Feature Selection for Linear SVM with Provable Guarantees

Saurabh Paul * Malik Magdon-Ismail † Petros Drineas ‡

June 4, 2018

Abstract

We give two provably accurate feature-selection techniques for the linear SVM. The algorithms run in deterministic and randomized time respectively. Our algorithms can be used in an unsupervised or supervised setting. The supervised approach is based on sampling features from support vectors. We prove that the margin in the feature space is preserved to within ϵ -relative error of the margin in the full feature space in the worst-case. In the unsupervised setting, we also provide worst-case guarantees of the radius of the minimum enclosing ball, thereby ensuring comparable generalization as in the full feature space and resolving an open problem posed in [1]. We present extensive experiments on real-world datasets to support our theory and to demonstrate that our method is competitive and often better than prior state-of-the-art, for which there are no known provable guarantees.

1 Introduction

The linear Support Vector Machine (SVM) is a popular classification method [2]. Few theoretical results exist for feature selection with SVMs. Empirically, numerous feature selection techniques work well (e.g. [3, 4]). We present a deterministic and a randomized feature selection technique for the linear SVM with a provable worst-case performance guarantee on the margin. The feature selection is unsupervised if features are selected obliviously to the data labels; otherwise, it is supervised. Our algorithms can be used in an unsupervised or supervised setting. In the unsupervised setting, our algorithm selects a number of features proportional to the rank of the data and preserves both the margin and radius of minimum enclosing ball to within ϵ -relative error in the worst-case, thus resolving an open problem posed in [1]. In the supervised setting, our algorithm selects O(#support vectors) features using only the set of support vectors, and preserves the margin for the support vectors to within ϵ -relative error in the worst-case.

SVM basics. The training data has n points $\mathbf{x}_i \in \mathbb{R}^d$, with respective labels $y_i \in \{-1, +1\}$ for $i = 1 \dots n$. For linearly separable data, the primal SVM learning problem constructs a hyperplane \mathbf{w}^* which maximizes the geometric margin (the minimum distance of a data point to the hyperplane), while separating the data. For non-separable data the "soft"

^{*}Computer Science Department, Rensselaer Polytechnic Institute, Troy, NY, USA, pauls2@rpi.edu

[†]Computer Science Department, Rensselaer Polytechnic Institute, Troy, NY, USA magdon@cs.rpi.edu

[‡]Computer Science Department, Rensselaer Polytechnic Institute, Troy, NY, USA, drinep@cs.rpi.edu

1-norm margin is maximized. The dual lagrangian formulation of the soft 1-norm SVM reduces to the following quadratic program:

$$\max_{\alpha_i} : \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

subject to:
$$\sum_{i=1}^n y_i \alpha_i = 0; \quad 0 \le \alpha_i \le C, \quad i = 1 \dots n.$$
 (1)

The regularizer C is part of the input and the hyperplane classifier can be constructed from the α_i . The out-of-sample performance is related to the VC-dimension of the resulting "fat"-separator. Assuming that the data lie in a ball of radius B, and that the hypothesis set consists of hyperplanes of width γ (the margin), then the VC-dimension is $O(B^2/\gamma^2)$ ([5]). Thus, by the VC-bound ([6]), the out-of-sample error is bounded by the in-sample error and a term monotonic in B^2/γ^2 .

Our Contributions. We give two provably accurate feature selection techniques for linear SVM in both unsupervised and supervised settings with worst-case performance guarantees on the margin. We use the single set spectral sparsification technique from [7] as our deterministic algorithm (the algorithm runs in deterministic time, hence the name 'deterministic') and leverage-score sampling ([1]) as the randomized algorithm. We give a new simple method of extending unsupervised feature selection to supervised in the context of SVMs by running the unsupervised technique on the support vectors. This allows us to select only O(#support vectors) features for the deterministic algorithm ($\tilde{O}(\#\text{support vectors})$) features for the randomized algorithm, where \tilde{O} hides the logarithmic factors) while still preserving the margin on the support vectors. Since the support vectors are a sufficient statistic for learning a linear SVM, preserving the margin on the support vectors should be enough for learning on all the data with the sampled feature set.

More formally, let γ^* be the optimal margin for the support vector set (which is also the optimal margin for all the data). The optimal margin γ^* is obtained by solving the SVM optimization problem using all the features. For a suitably chosen number of features r, let $\tilde{\gamma}^*$ be the optimal margin obtained by solving the SVM optimization problem using the support vectors in the sampled feature space. We prove that the margin is preserved to within ϵ -relative error: $\tilde{\gamma}^{*2} \geq (1 - \epsilon) \gamma^{*2}$. For the deterministic algorithm, the number of features required is $r = O(\#\text{support vectors}/\epsilon^2)$, whereas the randomized algorithm requires $r = \tilde{O}(\#\text{support vectors}/\epsilon^2)$ features to be selected.

In the unsupervised setting, by running our algorithm on all the data, instead of only the support vectors, we get a stronger result statistically, but using more features. The deterministic algorithm requires $O(\rho/\epsilon^2)$ features to be selected, while the randomized algorithm requires $O(\rho/\epsilon^2\log(\rho/\epsilon^2))$ features to be selected. Again, defining $\tilde{\gamma}^*$ as the optimal margin obtained by solving the SVM optimization problem using all the data in the sampled feature space, we prove that $\tilde{\gamma}^{*2} \geq (1-\epsilon)\,\gamma^{*2}$. We can now also prove that the data radius is preserved, $\tilde{B}^2 \leq (1+\epsilon)\,B^2$. This means that B^2/γ^{*2} is preserved to within ϵ -relative error, which means that the generalisation error is also preserved to within ϵ -relative error. The rank of the data is the effective dimension of the data, and one can construct this many combinations of pure features that preserve the geometry of the SVM exactly. What makes our result non-trivial is that we select this many pure features and preserve the geometry of the SVM to within ϵ -relative error.

On the practical side, we provide an efficient heuristic for our supervised feature selec-

tion using BSS which allows our algorithm to scale-up to large datasets. While the main focus of this paper is theoretical, we compare both supervised and unsupervised versions of feature selection using single-set spectral sparsification and leverage-score sampling with the corresponding supervised and unsupervised forms of Recursive Feature Elimination (RFE) ([3]), LPSVM ([4]), uniform sampling and rank-revealing QR factorization (RRQR) based method of column selection. Feature selection based on the single-set spectral sparsification and leverage-score sampling technique is competitive and often better than RFE and LPSVM, and none of the prior art comes with provable performance guarantees in either the supervised or unsupervised setting.

