

Automated Retraining Methods for Document Classification and Their Parameter Tuning

Stefan Siersdorfer and Gerhard Weikum

Max-Planck-Institute for Computer Science, Germany
{stesi, weikum}@mpi-sb.mpg.de

Abstract. This paper addresses the problem of semi-supervised classification on document collections using retraining (also called self-training). A possible application is focused Web crawling which may start with very few, manually selected, training documents but can be enhanced by automatically adding initially unlabeled, positively classified Web pages for retraining. Such an approach is by itself not robust and faces tuning problems regarding parameters like the number of selected documents, the number of retraining iterations, and the ratio of positive and negative classified samples used for retraining. The paper develops methods for automatically tuning these parameters, based on predicting the leave-one-out error for a re-trained classifier and avoiding that the classifier is diluted by selecting too many or weak documents for retraining. Our experiments with three different datasets confirm the practical viability of the approach.

1 Introduction

Automatic document classification is useful for a wide range of applications such as organizing web, intranet, or portal pages into topic directories, filtering news feeds or mail, focused crawling on the web or in intranets, and many more. In some applications, the availability of good training data for the classifier is the key bottleneck. As an example, consider a personalized or community information tool that uses thematically focused crawling [10, 29] to build and maintain a directory or index for browsing, search, or recommendations.

To overcome the training bottleneck, semi-supervised learning techniques could be applied. In our given setting, the classifier could be bootstrapped by training it with whatever explicitly class-labeled training data are available and used for making decisions about the classes of previously unseen, unlabeled test documents retrieved by the crawler. These decisions would have a certain degree of uncertainty, depending on the classifier's statistical learning model. However, some of the test documents are usually accepted for their corresponding classes with high statistical confidence, and these could then be selected for retraining the classifier, now with considerably more training documents. Obviously, this simple idea does not provide a robust solution, for the automatically selected, additional training data may also increase the classifier's uncertainty and may eventually lead to an unintended topic drift. In this paper we address the issue of how to make such a semi-supervised classifier robust and practically viable.

There are various approaches, like Transductive SVM, EM-iterated Bayesian Classifiers, or Co-Training [30, 23, 7], that successfully use information from ini-

tially *unlabeled* documents to improve classification results. However these methods come with parameters, which have a crucial influence on the quality of the classification results and need to be *tuned manually* on a per application basis.

To this end we propose a retraining algorithm that performs *automatic parameter tuning*. When our method considers adding a batch of initially unlabeled documents to the training set, it predicts the resulting improvement (or degeneration) of the classifier’s accuracy by performing a leave-one-out validation. The training set is extended, by selecting the unlabeled documents with highest classification confidence and then retraining the classifier, only as long as the predictor expects an improvement. To avoid extensive leave-one-out validations, which are resource-intensive, the predictor invokes the validation merely after a certain number of iterations and rather uses a spline interpolation technique for less expensive estimation.

A particularly subtle but important point in this procedure is that one should often select different numbers of positive and negative samples for retraining, depending on the ratio in the underlying corpus. In the case of a focused Web crawler, usually a much larger fraction of negative (i.e., thematically uninteresting) documents is seen and may lead to a wrong classifier bias unless such corrective steps are taken.

The novel contributions of this paper are the following:

1. We develop a robust, practically viable procedure for automated retraining of classifiers with careful selection of initially unlabeled documents.
2. We perform comprehensive experiments that evaluate our retraining procedure against state-of-the-art semi-supervised classification methods like EM-iterated Bayesian classifiers, Transductive SVMs, and Spectral Graph Transduction.

