

Fast Multi-task Learning for Query Spelling Correction

Xu Sun
Dept. of Statistical Science
Cornell University
Ithaca, NY 14853
xusun@cornell.edu

Anshumali Shrivastava
Dept. of Computer Science
Cornell University
Ithaca, NY 14853
anshu@cs.cornell.edu

Ping Li
Dept. of Statistical Science
Cornell University
Ithaca, NY 14853
pingli@cornell.edu

ABSTRACT

In this paper, we explore the use of a novel online multi-task learning framework for the task of search query spelling correction. In our procedure, correction candidates are initially generated by a ranker-based system and then re-ranked by our multi-task learning algorithm. With the proposed multi-task learning method, we are able to effectively transfer information from different and highly biased training datasets, for improving spelling correction on all datasets. Our experiments are conducted on three query spelling correction datasets including the well-known TREC benchmark dataset. The experimental results demonstrate that our proposed method considerably outperforms the existing baseline systems in terms of accuracy. Importantly, the proposed method is about one order of magnitude faster than baseline systems in terms of training speed. Compared to the commonly used online learning methods which typically require more than (e.g.,) 60 training passes, our proposed method is able to closely reach the empirical optimum in about 5 passes.

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Query Alteration*

General Terms

Algorithms, Performance, Experimentation

Keywords

Query Spelling Correction, Multi-task Learning

1. INTRODUCTION

Search queries present a particular challenge for traditional spelling correction methods, for at least three reasons [3]. Firstly, spelling errors are more common in search queries than in regular written texts. For example, [12]

showed that roughly 10-15% of queries contain misspelled terms. Secondly, most search queries consist of only a few key words rather than grammatical sentences, making a grammar-based approach inappropriate. Most importantly, many queries contain search terms, such as proper nouns and names, which are not well established in the language. For example, Chen *et al.* [10] reported that 16.5% of valid search terms do not occur in their spelling lexicon, which contains more than 200,000 entries.

Due to the practical importance, query spelling correction has received much attention and a variety of datasets have been developed (from different domains). As discussed in [17], queries in one specific dataset can be very biased from another dataset, and combining such biased datasets for query spelling correction is difficult. In fact, in this paper, we will also demonstrate an empirical phenomenon that simply merging biased datasets for training a unified speller may not bring improvement at all.

The difficulty in effectively combining biased datasets for spelling correction lies in at least three aspects: (i) their error-patterns can be high biased; (ii) their distributions of misspelled queries can be quite different; (iii) their domains can be quite different (hence, a domain adaptation problem exists). Thus, a natural question arises: *how can we effectively integrate such highly biased datasets for improving query spelling correction?* To the best of our knowledge, there are still no satisfactory solutions in the literature.

Our goal in this paper is to solve this well-known difficult problem with high accuracy as well as high efficiency. The training speed can be crucial because industrial query spelling correction tasks may have to deal with very large-scale training data (e.g., click-through query spelling logs which contain more than millions of entries, see [34]). Our proposed multi-task learning framework for query spelling correction is capable of adaptively integrating highly biased training datasets via automatically learning the task-relationships (data-similarities). This allows us to “softly” merge biased datasets for learning a unified speller.

In our procedure, the first step is to generate a large number of spelling correction candidates using a candidate generation module. Then the top candidates are accurately identified through our novel multi-task learning algorithm. Encouragingly, with the ability to effectively transfer information among those biased training datasets, our proposed method is capable of improving spelling correction accuracies on *all* related datasets. The prior standard practice for combining different datasets is to simply merge them and feed them to a single-task learner. Compared to this

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CIKM'12, October 29–November 2, 2012, Maui, HI, USA.

Copyright 2012 ACM 978-1-4503-1156-4/12/10 ...\$15.00.

baseline, our experimental results will demonstrate that our new approach will result in significantly better performance.

Our **major contributions** can be summarized as follows:

- To the best of our knowledge, this is the first study of multi-task learning for query spelling correction. The proposed method outperforms the state-of-the-art results on the TREC dataset.
- Our novel *adaptive* multi-task learning framework can effectively utilize *2nd-order gradient information*.
- Our method is of fast convergence speed. It is able to approach close to the empirical optimum in only a few (e.g., ≤ 5) passes. This is about one order of magnitude faster than traditional training (e.g., batch or online) methods. Fast training speed is crucial for large-scale query spelling correction tasks (e.g., [34]).

The rest of this paper is structured as follows. Section 2 describes the architecture of our spelling correction system. Section 3 presents the proposed multi-task learning framework for query spelling correction. Section 4 presents the extensive experiments on query spelling correction. Section 5 reviews related work, and Section 6 concludes the paper.

2. SYSTEM ARCHITECTURE

2.1 Candidates Generation

We first describe our procedure for generating correction candidates. Following Ganjisaffar *et al.* [17], we implement three candidate generators in our system to produce about 2,000 candidates (on average) for each query. Initially, a character-based candidate generator is adopted for producing all possible candidates within an edit distance 1. It considers replacing each character with all possible characters in the alphabet, transposing each pair of adjacent characters, deleting each character, inserting all possible characters after each character, and so on.

The AOL query logs showed that 16% of the query corrections are different from the original query only in adding/removing spaces [17]. For example, “ebayauction” is a query which should be corrected to “ebay auction”. As a more complicated example, the long query “broccoliandcheesebake” actually indicates “broccoli and cheese bake”. In order to handle this class of queries, we implement the word segmentation algorithm based on the Microsoft word breaker¹, with some minor modifications. For each possible segmentation of the character sequence, we use a language model to compute the probability of the segmentation. The most probable segmentation candidates are then added to the candidate list.

Neither of the above two candidate generators is able to produce candidates which have more than 1 edit distance from the original query (in spite of adding or removing multiple spaces). For example, for the query “washton university”, we might want to have “washington university” as a candidate, which however can not be produced by the above two candidate generators. To ease this problem, we perform a fuzzy search process based on a lexicon containing most frequent words, to quickly find known unigrams with a small

¹<http://web-ngram.research.microsoft.com>

edit distance to unigrams and bigrams located in the original query. For the lexicon containing frequent-word, we use the top 100K words based on their frequency on the web. With the three candidate generators, we can generate 2,000 candidates (on average) for each query.

2.2 Naive Ranking

An important factor in selecting and ranking the correction candidates is the prior probability of a correction phrase. It represents our prior belief about how likely a query will be chosen by the user without seeing any input from the user. In this work we make use of the Web n-gram service provided by Microsoft. Web n-gram model intends to model the n-gram probability of English phrases with the parameters estimated from the entire Web data. It also differentiates the data sources to build different language models from the title, anchor text and body of Web pages, as well as the queries from query log. To build our spelling system, we make use of the tri-gram language model. Note that, even though the Web n-gram model contains a large amount of entries, it may still suffer from data sparseness in higher-order language models.

2.3 Re-Ranking

Re-ranking of top search results has shown to improve the quality of rankings in query spelling correction [17, 28]. Since the main focus of spelling correction is on the top ranked list of candidates, we add a re-ranker on top of the naive ranker. This step significantly improved the results of the naive ranker, and the cost is tractable. This is because re-ranking only needs to deal with top- k candidates, and k is typically small (in our case, following [28], we set $k = 40$). Hence, the re-ranker is solving a much easier problem, compared to the original naive ranker which needs to consider about 2,000 candidates for each query.