Related Work. All the prior art is heuristic in that there are no performance guarantees; nevertheless, they have been empirically tested against each other. Our algorithm comes with provable bounds, and performs comparably or better in empirical tests. We give a short survey of the prior art: Guyon et al. [3] and Rakotomamonjy [8] proposed SVM based criteria to rank features based on the weights. Weston et al. [9] formulated a combinatorial optimization problem to select features by minimising B^2/γ^2 . Weston et al. [10] used the zero norm to perform error minimization and feature selection in one step. A Newton based method of feature selection using linear programming was given in [4]. Tan et al. [11] formulated the ℓ_0 -norm Sparse SVM using mixed integer programming. Do et al. [12] proposed R-SVM which performs feature selection and ranking by optimizing the radius-margin bound with a scaling factor, and extend this work in [13] using a quadratic optimization problem with quadratic constraints. Another line of work includes the doubly regularised Support Vector Machine (DrSVM) [14] which uses a mixture of ℓ_2 -norm and ℓ_1 -norm penalties to solve the SVM optimization problem and perform variable selection. Subsequent works on DrSVM involve reducing the computational bottleneck ([15],[16]). Gilad-Bachrach et al. [17] formulate the margin as a function of set of features and score to sets of features according to the margin induced. Park et al. [18] studied the Fisher consistency and oracle property of penalized SVM where the dimension of inputs is fixed and showed that their method is able to identify the right model in most cases.

Paul et al. [19, 20] used random projections on linear SVM and showed that the margin and data-radius are preserved. However, this is different from our work, since they used linear combinations of features and we select pure features.

BSS and leverage-score sampling have been used to select features for k-means ([21, 22]), regularized least-squares classifier ([1, 23]). Our work further expands research into sparsification algorithms for machine learning.

2 Background

Notation: $\mathbf{A}, \mathbf{B}, \ldots$ denote matrices and $\boldsymbol{\alpha}, \mathbf{b}, \ldots$ denote column vectors; \mathbf{e}_i (for all $i = 1 \ldots n$) is the standard basis, whose dimensionality will be clear from context; and \mathbf{I}_n is the $n \times n$ identity matrix. The Singular Value Decomposition (SVD) of a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ of rank $\rho \leq \min\{n,d\}$ is equal to $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{n \times \rho}$ is an orthogonal matrix containing the left singular vectors, $\mathbf{\Sigma} \in \mathbb{R}^{\rho \times \rho}$ is a diagonal matrix containing the singular values $\sigma_1 \geq \sigma_2 \geq \ldots \sigma_{\rho} > 0$, and $\mathbf{V} \in \mathbb{R}^{d \times \rho}$ is a matrix containing the right singular vectors. The spectral norm of \mathbf{A} is $\|\mathbf{A}\|_2 = \sigma_1$.

Matrix Sampling Formalism: Let **A** be the data matrix consisting of n points and d dimensions, $\mathbf{S} \in \mathbb{R}^{d \times r}$ be a matrix such that $\mathbf{AS} \in \mathbb{R}^{n \times r}$ contains r columns of **A** (**S** is a

sampling matrix as it samples r columns of \mathbf{A}). Let $\mathbf{D} \in \mathbb{R}^{r \times r}$ be the diagonal matrix such that $\mathbf{ASD} \in \mathbb{R}^{n \times r}$ rescales the columns of \mathbf{A} that are in \mathbf{AS} . We will replace the sampling and re-scaling matrices by a single matrix $\mathbf{R} \in \mathbb{R}^{d \times r}$, where $\mathbf{R} = \mathbf{SD}$ first samples and then rescales r columns of \mathbf{A} .

Let **X** be a generic data matrix in d dimensions whose rows are data vectors \mathbf{x}_i^T , and let **Y** be the diagonal label matrix whose diagonal entries are the labels, $\mathbf{Y}_{ii} = y_i$. Let $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_n] \in \mathbb{R}^n$ be the vector of lagrange multipliers to be determined by solving eqn. (2). In matrix form, the SVM dual optimization problem is

$$\max_{\alpha} : \mathbf{1}^{T} \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^{T} \mathbf{Y} \mathbf{X} \mathbf{X}^{T} \mathbf{Y} \boldsymbol{\alpha}$$
subject to:
$$\mathbf{1}^{T} \mathbf{Y} \boldsymbol{\alpha} = 0; \qquad \mathbf{0} \le \boldsymbol{\alpha} \le \boldsymbol{C}.$$
 (2)

(In the above, **1**, **0**, C are vectors with the implied constant entry.) When the data and label matrices contain all the data, we will emphasize this using the notation $\mathbf{X}^{\mathbf{tr}} \in \mathbb{R}^{n \times d}$, $\mathbf{Y}^{\mathbf{tr}} \in \mathbb{R}^{n \times n}$. Solving (2) with these full data matrices gives a solution $\dot{\boldsymbol{\alpha}}^*$. The data \mathbf{x}_i for which $\dot{\alpha}_i^* > 0$ are support vectors and we denote by $\mathbf{X}^{\mathbf{sv}} \in \mathbb{R}^{p \times d}$, $\mathbf{Y}^{\mathbf{sv}} \in \mathbb{R}^{p \times p}$ the data and label matrices containing only the p support vectors. Solving (2) with $(\mathbf{X}^{\mathbf{tr}}, \mathbf{Y}^{\mathbf{tr}})$ or $(\mathbf{X}^{\mathbf{sv}}, \mathbf{Y}^{\mathbf{sv}})$ result in the same classifier. Let $\boldsymbol{\alpha}^*$ be the solution to (2) for the support vector data. The optimal separating hyperplane is $\mathbf{w}^* = (\mathbf{X}^{\mathbf{tr}})^T \mathbf{Y}^{\mathbf{tr}} \dot{\boldsymbol{\alpha}}^* = (\mathbf{X}^{\mathbf{sv}})^T \mathbf{Y}^{\mathbf{sv}} \boldsymbol{\alpha}^*$, where $\mathbf{X}^{\mathbf{sv}}$ is the support vector matrix. The geometric margin is $\gamma^* = 1/\|\mathbf{w}^*\|_2$, where $\|\mathbf{w}^*\|_2^2 = \sum_{i=1}^n \alpha_i^*$. The data radius is $B = \min_{\mathbf{x}^*} \max_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{x}^*\|_2$.

Our goal is to study how the SVM performs when run in the sampled feature space. Let \mathbf{X} , \mathbf{Y} be data and label matrices (such as those above) and $\mathbf{R} \in \mathbb{R}^{d \times r}$ a sampling and rescaling matrix which selects r columns of \mathbf{X} . The transformed dataset into the r selected features is $\tilde{\mathbf{X}} = \mathbf{X}\mathbf{R}$, and the SVM optimization problem in this feature space becomes

$$\max_{\hat{\boldsymbol{\alpha}}} : \quad \mathbf{1}^T \hat{\boldsymbol{\alpha}} - \frac{1}{2} \hat{\boldsymbol{\alpha}}^T \mathbf{Y} \mathbf{X} \mathbf{R} \mathbf{R}^T \mathbf{X}^T \mathbf{Y} \hat{\boldsymbol{\alpha}},$$

subject to:
$$\mathbf{1}^T \mathbf{Y} \hat{\boldsymbol{\alpha}} = 0: \qquad \mathbf{0} < \hat{\boldsymbol{\alpha}} < \boldsymbol{C}.$$
 (3)

For the supervised setting, we select features from the support vector matrix and so we set $\mathbf{X} = \mathbf{X^{sv}}$ and $\mathbf{Y} = \mathbf{Y^{sv}}$ and we select $r_1 \ll d$ features using \mathbf{R} . For the unsupervised setting, we select features from the full data matrix and so we set $\mathbf{X} = \mathbf{X^{tr}}$ and $\mathbf{Y} = \mathbf{Y^{tr}}$ and we select $r_2 \ll d$ features using \mathbf{R} .

