

Use of Maximum Entropy in Back-off Modeling for a Named Entity Tagger

Rohini K. Srihari, Munirathnam Srikanth , Cheng Niu, Wei Li

Abstract—

In this paper, we describe improvements made to a statistical technique for “Named Entity ” (NE) Tagging. NE Tagging is the process of detecting and classifying proper noun sequences in text as one of name, location, organization, product, etc. For example, in the sentence “Tom Gates is the president of The Macrosoft Group”, Tom Gates is classified as a person name and The Macrosoft Group is classified as an organization name. It is an important step in both information extraction systems as well as information retrieval systems. The classifier is based on a Hidden Markov Model that uses various probabilities such as word co-occurrence, word features, and class probabilities. Since these probabilities are estimated based on observations seen in a corpus, “back-off models” are used to reflect the strength of support for a given statistic. For example, a word trigram uses a word bigram as an initial (first-order) back-off statistic. Typically, the back-off model consists of a linear combination of the n back-off statistics. The weights used in the linear combination are determined by adhoc procedures. Using maximum entropy techniques, it is possible to eliminate the use of back-off models, since the probabilities may be estimated more reliably. Experimental results using the NE tagger that is part of the information extraction system *Textract* are discussed in this paper.

I. INTRODUCTION

The last decade has seen great advances in the area of IE (Information Extraction). In US, the DARPA sponsored Tipster Text Program [Grishman 1997] and the Message Understanding Conferences (MUC) [1], [2] have been the driving force for developing this technology. Information retrieval technology, i.e., *search engines*, are able to direct users to documents that may contain the required information. It is necessary to wade through several documents before the information is actually found; this can be a very tedious process. Information extraction goes one step further by attempting to “understand” the document and summarize the salient information in a user-friendly manner. One of the goals of IE systems is general question-answering, e.g., *Who is the premier of Ontario and what year did he take office?*

Named entity tagging is a critical first step in any information extraction system. NE tagging refers to the task of identifying word groups representing entities such as people names, location names, organization names, dates and times, currency, etc. The number and type of NE classes may vary depending on the domain. In a medical domain, for example, NE classes may include drug names and disease names. NE tagging reflects the intuition that most information worth pursuing, concerns one or more of these named entities; thus, named entities serve as anchors for

mining valuable information. Depending on the granularity of the tagging, NE tagging alone can suffice in many IE applications. NE taggers typically add “mark-up”, in the form of special SGML tags, which denote the extent and type (person, place, etc.) of entity detected. Consider the following sentences which illustrate the output of the NE tagger which is part of the *Textract* system developed by our organization:

The British balloon, called the Virgin Global Challenger, is to be flown by
 <ENAMEX TYPE="PERSON"> Richard Branson</ENAMEX>, chairman of Virgin <ENAMEX TYPE="ORGANIZATION"> Atlantic Airways</ENAMEX>;
 <ENAMEX TYPE="PERSON">Per Lindstrand</ENAMEX>, chairman of <ENAMEX TYPE="ORGANIZATION"> Lindstrand Balloons Ltd. of Oswestry</ENAMEX>, <ENAMEX TYPE="LOCATION">England</ENAMEX>, and an Irish balloonist, <ENAMEX TYPE="ORGANIZATION"> Rory McCarthy</ENAMEX>.

<ENAMEX TYPE="PERSON">Branson</ENAMEX> and <ENAMEX TYPE="PERSON">Lindstrand</ENAMEX>, who have set several ballooning records, were the first pilots of hot-air balloons to cross both the <ENAMEX TYPE="LOCATION">Atlantic Ocean</ENAMEX>, in <TIMEX TYPE="DATE">1987</TIMEX>, and the <ENAMEX TYPE="LOCATION">Pacific</ENAMEX>, in <TIMEX TYPE="DATE">1991</TIMEX>.