Related Work. There is a considerable prior of work on classification using unlabeled data (also called semi-supervised learning), see [26] for an overview. Naive Retraining where new documents with highest classification confidence are iteratively added to the training set, is, e.g., described in [10]; but these methods perform often worse than the underlying base learning method. A more enhanced EM (Expectation Maximization)-based variant for Bayesian Classifiers is proposed in [23] and applied to text classification. For Transductive SVM [15, 30] and Semi-Supervised SVM [5] unlabeled samples are taken into account (opposite to standard SVM) in a modified optimization problem (standard SVM cannot use unlabeled samples at all). Co-training [7] splits the feature space into conditionally independent dimensions and performs retraining on the corresponding classifiers. Recent graph-based semi-supervised learning algorithms work by formulating the assumption that nearby points, and points in the same structure should have similar labels [18, 31, 16]. In [6] semi-supervised learning is combined with ensemble classification methods. An approach for the case that only positive (and no negative) training data plus unlabeled data are available is described in [20]. In [3] semi-supervised learning is used for text summarization; in [32] a retraining method with user feedback as a stopping criterion is used for

image retrieval. However, to our knowledge, none of these methods deals with the problem of automatically tuning their parameters.

The issue of asymmetric distribution of documents among different classes is addressed, e.g., in [19, 8, 13], and the problem of automated parameter tuning has been considered in the field of machine learning, e.g., in [17], but, to our knowledge, not in the context of retraining.

2 Technical Basics

Classifying text documents into thematic categories usually follows a supervised learning paradigm and is based on training documents that need to be provided for each topic. Both training documents and test documents, which are later given to the classifier, are represented as multidimensional feature vectors. In the prevalent bag-of-words model the features are derived from word occurrence frequencies, e.g. based on *tf*idf* feature weights [4, 22]. Often feature selection algorithms are applied to reduce the dimensionality of the feature space and eliminate “noisy”, non-characteristic features, based on information-theoretic measures for feature ordering (e.g., relative entropy or information gain).

Feature vectors of topic labeled text documents (e.g., capturing *tf*idf* weights of terms) are used to train a classification model for each topic, using probabilistic (e.g., Naive Bayes) or discriminative models (e.g., SVM). Linear support vector machines (SVMs) construct a hyperplane $\mathbf{w} \cdot \mathbf{x} + b = 0$ that separates the set of positive training examples from a set of negative examples with maximum margin δ . This training requires solving a quadratic optimization problem whose empirical performance is somewhere between quadratic and cubic in the number of training documents [9]. For a new, previously unseen, (test) document \mathbf{d} the SVM merely needs to test whether the document lies on the “positive” side or the “negative” side of the separating hyperplane. The decision simply requires computing a scalar product of the vectors \mathbf{w} and \mathbf{d} . SVMs have been shown to perform very well for text classification (see, e.g., [12, 14]).

Unlike the inductive SVM setting, for Transductive SVM (TSVM) [15, 30], a hyperplane is computed that separates *both* training and (unlabeled) test data with maximum margin.

The most widely used technique for empirically estimating the classifier quality is *cross-validation* [22] on a set of independent data samples with known topic memberships (aka. class labels). The partitioning is systematically varied by dividing the overall data into k groups and investigating each of the k choices for using one group as test data and the other $k - 1$ groups for training; the empirical results are finally averaged over all choices. An important variation is *leave-one-out validation* [22]. Here the n documents of a data collection are divided by the ratio $(n - 1) : 1$. Both methods are also popular for predicting a classifier’s quality. Leave-one-out prediction is more accurate than prediction based on cross-validation but requires training the classifier n times, unless special properties of the classifier’s underlying model could be exploited.

In this paper we consider only binary classifiers that make a decision for a single topic, based on positive and negative training examples.

3 Retraining and Its Parameter Tuning

3.1 A Simple Base Algorithm

Consider a training set T and a set of unlabeled data U . We can perform retraining by iteratively building a classifier C on T , classifying the documents in U and adding the documents with the highest classification confidence, p positively and n negatively classified documents in one iteration. Classification confidence could be estimated, e.g., by the distance from the separating hyperplane in the SVM case or by the probability of accepting a document for a class.