3 Our main tools

In this section, we describe our main tools of feature selection from the numerical linear algebra literature, namely single-set spectral sparsification and leverage-score sampling.

Single-set Spectral Sparsification. The BSS algorithm ([7]) is a deterministic greedy technique that selects columns one at a time. The algorithm samples r columns in deterministic time, hence the name deterministic sampling. Consider the input matrix as a set of d column vectors $\mathbf{V}^T = [\mathbf{v}_1, \mathbf{v}_2,, \mathbf{v}_d]$, with $\mathbf{v}_i \in \mathbb{R}^\ell$ (i = 1, ..., d). Given ℓ and $r > \ell$, we iterate over $\tau = 0, 1, 2, ... - 1$. Define the parameters $L_\tau = \tau - \sqrt{r\ell}$, $\delta_L = 1$, $\delta_U = \left(1 + \sqrt{\ell/r}\right) / \left(1 - \sqrt{\ell/r}\right)$ and $U_\tau = \delta_U \left(\tau + \sqrt{\ell r}\right)$. For $U, L \in \mathbb{R}$ and $\mathbf{A} \in \mathbb{R}^{\ell \times \ell}$ a symmetric positive definite matrix with eigenvalues $\lambda_1, \lambda_2, ..., \lambda_\ell$, define $\Phi(L, \mathbf{A}) = \sum_{i=1}^\ell \frac{1}{\lambda_i - L}$ and $\hat{\Phi}(U, \mathbf{A}) = \sum_{i=1}^\ell \frac{1}{U - \lambda_i}$ as the lower and upper potentials respectively. These potential

functions measure how far the eigenvalues of **A** are from the upper and lower barriers U and L respectively. We define $\mathcal{L}(\mathbf{v}, \delta_L, \mathbf{A}, L)$ and $\mathcal{U}(\mathbf{v}, \delta_U, \mathbf{A}, U)$ as follows:

$$\mathcal{L}\left(\mathbf{v}, \delta_{L}, \mathbf{A}, L\right) = \frac{\mathbf{v}^{T} \left(\mathbf{A} - (L + \delta_{L}) \mathbf{I}_{\ell}\right)^{-2} \mathbf{v}}{\Phi \left(L + \delta_{L}, \mathbf{A}\right) - \Phi \left(L, \mathbf{A}\right)} - \mathbf{v}^{T} \left(\mathbf{A} - (L + \delta_{L}) \mathbf{I}_{\ell}\right)^{-1} \mathbf{v}$$

$$\mathcal{U}\left(\mathbf{v}, \delta_{U}, \mathbf{A}, U\right) = \frac{\mathbf{v}^{T} \left(\left(U + \delta_{U}\right) \mathbf{I}_{\ell} - \mathbf{A}\right)^{-2} \mathbf{v}}{\hat{\Phi}\left(U, \mathbf{A}\right) - \hat{\Phi}\left(U + \delta_{U}, \mathbf{A}\right)} + \mathbf{v}^{T} \left(\left(U + \delta_{U}\right) \mathbf{I}_{\ell} - \mathbf{A}\right)^{-1} \mathbf{v}.$$

At every iteration, there exists an index i_{τ} and a weight $t_{\tau} > 0$ such that, $t_{\tau}^{-1} \leq \mathcal{L}(\mathbf{v}_i, \delta_L, \mathbf{A}, L)$ and $t_{\tau}^{-1} \geq \mathcal{U}(\mathbf{v}_i, \delta_U, \mathbf{A}, U)$. Thus, there will be at most r columns selected after τ iterations. The running time of the algorithm is dominated by the search for an index i_{τ} satisfying $\mathcal{U}(\mathbf{v}_i, \delta_U, \mathbf{A}_{\tau}, U_{\tau}) \leq \mathcal{L}(\mathbf{v}_i, \delta_L, \mathbf{A}_{\tau}, L_{\tau})$ and computing the weight t_{τ} . One needs to compute the upper and lower potentials $\hat{\Phi}(U, \mathbf{A})$ and $\Phi(L, \mathbf{A})$ and hence the eigenvalues of \mathbf{A} . Cost per iteration is $O(\ell^3)$ and the total cost is $O(r\ell^3)$. For i = 1, ..., d, we need to compute \mathcal{L} and \mathcal{U} for every \mathbf{v}_i which can be done in $O(d\ell^2)$ for every iteration, for a total of $O(rd\ell^2)$. Thus total running time of the algorithm is $O(rd\ell^2)$. We include the algorithm as Algorithm 1.

Input: $\mathbf{V}^T = [\mathbf{v}_1, \mathbf{v}_2, ... \mathbf{v}_d] \in \mathbb{R}^{\ell \times d}$ with $\mathbf{v}_i \in \mathbb{R}^{\ell}$ and $r > \ell$. Output: Matrices $\mathbf{S} \in \mathbb{R}^{d \times r}$, $\mathbf{D} \in \mathbb{R}^{r \times r}$.

- 1. Initialize $\mathbf{A}_0 = \mathbf{0}_{\ell \times \ell}$, $\mathbf{S} = \mathbf{0}_{d \times r}$, $\mathbf{D} = \mathbf{0}_{r \times r}$.
- 2. Set constants $\delta_L = 1$ and $\delta_U = \left(1 + \sqrt{\ell/r}\right) / \left(1 \sqrt{\ell/r}\right)$.
- 3. for $\tau = 0$ to r 1 do
 - Let $L_{\tau} = \tau \sqrt{r\ell}$; $U_{\tau} = \delta_U \left(\tau + \sqrt{\ell r} \right)$.
 - Pick index $i \in \{1, 2, ...d\}$ and number $t_{\tau} > 0$, such that

$$\mathcal{U}\left(\mathbf{v}_{i}, \delta_{U}, \mathbf{A}_{\tau}, U_{\tau}\right) < \mathcal{L}\left(\mathbf{v}_{i}, \delta_{L}, \mathbf{A}_{\tau}, L_{\tau}\right).$$

- Let $t_{\tau}^{-1} = \frac{1}{2} \left(\mathcal{U} \left(\mathbf{v}_{i}, \delta_{U}, \mathbf{A}_{\tau}, U_{\tau} \right) + \mathcal{L} \left(\mathbf{v}_{i}, \delta_{L}, \mathbf{A}_{\tau}, L_{\tau} \right) \right)$
- Update $\mathbf{A}_{\tau+1} = \mathbf{A}_{\tau} + t_{\tau} \mathbf{v}_i \mathbf{v}_i^T$; set $\mathbf{S}_{i_{\tau}, \tau+1} = 1$ and $\mathbf{D}_{\tau+1, \tau+1} = 1/\sqrt{t_{\tau}}$.
- 4. end for
- 5. Multiply all the weights in **D** by $\sqrt{r^{-1}\left(1-\sqrt{(\ell/r)}\right)}$.
- 6. Return **S** and **D**.

Algorithm 1: Single-set Spectral Sparsification

We present the following lemma for the single-set spectral sparsification algorithm.