As the above output indicates, some tagging errors are encountered. For example, the word “Virgin” is not included in the organization name “Virgin Atlantic Airways”. Initially, we implemented a Hidden Markov Model (HMM) for NE (Named Entity) to identify/classify proper names (person, organization, etc.) and other entities (time, date, percent, etc.). Due to the sparse data problem, although temporal and numerical NEs are expressed in English in very predictable patterns, there was a considerable amount of mistagging. We have identified cases where a statistical method is ineffective while simple pattern rules work exceptionally well. We then complemented the statistical model by a rule-based NE tagger. This tagger contains pattern match rules, or templates, for typical time and date expressions, as well as commonly seen name and organization expressions (e.g. names beginning with *Mr.* and organizations ending in *Inc.*). More details concerning the rule-based tagger are presented later in this paper. The pattern matcher is based on a Finite State Transducer (FST) toolkit developed in our organization. The rule-based tagger and the stochastic tagger are combined in a

pipeline architecture, with the rule-based tagger applied first. architecture. Using the taggers in tandem has resulted in a significant performance enhancement. Various corpus-based methods (e.g. [3] and rule-based methods (e.g. [4], [5]) have been applied to IE. But how to combine the strength of the two remains a topic to be further studied.

This paper focuses on the stochastic NE tagger, specifically, attempts to improve the performance of the tagger by refining the statistical modeling. Three different techniques are discussed and compared, including the use of ME in improving the estimation of the required probabilities.

II. STOCHASTIC MODEL

We follow the basic approach presented in [6], which is to consider the raw text encountered when decoding as though it had passed through a noisy channel. The assumption is that the original text had been correctly marked with named entities; the noisy channel corrupts this tagging. The task of the generative model is then to model the original process that generated the name-class annotated words, before they went through the noisy channel.

In another words, the goal is to find the most likely sequence of name-class (NC) given a sequence of words (W), that is,

$$Pr(NC | W) \quad (1)$$

In order to treat this as a generative model, we use Bayes' Rule:

$$Pr(NC | W) = \frac{Pr(W, NC)}{Pr(W)} \quad (2)$$

Since the *a priori* probability of the word sequence, i.e., the denominator is constant for any given sentence, equation(2) can be maximized by maximizing the numerator alone.

The model used here is an ergodic HMM, with eight internal states representing the five name classes, a NOT-A-NAME class, and tow special states representing START- and END-OF-SENTENCE states. Within each of the NC states, a statistical bigram language model is employed, with the usual one-word-per-state emission. This implies that the number of states in each of the NC states is equal to the vocabulary size, $|V|$.

The generation of words and name-classes proceeds in three steps:

1. Select a name-class NC , conditioning on the previous name-class and previous word.
2. Generate the first word inside that name-class, conditioning on the current and previous name-classes.
3. Generate all subsequent words inside the current name-class, where each subsequent word is conditioned on its immediate predecessor.

These three steps are repeated until the entire observed word sequence is generated. Using the Viterbi algorithm, it is possible to efficiently search the entire space of all possible name-class assignments, maximizing the numerator of Equation (2).

A. Words and Word-features

In this model, each individual term is considered as a two-element vector *word/feature*, denoted $\langle w, f \rangle$. The word feature is a simple, deterministic computation of each word as it is added to, or looked up in the vocabulary.

Table I illustrates the word features being used along with examples. Since it is not necessary to distinguish between two instances of certain features, e.g. *twoDigitNum*, *fourDigitNum*, a special word is used to represent all words having that feature. Thus statistics about the words *92* and *93* would be combined in the same word class.

Word Feature	Example	Intuition
twoDigitNum	99	Two-digit year
fourDigitNum	1998	Four-digit year
hasDollar	\$5,000.00	Monetary amount
hasPercent	5.75%	percentage
hasDigitAndAlpha	A89-23	Product code
hasDigitAndDash	01-98	Date
hasDigitAndSlash	12/31/97	Date
hasDigitAndComma	38,000.00	Monetary Amt
hasDigitandPeriod	1.00	Monetary Amt, Pct
otherNum	123456	Other Number
allCaps	IBM	Organization
capPeriod	J.	Person name initial
firstWord	first word of sent	No useful cap info
initCap	Jim	Capitalized word
lowerCase	winter	uncapitalized word
other	,	Punct,other words

TABLE I
WORD FEATURES, EXAMPLES AND REASON FOR USE

B. Calculation of Probabilities

The probability for generating the first word of a name-class is factored into two parts:

$$Pr(NC | NC_{-1}, w_{-1}) \times Pr(\langle w, f \rangle_{first} | NC, NC_{-1}) \quad (3)$$

The top level model for generating all but the first word in a name-class is

$$Pr(\langle w, f \rangle | \langle w, f \rangle_{-1}, NC) \quad (4)$$

A special *+end+* word is used to compute any current word to be the final word of its name-class. That is, when $NC_{-1} = START - OF - SENTENCE$ in $Pr(NC | NC_{-1}, w_{-1}), w_{-1} = +end+$.