This algorithm provides us with a tradeoff. On one hand, a higher number of training examples could potentially improve the classification accuracy; on the other hand, there are potentially incorrectly labeled documents among the automatically labeled training docs U_{pos} and U_{neg} , which can dilute the training set. The algorithm thus has two important tuning parameters:

1. the number m of iterations
2. the ratio p/n between new positively classified and negatively classified docs used for retraining

In the following we show how we can automatically tune these parameters. Note that the total number of selected documents for each retraining step, $r := p + n$ could be considered as an additional tuning parameter. However, we can simply choose it sufficiently small to be on the conservative side.

3.2 Tuning the Number of Iterations

Because of the tradeoffs mentioned above, a higher number of iterations do not necessarily imply a lower error. Our idea now is to approximate this error curve on the test set U by an estimated error curve.

For a retraining step we can build an error estimator by performing leave-one-out validation of the current classifier C on the original training set T_0 , i.e., the part of the training set that consists of the manually labeled documents (which are correct with perfect confidence).

For a set of sample estimates

$$\{(i_0, estError(i_0)), \dots, (i_l, estError(i_l))\}, \quad (1)$$

where the i_j values are the iteration numbers and $estError(i_j)$ is the estimated error, we can now approximate the overall error curve by fitting the sample estimates.

There are various approaches to this curve fitting. In our experiments we obtained good performance using cubic splines. Cubic splines are used in many areas, e.g., bio medicine, signal processing, and computer graphics [11, 24, 27]. In our experiments we also tested other approaches like linear splines, and error estimation by the less time consuming k-fold-cross-validation instead of leave-one-out.

Having approximated the error estimation curve $S(x)$, we choose the retraining classifier C in the iteration i with minimum $S(i)$. Choosing the number of

supporting points for the fitting is an efficiency issue. The more supporting points the better the approximation but the higher the overall cost for computing the estimator values.

The classifier can be optimized in the same way for other quality measures like the F-measure (the harmonic mean of precision and recall).

3.3 Tuning the Ratio of Positive and Negative Samples

For an effective classification the training set should be an appropriate representation of the test set. For binary classification, it is especially helpful if the ratio between positive and negative documents is approximately the same for the test and the training set. For example, Bayesian classifiers take the prior class probabilities explicitly into account. For SVM a badly proportionalized training set can also lead to a disadvantageous bias [8]. The assumption of having a training set with the same ratio of positive and negative documents as a test set is not at all self-guaranteed or easy to satisfy in practice. Typically a human, collecting training documents, would rather choose roughly the same number of documents for each class, even if there are significant (but a priori unknown) differences in the real world.

The idea is to overcome this problem by adjusting the training set such that it better represents the test set. To do so, in each iteration of our retraining algorithm we approximate the ratio between positive and negative documents by applying the current classifier to the set of initially unlabeled data U_0 (test data). Among a small number r of new retraining documents we choose the number of positive and negative documents, n and p , such that the difference between the overall ratio of positive and negative training docs and the estimated ratio on the unlabeled data is minimized.

More formally let t_{pos} be the number of positive, t_{neg} be the number of negative training documents in the current iteration, v_{pos} be the number of unlabeled documents classified as positive by the current classifier C , and v_{neg} be the number of documents classified as negative. Then we choose the number of *newly* added positive and negative documents for retraining, p and n , such that the ratio $(t_{pos} + p) : (t_{neg} + n)$ between the overall number of positive and negative training documents provides the best approximation for the ratio $v_{pos} : v_{neg}$ of positive and negative test documents estimated by the current classifier:

$$p = \arg \min_{x \in \{0, \dots, r\}} \left| \frac{t_{pos} + x}{t_{neg} + r - x} - \frac{v_{pos}}{v_{neg}} \right| \quad (2)$$