Lemma 1. BSS ([7]): Given $\mathbf{V} \in \mathbb{R}^{d \times \ell}$ satisfying $\mathbf{V}^T \mathbf{V} = \mathbf{I}_{\ell}$ and $r > \ell$, we can deterministically construct sampling and rescaling matrices $\mathbf{S} \in \mathbb{R}^{d \times r}$ and $\mathbf{D} \in \mathbb{R}^{r \times r}$ with $\mathbf{R} = \mathbf{SD}$, such that, for all $\mathbf{y} \in \mathbb{R}^{\ell}$: $\left(1 - \sqrt{\ell/r}\right)^2 \|\mathbf{V}\mathbf{y}\|_2^2 \le \left\|\mathbf{V}^T \mathbf{R}\mathbf{y}\right\|_2^2 \le \left(1 + \sqrt{\ell/r}\right)^2 \|\mathbf{V}\mathbf{y}\|_2^2$.

We now present a slightly modified version of Lemma 1 for our theorems.

Lemma 2. Given $\mathbf{V} \in \mathbb{R}^{d \times \ell}$ satisfying $\mathbf{V}^T \mathbf{V} = \mathbf{I}_{\ell}$ and $r > \ell$, we can deterministically construct sampling and rescaling matrices $\mathbf{S} \in \mathbb{R}^{d \times r}$ and $\mathbf{D} \in \mathbb{R}^{r \times r}$ such that for $\mathbf{R} = \mathbf{S}\mathbf{D}$, $\|\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}\|_2 \leq 3\sqrt{\ell/r}$

Proof. From Lemma 1, it follows, $\sigma_{\ell} \left(\mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V} \right) \geq \left(1 - \sqrt{\ell/r} \right)^2$, $\sigma_1 \left(\mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V} \right) \leq \left(1 + \sqrt{\ell/r} \right)^2$. Thus, $\lambda_{max} \left(\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V} \right) \leq \left(1 - \left(1 - \sqrt{\ell/r} \right)^2 \right) \leq 2\sqrt{\ell/r}$. Similarly, $\lambda_{min} \left(\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V} \right) \geq \left(1 - \left(1 + \sqrt{\ell/r} \right)^2 \right) \geq 3\sqrt{\ell/r}$. Combining these two results, we have $\| \mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V} \|_2 \leq 3\sqrt{\ell/r}$.

Leverage-Score Sampling. Our randomized feature selection method is based on importance sampling or the so-called leverage-score sampling of [1]. Let V be the top- ρ right singular vectors of the training set X. A carefully chosen probability distribution of the form

$$p_i = \frac{\|\mathbf{V}_i\|_2^2}{n}$$
, for $i = 1, 2, ..., d$, (4)

i.e. proportional to the squared Euclidean norms of the rows of the right-singular vectors is constructed. Select r rows of \mathbf{V} in i.i.d trials and re-scale the rows with $1/\sqrt{p_i}$. The time complexity is dominated by the time to compute the SVD of \mathbf{X} .

Lemma 3. Let $\epsilon \in (0, 1/2]$ be an accuracy parameter and $\delta \in (0, 1)$ be the failure probability. Given $\mathbf{V} \in \mathbb{R}^{d \times \ell}$ satisfying $\mathbf{V}^T \mathbf{V} = \mathbf{I}_{\ell}$. Let $\tilde{p} = min\{1, rp_i\}$, let p_i be as Eqn. 4 and let $r = O\left(\frac{n}{\epsilon^2}\log\left(\frac{n}{\epsilon^2\sqrt{\delta}}\right)\right)$. Construct the sampling and rescaling matrix \mathbf{R} . Then with probability at least 0.99, $\|\mathbf{V}^T\mathbf{V} - \mathbf{V}^T\mathbf{R}^T\mathbf{R}\mathbf{V}\|_2 \le \epsilon$.

4 Theoretical Analysis

Though our feature selection algorithms are relatively simple, we show that running the linear SVM in the feature space results in a classifier with provably comparable margin to the SVM classifier obtained from the full feature space. Our main results are in Theorems 1 & 3. We state the theorems for BSS, but similar theorems can be stated for leverage-score sampling.

4.1 Margin is preserved by Supervised Feature Selection

Theorem 1 says that you get comparable margin from solving the SVM on the support vectors (equivalently all the data) and from solving the SVM on support vectors in a feature space with only O(#support vectors) features.

Theorem 1. Given $\epsilon \in (0,1)$, run supervised BSS-feature selection on \mathbf{X}^{sv} with $r_1 = 36p/\epsilon^2$, to obtain the feature sampling and rescaling matrix \mathbf{R} . Let γ^* and $\tilde{\gamma}^*$ be the margins obtained by solving the SVM dual (2) with $(\mathbf{X}^{sv}, \mathbf{Y}^{sv})$ and $(\mathbf{X}^{sv}\mathbf{R}, \mathbf{Y}^{sv})$ respectively. Then, $\tilde{\gamma}^{*2} \geq (1-\epsilon) \gamma^{*2}$.

Proof. Let $\mathbf{X^{tr}} \in \mathbb{R}^{n \times d}$, $\mathbf{Y^{tr}} \in \mathbb{R}^{n \times n}$ be the feature matrix and class labels of the training set (as defined in Section 2) and let $\dot{\boldsymbol{\alpha}}^* = [\alpha_1^*, \alpha_2^*, \dots, \alpha_n^*]^T \in \mathbb{R}^n$ be the vector achieving the optimal solution for the problem of eqn. (2). Then,

$$Z_{opt} = \sum_{j=1}^{n} \dot{\alpha}_{j}^{*} - \frac{1}{2} \dot{\alpha}^{*T} \mathbf{Y}^{tr} \mathbf{X}^{tr} \left(\mathbf{X}^{tr} \right)^{T} \mathbf{Y}^{tr} \dot{\alpha}^{*}$$
 (5)

Let $p \leq n$ be the support vectors with $\dot{\alpha}_j > 0$. Let $\boldsymbol{\alpha}^* = \left[\alpha_1^*, \alpha_2^*, \dots, \alpha_p^*\right]^T \in \mathbb{R}^p$ be the vector achieving the optimal solution for the problem of eqn. (5). Let $\mathbf{X}^{\mathbf{sv}} \in \mathbb{R}^{p \times d}$, $\mathbf{Y}^{\mathbf{sv}} \in \mathbb{R}^{p \times p}$ be the support vector matrix and the corresponding labels respectively. Let $\mathbf{E} = \mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}$. Then, we can write eqn (5) in terms of support vectors as,

$$Z_{opt} = \sum_{i=1}^{p} \alpha_{i}^{*} - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y}^{sv} \mathbf{X}^{sv} (\mathbf{X}^{sv})^{T} \mathbf{Y}^{sv} \boldsymbol{\alpha}^{*}$$

$$= \sum_{i=1}^{p} \alpha_{i}^{*} - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^{T} \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \boldsymbol{\alpha}^{*}$$

$$= \sum_{i=1}^{p} \alpha_{i}^{*} - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^{T} \mathbf{R} \mathbf{R}^{T} \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \boldsymbol{\alpha}^{*}$$

$$- \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{E} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \boldsymbol{\alpha}^{*}. \tag{6}$$

