

Multi-class Feature Selection with Support Vector Machines

Olivier Chapelle
Yahoo! Research
Santa Clara, CA
chap@yahoo-inc.com

Sathiya Keerthi
Yahoo! Research
Santa Clara, CA
selvarak@yahoo-inc.com

ABSTRACT

Feature selection is an important component of text categorization that has mostly been addressed by filter methods. In the context of SVM classifiers we propose in this paper a general framework where feature selection is embedded as a part of the learning algorithms. As special cases of this framework we derive two algorithms, one based on the L_1 regularization of the weight vector and the other being an extension of the Recursive Feature Elimination algorithm. We also derive a new method that performs even better than these two methods. Our framework and methods are developed in a multi-class setting where the goal is to find a small set of features for all the classes simultaneously. On some datasets the size of the feature set found by our new method is one order of magnitude smaller than the one found by a filter method based on information gain. Finally, we devise a computationally efficient optimization technique for this method.

Categories and Subject Descriptors

I.2.6 [Learning]: *Support Vector Machines*; I.5.2 [Design Methodology]: *Feature Selection*; H.3.3 [Information Search and Retrieval]: *Text Classification*

General Terms

Algorithms, Performance, Experimentation

Keywords

Support Vector Machines, Feature selection, Text categorization, Scaling factors

1. INTRODUCTION

Feature selection is an important component of text classification. It is used to help reduce the load on computational resources and, in cases where there are many noisy features, to help in lifting the performance by eliminating

such features. Several feature selection methods have been suggested in the literature, particularly with respect to binary classification.

Guyon and Elisseeff [11] classify feature selection methods into three types: *Filter*, *Wrapper* and *Embedded* methods. Filter methods select features as a pre-processing step, independently of the prediction method. Because text classification involves a large number of features and filter methods are computationally very efficient they have been popularly used in text classification. Yang and Pedersen [23] and Forman [7] compared a number of filter methods for text classification. These studies show information gain, Chi-squared and Bi-normal separation as the leading filter measures. Wrapper methods use the prediction method as a black box to score subsets of features. In text classification they have not been tried because of their expensive need to try out a very large number of subset selections. Embedded methods perform feature selection as part of the training process of the prediction method. Linear classifiers that use L_1 regularization on the weights [10, 18] fall in this class. Recursive Feature Elimination (RFE) [12], a backward elimination method that uses smallness of weights to decide feature removal, is another effective embedded method.

A large number of text classification problems occurring in practice involve many categories. They are either of a *multi-class* type (assigning exactly one class to each document) or of a *multi-labeled* type (assigning a variable number of classes to each document). In such multi-class and multi-labeled settings it is natural to look for a small common set of features that works well for all the classes. The need for selecting a small common set of features is also motivated by computational resource constraints in online and active learning settings. When the average number of features occurring in a document is large (this can happen, for example when the document representation uses rich data types such as multiple text fields and phrases and/or when the documents themselves are large) it is expensive to process, store and transport over a network, an individual set of different features for each class.

It is easy to extend filter methods for doing simultaneous feature subset selection [8, 23]. For example, for the multi-labeled case one can take the information gain values of the individual binary classifiers and combine them (say, via an averaging or max operation) to form a single measure, using which the various features can be ordered. But we are unaware of any embedded method that has been developed for the same purpose. The chief aim of this paper is to fill this gap. We consider the regularized linear classification setting

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.

SIGIR '07 Amsterdam

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

(with SVMs being the particular model taken for implementation) and develop new embedded methods that address the simultaneous feature subset selection problem. These methods include suitable adaptations of L_1 regularization and RFE to deal with simultaneous feature subset selection, as well as a new method that is even more effective than those adaptations. A nice feature is that all these methods are related to a common linear classifier model that employs scaling factors to address feature selection. We also develop scalable algorithms for the various methods. Comparative evaluation of the new embedded methods with respect to a baseline filter method that uses information gain shows the new methods to be very effective.

