

# BOOSTEXTER FOR TEXT CATEGORIZATION IN SPOKEN LANGUAGE DIALOGUE

Marie Rochery, Robert Schapire, Mazin Rahim, Narendra Gupta

AT&T Labs - Research, 180 Park Avenue, Florham Park, NJ 07932

## ABSTRACT

Data collection plays a critical role in the development of accurate syntactic and semantic models for natural-language dialogue systems. This task, which is known to be both labor intensive and financially expensive, is the bottleneck when rapidly prototyping voice-enabled services. In this paper, we propose a mathematical framework for minimizing the reliance of data by refining human expertise with the statistics of any available training data. In particular, we extend BoosTexter, a member of the boosting family of algorithms, to combine and balance hand-crafted rules with the statistics of the data in spoken language understanding. Experiments on two voice-enabled applications for customer care and help desk are presented.

## 1. INTRODUCTION

Building spoken natural-language dialogue systems for automated customer care and help desk applications presents several technical challenges: (1) the need for large vocabulary recognition to accommodate for the variety of input requests, (2) parsing and understanding users' requests, and (3) supporting mixed-initiative and conversational dialogue. The open-format natural language input for these sets of applications is significantly more complex than travel reservation systems [4] for example, and presents a major challenge to both speech recognition and language understanding.

In this paper, we describe one component of our dialogue system, namely the module for identifying a user's request. We consider this as a multi-class multi-label text categorization problem. We apply a machine-learning system called BoosTexter [7] which was based on the boosting family of algorithms first proposed by Freund and Schapire [2, 5]. BoosTexter is strictly data driven. It combines many simple and moderately accurate categorization rules that are trained sequentially into a single, highly accurate rule that can accurately predict a class. It has been shown to outperform traditional methods for text categorization [6].

This paper presents an extension to BoosTexter that combines and balances human expertise with the statistics of the training data. Human expertise ranges from pragmatic knowledge to application specific rules. Our proposed algorithm provides several advantages: (1) it enables rapid prototyping of spoken natural-language services when the amount of data available is severely limited, (2) it provides

a framework for supporting new service requests for which no data is available.

We present two sets of experiments, one for a customer care application and the other for a Help Desk application. Our results show that including expert knowledge when data is limited can help to reduce the amount of data necessary for building these applications by factors of 2 to 4.

## 2. BOOSTEXTER

We begin by describing the machine-learning system BoosTexter.

We assume that we are given a set of training examples  $(x_1, y_1), \dots, (x_m, y_m)$ . Each  $x_i$  is called an *instance*. In this paper, each  $x_i$  will generally be the text of a transcribed or recognized utterance; however, in general,  $x_i$  may incorporate other information about what was spoken. Each  $y_i$  is the *label* or *class* assigned to the instance  $x_i$ ; for instance,  $y_i$  may indicate call type. For simplicity, we assume for now that there are only two classes,  $-1$  and  $+1$ .

The goal of a learning algorithm is to use the training data to derive a rule that accurately predicts the class of any new instance  $x$ ; such a prediction rule is called a *classifier*. The approach that we take is based on a machine-learning method called *boosting* [2, 5]. In particular, we use a variant of Schapire and Singer's BoosTexter system [7].

The basic idea of boosting is to build a highly accurate classifier by combining many "weak" or "simple" *base classifiers*, each one of which may only be moderately accurate. To obtain these base classifiers, we assume we have access to a *base learning algorithm* that we use as a black-box subroutine.

The collection of base classifiers is constructed in rounds. On each round  $t$ , the base learner is used to generate a base classifier  $h_t$ . Besides supplying the base learner with training data, the boosting algorithm also provides a set of non-negative weights  $w_t$  over the training examples. Intuitively, the weights encode how important it is that  $h_t$  correctly classify each training example. Generally, the examples that were most often misclassified by the preceding base classifiers will be given the most weight so as to force the base learner to focus on the "hardest" examples.

Following Schapire and Singer [6], we use *confidence-rated* classifiers  $h$  that, rather than outputting simply  $-1$  or  $+1$ , output a real number  $h(x)$  whose sign ( $-1$  or  $+1$ ) is