Let $\tilde{\boldsymbol{\alpha}}^* = \left[\tilde{\alpha}_1^*, \tilde{\alpha}_2^*, \dots, \tilde{\alpha}_p^*\right]^T \in \mathbb{R}^p$ be the vector achieving the optimal solution for the dimensionally-reduced SVM problem of eqn. (6) using $\tilde{\mathbf{X}}^{\mathbf{sv}} = \mathbf{X}^{\mathbf{sv}}\mathbf{R}$. Using the SVD of $\mathbf{X}^{\mathbf{sv}}$,

$$\tilde{Z}_{opt} = \sum_{i=1}^{p} \tilde{\alpha}_{i}^{*} - \frac{1}{2} \tilde{\alpha}^{*T} \mathbf{Y}^{sv} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{T} \mathbf{R} \mathbf{R}^{T} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \tilde{\alpha}^{*}.$$
 (7)

Since the constraints on α^* , $\tilde{\alpha}^*$ do not depend on the data it is clear that $\tilde{\alpha}^*$ is a feasible solution for the problem of eqn. (6). Thus, from the optimality of α^* , and using eqn. (7), it follows that

$$Z_{opt} = \sum_{i=1}^{p} \alpha_{i}^{*} - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^{T} \mathbf{R} \mathbf{R}^{T} \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \boldsymbol{\alpha}^{*}$$

$$- \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{E} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \boldsymbol{\alpha}^{*}$$

$$\geq \sum_{i=1}^{p} \tilde{\alpha}_{i}^{*} - \frac{1}{2} \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^{T} \mathbf{R} \mathbf{R}^{T} \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \tilde{\boldsymbol{\alpha}}^{*}$$

$$- \frac{1}{2} \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{E} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \tilde{\boldsymbol{\alpha}}^{*}$$

$$= \tilde{Z}_{opt} - \frac{1}{2} \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma} \mathbf{E} \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{sv} \tilde{\boldsymbol{\alpha}}^{*}. \tag{8}$$

We now analyze the second term using standard submultiplicativity properties and $\mathbf{V}^T\mathbf{V}=$

I. Taking $\mathbf{Q} = \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{sv} \mathbf{U} \boldsymbol{\Sigma}$,

$$\frac{1}{2}\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{Y}^{\mathbf{s}\mathbf{v}}\mathbf{U}\boldsymbol{\Sigma}\mathbf{E}\boldsymbol{\Sigma}\mathbf{U}^{T}\mathbf{Y}^{\mathbf{s}\mathbf{v}}\tilde{\boldsymbol{\alpha}}^{*}$$

$$\leq \frac{1}{2}\|\mathbf{Q}\|_{2}\|\mathbf{E}\|_{2}\|\mathbf{Q}^{T}\|_{2}$$

$$= \frac{1}{2}\|\mathbf{E}\|_{2}\|\tilde{\mathbf{Q}}^{*T}\mathbf{Y}^{\mathbf{s}\mathbf{v}}\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{T}\|_{2}^{2}$$

$$= \frac{1}{2}\|\mathbf{E}\|_{2}\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{Y}^{\mathbf{s}\mathbf{v}}\mathbf{X}^{\mathbf{s}\mathbf{v}}\|_{2}^{2}.$$
(9)

Combining eqns. (8) and (9), we get

$$Z_{opt} \geq \tilde{Z}_{opt} - \frac{1}{2} \|\mathbf{E}\|_2 \|\tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{sv}} \mathbf{X}^{\mathbf{sv}}\|_2^2.$$
 (10)

We now proceed to bound the second term in the right-hand side of the above equation. Towards that end, we bound the difference:

$$\begin{aligned} & \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{U} \boldsymbol{\Sigma} \left(\mathbf{V}^{T} \mathbf{R} \mathbf{R}^{T} \mathbf{V} - \mathbf{V}^{T} \mathbf{V} \right) \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \tilde{\boldsymbol{\alpha}}^{*} \right| \\ &= & \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{U} \boldsymbol{\Sigma} \left(-\mathbf{E} \right) \boldsymbol{\Sigma} \mathbf{U}^{T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \tilde{\boldsymbol{\alpha}}^{*} \right| \\ &\leq & \left\| \mathbf{E} \right\|_{2} \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^{T} \right\|_{2}^{2} \\ &= & \left\| \mathbf{E} \right\|_{2} \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}} \right\|_{2}^{2} \\ &= & \left\| \mathbf{E} \right\|_{2} \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}} \right\|_{2}^{2}. \end{aligned}$$

We can rewrite the above inequality as

$$\begin{aligned} & \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}} \mathbf{R} \right\|_{2}^{2} - \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}} \right\|_{2}^{2} \\ & \leq & \left\| \mathbf{E} \right\|_{2} \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}} \right\|_{2}^{2} \\ & \leq & \frac{\left\| \mathbf{E} \right\|_{2}}{1 - \left\| \mathbf{E} \right\|_{2}} \left\| \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}} \mathbf{R} \right\|_{2}^{2}. \end{aligned}$$

Combining with eqn. (10), we get

$$Z_{opt} \ge \tilde{Z}_{opt} - \frac{1}{2} \left(\frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2} \right) \|\tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{sv} \mathbf{X}^{sv} \mathbf{R} \|_2^2.$$
 (11)

Now recall that $\mathbf{w}^{*T} = \boldsymbol{\alpha}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}}, \ \tilde{\mathbf{w}}^{*T} = \tilde{\boldsymbol{\alpha}}^{*T} \mathbf{Y}^{\mathbf{s}\mathbf{v}} \mathbf{X}^{\mathbf{s}\mathbf{v}} \mathbf{R}, \ \|\mathbf{w}^*\|_2^2 = \sum_{i=1}^p \alpha_i^*, \text{ and } \|\tilde{\mathbf{w}}^*\|_2^2 = \sum_{i=1}^p \tilde{\alpha}_i^*.$ Then, the optimal solutions Z_{opt} and \tilde{Z}_{opt} can be expressed as follows:

$$Z_{opt} = \|\mathbf{w}^*\|_2^2 - \frac{1}{2} \|\mathbf{w}^*\|_2^2 = \frac{1}{2} \|\mathbf{w}^*\|_2^2.$$
 (12)

$$\tilde{Z}_{opt} = \|\tilde{\mathbf{w}}^*\|_2^2 - \frac{\tilde{\mathbf{1}}}{2} \|\tilde{\mathbf{w}}^*\|_2^2 = \frac{\tilde{\mathbf{1}}}{2} \|\tilde{\mathbf{w}}^*\|_2^2.$$
(13)

Combining eqns. (11), (12) and (13), we get $\|\mathbf{w}^*\|_2^2 \ge \|\tilde{\mathbf{w}}^*\|_2^2 - \left(\frac{\|\mathbf{E}\|_2}{1-\|\mathbf{E}\|_2}\right) \|\tilde{\mathbf{w}}^*\|_2^2 = \left(1 - \frac{\|\mathbf{E}\|_2}{1-\|\mathbf{E}\|_2}\right) \|\tilde{\mathbf{w}}^*\|_2^2$. Let $\gamma^* = \|\mathbf{w}^*\|_2^{-1}$ be the geometric margin of the problem of eqn. (6)

and let $\tilde{\gamma}^* = \|\tilde{\mathbf{w}}^*\|_2^{-1}$ be the geometric margin of the problem of eqn. (7). Then, the above equation implies: $\gamma^{*2} \leq \left(1 - \frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2}\right)^{-1} \tilde{\gamma}^{*2} \Rightarrow \tilde{\gamma}^{*2} \geq \left(1 - \frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2}\right) \gamma^{*2}$. The result follows because $\|\mathbf{E}\|_2 \leq \epsilon/2$ by our choice of r, and so $\|\mathbf{E}\|_2/(1 - \|\mathbf{E}\|_2) \leq \epsilon$.