2. FRAMEWORK

In this section we present our main ideas for multi-class feature selection. Before doing this we first review Support Vector Machines and two related embedded methods for feature selection, viz. RFE and L_1 -SVMs.

2.1 Support Vector Machines

Let us first introduce some notations. In the following, the index i will always run over the documents $1 \dots n$, j over the features $1 \dots d$, and k over the classes $1 \dots c$. A training set $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{1 \leq i \leq n}$ is given, where $\mathbf{x}_i \equiv (x_{i1} \dots x_{id})^\top$ is the d dimensional vector representation of the i -th example and $\mathbf{y}_i \equiv (y_{i1} \dots y_{ic})^\top$ its label vector. $y_{ik} = 1$ if that example belongs to the k -th category and -1 otherwise. The linear classifier for the k -th class uses a d dimensional weight vector, \mathbf{w}_k . The j -th element of \mathbf{w}_k will be written as w_{jk} . We will use \mathbf{w}_k^2 to denote the square of the Euclidean norm of \mathbf{w}_k .

Consider the design of \mathbf{w}_k , the weight vector for class k . Support Vector Machines [2] minimize the following objective function¹

$$\frac{1}{2} \mathbf{w}_k^2 + \frac{C}{2} \sum_{i=1}^n \ell(y_{ik}(\mathbf{w}_k \cdot \mathbf{x}_i)). \quad (1)$$

The loss function ℓ is $\ell(t) = \max(0, 1-t)^p$. We take $p = 2$ in the rest of this paper. SVMs are often used to find non-linear decision boundaries through a *kernel* function, but in text classification this additional step is not usually necessary [13]. Note that in the linear case, fast methods exist to train SVMs [14]. For instance in that paper, a training time of 10 minutes has been reported for a training set of one million documents.

2.2 Recursive Feature Elimination

One of the standard embedded methods for doing feature selection with SVMs is *Recursive Feature Elimination* (RFE) [12]. It was originally suggested for binary classification. So to begin with let us describe RFE just for one binary classifier, say for the k -th class with weight vector \mathbf{w}_k . RFE is based on the idea that the importance of a feature j should be related to the magnitude of its weight $|w_{jk}|$. Combining this idea with a backward elimination gives the following algorithm:

1. Train SVM on the active features.
2. Remove the feature with the smallest $|w_{jk}|$.

¹A threshold is not included, but can easily be learned by adding a constant component 1 to the examples.

3. Go back to step 1 or stop if there are only m features left.

Of course removing one feature at a time in step 2 is time consuming and in practice it is common to remove as much as half of the active features in each iteration. This algorithm is very simple to implement. Also it can be made more efficient by using *seeding*, i.e. when retraining the SVM in step 1, use the weight vector from the previous step as a starting point.

One can think of extending RFE to the multi-class simultaneous feature subset selection problem in various ways. For step 2 of RFE, Chen et al [4] propose using the smallness of $\max_k |w_{jk}|$ as the criterion. One could instead use the smallness of $\sum_k w_{jk}^2$. As we will see in section 2.6 below, this measure has good theoretical support.

2.3 L_1 Support Vector Machines

Another way of obtaining a sparse solution (i.e. with a few non-zero weight components) is to change the L_2 norm regularizer in (1) by an L_1 norm, $\sum_j |w_{jk}|$. We refer to this model as L_1 -Support Vector Machines (L_1 -SVM).² The use of the L_1 norm tends to give sparse solution. For binary classification this approach has been pioneered by Mangasarian and his colleagues (see for instance [19]). One difficulty with L_1 -SVM is that one cannot use standard unconstrained optimization techniques to solve this problem (because of the non differentiability). One can resort to LP solvers [1]. See also [24] for a path tracking method and [10] for a coordinate-wise optimization method.