We use a conditional log-linear method, the *Maximum Entropy* model, for re-ranking. The maximum entropy model produces a probability distribution over multiple classes and have the advantage of handling large numbers of overlapping features. Assuming a feature function that maps a pair of observation sequence \mathbf{x} and a classification label y to a feature vector \mathbf{f} , the probability of y conditioned on the observation sequence \mathbf{x} is modeled as follows [23]:

$$P(y|\mathbf{x}, \mathbf{w}) = \frac{\exp[\mathbf{w}^\top \mathbf{f}(y, \mathbf{x})]}{\sum_{y'} \exp[\mathbf{w}^\top \mathbf{f}(y', \mathbf{x})]}, \quad (1)$$

where \mathbf{w} is a weight vector. Given a training set consisting of n labeled samples, (\mathbf{x}_i, y_i) , for $i = 1 \dots n$, weights are gained via maximizing the objective function,

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n \log P(y_i|\mathbf{x}_i, \mathbf{w}) - R(\mathbf{w}). \quad (2)$$

The first term of this equation represents a conditional log-likelihood of the training data. The second term is a regularizer for reducing overfitting. We employ an L_2 prior. In what follows, we denote the conditional log-likelihood of each sample, $\log P(y_i|\mathbf{x}_i, \mathbf{w})$, as $\ell(i, \mathbf{w})$. Then, we have

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n \ell(i, \mathbf{w}) - \frac{\|\mathbf{w}\|^2}{2\sigma^2}. \quad (3)$$

Since the evaluation is based on expected F-score, probability information is required in computing expected F-score.

The original conditional probability produced by maximum entropy models is not proper for computing expected F-score, because the sum of conditional probabilities of top- k candidates are not 1. To deal with this problem, we recompute the probability of a candidate as follows:

$$P'(y) = \exp[cP(y)] / \sum_{k'=1, \dots, k} \exp[cP(y_{k'})],$$

where c is a scalar to control the density of the distribution. Since essentially c has a similar function like the σ in regularization, to avoid introducing new hyper-parameters, we simply set c based on σ .

2.4 Re-Ranking Features

In the re-ranking phase we add features which are extracted from top- k query-candidate pairs. In addition, other than the features which based on transformations between query-candidate pairs, we also design features which consider the top candidates of the query (i.e., comparisons with other candidates in the top- k list). In this way, the features can include valuable information that may help in the ranking process.

We have about 30 feature templates for each query-candidate pair. These features include: error model features (e.g., edit distance), candidate language model features, query language model features, surface features capturing differences of the query and the candidate, frequency of the query and the candidate, the rank of the candidate in the naive ranking phase, and so on. Although we have used language model scores for naive ranking for producing top candidates, we keep using Web-scale n-gram language model features for re-ranking the top candidates, and we find such Web-scale n-gram language model features are still quite useful in the re-ranking phase.

In the Web-scale n-gram language model features, the log of n-gram language model probabilities of an original query and its candidate corrections are used for re-ranking. In addition, the average, max, min, and standard deviation of the language model probabilities of the top- k candidates are employed as features.

2.5 Online Learning

There are two major approaches for training a log-linear model: batch training and online training. Batch training methods include, for example, steepest gradient descent, conjugate gradient descent (CG), and limited-memory BFGS (LBFGS) [30]. In such training methods, gradients are computed by using all training instances. Typically, the training process is quite slow in practice.

To speed up the training process, online algorithms have become increasingly popular. A representative online learning method is the stochastic gradient descent (SGD) [7]. The SGD uses a small randomly-drawn subset of the training samples to approximate the gradient of the objective function, which allows one to update the model weights much more frequently, and consequently, to speed up the convergence. Suppose $\hat{\mathcal{S}}$ is a randomly drawn subset of the full training set \mathcal{S} , the stochastic objective function is then given by

$$\mathcal{L}_{stoch}(\mathbf{w}, \hat{\mathcal{S}}) = \sum_{i \in \hat{\mathcal{S}}} \ell(i, \mathbf{w}) - \frac{|\hat{\mathcal{S}}|}{|\mathcal{S}|} \frac{\|\mathbf{w}\|^2}{2\sigma^2}.$$

The extreme case is a batch size of 1, and it gives the maximum frequency of updates, which we adopt in this work. In this case, $|\hat{\mathcal{S}}| = 1$ and $|\mathcal{S}| = n$ (suppose the full training set contains n samples). In this case, we have

$$\mathcal{L}_{stoch}(\mathbf{w}, \hat{\mathcal{S}}) = \ell(i, \mathbf{w}) - \frac{1}{n} \frac{\|\mathbf{w}\|^2}{2\sigma^2}, \quad (4)$$

where $\hat{\mathcal{S}} = \{i\}$. Model weights are updated like this:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \gamma_k \nabla_{\mathbf{w}_k} \mathcal{L}_{stoch}(\mathbf{w}, \hat{\mathcal{S}}), \quad (5)$$

where k is the update counter, γ_k is the learning rate.

3. OUR PROPOSAL

In this section, we introduce our adaptive (and online) multi-task learning framework (MTL below). For every positive integer q , we define $\mathcal{N}_q = \{1, \dots, q\}$. Let T be the number of tasks which we want to simultaneously learn. For each task $t \in \mathcal{N}_T$, there are n data examples $\{(\mathbf{x}_{t,i}, \mathbf{y}_{t,i}) : i \in \mathcal{N}_n\}$ available. In practice, the number of examples per task may vary but we keep it constant for simplicity of notation. We use \mathbf{D} to denote the $n \times T$ matrix whose t -th column is given by the vector \mathbf{d}_t of data examples.

3.1 Adaptive Multi-Task Learning Model

Our goal is to learn the weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_T$ from the data \mathbf{D} . For simplicity of notation, we assume that each of the weight vectors is of the same size f (this is also the feature dimension), and corresponds to the same ordering of features. We use \mathbf{W} to denote the $f \times T$ matrix whose t -th column is given by the vector \mathbf{w}_t . We learn \mathbf{W} by maximizing the objective function,

$$\text{Obj}(\mathbf{W}, \mathbf{D}) \triangleq \text{Likelihood}(\mathbf{W}, \mathbf{D}) - R(\mathbf{W}), \quad (6)$$

where $\text{Likelihood}(\mathbf{W}, \mathbf{D})$ is the accumulative likelihood over all interactive tasks, namely,

$$\text{Likelihood}(\mathbf{W}, \mathbf{D}) = \sum_{t \in \mathcal{N}_T} \mathcal{L}(\mathbf{w}_t, \mathbf{D}), \quad (7)$$

and $\mathcal{L}(\mathbf{w}_t, \mathbf{D})$ is defined as follows:

$$\mathcal{L}(\mathbf{w}_t, \mathbf{D}) \triangleq \sum_{t' \in \mathcal{N}_T} [\alpha_{t,t'} \mathcal{L}(\mathbf{w}_t, \mathbf{d}_{t'})]. \quad (8)$$