We now state a similar theorem for leverage-score sampling.

Theorem 2. Given $\epsilon \in (0,1)$, run supervised Leverage-score sampling based feature selection on \mathbf{X}^{sv} with $r_1 = \tilde{O}(p/\epsilon^2)$, to obtain the feature sampling and rescaling matrix \mathbf{R} . Let γ^* and $\tilde{\gamma}^*$ be the margins obtained by solving the SVM dual (2) with $(\mathbf{X}^{sv}, \mathbf{Y}^{sv})$ and $(\mathbf{X}^{sv}\mathbf{R}, \mathbf{Y}^{sv})$ respectively. Then with probability at least 0.99, $\tilde{\gamma}^{*2} > (1 - \epsilon) \gamma^{*2}$.

4.2 Geometry is preserved by Unsupervised Feature Selection

In the unsupervised setting, the next theorem says that with a number of features proportional to the rank of the training data (which is at most the number of data points), we preserve B^2/γ^2 , thus ensuring comparable generalization error bounds (B is the radius of the minimum enclosing ball).

Theorem 3. Given $\epsilon \in (0,1)$, run unsupervised BSS-feature selection on the full data \mathbf{X}^{tr} with $r_2 = O\left(\rho/\epsilon^2\right)$, where $\rho = \operatorname{rank}(\mathbf{X}^{tr})$, to obtain the feature sampling and rescaling matrix \mathbf{R} . Let γ^* and $\tilde{\gamma}^*$ be the margins obtained by solving the SVM dual (2) with $(\mathbf{X}^{tr}, \mathbf{Y}^{tr})$ and $(\mathbf{X}^{tr}\mathbf{R}, \mathbf{Y}^{tr})$ respectively; and, let B and \tilde{B} be the radii for the data matrices \mathbf{X}^{tr} and $\mathbf{X}^{tr}\mathbf{R}$ respectively. Then,

$$\frac{\tilde{B}^2}{\tilde{\gamma}^{*2}} \le \frac{(1+\epsilon)}{(1-\epsilon)} \frac{B^2}{\gamma^{*2}} = (1+O(\epsilon)) \frac{B^2}{\gamma^{*2}}.$$

Proof. (sketch) The proof has two parts. First, as in Theorem 1 we prove that $\tilde{\gamma}^{*2} \geq (1-\epsilon) \cdot \gamma^{*2}$. This proof is almost identical to the proof of Theorem 1 (replacing $(\mathbf{X^{sv}}, \mathbf{Y^{sv}})$ with $(\mathbf{X^{tr}}, \mathbf{Y^{tr}})$), and so we omit it. Second, we prove that $\tilde{B}^2 \leq (1+\epsilon)B^2$. We give the result (with proof) as Theorem 5. The theorem follows by combining these two results. \square

We now state a similar theorem for leverage-score sampling.

Theorem 4. Given $\epsilon \in (0,1)$, run unsupervised Leverage-score feature selection on the full data \mathbf{X}^{tr} with $r_2 = \tilde{O}\left(\rho/\epsilon^2\right)$, where $\rho = \operatorname{rank}(\mathbf{X}^{tr})$, to obtain the feature sampling and rescaling matrix \mathbf{R} . Let γ^* and $\tilde{\gamma}^*$ be the margins obtained by solving the SVM dual (2) with $(\mathbf{X}^{tr}, \mathbf{Y}^{tr})$ and $(\mathbf{X}^{tr}\mathbf{R}, \mathbf{Y}^{tr})$ respectively; and, let B and \tilde{B} be the radii for the data matrices \mathbf{X}^{tr} and $\mathbf{X}^{tr}\mathbf{R}$ respectively. Then with probability at least 0.99,

$$\frac{\tilde{B}^2}{\tilde{\gamma}^{*2}} \le \frac{(1+\epsilon)}{(1-\epsilon)} \frac{B^2}{\gamma^{*2}} = (1+O(\epsilon)) \frac{B^2}{\gamma^{*2}}.$$

4.3 Proof That the Data Radius is preserved by Unsupervised BSS-Feature Selection.

Theorem 5. Let $r_2 = O(n/\epsilon^2)$, where $\epsilon > 0$ is an accuracy parameter, n is the number of training points and r_2 is the number of features selected. Let B be the radius of the minimum ball enclosing all points in the full-dimensional space, and let \tilde{B} be the radius of the ball enclosing all points in the sampled subspace obtained by using BSS in an unsupervised manner. For \mathbf{R} as in Lemma 2, $\tilde{B}^2 \leq (1+\epsilon)B^2$.

Proof. We consider the matrix $\mathbf{X}_B \in \mathbb{R}^{(n+1)\times d}$ whose first n rows are the rows of $\mathbf{X}^{\mathbf{tr}}$ and whose last row is the vector \mathbf{x}_B^T ; here \mathbf{x}_B denotes the center of the minimum radius ball enclosing all n points. Then, the SVD of \mathbf{X}_B is equal to $\mathbf{X}_B = \mathbf{U}_B \mathbf{\Sigma}_B \mathbf{V}_B^T$, where $\mathbf{U}_B \in \mathbb{R}^{(n+1)\times \rho_B}$, $\mathbf{\Sigma}_B \in \mathbb{R}^{\rho_B \times \rho_B}$, and $\mathbf{V} \in \mathbb{R}^{d \times \rho_B}$. Here ρ_B is the rank of the matrix \mathbf{X}_B and clearly $\rho_B \leq \rho + 1$. (Recall that ρ is the rank of the matrix $\mathbf{X}^{\mathbf{tr}}$.) Let B be the radius of the minimal radius ball enclosing all n points in the original space. Then, for any $i = 1, \ldots, n$,

 $B^{2} \ge \|\mathbf{x}_{i} - \mathbf{x}_{B}\|_{2}^{2} = \|(\mathbf{e}_{i} - \mathbf{e}_{n+1})^{T} \mathbf{X}_{B}\|_{2}^{2}.$ (14)

Now consider the matrix $\mathbf{X}_{B}\mathbf{R}$ and notice that

$$\begin{aligned} & \left\| \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right)^{T} \mathbf{X}_{B} \right\|_{2}^{2} - \left\| \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right)^{T} \mathbf{X}_{B} \mathbf{R} \right\|_{2}^{2} \right\| \\ &= \left\| \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right)^{T} \left(\mathbf{X}_{B} \mathbf{X}_{B}^{T} - \mathbf{X}_{B} \mathbf{R} \mathbf{R}^{T} \mathbf{X}_{B}^{T} \right) \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right) \right\| \\ &= \left\| \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right)^{T} \mathbf{U}_{B} \mathbf{\Sigma}_{B} \mathbf{E}_{B} \mathbf{\Sigma}_{B} \mathbf{U}_{B}^{T} \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right) \right\| \\ &\leq \left\| \mathbf{E}_{B} \right\|_{2} \left\| \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right)^{T} \mathbf{U}_{B} \mathbf{\Sigma}_{B} \right\|_{2}^{2} \\ &= \left\| \mathbf{E}_{B} \right\|_{2} \left\| \left(\mathbf{e}_{i} - \mathbf{e}_{n+1} \right)^{T} \mathbf{X}_{B} \right\|_{2}^{2}. \end{aligned}$$