It is not obvious how the L_1 -SVM model can be extended to deal with the simultaneous feature subset selection problem in multi-class settings. Note that replacing the L_2 norm regularizer in (1) by an L_1 norm for each k separately will only lead to the selection of an independent set of different features for that class. Below, in section 2.5 we derive a new L_1 -SVM model that addresses the simultaneous feature subset selection problem.

2.4 Feature selection via scaling factors

We now come back to the multi-class case and introduce the framework of scaling factors for doing feature selection.

The goal is to learn c classification functions $f_k(\mathbf{x}) = \mathbf{w}_k \cdot \mathbf{x}$, $1 \leq k \leq c$ and to do simultaneous feature selection, i.e. find a small set of features which are good for all the classifiers. Let's say that we want to find m features. A natural optimization problem to solve is (see also [16]):

$$\min \frac{1}{2} \sum_{k=1}^c \mathbf{w}_k^2 + \frac{C}{2} \sum_{i=1}^n \ell(y_{ik} \sum_{j=1}^d \sqrt{\sigma_j} w_{jk} x_{ij}) \quad (2)$$

subject to the constraints $\sigma_j \in \{0, 1\}$, $\sum_j \sigma_j = m$.

The reason for using the squared root in (2) will become clear later, but for now it does not change anything since $\sigma_j \in \{0, 1\}$. Minimizing this objective function would indeed select m features (those corresponding to $\sigma_j = 1$) and effectively discard the others.

The problem is that this optimization problem is combinatorial and thus difficult to solve. But we can relax the

²In " L_1 -SVM", L_1 refers to the norm of the regularizer. In the literature the same acronym has also been used in a different context to refer to the norm of the training errors.

constraint $\sigma_j \in \{0, 1\}$ by $\sigma_j \geq 0$. Also, by making the change of variables $w_{jk} \leftarrow w_{jk} \sqrt{\sigma_j}$, we can then rewrite the relaxed version of (2) as

$$\min \frac{1}{2} \sum_{k=1}^c \left(\sum_{j=1}^d \frac{w_{jk}^2}{\sigma_j} + C \sum_{i=1}^n \ell(y_{ik}(\mathbf{w}_k \cdot \mathbf{x}_i)) \right)$$

subject to the constraints $\sigma_j \geq 0, \sum \sigma_j = m$.

(In the above, the choice $\sigma_j = 0$ will correspond to $w_{jk} = 0 \forall k$ and feature j being eliminated.) This is a convex problem because the function $(x, y) \mapsto x^2/y$ is jointly convex.

Instead of having the constraint $\sum \sigma_j = m$, we can introduce a Lagrange multiplier and add $\frac{\lambda}{2} \sum \sigma_j$ in the objective function and get:

$$\min \frac{1}{2} \sum_{k=1}^c \left(\sum_{j=1}^d \frac{w_{jk}^2}{\sigma_j} + \lambda \sum \sigma_j + C \sum_{i=1}^n \ell(y_{ik}(\mathbf{w}_k \cdot \mathbf{x}_i)) \right)$$

subject to the constraints $\sigma_j \geq 0$.

Note that the two problems given above are equivalent: for every m , there exists a λ such that the solutions are identical (and vice versa). Dividing this objective function by $\sqrt{\lambda}$, making the change of variable $\sigma_j \leftarrow \sqrt{\lambda} \sigma_j$ and introducing $\tilde{C} = C/\sqrt{\lambda}$, we obtain the following optimization problem:

$$\min \frac{1}{2} \sum_{k=1}^c \left(\sum_{j=1}^d \frac{w_{jk}^2}{\sigma_j} + \sum \sigma_j + \tilde{C} \sum_{i=1}^n \ell(y_{ik}(\mathbf{w}_k \cdot \mathbf{x}_i)) \right) \quad (3)$$

subject to the constraints $\sigma_j \geq 0$.

The sparsity is not controlled by m anymore but by \tilde{C} . Small values of \tilde{C} will yield sparse solutions.