$\alpha_{t,t'}$ is a real-valued *task-similarity*, with $\alpha_{t,t'} = \alpha_{t',t}$ (symmetric). Intuitively, a task-similarity $\alpha_{t,t'}$ measures the *similarity of patterns* between the t -th task and the t' -th task. $\mathcal{L}(\mathbf{w}_t, \mathbf{d}_{t'})$ is defined as follows:

$$\begin{aligned} \mathcal{L}(\mathbf{w}_t, \mathbf{d}_{t'}) &\triangleq \sum_{i \in \mathcal{N}_n} \log P(\mathbf{y}_{t',i} | \mathbf{x}_{t',i}, \mathbf{w}_t) \\ &= \sum_{i \in \mathcal{N}_n} \ell_{t'}(i, \mathbf{w}_t), \end{aligned} \quad (9)$$

where $P(\cdot)$ is a prescribed probability function. We can flexibly use any prescribed probability function. This makes our MTL method a flexible and general framework for no matter structured or non-structured classification tasks. In this paper, we will adopt the maximum entropy probability function (Eq. 1), which works well for spelling correction.

Finally, $R(\mathbf{W})$ is a regularization term for dealing with overfitting. In this paper, we simply use L_2 regularization:

$$R(\mathbf{W}) = \sum_{t \in \mathcal{N}_T} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2}. \quad (10)$$

MTL with fixed task-similarities (MTL-F)

Input: Initialize $\mathbf{W}^{(0)}$; given $\mathbf{D}, \mathbf{A}^*, \beta$; $k \leftarrow 0$
for $t \leftarrow 1$ to T
 . Initialize $\boldsymbol{\eta}^{(0)}$
 . **Repeat** until convergence
 . . $\mathbf{g}_t \leftarrow -\frac{1}{n} \nabla_{\mathbf{w}_t} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2}$
 . . **for** $t' \leftarrow 1$ to T
 . . . Draw $i \in \mathcal{N}_n$ at random
 . . . $\mathbf{g}_t \leftarrow \mathbf{g}_t + \mathbf{A}_{t,t'}^* \nabla_{\mathbf{w}_t} \ell_{t'}(i, \mathbf{w}_t)$
 . . $\mathbf{w}_t^{(k+1)} \leftarrow \mathbf{w}_t^{(k)} + \boldsymbol{\eta}^{(k)} \cdot \mathbf{g}_t$
 . . **if** $k+1 \bmod 2 = 0$
 . . . $\mathbf{v}_i \leftarrow \frac{\mathbf{w}_t^{(k+1)}(i) - \mathbf{w}_t^{(k)}(i)}{\mathbf{w}_t^{(k)}(i) - \mathbf{w}_t^{(k-1)}(i)}$
 . . . Lower-bounds \mathbf{v}_i with β
 . . . $\boldsymbol{\eta}^{(k+1)} \leftarrow \mathbf{v} \cdot \boldsymbol{\eta}^{(k)}$
 . . **else**
 . . . $\boldsymbol{\eta}^{(k+1)} \leftarrow \boldsymbol{\eta}^{(k)}$
 . . $k \leftarrow k+1$
Output: $\forall t, \mathbf{w}_t^{(k)}$ converges to \mathbf{w}_t^* ; $\mathbf{W}^{(k)}$ converges.

MTL with unknown task-similarities (MTL)

Input: Initialize $\mathbf{W}^{(0)}, \mathbf{A}^{(0)}$; given \mathbf{D} ; $k \leftarrow 0$
Repeat until convergence
 . $\mathbf{W}^{(k+1)} \leftarrow \text{MTL-F}(\mathbf{W}^{(k)}, \mathbf{A}^{(k)}, \mathbf{D})$
 . **for** $t \leftarrow 1$ to T
 . . **for** $t' \leftarrow 1$ to T
 . . . Update $\mathbf{A}_{t,t'}^{(k+1)}$ with Eq.20/21
 . $k \leftarrow k+1$
Output: $\mathbf{A}^{(k)}$ converges to $\hat{\mathbf{A}}$; $\mathbf{W}^{(k)}$ converges to $\hat{\mathbf{W}}$.

Figure 1: MTL algorithms (using batch size of 1). The derivation of $\frac{1}{n}$ before the regularization term was explained in Eq. 4. Lower-bounding \mathbf{v}_i with β is for stability consideration in the online setting.

To summarize, the overall objective function is as follows:

$$\text{Obj}(\mathbf{W}, \mathbf{D}) = \sum_{t, t' \in \mathcal{N}_T} \left[\alpha_{t,t'} \sum_{i \in \mathcal{N}_n} \ell_{t'}(i, \mathbf{w}_t) \right] - \sum_{t \in \mathcal{N}_T} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2}.$$

To simplify denotations, we introduce a $T \times T$ matrix \mathbf{A} , such that $\mathbf{A}_{t,t'} \triangleq \alpha_{t,t'}$. We also introduce a $T \times T$ functional matrix Φ , such that $\Phi_{t,t'} \triangleq \mathcal{L}(\mathbf{w}_t, \mathbf{d}_{t'})$. Then, the objective function can be compactly expressed as follows:

$$\text{Obj}(\mathbf{W}, \mathbf{D}) = \text{tr}(\mathbf{A}\Phi^\top) - \sum_{t \in \mathcal{N}_T} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2}, \quad (11)$$

In the following content, we will first discuss a simple case that the task-similarity matrix \mathbf{A} is fixed. After that, we will focus on the case that \mathbf{A} is unknown, because the task-relationships are unknown among the three query spelling correction tasks that we will focus on.

3.2 MTL with Fixed Task-Similarities

Although the task-similarities are unknown for query spelling correction, we will present a learning algorithm that iteratively reduce the problem to a case of fixed task-similarities. Hence, it is important to discuss the case of fixed task-similarities. With fixed task-similarities, the op-

timization problem is as follows:

$$\mathbf{W}^* = \underset{\mathbf{W}}{\text{argmax}} \left[\text{tr}(\mathbf{A}^* \Phi^\top) - \sum_{t \in \mathcal{N}_T} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2} \right]. \quad (12)$$

It is clear to see that we can independently optimize \mathbf{w}_t and $\mathbf{w}_{t'}$ ($t \neq t'$) given fixed task-similarities. Hence, we can independently optimize each column of \mathbf{W} and derive \mathbf{W}^* :

$$\mathbf{w}_t^* = \underset{\mathbf{w}_t}{\text{argmax}} \psi(\mathbf{w}_t, \mathbf{D}), \quad (13)$$

where $\psi(\mathbf{w}_t, \mathbf{D})$ has the form as follows:

$$\psi(\mathbf{w}_t, \mathbf{D}) = \sum_{t' \in \mathcal{N}_T} \left[\alpha_{t,t'}^* \mathcal{L}(\mathbf{w}_t, \mathbf{d}_{t'}) \right] - \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2}. \quad (14)$$

3.2.1 2nd-Order Gradient Information

For high convergence speed, an important issue of MTL-F is to effectively and efficiently approximate the Hessian matrix. Following the work of [20] on single-task learning, we present a simple yet effective method to approximate the eigenvalues of the Jacobian matrix of a fixed point iterative mapping. In MTL-F, the update formula is as follows:

$$\mathbf{w}_t^{(k+1)} = \mathbf{w}_t^{(k)} + \boldsymbol{\eta}_t \cdot \mathbf{g}_t, \quad (15)$$