In the above, we let $\mathbf{E}_B \in \mathbb{R}^{\rho_B \times \rho_B}$ be the matrix that satisfies $\mathbf{V}_B^T \mathbf{V}_B = \mathbf{V}_B^T \mathbf{R} \mathbf{R}^T \mathbf{V}_B + \mathbf{E}_B$, and we also used $\mathbf{V}_B^T \mathbf{V}_B = \mathbf{I}$. Now consider the ball whose center is the (n+1)-th row of the matrix $\mathbf{X}_B \mathbf{R}$ (essentially, the center of the minimal radius enclosing ball for the original points in the sampled space). Let $\tilde{i} = \arg\max_{i=1...n} \left\| (\mathbf{e}_i - \mathbf{e}_{n+1})^T \mathbf{X}_B \mathbf{R} \right\|_2^2$; then, using the above bound and eqn. (14), we get $\left\| (\mathbf{e}_{\tilde{i}} - \mathbf{e}_{n+1})^T \mathbf{X}_B \mathbf{R} \right\|_2^2 \leq (1 + \|\mathbf{E}_B\|_2) \left\| (\mathbf{e}_{\tilde{i}} - \mathbf{e}_{n+1})^T \mathbf{X}_B \right\|_2^2 \leq (1 + \|\mathbf{E}_B\|_2) B^2$. Thus, there exists a ball centered at $\mathbf{e}_{n+1}^T \mathbf{X}_B \mathbf{R}$ (the projected center of the minimal radius ball in the original space) with radius at most $\sqrt{1 + \|\mathbf{E}_B\|_2} B$ that encloses all the points in the sampled space. Recall that \tilde{B} is defined as the radius of the minimal radius ball that encloses all points in sampled subspace; clearly, $\tilde{B}^2 \leq (1 + \|\mathbf{E}_B\|_2) B^2$. We can now use Lemma 2 on \mathbf{V}_B to conclude that (using $\rho_B \leq \rho + 1$) $\|\mathbf{E}_B\|_2 \leq \epsilon$.

5 Experiments

We compared BSS and leverage-score sampling with RFE ([3]), LPSVM ([4]), rank-revealing QR factorization (RRQR), random feature selection and full-data without feature selection on synthetic and real-world datasets. For the supervised case, we first run SVM on the training set, then run a feature selection method on the support-vector set and recalibrate the model using the support vector-set. For unsupervised feature selection, we perform feature selection on the training set. For LPSVM, we were not able to control the number of features and report the out-of-sample error using the features output by the algorithm. We did not extrapolate the values of out-of-sample error for LPSVM. We repeated random feature selection and leverage-score sampling five times. We performed ten-fold cross-validation and repeated it ten times. For medium-scale datasets like TechTC-300 we do not perform approximate BSS. For large-scale datasets like Reuters-CCAT ([24]) we use the approximate BSS method as described in Section 5.5. We used LIBSVM ([25]) as our

SVM solver for medium-scale datasets and LIBLINEAR ([26]) for large-scale datasets. We do not report running times in our experiments, since feature selection is an offline-task. We implemented all our algorithms in MATLAB R2013b on an Intel i-7 processor with 16GB RAM. BSS and leverage-score sampling are better than LPSVM and RRQR and comparable to RFE on 49 TechTC-300 datasets.

5.1 Other Feature Selection Methods

In this section, we describe other feature-selection methods with which we compare BSS and Leverage-score sampling.

Rank-Revealing QR Factorization (RRQR): Within the numerical linear algebra community, subset selection algorithms use the so-called Rank Revealing QR (RRQR) factorization. Let **A** be a $n \times d$ matrix with (n < d) and an integer k (k < d) and assume partial QR factorizations of the form

$$\mathbf{AP} = \mathbf{Q} egin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{R}_{22} \end{pmatrix},$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix, $\mathbf{P} \in \mathbb{R}^{d \times d}$ is a permutation matrix, $\mathbf{R}_{11} \in \mathbb{R}^{k \times k}$, $\mathbf{R}_{12} \in \mathbb{R}^{k \times (d-k)}$, $\mathbf{R}_{22} \in \mathbb{R}^{(d-k) \times (d-k)}$ The above factorization is called a RRQR factorization if $\sigma_{min}\left(\mathbf{R}_{11}\right) \geq \sigma_{k}\left(\mathbf{A}\right)/p(k,d)$, $\sigma_{max}\left(\mathbf{R}_{22}\right) \leq \sigma_{min}(\mathbf{A})p(k,d)$, where p(k,d) is a function bounded by a low-degree polynomial in k and d. The important columns are given by $\mathbf{A}_{1} = \mathbf{Q}\begin{pmatrix}\mathbf{R}_{11}\\\mathbf{0}\end{pmatrix}$ and $\sigma_{i}\left(\mathbf{A}_{1}\right) = \sigma_{i}\left(\mathbf{R}_{11}\right)$ with $1 \leq i \leq k$. We perform feature selection

Random Feature Selection: We select features uniformly at random without replacement which serves as a baseline method. To get around the randomness, we repeat the sampling process five times.

using RRQR by picking the important columns which preserve the rank of the matrix.

Recursive Feature Elimination: Recursive Feature Elimination (RFE), [3] tries to find the best subset of features which leads to the largest margin of class separation using SVM. At each iteration, the algorithm greedily removes the feature that decreases the margin the least, until the required number of features remain. At each step, it computes the weight vector and removes the feature with smallest weight. RFE is computationally expensive for high-dimensional datasets. Therefore, at each iteration, multiple features are removed to avoid the computational bottleneck.

LPSVM: The feature selection problem for SVM can be formulated in the form of a linear program. LPSVM [4] uses a fast Newton method to solve this problem and obtains a sparse solution of the weight vector, which is used to select the features.

5.2 BSS Implementation Issues

At every iteration, there can be multiple columns which satisfy the condition $\mathcal{U}(\mathbf{v}_i, \delta_U, \mathbf{A}_{\tau}, U_{\tau}) \leq \mathcal{L}(\mathbf{v}_i, \delta_L, \mathbf{A}_{\tau}, L_{\tau})$. [7] suggest picking any column which satisfies this constraint. Selecting a column naively leaves out important features required for classification. Therefore, we choose the column \mathbf{v}_i which has not been selected in previous iterations and whose Euclidean-norm is highest among the candidate set. Columns with zero Euclidean norm never get selected by the algorithm.

In our implementation, we do not use the data center as one of the inputs (since computing the center involves solving a quadratic program).

Table 1: Most frequently selected features using the synthetic dataset.