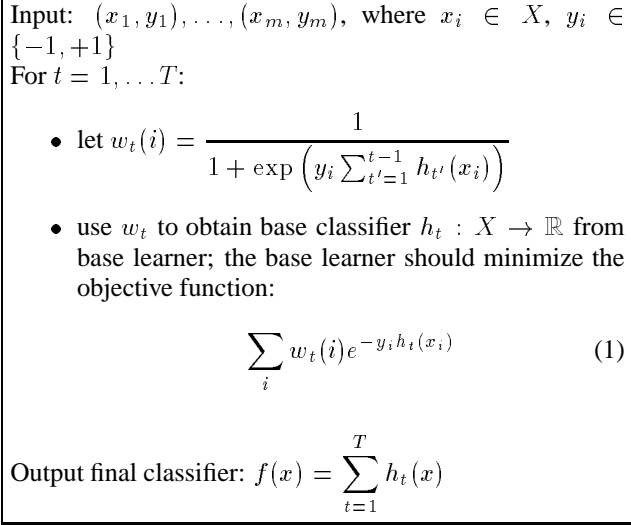


Figure 1: A binary boosting algorithm.

interpreted as a prediction, and whose magnitude  $|h(x)|$  is a measure of “confidence.”

The pseudo-code of the boosting algorithm is given in Figure 1. This is a variant of Freund and Schapire’s original AdaBoost algorithm [2] as modified by Collins, Schapire and Singer [1] to minimize logistic loss rather than exponential loss.

The real-valued predictions of the final classifier  $f$  can be converted into probabilities by passing them through a logistic function; that is, we can regard the quantity

$$\frac{1}{1 + e^{-f(x)}}$$

as an estimate of the probability that  $x$  belongs to class  $+1$ . In fact, the boosting procedure here described is designed to minimize the negative conditional log likelihood of the data under this model, namely,

$$\sum_i \ln(1 + e^{-y_i f(x_i)}). \quad (2)$$

The base learning algorithm that we use is the same as in Schapire and Singer’s BoosTexter system [7]. In particular, each base classifier tests for the presence or absence of a particular word, short phrase or other simple pattern, henceforth referred to simply as a *term*. If the term is present, then one prediction is made; otherwise, some other prediction is made. For instance, the base classifier might be: “If the word ‘yes’ occurs in the utterance, then predict  $+1.731$ , else predict  $-2.171$ .” Schapire and Singer [7] describe a base learning algorithm that efficiently finds the best base classifier of this form, i.e., the one minimizing criterion (1).

Schapire and Singer [6, 7] also describe in detail how to extend this approach to multiclass problems in which more than two classes are allowed and furthermore in which each

example may belong to multiple classes. The intuitive idea is to reduce to binary questions which ask if each example is or is not in each of the classes.

## 2.1. Incorporating human knowledge

Boosting, like many machine-learning methods, is entirely data-driven in the sense that the classifier it generates is derived exclusively from the evidence present in the training data itself. When data is abundant, this approach makes sense. However, in some applications, data may be severely limited, but there may be human knowledge that, in principle, might compensate for the lack of data.

In its standard form, boosting does not allow for the direct incorporation of such prior knowledge. In this section, we describe a modification of boosting that combines and balances human expertise with available training data. We aim for an approach that allows the rough human judgments to be refined, reinforced and adjusted by the statistics of the training data in a well controlled manner.

As before, we limit our attention to binary classification; extension to multiclass problems can be made along the lines of Schapire and Singer [6, 7]. In our approach, a human expert must begin by constructing a rule  $p$  mapping each instance  $x$  to an estimated probability  $p(x) \in [0, 1]$  which is interpreted as the guessed probability that instance  $x$  belongs to class  $+1$ . We discuss below some methods for constructing such a function  $p$ .

To apply boosting using  $p$  and a training set, we create a new *weighted* training set. This new set includes all of the original training examples  $(x_i, y_i)$ , each with unit weight. In addition, for each training example  $(x_i, y_i)$ , we create two new training examples  $(x_i, +1)$  and  $(x_i, -1)$  with weights  $\eta p(x_i)$  and  $\eta(1 - p(x_i))$ , respectively, where  $\eta$  is a parameter of the algorithm controlling the confidence in encoding knowledge. During training, these weights  $w_0$  are used in computing  $w_t$  so that

$$w_t(i) = \frac{w_0(i)}{1 + \exp\left(y_i \sum_{t'=0}^{t-1} h_{t'}(x_i)\right)}$$

(here,  $i$  ranges over all of the examples in the *new* training set).

One final modification that we make is to add a 0-th base classifier  $h_0$  that is based on  $p$  so as to incorporate  $p$  right from the start. In particular, we take

$$h_0(x) = \ln\left(\frac{p(x)}{1 - p(x)}\right)$$

and include  $h_0$  in computing the final classifier  $f$ .