Once (3) has been optimized, there are two possibilities:

1. Just take the solution (i.e the \mathbf{w}_k) as it is. We will refer to this as the L_1 -SVM solution, for reasons that will become obvious below in section 2.5.
2. Go back to the original (non relaxed) formulation (2) by fixing $\sigma_j \leftarrow 1_{\sigma_j > 0}$ and optimize (2). This is equivalent to finding the features through the L_1 -SVM algorithm and then train a standard SVM on these features. We call this method SSVM (Sparse Support Vector Machines) as it uses the standard SVM formulation but is constrained to give sparse solutions.

The details of the optimization technique used to solve (3) are given in the appendix. The technique uses a combination of a Newton-type algorithm and a path tracking algorithm. In terms of training time, it is very efficient because it only involves sparse matrix vector multiplications, where the matrix contains the active features of the support vectors.

2.5 Multi-Class L_1 -SVM

For fixed w_{jk} the objective function in (3) can be minimized in closed form with respect to σ_j . Indeed, the optimal value of σ_j is given by $\sigma_j = \sqrt{\sum w_{jk}^2}$. Plugging this value in (3), we obtain

$$\sum_{j=1}^d \sqrt{\sum_{k=1}^c w_{jk}^2} + \frac{\tilde{C}}{2} \sum_{i=1}^n \ell(y_{ik}(\mathbf{w}_k \cdot \mathbf{x}_i)), \quad (4)$$

which can be seen as a multiclass extension of L_1 -SVM. Indeed, when there is only a single binary classification problem, the regularizer in (4) is just the L_1 norm of the weight vector. In the multi-class case (4) the optimization is more worrisome, not only because the objective function is not differentiable, but also because linear programming techniques cannot be employed.

2.6 Multi-Class RFE

Another possibility is the following. Instead of optimizing on σ , fix all of them to the value 1 and decide which one to keep based on the gradient. More precisely, let $T(w, \sigma)$ denote the objective function in (3) (with $\tilde{C} = C$) and define $V(\sigma) = \min_w T(w, \sigma)$. Solving problem (2) is actually equivalent to minimizing V over $\sigma_j \in \{0, 1\}, \sum \sigma_j = m$. One can do something similar to RFE: start with all the features (i.e. all $\sigma_j = 1$) and put one of the σ_j to 0 such that V is minimized. Keep doing that until m features are left. An approximate way of selecting the feature \hat{j} to suppress is to make a linear approximation of V and select the smallest component of the gradient of V :

$$\hat{j} = \arg \max \left. \frac{\partial V}{\partial \sigma_j} \right|_{\sigma=1}.$$

From the definition of V we get,

$$\frac{\partial V}{\partial \sigma_j} = \underbrace{\frac{\partial T}{\partial w}}_{=0} \frac{\partial w}{\partial \sigma} + \frac{\partial T}{\partial \sigma_j} = \frac{1}{2} \left(- \sum_k w_{jk}^2 + 1 \right).$$

So

$$\hat{j} = \arg \min \sum_k w_{jk}^2 \quad (5)$$

The criterion (5) is very intuitive: it removes features for which the weights are small.

3. EXPERIMENTS

The performance evaluation was done in a multi-labeled setting: each class was considered as an independent binary problem. This is in contrast to a multi-class setting where a test point is only assigned one label based on the maximum output of the binary classifiers. We considered the multi-labeled case because several of our datasets are of this type. But note that our algorithms can easily be modified for the multi-class case.