The update term \mathbf{g}_t is derived by weighted sampling over different tasks. The weighted sampling is based on fixed task-similarities, \mathbf{A}^* . \mathbf{g}_t has a form as follows:

$$\mathbf{g}_t = \sum_{t' \in \mathcal{N}_T} \left[\alpha_{t,t'}^* \nabla_{\mathbf{w}_t} \ell_{t'}(i_{t'}, \mathbf{w}_t) \right] - \frac{1}{n} \nabla_{\mathbf{w}_t} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2}, \quad (16)$$

where $\alpha_{t,t'}^* = \mathbf{A}_{t,t'}^*$ and $i_{t'}$ indexes a random sample selected from $\mathbf{d}_{t'}$. Then, the expectation (over distribution of data) of the update term is as follows:

$$\begin{aligned} \mathbb{E}(\mathbf{g}_t) &= \sum_{t' \in \mathcal{N}_T} \left\{ \alpha_{t,t'}^* \left[\frac{1}{n} \nabla_{\mathbf{w}_t} \mathcal{L}(\mathbf{w}_t, \mathbf{d}_{t'}) \right] \right\} - \frac{1}{n} \nabla_{\mathbf{w}_t} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2} \\ &= \frac{1}{n} \left\{ \sum_{t' \in \mathcal{N}_T} \left[\alpha_{t,t'}^* \nabla_{\mathbf{w}_t} \mathcal{L}(\mathbf{w}_t, \mathbf{d}_{t'}) \right] - \nabla_{\mathbf{w}_t} \frac{\|\mathbf{w}_t\|^2}{2\sigma_t^2} \right\} \\ &= \frac{1}{n} \nabla_{\mathbf{w}_t} \psi(\mathbf{w}_t, \mathbf{D}). \end{aligned}$$

In addition, $\boldsymbol{\eta}_t \in \mathbb{R}_+^f$ is a positive vector-valued step size and \cdot denotes component-wise (Hadamard) product of two vectors. As presented in [6], the optimal step size is the one that asymptotically approaches to \mathbf{H}_t^{-1} , the inverse Hessian matrix of $\psi(\mathbf{w}_t, \mathbf{D})$ in our setting. To avoid actually evaluating \mathbf{H}_t^{-1} , we can approximate \mathbf{H}_t^{-1} with its eigenvalues. Following [20], we consider an update iterate as a fixed-point iterative mapping (though a stochastic one) \mathcal{M} . Taking partial derivative of \mathcal{M} with respect to \mathbf{w}_t , we have

$$\mathbf{J}_t = \mathcal{M}' = \mathbf{I} - \text{diag}(\boldsymbol{\eta}_t) \mathbf{H}_t. \quad (17)$$

By exploiting this linear relation between Jacobian and Hessian, we can obtain approximate eigenvalues of inverse Hessian using eigenvalues of Jacobian:

$$\text{eigen}_i(\mathbf{H}_t^{-1}) = \frac{\boldsymbol{\eta}_t(i)}{1 - \text{eigen}_i(\mathbf{J}_t)}. \quad (18)$$

In addition, $\text{eigen}_i(\mathbf{J}_t)$ can be asymptotically approximated:

$$\text{eigen}_i(\mathbf{J}_t) \approx \lambda_i \triangleq \frac{\mathbf{w}_t^{(k+1)}(i) - \mathbf{w}_t^{(k)}(i)}{\mathbf{w}_t^{(k)}(i) - \mathbf{w}_t^{(k-1)}(i)}. \quad (19)$$

When k is sufficiently large, λ_i will be sufficiently close to $eigen_i(\mathbf{J}_t)$. Therefore, we can asymptotically approximate the inverse of the Hessian matrix via efficient estimation of the Jacobian matrix of fixed-point mapping.

In multi-task setting, this optimization problem is a cost-sensitive optimization problem. To summarize our discussion, we present the *2nd-order gradient based MTL algorithm with fixed task-similarities* (MTL-F below) in Figure 1 (upper). As we can see, the algorithm iteratively approximates the eigenvalues of the inverse Hessian matrix via Eq. 19. This approximation is based on data points that are cost-sensitively sampled from multiple tasks. Then, it updates the vector-valued step size based on the inverse Hessian information, and the multi-task model weights are updated accordingly. This iterative process continues until convergence.

3.3 MTL with Unknown Task-Similarities

For our query spelling correction tasks, the task-similarities are unknown. To solve this problem, we present an algorithm to learn the task-similarities and model weights in an alternating optimization manner. Our alternating learning algorithm with unknown task-similarities, called *MTL*, is presented in Figure 1 (bottom).

In the MTL learning, the MTL-F algorithm is employed as a subroutine. In the beginning of the MTL, model weights \mathbf{W} and task-similarities \mathbf{A} are initialized. \mathbf{W} is then optimized to $\hat{\mathbf{W}}$ by using the MTL-F algorithm, based on the fixed \mathbf{A} . Then, in an alternative way, \mathbf{A} is updated based on the optimized weights $\hat{\mathbf{W}}$. After that, \mathbf{W} are optimized based on updated (and fixed) task-similarities. This iterative process continues until empirical convergence of \mathbf{A} and \mathbf{W} .

In updating task-similarities \mathbf{A} based on \mathbf{W} , a natural idea is to estimate a task-similarity $\alpha_{t,t'}$ based on the similarity between weight vectors, \mathbf{w}_t and $\mathbf{w}_{t'}$. It is unclear which similarity measure best fits the tasks of query spelling correction. To study this, we propose two candidate similarity measures for query spelling correction:

Polynomial kernel (Poly): We can use (normalized) polynomial kernel to estimate similarities:

$$\alpha_{t,t'} \triangleq \frac{1}{C} \frac{\langle \mathbf{w}_t, \mathbf{w}_{t'} \rangle^d}{\|\mathbf{w}_t\|^d \cdot \|\mathbf{w}_{t'}\|^d}, \quad (20)$$

where $\langle \mathbf{w}_t, \mathbf{w}_{t'} \rangle$ means inner product between the two vectors; d is the degree of the polynomial kernel; $\|\mathbf{w}_t\|^d \cdot \|\mathbf{w}_{t'}\|^d$ is the normalizer. C is a real-valued constant for tuning the magnitude of task-similarities. Intuitively, a big value of C will result in “weak multi-tasking” and a small value of C will make “strong multi-tasking”. For example, when $d = 1$, the normalized kernel has exactly the form $\frac{1}{C} \cos \theta$, where θ is the angle between \mathbf{w}_t and $\mathbf{w}_{t'}$ in the Euclidean space.