and

$$n = r - p \quad (3)$$

3.4 The Enhanced Retraining Algorithm

With the parameter tuning methods described above, our retraining algorithm now works as follows: We retrain as long as documents for retraining are available. In each retraining iteration we add a small number r of documents to the

```

Input: training set  $T = T_0$ ; set of unlabeled Data  $U = U_0$ ; stepsize

set of classifiers C-Set = empty
set of supporting points Support-Set = empty
iteration number  $i = 0$ ;
while (U is not empty) do
  build classifier C on T; add (i,C) to C-Set
  estimate  $p$  and  $n$  \\\as described above; classify U
   $U_{\text{pos}} :=$  top- $p$  positively classified docs
   $U_{\text{neg}} :=$  top- $n$  negatively classified docs
   $T = T + U_{\text{pos}} + U_{\text{neg}}$  ;  $U = U - U_{\text{pos}} - U_{\text{neg}}$ 
  if ( $i \bmod \textit{stepsize} = 0$ )
    estimate error  $\text{estError}$  of C by leave-one-out on  $T_0$ 
    add (i, $\text{estError}$ ) to Support-Set
   $i++$ 

compute interpolating curve S on Support-Set \\\as described above
choose  $j$  which minimizes  $S(i)$ 
return Classifier  $c$  from C-Set with iteration number =  $j$ 

```

Fig. 1. Enhanced Retraining Algorithm

training set, determining the ratio between new positive and negative training documents as described in Section 3.3. Every *stepsize* iterations we compute and save an error estimator. We apply curve fitting to the estimated error, and choose the classifier corresponding to the minimum estimated error (see Section 3.2.).

The pseudo code in Figure 1 summarizes our modified retraining algorithm.

4 Experiments

Setup. We performed a series of experiments with real-life data from the following sources: 1) The Newsgroups collection at [1] with 17,847 postings collected from 20 Usenet newsgroups such as 'rec.autos', 'sci.space', etc. 2) The Reuters collection [21] with 21,578 newswire articles; 12,904 of them are subdivided into categories ('earn', 'grain', 'trade', etc.). 3) The Internet Movie Database (IMDB) at [2] with short movie descriptions from 20 topics according to particular movie genres ('drama', 'horror' etc.). Only 34,681 movies were considered that have a unique genre.

For every data collection we considered each class with at least 300 documents. We obtained 20 classes for Newsgroups, 8 for Reuters and 9 for IMDB. For each class we randomly chose 100 documents as positive training examples and 100 negative examples from all other classes. For testing we considered two cases: 1) the symmetric case: we chose equal numbers of positive and negative test documents for each class (200 per class), and 2) the asymmetric case: we

chose the number of positive and negative test documents in a ratio of 1 : 6 (i.e., 200:1200).

In all experiments, the standard bag-of-words model [4] (using term frequencies to build L1-normalized feature vectors, stemming with the algorithm of Porter [25], and deletion of stopwords) was used for document representation. We used binary classifiers so as to recognize documents from one specific topic against all other topics; this setup was repeated for every topic.

For each data collection we computed the macro-averaged error (i.e., the average ratio of incorrectly classified documents to the number of test documents) along with the 95 percent confidence interval and the macro-averaged F1 value (the harmonic mean of precision and recall).

Results. We compared the following classification methods:

1. Standard linear SVM (**SVM**)
2. Standard linear TSVM. Here the fraction f of unlabeled examples to be classified into the positive class is a selectable parameter. As default setting we used the ratio between the positive and the negative examples in the training data. (**TSVM**)
3. Linear TSVM where the ratio f between positive and negative test documents was set according to the SVM classification (Method 1) on the test documents. (**TSVM+est**)
4. The augmented EM-iterated Bayesian classifier with weighting of the unlabeled data as described in [23]. Here we determined the weighting parameter λ by leave-one-out validation (considering the values between 0 and 1 with a step width of 0.2), choosing the λ with the lowest estimated error. (**EM-Bayes**)
5. Spectral Graph Transduction as described in [16] (**SGT**)
6. Our retraining approach with linear SVM (Method 1) as the underlying base classifier and 10 new retraining documents per iteration and
 - (a) error/F1 prediction by leave-one-out estimation invoked after every 10 iterations and cubic spline interpolation (**RetCspL1o**)
 - (b) error/F1 prediction by leave-one-out estimation invoked after every 10 iterations, linear spline interpolation (**RetLspL1o**)