Essentially, these modifications have the effect of changing the objective function in (2) to one that incorporates prior knowledge, namely,

$$\sum_i \ln(1 + e^{-y_i f(x_i)}) + \eta \sum_i \text{RE}\left(p(x_i) \parallel \frac{1}{1 + e^{-f(x_i)}}\right)$$

where  $\text{RE}(p \parallel q) = p \ln(p/q) + (1-p) \ln((1-p)/(1-q))$  is binary relative entropy. Thus, we balance the conditional likelihood of the data against the distance of the data-generated model from the model provided by the human. The relative importance of the two terms is controlled by the parameter  $\eta$ .

### 3. PRIOR KNOWLEDGE

Prior knowledge may be acquired from several sources, e.g. human judgment, application guidelines and manuals, world knowledge, and in-domain website. In fact while developing a spoken dialogue system designers do have access to one or more such sources of knowledge. Designers use these sources of knowledge to deduce information crucial for the development of the dialogue system, i.e. the functionalities to support, and a basic understanding of how users may interact with the application. It would be only prudent, therefore, to also use these sources of knowledge for bootstrapping the text categorization module needed for the natural language understanding, especially when data is limited.

As an example, prior knowledge allows us to encode rules that can classify user responses to confirmation questions like: “So you want to fly from Boston to New York on Sunday evening?” A user response containing “yes”, “okay”, “correct”, “all right”, or “fine”, etc. is highly indicative of a positive confirmation. We can formally express this by a rule with an estimated probability, of say 0.9:

- `yes|okay|correct|all right|fine → class(Yes, .9).`

Another example of applying prior knowledge is for classifying users requests to be connected to an operator/service agent. This can be expressed by the following rule:

- `speak & (human|operator|(service & agent))  
→ class(Agent, .95).`

In general to incorporate prior knowledge the three logical operators, OR, AND and NOT must be supported. More specifically given such rules (with logical connectives), we require a method to map them onto the BoosTexter 0-th base classifier  $h_0$ . To do this, disjunctive terms in the rules (prior knowledge) are represented as separate classifiers. Each of these classifiers are assigned a  $p(x)$  value equal to the subjective probability assigned to rules in prior knowledge. The values of  $p(x)$  are then used to compute  $h_0(x)$  as explained in section 2.1.

## 4. EXPERIMENTS

In this section, we describe and analyze the experiments we performed using BoosTexter for text categorization with in-domain knowledge provided by the application guidelines.

### 4.1. Test databases

We ran experiments with data from two different applications. The first database used is the *How May I Help You?* database [3]. In this task, there are 15 different classes. We did experiments with 50 to 1600 sentences in the training set and 2991 sentences in the test set. The second database is for an application called “*TTS Help Desk*”. This application provides information about AT&T Text-to-Speech engine. There are 22 different classes. We trained models on 100 to 2675 sentences and tested on 2000 sentences.

### 4.2. Results

In the first experiment that was performed on the *How May I Help You?* database, we measured the classification accuracy as a function of the number of examples used during training. The classification accuracy is the percentage of sentences with a correctly predicted label. There is no rejection and the label with the highest score is kept as the predicted label. We compared models built only with some training examples and models built with both hand-crafted rules (prior knowledge) and training examples. We used approximately one rule per class where a rule is a combination of many words or phrases. We trained the models on 50, 100, 200, 300 rounds when the number of available training examples was respectively 50, 100, 200, 400 and up. The parameter  $\eta$  was selected empirically based on the number of available training examples. We set  $\eta$  to 1 when the number of training examples was less than or equal to 200, 0.1 when it was between 400 and 800, and 0.01 when it was greater.

The dashed line in Figure 2 shows the classification accuracy for models built on hand-crafted rules and training examples whereas the solid lines show the classification accuracy for models built either on training examples only or on hand-crafted rules only. An improvement in accuracy is observed when using hand-crafted rules and training examples together. This comes from the fact that some patterns of the hand-crafted rules are not in the data at all or are not in a sufficient number of sentences to have a statistical impact when training only on the data.

In this experiment, when fewer training examples were available (<100 examples) exploiting human expertise provided classification accuracy levels that are equivalent to models trained on four times the amount of training data. When the number of training examples is larger (> 100), accuracy levels becomes equivalent to two times the amount of training data. When larger than 6000 sentences were available, both models were found to converge to similar classification accuracy.