Correlation (Cor): Since the *covariance* of task weight vectors is a natural way to estimate inter-task interactions, we consider using covariance information for estimating task-similarities. However, we find directly using a covariance matrix (to estimate task-similarities) faces the problem of stability in our online setting. Hence, we use the correlation matrix via normalizing the covariance matrix.

$$\alpha_{t,t'} \triangleq \frac{1}{C} \frac{cor(\mathbf{w}_t, \mathbf{w}_{t'})}{std(\mathbf{w}_t)std(\mathbf{w}_{t'})}, \quad (21)$$

Table 1: Statistics of the three spelling correction datasets, and the top-1 accuracy of naive ranking (without re-ranking). #Candidates is the number of correction candidates generated for re-ranking.

Data	#Queries	#Candidates	Top-1 Acc. of NR (%)
JDB	1.2×10^4	4.8×10^5	45.7
MSN	5.0×10^3	2.0×10^5	35.5
TREC	6.0×10^3	2.4×10^5	31.0

where $cov(\mathbf{w}_t, \mathbf{w}_{t'})$ is the covariance between \mathbf{w}_t and $\mathbf{w}_{t'}$. $std(\cdot)$ is standard deviation.

3.4 Accelerated MTL Learning

The MTL learning algorithm can be further accelerated. The naive 2MTL learning algorithm waits for the convergence of the model weights \mathbf{W} (in the MTL-F step) before updating the task-similarities \mathbf{A} . In practice, we can update task-similarities \mathbf{A} before the convergence of the model weights \mathbf{W} .

For example, we can update task-similarities \mathbf{A} after running the MTL-F step over a small number of training passes. We will adopt this accelerated version of the MTL learning for experiments so that the task-similarities will not be repeatedly updated. In the experiment section, we will compare the accelerated MTL method with a variety of strong baseline methods.

4. EXPERIMENTS ON TARGET DATA

We will test the proposed method on three query spelling correction tasks. We will study the MTL method with different similarity measures for spelling correction, and compare them against a number of strong baselines, including traditional batch and online learning methods. The results have been averaged over 5 runs for random permutations of the training data order, and standard deviations are given.

4.1 Datasets

Since the **TREC** dataset (released by Microsoft Speller Challenge 2011) is a new and well-known benchmark dataset for query spelling correction, we will use this dataset for experiments. The TREC dataset is based on the publicly available TREC queries (2008 Million Query Track). This dataset contains 5,892 queries and corrections annotated by the Speller Challenge organizers. There could be more than one plausible corrections for a query. In this dataset only 5.3% of queries are judged as misspelled. We use the same split of the training and testing data.

To improve the performance of query spelling correction, we use two auxiliary datasets for performing multi-task learning. One auxiliary dataset is the **JDB** dataset [17] collected from the publicly available AOL query sets, with a total of 12,000 samples, and 2,000 of them (16.7%) are different from the original query. The second auxiliary dataset is collected from **MSN** queries, and for each query there is at most only one correction. In this dataset, about 11% of queries are judged as misspelled. Also, we evenly split the two auxiliary datasets into training and testing data. We show the statistics of the three datasets in Table 1. In the table, we also show the top-1 accuracy (accuracy of the top-1 candidate) of naive ranking using language models (without

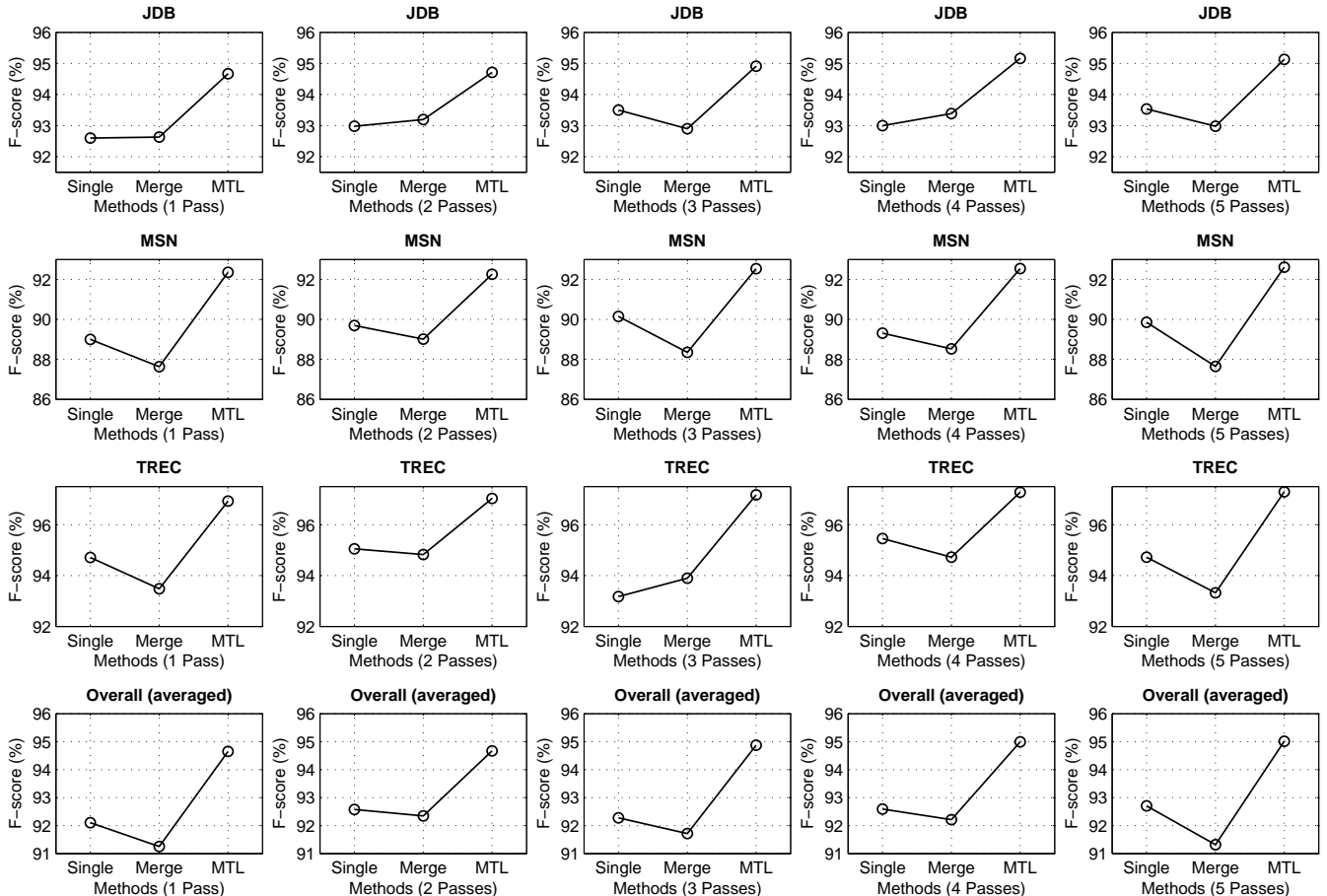


Figure 2: F-scores of different methods in 5 passes. *Single* and *Merge* are based on SGD training.

re-ranking). As we can see, the naive ranking has low accuracies. This reflects the importance of the re-ranking.

The two auxiliary datasets are biased from the TREC dataset in three aspects. First, the queries in those datasets are from three different domains: TREC, AOL, and MSN domains. Second, the distribution of misspelled queries are very different: 5.3% in TREC dataset, 16.7% in JDB dataset, and 11% in MSN dataset. Finally, the number of correct spelling suggestions are different: For the MSN dataset, for each query there is only one correction. On the other hand, for the TREC and JDB datasets, for each query there can be multiple corrections. Hence, it is expected to be a big challenge to integrate those three highly biased datasets for improving spelling correction.