3.1 Algorithms

We evaluated the different algorithms described in the previous section, i.e., RFE, L_1 -SVM, and SSVM and compared them to a popular filter method for feature selection, Information Gain (IG) [23]. Even though IG has a natural extension to the multi-class scenario, there are different ways of using it for the multi-labeled scenario. We selected the features using the sum of the information gains of all the binary problems as the overall IG criterion. Thus, the score of feature j is:

$$\sum_{k=1}^c -H(P(y_{ik} = 1)) + P(x_{ij} > 0)H(P(y_{ik} = 1|x_{ij} > 0)) + P(x_{ij} = 0)H(P(y_{ik} = 1|x_{ij} = 0)), \quad (6)$$

where H is the binary entropy function, $H(t) = t \log(t) + (1-t) \log(1-t)$ and P denotes the empirical frequency over $i = 1 \dots n$.

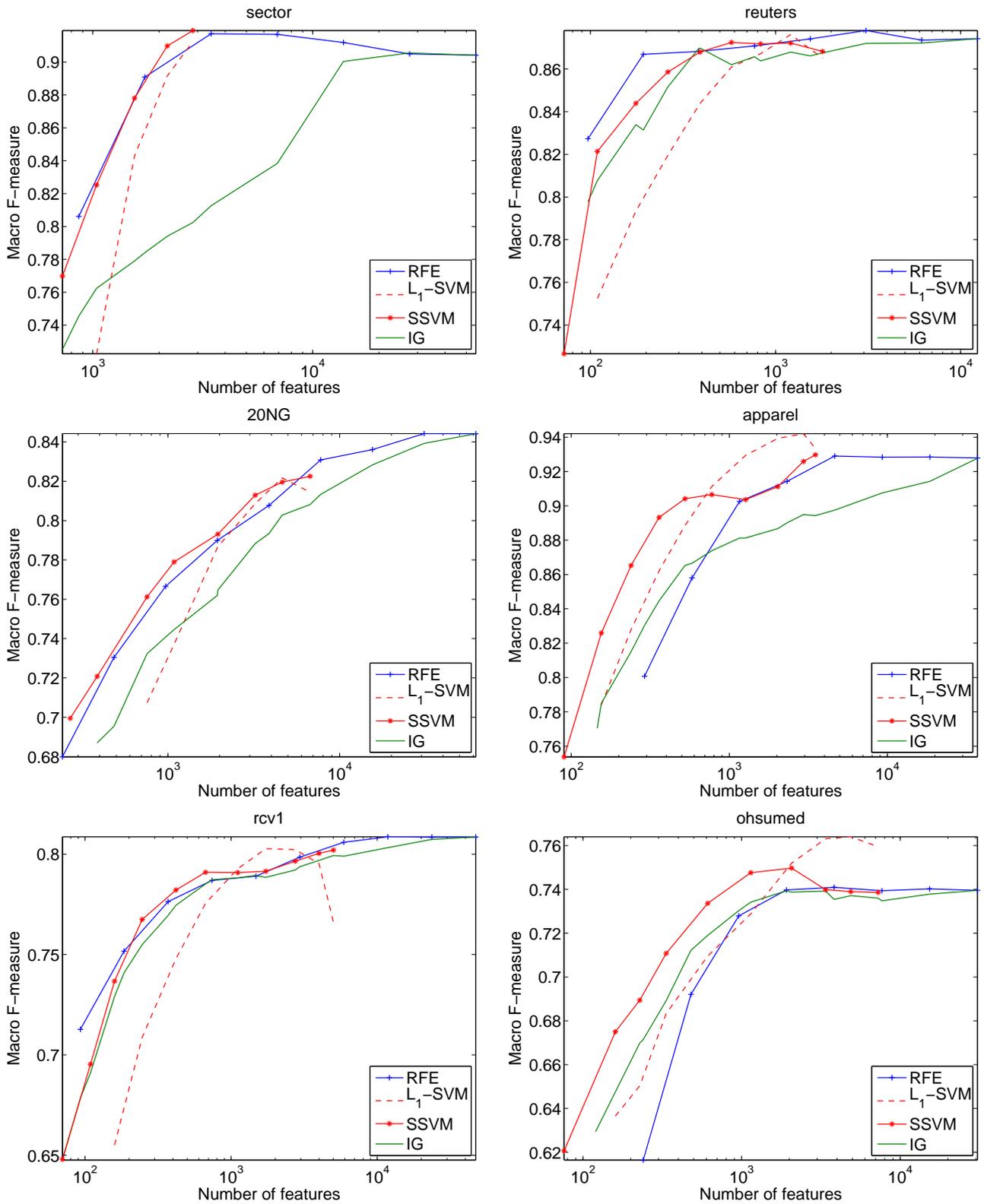


Figure 1: Microaverage F-measure as a function of the number of selected features on different datasets. webkb is shown on figure 2

Table 1: Properties of the datasets used in this paper, viz. the number of training documents (n), the number of features (d) and the number of classes (c). For these datasets marked *, only the c largest classes were retained.

	n	d	c	Description
sector	6412	55197	105	Industry sector [20]
20NG	15935	62061	20	20 Newsgroups [20]
apparel	22148	37144	22	Yahoo! Shopping
webkb	5505	3000	7	[20]
reuters*	9228	12317	45	ModApte split [20, 9]
rcv1*	11575	47236	52	Half size [17]
ohsumed	22926	30689	23	[9]

For all algorithms, the thresholds were further optimized after training to maximize the F-measure on the training set. This is the SCut strategy described in [22] except that we did not use a validation set. For evaluation we used the *microaveraged* F-measure [17, section 5.3].

For RFE, SSVM and IG, the C parameter of the SVM was fixed to 10. This large value corresponds closely to the hard margin SVM. For SSVM and L_1 -SVM, \tilde{C} was gradually increased from a small value (corresponding to a small number of features) to higher values until 5000 features are selected. For RFE, at each step of the algorithm, half of the features were discarded.