The second experiment was performed on the Help Desk task. Figure 3 shows the comparison between models built with hand-crafted rules and training examples and models built only with training examples. In this experiment there is also about 1 rule per class where a rule is again a combination of many words and phrases. We trained the models on

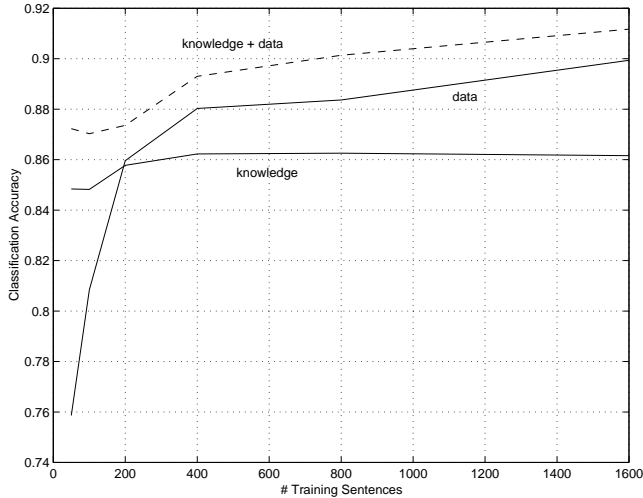


Figure 2: Comparison of performance using data and knowledge separately or together on the How May I Help You task.

100, 200, 400, 600, 800, 1000 when the number of training examples was respectively 100, 200, 400, 800, 1600, and up. We set  $\eta$  to 0.1 when the number of training examples was less than or equal to 1600 and to 0.01 otherwise.

Figure 3 shows an improvement in classification accuracy when hand-crafted rules are being used. This improvement is up to 9% absolute with 100 training examples and drops to 0.5% when more data becomes available.

We can notice from the experiments that the number of rounds and the choice of the parameter  $\eta$  are dependent on the number of training examples available. We can also notice that the knowledge-only curves are not perfectly flat. This comes from the fact that the models from the knowledge take into account the class distribution of the available training examples as explained in the section 3.

## 5. SUMMARY

The use of BoosTexter for text categorization in natural-language understanding was described in this paper. We presented an extension to BoosTexter that incorporates human knowledge of the application in the form of hand-crafted rules. Each set of rules, associated with estimated probabilities of the distribution of the class, is refined and enhanced by the statistics of the training data. A new objective function was proposed which has an additional relative entropy term that balances the conditional likelihood of the data against the distance of the data-generated model from that obtained using human knowledge.

Two experiments were conducted on speech data collected from natural-language dialogue applications. Both experiments demonstrate that introducing prior knowledge of the domain can significantly cut the amount of data needed for building the application. For both applications, exploit-

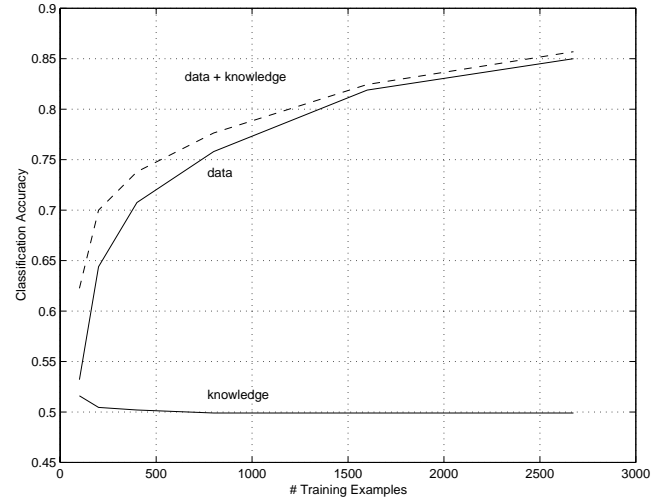


Figure 3: Comparison of performance using data and knowledge separately or together on the Help Desk task.

ing human expertise helped to reduce reliance on the data by factors of two to four.

## Acknowledgment

The authors would like to acknowledge the technical contribution of H. Alshawi, S. Bangalore, S. Douglas, Giuseppe Di Fabrizio, G. Riccardi and Y. Singer.

## References

- [1] Michael Collins, Robert E. Schapire, and Yoram Singer. Logistic regression, AdaBoost and Bregman distances. In *Proceedings of the Thirteenth Annual Conference on Computational Learning Theory*, 2000.
- [2] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, August 1997.
- [3] A.L. Gorin, G. Riccardi, and J.H. Wright. How May I Help You? *Speech Communication*, 23:113–127, 1997.
- [4] E. Levin, S. Narayanan, R. Pieraccini, K. Biatov, E. Bocchieri, G. Fabbriozio, W. Eckert, S. Lee, A. Pokrovsky, M. Rahim, P. Ruscitti, and M. Walker. The AT&T darpa communicator mixed-initiative spoken dialogue system. In *ICSLP*, 2000.
- [5] Robert E. Schapire. A brief introduction to boosting. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence*, 1999.
- [6] Robert E. Schapire and Yoram Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37(3):297–336, December 1999.
- [7] Robert E. Schapire and Yoram Singer. BoosTexter: A boosting-based system for text categorization. *Machine Learning*, to appear.