4.2 Settings

Four baselines are adopted to make a comparison with the proposed multi-task learning method, including the limited-memory BFGS batch training method [30] with *single-task setting* (LBFGS-Single), the limited-memory BFGS batch training method with *merged setting* (LBFGS-Merge), the stochastic gradient descent method with single-task setting (SGD-Single), and the SGD method with merged setting (SGD-Merge). For the *single-task setting*, it uses only the task’s data to train the re-ranker of the speller (i.e., no data from other spelling correction tasks). For the *merged setting*,

it merges all of the training datasets of different tasks to train a unified re-ranker for the speller.

For the multi-task learning method, its hyper-parameters of similarity kernels are tuned in preliminary experiments, and we find using $d = 1$ worked well. In practice, we break the symmetric setting of C (in similarity kernels) and set different values of C for different tasks. We set $\beta = 0.99$ for the proposed MTL method, following the prior work on single-task learning [20]. The proposed method and the four baseline methods use exactly the same features, which is presented before. For the LBFGS method, we use the OWLQN software. The hyper-parameters of LBFGS were left unchanged from the default settings of the OWLQN software.

4.3 Results

We use the expected F1 score as the evaluation metric. This is the same evaluation metrics as the Speller Challenge 2011². Following the previous work of [28], top 40 corrections are used as the default setting. Using this setting, we can compare our results with the results of [28].

The performance comparisons from 1-pass to 5-pass settings are highlighted in Figure 2. In this figure, the *MTL* represents MTL-Poly; *Single* represents the *SGD-Single* baseline. Similarly, *Merge* represents *SGD-Merge*. We can

²<http://web-ngram.research.microsoft.com/spellerchallenge/Rules.aspx>

Table 2: Results (expected F1-score, standard deviation, and training time) of MTL and baselines in 5 passes. *MTL-Poly* and *MTL-Cor* represent the MTL method with polynomial kernel and correlation kernel, respectively. The standard deviations have been derived over 5 runs for random permutations of the training data order.

Method (#passes)	JDB F1 (%)	MSN F1 (%)	TREC F1 (%)	Overall F1 (%)	Overall Time (sec)
LBFGS-Single (1) (batch)	19.3	15.8	10.8	15.3 (± 0.0)	6.0
LBFGS-Single (5) (batch)	93.6	91.1	77.7	87.4 (± 0.0)	36.4
LBFGS-Merge (1) (batch)	19.5	15.7	10.7	15.3 (± 0.0)	5.5
LBFGS-Merge (5) (batch)	93.5	90.9	96.4	93.6 (± 0.0)	35.5
SGD-Single (1) (online)	92.5	88.9	94.7	92.1 (± 0.4)	22.3
SGD-Single (5) (online)	93.5	89.8	94.7	92.7 (± 0.5)	111.3
SGD-Merge (1) (online)	92.6	87.6	93.4	91.2 (± 1.1)	24.2
SGD-Merge (5) (online)	92.9	87.6	93.3	91.3 (± 1.6)	116.4
MTL-Poly (1) (new)	94.6	92.3	96.9	94.6 (± 0.1)	23.1
MTL-Poly (5) (new)	95.1	92.6	97.2	95.0 (± 0.1)	226.7
MTL-Cor (1) (new)	94.7	92.1	96.9	94.6 (± 0.1)	22.1
MTL-Cor (5) (new)	95.2	92.6	97.3	95.0 (± 0.1)	226.7

Table 3: Results of different methods on their convergence.

Method (#passes)	JDB F1 (%)	MSN F1 (%)	TREC F1 (%)	Overall F1 (%)	Overall Time (sec)
LBFGS-Single (200) (batch)	95.0	92.3	97.0	94.8	1479.1
LBFGS-Merge (200) (batch)	95.0	91.9	97.2	94.7	1360.9
SGD-Single (60) (online)	95.1	92.2	96.9	94.7	1370.9
SGD-Merge (60) (online)	95.0	91.8	97.2	94.7	1393.4
MTL-Poly (30) (new)	95.3	92.5	97.3	95.1	1463.3
MTL-Cor (30) (new)	95.3	92.5	97.3	95.1	1468.1

see that the proposed method MTL achieves much better F-score than all of the baseline methods on the 1-pass setting. In a similar way, MTL outperforms all of the baselines on the 2-pass, 3-pass, 4-pass, and 5-pass settings, and most of the differences are statistically significant.

In general, we find SGD-Merge does not have considerable advantage over the SGD-Single method. This indicates that it is difficult to combine such highly biased datasets for spelling correction, and a simple merge of such biased datasets does not give considerable improvement.

The experimental results in 5 passes are summarized in Table 2 with more details, including the results of batch training, the training time of different methods, and the comparisons between different similarity kernels. In addition, the results of different methods on their convergence state are shown in Table 3.

As we can see, the batch training (LBFGS) has weak performance in 5 passes, compared with SGD and MTL online training methods. The two similarity kernels (polynomial one and correlation one) have very similar performance on the MTL setting. Comparing the 5-pass results of MTL with the baseline results on their convergence, we find the MTL (with 5 passes) outperforms the SGD with 60 passes and LBFGS with 200 passes. In addition, the MTL with 5 passes is quite close to its empirical optimum on convergence. Compared with no matter the online or batch baselines, the MTL method is about one-order magnitude faster in terms of the empirical convergence speed.

In addition, we summarize all F-score curves by changing the number of passes, so that we can check the training process. The curves are shown in Figure 3. We also show the F-score curves based on training time. As we can see, the

MTL method can achieve better performance than baselines, with fewer training passes and fewer training time.

5. RELATED WORK

5.1 Spelling Correction

Spelling correction for regular written text is a long standing research topic. Previous researches can be roughly grouped into two categories: correcting non-word errors and real-word errors. In non-word error spelling correction, any word that is not found in a pre-compiled lexicon is considered to be misspelled. Then, a list of lexical words that are similar to the misspelled word are proposed as candidate spelling corrections.

Most traditional systems use a manually tuned similarity function (e.g., edit distance function) to rank the candidates, as reviewed by [11, 22]. During the last two decades, statistical error models learned on training data (i.e., query-correction pairs) have become increasingly popular, and have proven more effective [21, 37].

Real-word spelling correction is also referred to as context sensitive spelling correction (CSSC). Real-word spelling correction tries to detect incorrect usages of a valid word based on its context. A common strategy in CSSC is as follows. First, a pre-defined confusion set is used to generate candidate corrections, then a scoring model, such as a tri-gram language model or naive Bayes classifier, is used to rank the candidates according to their context (e.g., [19, 29]).