3.2 Datasets

We considered 7 datasets, which are summarized in table 1. Most of them were randomly split into a training set and a test set in a 2:1 ratio. The exact number of training documents can be seen in table 1. Note that the first 4 datasets are multi-class, while the 3 others are multi-labeled. All documents were coded with TFIDF and normalized to norm 1. We used a version of the `webkb` dataset where 3000 features have already been preselected by information gain. This might have inserted a bias in the experimental results.

3.3 Results

Performances of the different methods as a function of the number of features are plotted in figure 1. To have a more quantitative evaluation, we also computed the number of features that each method needs to keep in order to achieve a performance 3% below the performance of an SVM trained on all the features (right end of the black and green curves in figure 1). This was done by computing the intersects of the different curves with a horizontal line at the 3% loss level. We summarize in table 2 the relative feature set size with respect to information gain. For instance a 7 means that a given method has been able to find a feature set size 7 times smaller than IG (for the same performance).

Overall, SSVM is the best performing method. On some datasets, especially `sector` and `apparel`, IG is quite weak. Despite its simplicity RFE achieved good results. It is noteworthy that L_1 -SVM achieves better results on some datasets than a standard SVM trained on all the features (e.g., `apparel` and `ohsumed`). However this result in favor L_1 -SVM has to be taken with care because it is possible that the value of $C = 10$ that we used for SVM is sub-optimal. Note also that, on `20NG` the peak performance of L_1 -SVM was worse

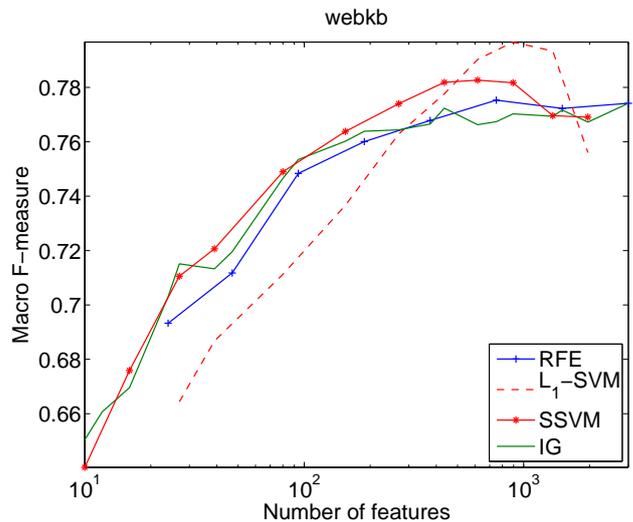


Figure 2: Continuation of figure 1.