The calculation of the above probabilities are as follows:

$$Pr(NC | NC_{-1}, w_{-1}) = \frac{c(NC, NC_{-1}, w_{-1})}{c(NC_{-1}, w_{-1})} \quad (5)$$

$$Pr(\langle w, f \rangle_{first} | NC, NC_{-1}) = \frac{c(\langle w, f \rangle_{first}, NC, NC_{-1})}{c(NC, NC_{-1})} \quad (6)$$

$$Pr(\langle w, f \rangle | \langle w, f \rangle_{-1}, NC) = \frac{c(\langle w, f \rangle, \langle w, f \rangle_{-1}, NC)}{c(\langle w, f \rangle_{-1}, NC)} \quad (7)$$

In the above equations, $c()$ represents the frequency of the events in the training data.

It is also necessary to train an *unknown word* model due to frequent occurrence of out-of-vocabulary words. This is achieved by partitioning the training data into two halves. Training for unknown words is computed on each half separately, and the results merged. It should be noted that the unknown word model can be significantly improved by utilizing features such as word suffixes and prefixes. The unknown word model can also be considered as the first level of back-off.

III. HYBRID APPROACH: COMBINE RULE-BASED AND STOCHASTIC TAGGER

The essential argument for this strategy is that by combining statistical models with an FST rule-based system, we will be able to exploit the best of both paradigms while overcoming their respective weaknesses. Traditional rule based systems have serious drawbacks such as knowledge acquisition bottleneck and lack of consistency/robustness, efficiency and portability. To overcome these weaknesses, the last decade has seen increasing efforts in two types of state-of-the-art technology applied to Natural Language Processing. They are Finite State Transducers and statistical models. However, they each have their own weaknesses. We discuss below their respective advantages and disadvantages, which led to our hybrid approach of combining the two for IE.

The most attractive feature of the FST formalism lies in its superior time and space efficiency [7] Applying a deterministic FST depends linearly only on the input size of the text. Our experiments also show that an FST rule system is extraordinarily robust. Besides, it has been verified by many research programs [4], [8], that FST is also a convenient tool for capturing linguistic phenomena, especially for idioms and semi-productive expressions. While FST overcomes the efficiency and robustness problems of the traditional rule based system, an FST rule system still suffers from the same disadvantage in knowledge acquisition as that for all hand-crafted systems.

The advantage of statistical corpus-based models lies in its automatic *learning* capabilities, *robustness*, and *portability*. A domain shift only requires the truthing of new data for re-training the model, with no need to change the program. However, a purely statistical model has serious limitation in sparse data, recognized as a major bottleneck for all corpus-based models [9]. In the present hybrid architecture of the NE tagger, the NE processing goes through two successive stages in accordance with the pipeline architecture just presented: first, the NE rules will identify the NEs by FST pattern matching; the remaining NEs will then be tagged by the HMM sub-module.

The rules which we have implemented currently include a grammar for temporal expressions (time, date, duration, frequency, age, etc.), a grammar for numerical expressions (money, percentage, length, weight, etc.), a small grammar for person in clear patterns, and a grammar for other non-MUC NEs (e.g. contact information such as address,

email).

IV. TRAINING, EVALUATION

As mentioned before, our work is modeled after MUC standards. This has enabled proper benchmarking of our NE performance using the MUC provided scorer. The *F - measure* statistic is used to evaluate performance; *f-measure* is a weighted combination of *precision*, and *recall*, the two most frequently used measures in information retrieval systems. They are defined as follows:

$$Prec = \frac{\text{number correct responses}}{\text{total number of responses}}$$

$$Recall = \frac{\text{num correct responses}}{\text{total num correct responses}}$$

$$F - \text{measure} = \frac{(b * b + 1) Recall * Prec}{(b * b * Recall + Prec)}$$

where $b * b$ represents a relative weight, typically 1.