Method	Newsg. avg(error)	IMDB avg(error)	Reuters avg(error)	Newsg. avg(F1)	IMDB avg(F1)	Reuters avg(F1)
SVM	0.097 ± 0.0035	0.246 ± 0.0075	0.075 ± 0.0049	0.726	0.481	0.783
TSVM	0.364 ± 0.0056	0.401 ± 0.0086	0.362 ± 0.0089	0.434	0.376	0.437
TSVM+est	0.096 ± 0.0035	0.249 ± 0.0075	0.076 ± 0.0049	0.728	0.475	0.78
EM-Bayes	0.202 ± 0.0047	0.267 ± 0.0077	0.093 ± 0.0054	0.596	0.498	0.75
SGT	0.216 ± 0.0048	0.329 ± 0.0082	0.167 ± 0.0069	0.543	0.402	0.606
RetCspL1o	0.077 ± 0.0031	0.207 ± 0.0071	0.058 ± 0.0043	0.749	0.497	0.818
RetCspL1o	0.08 ± 0.0032	0.211 ± 0.0071	0.059 ± 0.0044	0.749	0.496	0.817
RetLspL1o	0.081 ± 0.0032	0.212 ± 0.0071	0.058 ± 0.0043	0.744	0.49	0.813
RetLspL1o	0.083 ± 0.0032	0.209 ± 0.0071	0.06 ± 0.0044	0.744	0.491	0.812
RetCv	0.084 ± 0.0032	0.204 ± 0.007	0.059 ± 0.0044	0.745	0.499	0.816

Fig. 2. Macro-averaged Results for **Asymmetric** Test Set: Baseline and Retraining Methods

- (c) error/F1 prediction by 5-fold cross-validation invoked after every 10 iterations and cubic spline interpolation (**RetCsplCv**)
- (d) error/F1 prediction by 5-fold cross-validation invoked after every 10 iterations and linear spline interpolation (**RetLsplCv**)
- (e) error/F1 prediction by 5-fold cross-validation invoked after every iteration - and no interpolation (**RetCv**)

For SVM and TSVM we used the popular *SVMlight* implementation [14] with parameter $C = 1000$ (tradeoff between training error and margin). For the Spectral Graph Transductor we used the *SGTlight* implementation with parameterization as described in [16].

The average results for asymmetric test sets are shown in Figure 2 (best values in boldface). For lack of space, results for the symmetric case are omitted here; they can be found in [28]. The main observations are: In the asymmetric test case, our retraining algorithm clearly provides the best performance on all three datasets. For example, on the IMDB data, which is the hardest test case in terms of the absolute accuracy that was achievable, we reduce the error from approximately 25-27 percent (for SVM and TSVM with estimator and for EM-iterated Bayes) to 20.7 percent, quite a significant gain. The very bad performance of standard TSVM can be explained by the big gap between the parameter f , estimated on the training set, and the real ratio between positive and negative documents in the asymmetric test set.

As we regard the asymmetric test case, significantly more unacceptable test documents than acceptable ones, as the far more realistic setting (e.g. in focused crawling, news filtering, etc.), we conclude that the newly proposed retraining method is the clear winner and outperforms the previously known state-of-the-art algorithms by a significant margin.

An extended version of this paper is available as a technical report [28].

References

1. The 20 newsgroups data set. <http://www.ai.mit.edu/~jrennie/20Newsgroups/>.
2. Internet movie database. <http://www.imdb.com>.
3. M.-R. Amini and P. Gallinari. The use of unlabeled data to improve supervised learning for text summarization. In *SIGIR '02*, pages 105–112. ACM Press, 2002.
4. R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*. Addison Wesley, 1999.
5. K. P. Bennett and A. Demiriz. Semi-supervised support vector machines. In *NIPS 1999*, pages 368–374. MIT Press, 1999.
6. K. P. Bennett, A. Demiriz, and R. Maclin. Exploiting unlabeled data in ensemble methods. In *SIGKDD*, pages 289–296. ACM Press, 2002.
7. A. Blum and T. Mitchell. Combining labeled and unlabeled data with co-training. *Workshop on Computational Learning Theory*, 1998.
8. J. Brank, M. Grobelnik, N. Milic-Frayling, and D. Mladenic. Training text classifiers with SVM on very few positive examples. *Technical Report MSR-TR-2003-34*, Microsoft Corp., 2003.