Table 2: Relative reduction of the number of features in comparison to IG. The number of features was selected such that there is 3% loss in performance relative to the SVM trained on all the features. (Thus, higher the number above 1, better is the performance over IG.) Last column: number of features for SSVM.

	IG	RFE	L_1 -SVM	SSVM	# Feat
sector	1	7.08	5.92	7.29	1536
20NG	1	1.84	2.78	2.34	4528
apparel	1	5.19	8.76	13.36	429
webkb	1	0.78	0.42	1.06	83
reuters	1	1.68	0.57	1.31	189
rcv1	1	1.00	0.79	1.46	443
ohsumed	1	0.71	0.74	1.45	400
Mean	1	2.61	2.85	4.04	

than the SVM trained on all features.

We conclude this section with some running time analysis of the different algorithms. We take the `ohsumed` dataset for this purpose. As it can be seen from figure 3, running times of L_1 -SVM and SSVM are roughly linear with respect to the number of selected features. For selecting 5000 features the typical running times for L_1 -SVM and SSVM are in the order of a couple of hours. Note from the results in table 2 that, when SSVM is used, for most datasets we can expect accurate results with less than 1000 features. Therefore SSVM is very affordable. RFE is a much faster method since it requires only a few standard SVM trainings and, even when all the features are present, standard SVM training can be done very fast [14]. In general, RFE turns out to be an order of magnitude faster than SSVM and L_1 -SVM. Thus, if training time happens to be a bottleneck for SSVM, RFE could be the method of choice given its relatively good results.

4. DISCUSSION

The experiments of the previous section can be expanded to include other important investigations.

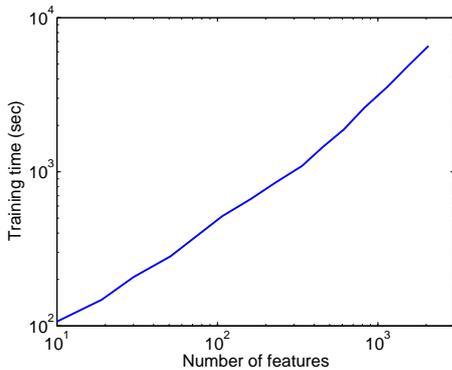


Figure 3: Total training time for L_1 -SVM as a function of the number of features on ohsumed. The relation is roughly linear. SSVM requires about the same time as L_1 -SVM since it uses the feature set found by L_1 -SVM and the final SVM training step is in comparison very fast. The total training time for RFE is 802 seconds.

In our experiments we used the microaveraged F-measure as the performance criterion of interest. In multi-class problems with imbalanced classes it is better to use the macroaveraged F-measure (mean of the F-measures of the various classes). Forman [8] points out the inadequacy of directly using filter methods (such as the information gain in (6)) for selecting a common set of features since they tend to be biased towards choosing features that work well for large and easy classes. It would be interesting to evaluate the usefulness of the embedded methods proposed in this paper for such situations.

Our experiments were only done in a multi-labeled setting. Many problems are in the pure multi-class form where exactly one class is assigned to each document. Even though the design process will develop the same binary classifiers as for the multi-labeled case, the decision process for multi-class problems uses the rule, $\arg \max_k \mathbf{w}_k \cdot \mathbf{x}$ to predict the class for document \mathbf{x} . For such pure multi-class problems a separate experimental study is needed for evaluating the new embedded methods.