The system has been trained using 312,712 words: these represent edited English text, such as that from the newswire organizations. As such, they are presumably free of spelling and grammatical errors. It should be noted that the text to which the NE tagger will be applied, exhibits these same characteristics.

In terms of speed, the total time required for NE tagging is roughly a simple addition of the two sub-modules as there is little overhead in combining the two NE modules in the pipeline architecture. Tested on PC Pentium II 450, the tagging speed is 25.5KB/second (or 91.9MB/hour). This speed is among the best in the NE world and is comparable to the speed of NetOwl and Nymble, 2 NE products used for industrial applications. It is expected that the on-going optimization of our system will further raise the speed of processing. In terms of performance, as shown in Table 1 below, our hybrid NE system using the proposed pipeline architecture obtains an overall F-scores 91.24F-scores of over 90% are generally considered as near human performance, certainly good enough for commercial application.

When tested alone, the FST sub-module obtains 34.37% F-score while the HMM NE sub-module achieves 84.25%. In other words, the addition of the pattern match rules gave us an increase of nearly 7 percentage points. Note that this score is achieved without the use of gazetteers, generally regarded as an important lexical information source for the NE task.

In general, for numerical and temporal NEs which are expressed in English in highly predictable patterns, the pattern matching rules achieve both high precision (numerical NE: 100%, temporal NE: 74% and 61%) and high recall (all above 94%). For other NE types (such as *person*), where there exists a high degree of trade-off between precision and recall, we only attempt to implement rules of high-precision at the expense of recall. This strategy is decided by our design: the FST NE module is supposed to overcome the weaknesses of the HMM module which never overrides any results from the FST NE. This leads to the expected NE result for *person* with high-precision (97%) and low-recall

(14%), leaving the bulk of the person tagging work to the statistical model.

Further overall error analysis shows the following:

- Many common location names were missed because there happened to be no obvious contextual cues captured by *Textract*.
- There is a significant amount of mistagging between organization entities and person entities because these two types are often used in similar contexts. Many such mistakes will be corrected once the gazetteers for location names, registered corporate names and frequently-used person names are incorporated into the system.

V. BACK-OFF MODELS

Due to the relatively few number of observations in the training data for the statistics being computed, the estimates can be very unreliable at times. To compensate for this, a back-off modeling strategy is necessary. Table II illustrates the various levels of back-off for each of the three statistics used in the model.

Name-Class bigrams	First-word Bigrams	Non-first-word Bigrams
$Pr(NC NC_{-1}, w_{-1})$	$Pr(\langle w, f \rangle_{first} NC, NC_{-1})$	$Pr(\langle w, f \rangle \langle w, f \rangle_{-1}, NC)$
$Pr(NC NC_{-1})$	$Pr(\langle w, f \rangle \langle begin, other \rangle, NC)$	$Pr(\langle w, f \rangle NC)$
$Pr(NC)$	$Pr(\langle w, f \rangle NC)$	$Pr(w NC)^*$ $Pr(f NC)$
$\frac{1}{\#name-classes}$	$Pr(w NC)^*$ $Pr(f NC)$	$\frac{1}{ v \#wordfeatures}$
	$\frac{1}{ V \#wordfeatures}$	

TABLE II
BACK-OFF MODELS

A linear interpolation strategy is employed for backing-off. Given k models $\{P_i(w | h)\}_{i=1\dots k}$, they can be combined linearly as follows:

$$P_{Combined}(w | h) = \sum_{i=1}^k \lambda_i P_i(w | h)$$

where $0 < \lambda_i \leq 1$ and $\sum_i \lambda_i = 1$.

Linear interpolation is both a way of combining knowledge sources and smoothing training data. Because it reconciliates the different information sources in a straightforward manner, it has the following advantages:

- Generality: Any language model may be used as a component.
- Easy to implement, experiment with and analyze.
- Guaranteed to be no worse than any of its componenets. This is true since each component may be viewed as a special case of the interpolation, with a weight of 1 for that component and 0 for all others.