9. C. Burges. A tutorial on Support Vector Machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2(2), 1998.
10. S. Chakrabarti. *Mining the Web: Discovering Knowledge from Hypertext Data*. Morgan-Kaufman, 2002.
11. E. Chen and C. Lam. Predictor-corrector with cubic spline method for spectrum estimation in compton scatter correction of spect. *Computers in biology and medicine*, 1994, vol. 24, no. 3, pp. 229, Ingenta.
12. S. Dumais and H. Chen. Hierarchical classification of Web content. *SIGIR*, 2000.
13. H. Guo and H. L. Viktor. Learning from imbalanced data sets with boosting and data generation: the databoost-im approach. *SIGKDD Explorations*, 30 - 39, 2004.
14. T. Joachims. Text categorization with Support Vector Machines: Learning with many relevant features. *ECML*, 1998.
15. T. Joachims. Transductive inference for text classification using support vector machines. In *ICML'99*, 200 - 209, 1999.
16. T. Joachims. Transductive learning via spectral graph partitioning. In *ICML*, pages 290–297, 2003.
17. R. Kohavi and G. John. Automatic parameter selection by minimizing estimated error. In *Machine Learning*, 1995.
18. B. Krishnapuram, D. Williams, Y. Xue, A. Hartemink, L. Carin, and M. Figueiredo. On semi-supervised classification. In *NIPS*. MIT Press, 2005.
19. M. Kubat and S. Matwin. Addressing the curse of imbalanced training sets: One-sided selection. In *ICML'97, Nashville, TN, U.S.A.*, 179-186, 1997.
20. W. S. Lee and B. Liu. Learning with positive and unlabeled examples using weighted logistic regression. In *ICML'03, Washington USA*, 2003.
21. D. D. Lewis. Evaluating text categorization. In *Proceedings of Speech and Natural Language Workshop*, pages 312–318. Defense Advanced Research Projects Agency, Morgan Kaufmann, Feb. 1991.
22. C. Manning and H. Schuetze. *Foundations of Statistical Natural Language Processing*. MIT Press, 1999.
23. K. Nigam, A. McCallum, S. Thrun, and T. Mitchell. Text classification from labeled and unlabeled documents using em. *Machine Intelligence*, 39(2/3), 2000.
24. E. Okanla and P. Gaydecki. A real-time audio frequency cubic spline interpolator. *Signal processing*, 1996, vol. 49, no. 1, pp. 45, Ingenta.
25. M. Porter. An algorithm for suffix stripping. *Automated Library and Information Systems*, 14(3).
26. M. Seeger. Learning with labeled and unlabeled data. *Tech. Rep., Institute for Adaptive and Neural Computation, University of Edinburgh, UK*, 2001.
27. C. Seymour and K. Unsworth. Interactive shape preserving interpolation by curvature continuous rational cubic splines. *Appl. Math.* 102 (1999), no. 1, 87–117.
28. S. Siersdorfer and G. Weikum. Automated retraining methods for document classification and their parameter tuning. In *Technical Report MPI-I-2005-5-002, Max-Planck-Institute for Computer Science, Germany*, <http://www.mpi-sb.mpg.de/~stesi/sources/2005/report05retr.pdf>, 2005.
29. S. Sizov, M. Biwer, J. Graupmann, S. Siersdorfer, M. Theobald, G. Weikum, and P. Zimmer. The BINGO! system for information portal generation and expert Web search. *Conference on Innovative Systems Research (CIDR)*, 2003.
30. V. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
31. D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf. Learning with local and global consistency. In *NIPS*. MIT Press, 2004.
32. Z. Zhou, K. Chen, and Y. Jiang. Exploiting unlabeled data in content-based image retrieval. In *ECML'03, Pisa, Italy*, 2004.