For both the scenarios mentioned above we expect our embedded methods to work well; nevertheless, detailed experimental investigations are necessary to verify that.

Madigan et al [18] found coordinate-wise optimization methods to be efficient for training multinomial logistic regression classifiers. While it is true that their model does not address the simultaneous feature subset selection problem, their coordinate-wise optimization strategy is applicable to the L_1 -SVM model in (4). It would be interesting to implement and evaluate the efficiency of such an optimization approach against the method that we have used in this paper. Conversely, [18] could take advantage of our scaling framework: they could have only one scaling factor per component in (3) and put a Laplace hyperprior on the scaling factor. As shown in [21], this would lead to a sparse solution and would be very similar in spirit to what we have done.

A more sophisticated multi-labeled classification algorithm based on SVMs and ranking has been proposed in [5]. Again this algorithm could be modified to include scaling factors in the aim of simultaneous feature selection.

For RFE we have not evaluated the goodness of the measure, $\max_k |w_{jk}|$ proposed by [4] for deciding the elimination of features; see also, section 2.2. This is another useful direction for experimental investigation.

5. CONCLUSION

In this paper we have proposed a new class of embedded methods for feature selection. In particular, for simultaneous feature selection in multi-class and multi-labeled text classification, these embedded methods are the first of their kind. The methods do effective feature selection: on some datasets, the feature set size was one order of magnitude smaller than the one selected by Information Gain (for the same performance level). They are also computationally efficient and the structure of the optimization method allows ample scope for enhancing efficiency using parallel processing techniques.

Finally it is noteworthy that the framework of scaling factors that was used for deriving the new embedded methods is powerful and new algorithms for feature selection can easily be derived from it.

Appendix

The optimization of (3) has been implemented as an interior point method: a log barrier on the scaling factors has been added to the objective: $-t \sum \log(\sigma_j)$. Depending on the value of t , this will force more or less the σ to be away from 0. We will describe later how to minimize this modified objective function for a given value of t . Usually [3], t is initialized to a large value, the new unconstrained objective function is minimized, then t is decreased and this process iterates until t has reached a sufficiently small value.

We could have done exactly this procedure, but given that the optimal solution is likely to be sparse, some time is wasted on optimizing weights and scaling features which will zero at the end anyway. Instead, we implemented a path tracking approach [24, 15]: \tilde{C} is increased from a small value to its final value. Given the solution and the set of active features for a given \tilde{C} , the new solution for a larger \tilde{C} can be found efficiently because the set of active features is likely not to change too much. Pseudo-code is given in algorithm 1. During this process, there is no need to vary t since the gradual change in \tilde{C} already makes the optimization well behaved. Thus t is fixed to a relatively small value, $t = 10^{-3}$.

Two important ingredients in this algorithm are the criterion to add and remove features. If t were 0, one could simply remove the features corresponding to $\sigma_j = 0$. But because of the log barrier, all the σ_j will be at least equal to $2t$. That is the reason why we decided to have a relative threshold with respect to t .

On the other hand, adding features is based on the link with L_1 -SVM (4). One can indeed show that adding an inactive feature j will decrease the objective function if and only if

$$\sum_k \left(\frac{\tilde{C}}{2} \frac{\partial}{\partial w_{jk}} \sum_{i=1}^n \ell(y_{ik}(\mathbf{w}_k \cdot \mathbf{x}_i)) \right)^2 > 1.$$

Let us now the study the unconstrained minimization of (3) with the log barrier. To do so we used the the Levenberg-

categorization. In *Proceedings of SIGIR*, pages 137–145, 2001.

- [23] Y. Yang and J. Pedersen. A comparative study on feature selection in text categorization. In *International Conference on Machine Learning*, 1997.
- [24] J. Zhu, S. Rosset, T. Hastie, and R. Tibshirani. 1-norm support vector machines. In *Advances in Neural Information Processing Systems*, 2003.