However, linear interpolation models also suffer from some weaknesses, chiefly their inconsistency. Each information source typically partitions the event space (h, w) and provides estimates based on the relative frequency of training

data within each class of the partition. Therefore, within each of the component models, the estimates are consistent with the training data. But this reasonable measure of consistency is in general violated by the interpolated model.

The next two sections discuss two different strategies for computing the weights λ_i .

VI. FIXED WEIGHT BACKOFF

In the first set of experiments, a fixed weight was assigned to each statistic, including each of the backoff statistics as well. These weights were determined empirically and were fixed, i.e., not trigram specific. Some sample weights that were used:

$$Pr(NC | w_{-1}, NC_{-1}) \quad \lambda = 0.5$$

$$Pr(NC | NC_{-1}) \quad \lambda = 0.4$$

$$Pr(NC) \quad \lambda = 0.09$$

For fixed weights, the best performance, i.e., F-measure obtained was 90.37%.

VII. DYNAMICALLY COMPUTED WEIGHTS

In this experiment, trigram-specific weights for the back-off models were computed as follows:

For each conditional probability $Pr(X | Y)$, assign a weight of λ to the direct computation and a weight of $(1-\lambda)$ to the back-off model where

$$\lambda = \frac{1}{1 + \frac{\text{uniqueoutcomesof}Y}{c(Y)}}$$

. In the ideal case, the weights for the original model, i.e., $Pr(X | Y)$ would be high, reflecting the high reliability based on sufficient samples in the training data.

Using the above scheme for computing weights for $Pr(NC | NC_{-1}, w_{-1})$ and $Pr(\langle w, f \rangle_{-1} | NC, NC_{-1})$, an F-measure of 85.91% was obtained. The remaining weights used the same fixed weighting scheme as described in the previous section. Computing weights for $Pr(NC | NC_{-1}, w_{-1})$ only, resulted in an F-measure of 90.29% was obtained. An analysis of the results, namely, decrease in performance using computed weights indicated that the original models were being weighted very heavily, usually $\lambda \geq .8$. This can at least in part be attributed to the insufficiency of the training data. Many words w_i are observed co-occurring with the same word w_j several times in the corpus; more important, the training corpus does not contain enough examples of other co-occurrences involving these words (which are present in a general corpus).

It is possible to refine the weight computation as in [6] by computing the weights dynamically. In this case, λ is computed as

$$\lambda = \left(1 - \frac{oldc(Y)}{c(Y)}\right) \star \frac{1}{1 + \frac{\text{uniqueoutcomesof}Y}{c(Y)}}$$

VIII. USE OF ME FOR ESTIMATING NAME CLASS: NO BACKOFF

In the final model to be discussed, no backoff model is employed. Instead, ME techniques are used to more accurately estimate the probability $P(NC_i|NC_{i-1}, w_{i-1})$. The following constraints were used in estimating the above probability.

In Shannon's seminal work [10], he established a quantitative measure of information, in terms of the uncertainty in a probability distribution. Shannon's measure of uncertainty is given by the following:

$$S(p-1, p_2, \dots, p_n) = -k \sum_{i=1}^n p_i \ln p_i$$

While Shannon provided the theoretical framework, Jaynes [11] provided a quantitative technique for deriving the required probability distributions. In particular, Jaynes discusses the maximum entropy principle (*MaxEnt*). MaxEnt allows the selection of the probability distribution (from the infinite number possible) that maximizes the Shannon entropy measure, and is simultaneously consistent with all the given constraints. MaxEnt states that the quantity to be maximized is

$$-\sum_{i=1}^n p_i \ln p_i$$

subject to

$$\sum_{i=1}^n p_i = 1,$$

$$\sum_{i=1}^n p_i g_r(x_i) = \sum_{i=1}^n p_i g_{ri} = a_r, \quad r = 1, 2, \dots, m,$$

and

$$p_i \geq 0, i = 1, 2, \dots, n.$$

The g_r s represent various expected values such as mean and variance. Recently, entropy optimization principles [12] have been used in various applications, ranging from image processing to statistical language modeling.

For the problem at hand, instead of using the constraints on the conditional probability $P(NC_i|NC_{i-1}, w_{i-1})$, the constraints on $P(NC_i, NC_{i-1}, w_{i-1})$ are used, since this corresponds to the observations (training data).