When designed to handle regular written text, both CSSC and non-word error speller systems rely on a pre-defined vocabulary (i.e., either a lexicon or a confusion set). However, in query spelling correction, it is impossible to compile such a vocabulary, and the boundary between the non-word

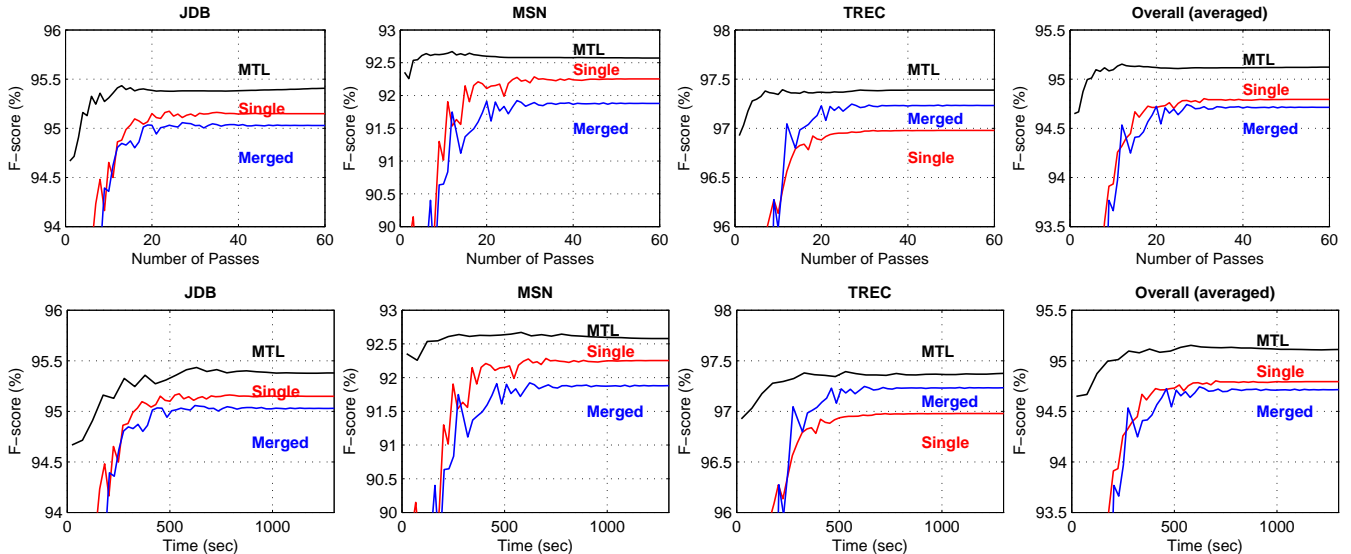


Figure 3: F-score curves of different methods. Upper panels are based on training passes. Bottom panels are based on training time (seconds). *Single* and *Merge* are based on SGD training.

and real-word errors is quite vague. Therefore, recent research on query spelling correction has focused on exploiting noisy Web data and query logs to infer knowledge about misspellings and word usage in search queries.

Cucerzan and Brill [12] discuss in detail the challenges of query spelling correction, and suggest the use of query logs. Ahmad and Kondrak [3] propose a method of estimating an error model from query logs using the EM algorithm. Li et al. [25] extend the error model by capturing word-level similarities learned from query logs. Chen *et al.* [10] suggest using web search results to improve spelling correction. Whitelaw *et al.* [39] present a query speller system in which both the error model and the language model are trained using Web data. Research in this direction also includes more recent work on utilizing large web corpora and query logs [10, 17, 28], employing large-scale n-gram models [17, 28], training phrase-based error model from clickthrough data [34, 18], and so on. Other related work includes [38, 15, 32, 36].

5.2 Multi-task Learning

Multi-task learning has been the focus of much interest in machine learning societies over the last decade. Traditional multi-task learning methods include: sharing hidden nodes in neural networks [8]; feature augmentation among interactive tasks [13]; producing a common prior in hierarchical Bayesian models [42, 43]; sharing parameters or common structures on the learning or predictor space [24, 4]; multi-task feature selection [41]; and matrix regularization based methods [5, 40], among others.

Recent development of multi-task learning is online multi-task learning, started from [14]. [14] assumes the tasks are related by a global loss function and the goal is to reduce the overall loss via online algorithm. With a similar but somewhat different motivation, [1] and [2] studied alternate formulations of online multi-task learning under traditional expert advice models. This is a formulation to exploit low dimensional common representations [16, 31]. Online multi-

task learning is also considered via reducing mistake bounds [9], and via perceptron-based online multi-task learning [33].

One of our target is adaptive online multi-task learning. Our adaptive (online) multi-task learning methods not only learn model weights, but also learn task relationships simultaneously from data [35]. More importantly, the proposed method can effectively estimate 2nd-order information, so that we can achieve very fast convergence of the multi-task learning. Finally, our proposal is a general framework which allows *non-structured* and *structured* classification.

6. CONCLUSIONS AND FUTURE WORK

In this paper, we proposed an adaptive (and online) multi-task learning method to integrate highly biased datasets for query spelling correction. We performed experiments on three different query spelling correction datasets, including the well-known TREC benchmark dataset. Experimental results demonstrated that the proposed method considerably outperformed the existing baseline systems in terms of accuracy.

Importantly, the proposed method was about one-order magnitude faster than baseline systems in terms of training speed. In contrast to the baseline methods which require more than 60 passes in training, the proposed method can approach very close to the empirical optimum in five passes. Since query spelling correction can have very large-scale training data in industrial applications (e.g., query spelling logs which contain more than millions of entries, see [34]), the proposed method’s ability to approach empirical optimum in 1-pass or a few passes will be of critical importance in such large-scale applications.

The proposed multi-task learning method is a general technique, and it can be easily applied to other classification tasks. As future work, we plan to apply this method to other large-scale web search and data mining tasks. Also, we will explore the possibility of integrating online learning with modern hashing algorithms [26, 27] to solve extremely large-scale problems.

7. ACKNOWLEDGMENTS

This work is partially supported by NSF (DMS-0808864, SES-1131848), ONR (YIP-N000140910911), and DARPA (FA-8650-11-1-7149). Xu Sun was a Postdoctoral Associate supported by ONR and NSF. Anshumali Shrivastava is a Ph.D. student supported by ONR and NSF. The authors thank Yanen Li for helpful discussions.

