Unobserving the word itself is not *a priori* a good idea. It could easily put too much pressure on the last state to explain the fixed state. This effect is even visible in this small example: the likelihood of the more typical PDT-DT tag sequence is even higher for MEMM* than the HMM.

These issues are quite important for NLP, since state-of-the-art statistical taggers are all based on one of these two models. In order to check which, if either, of label or observation bias is actually contributing to tagging error, we performed the following experiments with our simple HMM and MEMM taggers. First, we measured, on the training data, the entropy of the next-state distribution $P(s_{+1}|s)$ for each state s . For both the HMM and MEMM, we then measured the relative overproposal rate for each state: the number of errors where that state was incorrectly predicted in the test set, divided by the overall frequency of that state in the correct answers. The label bias hypothesis makes a concrete prediction: lower entropy states should have higher relative overproposal values, especially for the MEMM. Figure 6 shows that the trends, if any, are not clear. There does appear to be a slight tendency to have higher error on the low-entropy tags for the HMM, but if there is any superficial trend for the MEMM, it is the reverse.

On the other hand, if systematically unobserving unambiguous observations in the MEMM led to an increase in accuracy, then we would have evidence of observation bias. Figure 4 shows that this is exactly the case. The error rate of the MEMM drops when we unobserve these single-tag words (from 10.8% to 9.5%), and the error rate in positions before such words drops even more sharply (17.1% to 15.0%). The drop in overall error in fact cuts the gap between the HMM and the MEMM by about half.

The claim here is not that label bias is impossible for MEMMs nor that state-of-the-art maxent taggers would necessarily benefit from the unobserving of fixed-tag words – if there are already (tag,next-word) features in the model, this effect should be far weaker. The claim is that the independence assumptions embodied by the conditionally structured model were the primary root of the lower accuracy for this model. Label bias and observation bias are both explaining-away phenomena, and are both consequences of these assumptions. Explaining-away effects will be found quite

generally in conditionally-structured models, and should be carefully considered before such models are adopted. The effect can be good or bad: In the case of the NB-CL model, there was also an explaining-away effects among the words. This is exactly the cause for *flowers* being estimated as a less strong indicator than *transatlantic* in our example. However, in that case, we *wanted* certain features to be suppressed in the presence of more explanatory features. However, when some of the competing conditioned features are previous local decisions, ignoring them can be harmful.

4 Related Results

Johnson (2001) describes two parsing experiments. First, he examines a PCFG trained on the ATIS treebank, both with the RFES which maximize JL, and using a CG ascent to maximize what we have been calling CL*. He does not give results for maximizing the unconstrained CL, but even in the constrained case, the effects from section 2 occur. CL and parsing accuracy are both higher for the parser trained to maximize CL*. He also describes a conditional shift-reduce parsing model, but notes that it underperforms the simpler joint formulation. We take these two results not as contradictory, but as confirmation that conditional estimation, though often slow, generally improves classification accuracy, while conditional model structures must be used with caution. The conditional shift-reduce parsing model he describes can be expected to exhibit the same type of competing-explanation issues that we discovered in MEMM tagging. It should exhibit label bias: as an extreme case, once all the words have been shifted, the rest of the parser actions will be reductions with probability one.

Goodman (1996) describes algorithms for parse selection where the criterion being maximized in parse selection is the bracket-based accuracy measures that parses are scored by. He shows that there is benefit to optimizing as close as possible to the “accuracy” criterion in parse selection, not just in parameter estimation.

Finally, model structure and parameter estimation are not the entirety of factors which determine the behavior of a model. Model features are crucial, and the ability to incorporate richer features in a relatively sensible way can also lead to improved models. This is the main basis of the real world benefit which has been derived from maxent models.

5 Conclusions

We have argued that optimizing an objective that is as close to the task “accuracy” as possible is advantageous in NLP domains, even in data-poor cases where machine-learning results suggest discriminative approaches may not be reliable. We have also argued that model structure is a far more important issue, and that for POS tagging, the observation bias effect of the model’s independence assumptions far outweighs label bias as a source of error. Our results, combined with others in the literature, suggest that a conditional model structure is, in and of itself, undesirable, unless that structure gives rise to a method of incorporating better features, explaining why maximum-entropy taggers and parsers have had such success despite the inferior performance of their basic skeletal models.

References

- Alan Agresti. 1990. *Categorical Data Analysis*. John Wiley & Sons, New York.
- Adam L. Berger, Stephen A. Della Pietra, and Vincent J. Della Pietra. 1996. A maximum entropy approach to natural language processing. *Computational Linguistics*, 22:39–71.
- D. P. Bertsekas. 1995. *Nonlinear Programming*. Athena Scientific, Belmont, MA.
- Léon Bottou. 1991. *Une approche theorique de l'apprentissage connexioniste; applications a la reconnaissance de la parole*. Ph.D. thesis, Université de Paris XI.
- Thorsten Brants. 2000. TnT – a statistical part-of-speech tagger. In *ANLP 6*, pages 224–231.
- S. Chen and R. Rosenfeld. 1999. A gaussian prior for smoothing maximum entropy models. Technical Report CMU CS-99-108, Carnegie Mellon University.
- Jerome H. Friedman. 1997. On bias, variance, 0/1-loss, and the curse-of-dimensionality. *Data Mining and Knowledge Discovery*, 1(1):55–77.
- William A. Gale, Kenneth W. Church, and David Yarowsky. 1992. A method for disambiguating word senses in a large corpus. *Computers and the Humanities*, 26:415–439.
- Joshua Goodman. 1996. Parsing algorithms and metrics. In *ACL 34*, pages 177–183.
- Mark Johnson. 2001. Joint and conditional estimation of tagging and parsing models. In *ACL 39*, pages 314–321.
- A. Kilgarriff. 1998. Senseval: An exercise in evaluating word sense disambiguation programs.
- John Lafferty, Fernando Pereira, and Andrew McCallum. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML*.

- Claudia Leacock, Martin Chodorow, and George A. Miller. 1998. Using corpus statistics and Wordnet relations for sense identification. *Computational Linguistics*, 24:147.
- Andrew McCallum and Kamal Nigam. 1998. A comparison of event models for naive bayes text classification. In *Working Notes of the 1998 AAAI/ICML Workshop on Learning for Text Categorization*.
- Andrew Y. Ng and Michael Jordan. 2002. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. In *NIPS 14*.
- Judea Pearl. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, San Mateo, CA.
- W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling. 1988. *Numerical Recipes in C*. Cambridge University Press, Cambridge.
- Adwait Ratnaparkhi. 1998. *Maximum Entropy Models for Natural Language Ambiguity Resolution*. Ph.D. thesis, University of Pennsylvania.
- Scott M. Thede and Mary P. Harper. 1999. Second-order hidden Markov model for part-of-speech tagging. In *ACL 37*, pages 175–182.
- Kristina Toutanova and Christopher D. Manning. 2000. Enriching the knowledge sources used in a maximum entropy part-of-speech tagger. In *EMNLP/VLC 2000*, pages 63–70.