The conditional probability, $P(NC_i|NC_{i-1}, w_{i-1})$ is computed using

$$P(NC_i|NC_{i-1}, w_{i-1}) = \frac{P(NC_i, NC_{i-1}, w_{i-1})}{P(NC_{i-1}, w_{i-1})}. \quad (8)$$

The constraints are the expected values of the features in the training data. For each triple $(NC_i, NC_{i-1}, w_{i-1})$, the maximum likelihood estimate of its probability is taken as the expected value of the probability of the triple. Let $f(NC'_i, NC'_{i-1}, w'_{i-1})$ be the observed frequency of the current word, w'_i , belonging to name class NC'_i with the previous word w'_{i-1} from name class NC'_{i-1} . Then one of the

constraints is

$$\sum_{i=1}^n P(NC_i, NC_{i-1}, w_{i-1}) k(NC_i, NC_{i-1}, w_{i-1}) = f(NC'_i, NC'_{i-1}, w'_{i-1}) \quad (9)$$

where

$$k(NC_i, NC_{i-1}, w_{i-1}) = \begin{cases} 1 & \text{if } NC_i = NC'_i, NC_{i-1} = NC'_{i-1}, \text{ and } w_{i-1} = w'_{i-1} \\ 0 & \text{otherwise} \end{cases} \quad (10)$$

$k(NC_i, NC_{i-1}, w_{i-1})$ is an indicator function.

In addition to the constraint in equation 9, the following constraints on the bigram and unigram frequencies are used in the estimation of $P(NC_i, NC_{i-1}, w_{i-1})$.

$$\sum_{i=1}^n P(NC_i, NC_{i-1}, w_{i-1}) k_2(NC_i, NC_{i-1}, w_{i-1}) = f(NC'_i, NC'_{i-1}) \quad (11)$$

$$\sum_{i=1}^n P(NC_i, NC_{i-1}, w_{i-1}) k_1(NC_i, NC_{i-1}, w_{i-1}) = f(NC'_i) \quad (12)$$

Here, $k_1(\cdot)$ and $k_2(\cdot)$ are the corresponding indicator functions.

In the ME model, the probability $P(NC_i, NC_{i-1}, w_{i-1})$ is estimated by maximizing the Shannon entropy measure subject to the constraints in equation 9, equation 11, equation 12 and the natural constraint of probabilities

$$\sum_{i=1}^n P(NC_i, NC_{i-1}, w_{i-1}) = 1. \quad (13)$$

The maximum entropy probability distribution (MEPD) is computed using the Improved Iterative Scaling (IIS) method [13]. The IIS method is a special case of the Generalized Iterative Scaling (GIS) method [14] where the constraints are expected values of a 0,1-valued function. The indicator function satisfy this condition. Let m be the number of constraints. The IIS method starts with some initial value for the probability distribution, $P_0(\mathbf{x})$. In iteration n , the probability distribution $P_n(\mathbf{x})$ is obtained as a scaled value $P_{n-1}(\mathbf{x})$, i.e. $P_n(\mathbf{x}) = s_n P_{n-1}(\mathbf{x})$. The scaling factor, s_k , is estimated by constraining the probability distribution $P_k(\mathbf{x})$ to satisfy the k th constraint, where $k = n \pmod{m}$. Since the number of constraints can be enormous, computational issues may pose a problem. In this case, we have 8 named classes with, approximately, 30,000 unique words in the training set. This results in approximately $m = 8 * 8 * 30000 + 1$ number of constraints. Fortunately, the nature of the constraints caused the system to converge fairly soon.

Unfortunately, the strategy of eliminating backoff and using ME estimated conditional probability led to significant degradation in performance (on the same test document). There are several explanations for this:

- The entire event space was not used in estimating the probability distribution. In other words, to ..
- There were not sufficient training samples to reliably estimate the probability distribution.
- More thought needs to be given as to what constraints should be used; there may be some which are inconsistent.

It should be noted that performance evaluation was with respect to the end result of the tagger whereas we only changed part of the tagging process with maximum entropy estimated values. A fair comparison would be in terms of predicting the nameclass as against tagging words.

IX. FUTURE WORK

Much work remains, obviously. This includes using ME to estimate conditional probabilities by considering the *entire* event space. More consideration needs to be given to pruning out constraints that lead to inconsistencies. This is expected to lead to improved estimation.