8. REFERENCES

- [1] ABERNETHY, J., BARTLETT, P., AND RAKHLIN, A. Multitask learning with expert advice. In *COLT'07* (2007), vol. 4539 of *Lecture Notes in Computer Science*, Springer, pp. 484–498.
- [2] AGARWAL, A., RAKHLIN, A., AND BARTLETT, P. Matrix regularization techniques for online multitask learning. Tech. Rep. UCB/EECS-2008-138, University of California, Berkeley, Oct 2008.
- [3] AHMAD, F., AND KONDRAK, G. Learning a spelling error model from search query logs. In *HLT-EMNLP'05* (October 2005), pp. 955–962.
- [4] ANDO, R. K., AND ZHANG, T. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research* 6 (2005), 1817–1853.
- [5] ARGYRIOU, A., MICCHELLI, C. A., PONTIL, M., AND YING, Y. A spectral regularization framework for multi-task structure learning. In *Proceedings of NIPS'07* (2007), MIT Press.
- [6] BENVENISTE, A., METIVIER, M., AND PRIOURET, P. Adaptive algorithms and stochastic approximations. *Berlin: Springer* (1990).
- [7] BOTTOU, L. Online algorithms and stochastic approximations. *Online Learning and Neural Networks*. Saad, David. Cambridge University Press (1998).
- [8] CARUANA, R. Multitask learning. *Machine Learning* 28, 1 (1997), 41–75.
- [9] CAVALLANTI, G., CESA-BIANCHI, N., AND GENTILE, C. Linear algorithms for online multitask classification. In *COLT'08* (2008), Omnipress, pp. 251–262.
- [10] CHEN, Q., LI, M., AND ZHOU, M. Improving query spelling correction using web search results. In *EMNLP-CoNLL'07* (Prague, Czech Republic, June 2007), pp. 181–189.
- [11] CHURCH, K. W., AND GALE, W. A. Probability scoring for spelling correction. *Statistics and Computing volume 1* (1991), 93–103.
- [12] CUCERZAN, S., AND BRILL, E. Spelling correction as an iterative process that exploits the collective knowledge of web users. In *EMNLP'04* (Barcelona, Spain, July 2004), pp. 293–300.
- [13] DAUMÉ III, H. Frustratingly easy domain adaptation. In *ACL'07* (Prague, Czech Republic, June 2007), pp. 256–263.
- [14] DEKEL, O., LONG, P. M., AND SINGER, Y. Online multitask learning. In *COLT'06* (2006), vol. 4005 of *Lecture Notes in Computer Science*, Springer, pp. 453–467.
- [15] DUAN, H., AND HSU, B.-J. P. Online spelling correction for query completion. In *WWW'11* (2011), ACM, pp. 117–126.
- [16] EVGENIOU, T., MICCHELLI, C. A., AND PONTIL, M. Learning multiple tasks with kernel methods. *Journal of Machine Learning Research* 6 (2005), 615–637.
- [17] GANJISAFFAR, Y., ZILIO, A., JAVANMARDI, S., CETINDIL, I., SIKKA, M., KATUMALLA, S. P., KHATIB-ASTANEH, N., LI, C., AND LOPES, C. qSpell: Spelling correction of web search queries using ranking models and iterative correction. In *Spelling Alteration for Web Search Workshop* (July 2011).
- [18] GAO, J., LI, X., MICOL, D., QUIRK, C., AND SUN, X. A large scale ranker-based system for search query spelling correction. In *COLING'10* (2010), pp. 358–366.
- [19] GOLDING, A. R., AND ROTH, D. Applying winnow to context-sensitive spelling correction. In *ICML'96* (1996), Morgan Kaufmann, pp. 182–190.
- [20] HSU, C.-N., HUANG, H.-S., CHANG, Y.-M., AND LEE, Y.-J. Periodic step-size adaptation in second-order gradient descent for single-pass on-line structured learning. *Machine Learning* 77, 2-3 (2009), 195–224.
- [21] KERNIGHAN, M. D., CHURCH, K. W., AND GALE, W. A. A spelling correction program based on a noisy channel model. In *COLING'90* (1990), pp. 205–210.
- [22] KUKICH, K. Techniques for automatically correcting words in text. *ACM Comput. Surv.* 24, 4 (1992), 377–439.
- [23] LAFFERTY, J., MCCALLUM, A., AND PEREIRA, F. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML'01* (2001), pp. 282–289.
- [24] LAWRENCE, N. D., AND PLATT, J. C. Learning to learn with the informative vector machine. In *ICML'04* (2004), vol. 69, ACM.
- [25] LI, M., ZHU, M., ZHANG, Y., AND ZHOU, M. Exploring distributional similarity based models for query spelling correction. In *COLING-ACL'06* (Sydney, Australia, July 2006), pp. 1025–1032.
- [26] LI, P. AND KONIG, A.C. Theory and Applications of b-bit Minwise Hashing. *Communications of the ACM* 54 (2011), 101–109.
- [27] LI, P., SHRIVASTAVA, A., MOORE, J., AND KONIG, A.C. Hashing Algorithms for Large-Scale Learning. In *NIPS'11* (Granada, Spain, December 2011).
- [28] LI, Y., DUAN, H., AND ZHAI, C. Cloudspeller: Spelling correction for search queries by using a unified hidden markov model with web-scale resources. In *Spelling Alteration for Web Search Workshop* (July 2011).
- [29] MANGU, L., AND BRILL, E. Automatic rule acquisition for spelling correction. In *ICML'97* (Nashville, TN, 1997), Morgan Kaufmann, pp. 187–194.
- [30] NOCEDAL, J., AND WRIGHT, S. J. Numerical optimization. *Springer* (1999).
- [31] RAI, P., AND III, H. D. Infinite predictor subspace models for multitask learning. *Journal of Machine Learning Research - Proceedings Track 9* (2010), 613–620.

- [32] REYNAERT, M. Character confusion versus focus word-based correction of spelling and ocr variants in corpora. *IJDAR* 14, 2 (2011), 173–187.
- [33] SAHA, A., RAI, P., DAUMÉ III, H., AND VENKATASUBRAMANIAN, S. Online learning of multiple tasks and their relationships. In *AISTATS'10* (Ft. Lauderdale, Florida, 2011).
- [34] SUN, X., GAO, J., MICOL, D., AND QUIRK, C. Learning phrase-based spelling error models from clickthrough data. In *ACL'10* (Uppsala, Sweden, July 2010), pp. 266–274.
- [35] SUN, X., KASHIMA, H., TOMIOKA, R., UEDA, N., AND LI, P. A new multi-task learning method for personalized activity recognition. In *ICDM'11* (2011), IEEE.
- [36] SUN, X., SHRIVASTAVA, A., AND LI, P. Query spelling correction using multi-task learning. In *WWW (Companion Volume)* (2012), A. Mille, F. L. Gandon, J. Misselis, M. Rabinovich, and S. Staab, Eds., ACM, pp. 613–614.
- [37] TOUTANOVA, K., AND MOORE, R. Pronunciation modeling for improved spelling correction. In *ACL'02* (Philadelphia, USA, 2002), pp. 144–151.
- [38] WANG, Z., XU, G., LI, H., AND ZHANG, M. A fast and accurate method for approximate string search. In *ACL'11* (2011), pp. 52–61.
- [39] WHITELAW, C., HUTCHINSON, B., CHUNG, G., AND ELLIS, G. Using the web for language independent spellchecking and autocorrection. In *EMNLP'09* (2009), pp. 890–899.
- [40] XUE, Y., DUNSON, D., AND CARIN, L. The matrix stick-breaking process for flexible multi-task learning. In *ICML'07* (Corvalis, Oregon, 2007), ACM, pp. 1063–1070.
- [41] YANG, H., KING, I., AND LYU, M. R. Online learning for multi-task feature selection. In *CIKM'10* (2010), ACM, pp. 1693–1696.
- [42] YU, K., TRESP, V., AND SCHWAIGHOFER, A. Learning gaussian processes from multiple tasks. In *ICML'05* (2005), vol. 119, ACM, pp. 1012–1019.
- [43] ZHANG, J., GHAHRAMANI, Z., AND YANG, Y. Learning multiple related tasks using latent independent component analysis. In *NIPS'05* (2005).