We are also interested in using ME to more effectively combine our stochastic and NE taggers. Since they use orthogonal techniques, the combination should lead to more information. It is particularly noteworthy that New York University [15] has developed an NE system, MENE, whereby ME is used to optimally combine several independently developed NE taggers. The strength of their ME approach lies in its capability of evaluating the relative probabilities of different outputs provided by the external systems. The ME approach is superior to a pipeline architecture in cases where earlier modules in the pipeline have low precision; a pipeline model will propagate the errors. This problem can be avoided for applications like NE tagging by designing rules which capture only clear patterns (with high precision and possibly low recall) while more subtle cases are left to the statistical model.

X. CONCLUSION

This paper has discussed named entity tagging, a valuable first step in any information extraction system. A stochastic model for NE tagging has been presented along with improvements by introducing a rule-based NE tagger as well. The focus of the paper was on improvements to the stochastic model by considering three alternative strategies. The first two concerned the use of back-off models to compensate for unreliable probability estimates. The third method involved the use of ME techniques to more accurately estimate the conditional probabilities, thus eliminating the need for back-off models.

Results (as of writing this paper) indicate that the simplistic, fixed weight scheme is outperforming both the other, more sophisticated methods. This is expected to change as more refinement of both the back-off model, and the constraints used in the ME technique takes place.

REFERENCES

- [1] Beth Sundheim, Ed., *Proceedings of the Sixth Message Understanding Conference (MUC-6)*. Morgan Kaufmann, 1995.
- [2] *Proceedings of the Seventh Message Understanding Conference (MUC-7)*. Morgan Kaufmann, 1998.
- [3] S. Miller et al., "BBN: Description of the SIFT System as Used for MUC-7," in *Proceedings of the Seventh Message Understanding Conference (MUC-7)*, 1998.
- [4] Krupka G.R. and Hausman K., "IsoQuest Inc: Description of the NetOwl Text Extractor System as Used for MUC-7," in *Proceedings of the Seventh Message Understanding Conference (MUC-7)*, 1998.
- [5] D. E. Appelt et al., "SRI International: FASTUS System MUC-6 Test Results and Analysis," in *Proceedings of the Sixth Message Understanding Conference (MUC-6)*. 1995, Morgan Kaufmann Publishers.
- [6] D. M. Bikel et al., "Nymble: a high-performance learning name-finder," in *Proceedings of the Fifth Conference on Applied Natural Language Processing*. 1997, pp. 194–201, Morgan Kaufmann Publishers.
- [7] M. Mohri, "Finite-State Transducers in Language and Speech Processing," *Computational Linguistics*, vol. 23, no. 2, pp. 269–311, 1997.
- [8] J. R. Hobbs, "Fastus: A System for Extracting Information from Text," in *Proceedings of the DARPA workshop on Human Language Technology*, 1993, pp. 133–137.
- [9] Eugene Charniak, *Statistical Language Learning*, MIT Press, 1994.
- [10] C.E. Shannon, "A Mathematical Theory of Communication," *Bell System Technical Journal*, vol. 27, 1948.
- [11] E.T. Jaynes, "Information Theory and Statistical Mechanics," *Physical Reviews*, vol. 106, 1957.
- [12] J.N. Kapur and H.K. Kesavan, *Entropy Optimization Principles with Applications*, Academic Press, 1992.
- [13] Frederick Jelinek, *Statistical Methods in Speech Recognition*, MIT Press, 1997.
- [14] J. N. Darroch and D. Ratcliff, "Generalized iterative scaling for log-linear models," *The Annals of Mathematical Statistics*, pp. 1470–1480, 1972.
- [15] Borthwick A., Sterling J., Agichtein, and Grishman R., "NYU: Description of the MENE Named Entity System as used in MUC-7," in *Proceedings of the Fifth Conference on Applied Natural Language Processing*. 1998, Morgan Kaufmann Publishers, also <http://www.muc.saic.com/>.

XI. NOTES

* The work reported in this paper is supported by DoD SBIR Phase 1 grant (Contract No. F30602-97-C-0179) and Phase 2 grant (Contract No. F30602-98-C-0043).