	$r_1 = 30$		$r_1 = 40$		
	k = 40	k = 50	k = 40	k = 50	
BSS	40 , 39, 34, 36, 37	50 , 49, 48, 47, 44	40 , 39, 34, 37, 36	50 , 49, 48, 47, 44	
Lvg	40 , 39, 37, 36, 34	50 , 49, 48, 47, 46	40 , 39, 37, 31, 32	50 , 49, 48, 31, 47	
RFE	40 , 39, 38, 37, 36	50 , 49, 48, 47, 46	40 , 39, 38, 37, 36	50 , 49, 48, 47, 46	
LPSVM	40 , 39, 38, 37, 34	50 , 49, 48, 43, 40	40 , 39, 38, 37, 34	50 , 49, 48, 43, 40	
RRQR	40 , 30, 29, 28, 27	50 , 30, 29, 28, 27	40 , 39, 38, 37, 36	50 , 40, 39, 38, 37	

5.3 Experiments on Supervised Feature Selection

Synthetic Data: We generate synthetic data as described in [27], where we control the number of relevant features in the dataset. The dataset has n data-points and d features. The class label y_i of each data-point was randomly chosen to be 1 or -1 with equal probability. The first k features of each data-point \mathbf{x}_i are the relevant features and are drawn from $y_i \mathcal{N}\left(-j,1\right)$ distribution, where $\mathcal{N}\left(\mu,\sigma^2\right)$ is a random normal distribution with mean μ and variance σ^2 and j varies from 1 to k. The remaining (d-k) features are chosen from a $\mathcal{N}(0,1)$ distribution and are noisy features. By construction, among the first k features, the kth feature has the most discriminatory power, followed by (k-1)th feature and so on. We set n to 200 and d to 1000. We set k to 40 and 50 and ran two sets of experiments. We set the value of r_1 , i.e. the number of features selected, to 30 and 40 for all experiments. We performed ten-fold cross-validation and repeated it ten times. We used LIBSVM with default settings and set C=1. We compared with the other methods. The mean out-of-sample error was 0 for all methods for both k=40 and k=50. Table 1 shows the set of five most frequently selected features by the different methods for one such synthetic dataset. The top features picked up by the different methods are the relevant features by construction and also have good discriminatory power. This shows that supervised BSS and leverage-score sampling are as good as any other method in terms of feature selection. We repeated our experiments on ten different synthetic datasets and each time, the five most frequently selected features were from the set of relevant features. Thus, by selecting only 3% -4% of all features, we show that we are able to obtain the most discriminatory features along with good out-of-sample error using BSS and leverage-score sampling.

Table 2: A subset of the TechTC matrices of our study

	id1	id2				
(i)	Arts: Music: Styles: Opera	Arts: Education: Language: Reading Instructions				
(ii)	Arts: Music: Styles: Opera	US Navy: Decommisioned Attack Submarines				
(iii)	US: Michigan: Travel & Tourism	Recreation:Sailing Clubs: UK				
(iv)	US: Michigan: Travel & Tourism	Science: Chemistry: Analytical: Products				
(v)	US: Colorado: Localities: Boulder	Europe: Ireland: Dublin: Localities				

TechTC-300: For our first real dataset, we use 49 datasets of TechTC-300 ([28]) which contain binary classification tasks. Each data matrix consists of 150-280 documents (the rows of the data matrix), and each document is described with respect to 10,000-50,000 words (features are columns of the matrix). We removed all words with at most four letters

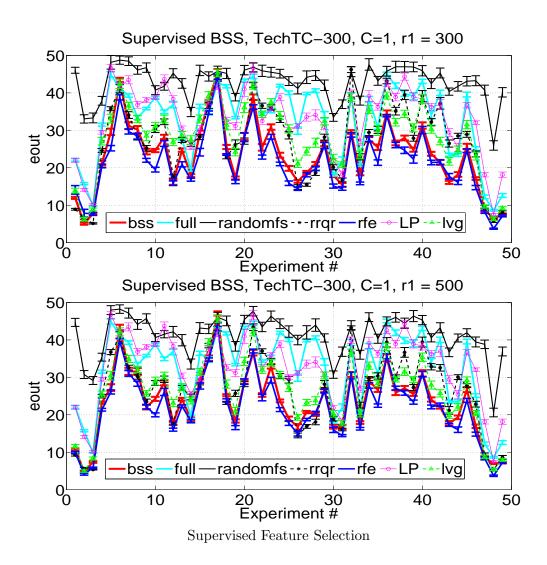


Figure 1: Plots of out-of-sample error of Supervised BSS and leverage-score compared with other methods for 49 TechTC-300 documents averaged over ten ten-fold cross validation experiments. Vertical bars represent standard deviation.

Table 3: Frequently occurring terms of the five TechTC-300 datasets of Table 2 selected by supervised BSS and Leverage-score sampling.

	BSS	Leverage-score Sampling		
(i)	reading, education, opera, frame	reading, opera, frame, spacer		
(ii)	submarine, hullnumber, opera, tickets	hullnumber, opera, music, tickets		
(iii)	michigan, vacation, yacht, sailing	sailing, yacht, michigan, vacation		
(iv)	chemical, michigan, environmental, asbestos	travel, vacation, michigan, services, environmental		
(v)	ireland, dublin, swords, boulder, colorado	ireland, boulder, swords, school, grade		

from the datasets. We set the parameter C=1 in LIBSVM and used default settings. We tried different values of C for the full-dataset and the out-of-sample error averaged over

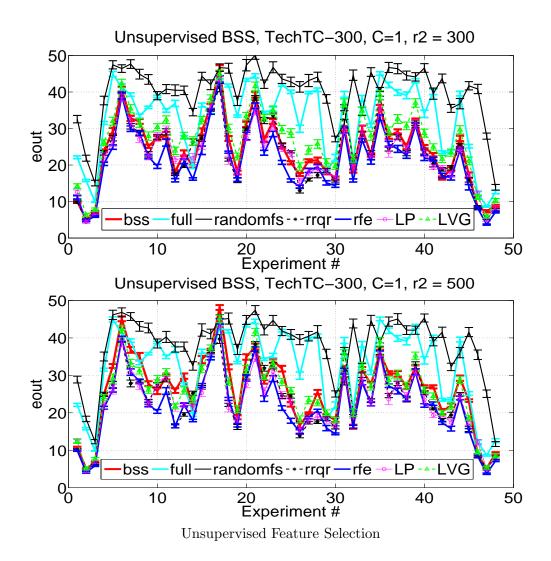


Figure 2: Plots of out-of-sample error of Unsupervised BSS and leverage-score compared with other methods for 49 TechTC-300 documents averaged over ten ten-fold cross validation experiments. Vertical bars represent standard deviation.

Table 4: Results of Approximate BSS. CCAT (train / test): (23149 / 781265), d=47236. Mean and standard deviation (in parenthesis) of out-of-sample error. Eout of full-data is 8.66 ± 0.54 .

Eout	r_1	BSS $(t = 128)$	BSS $(t = 256)$	RRQR	RFE	LPSVM
CCAT	1024	$10.53 \ (0.59)$	$10.35 \ (0.64)$	9.97(0.62)	8.92(0.57)	9.97(0.55)
CCAT	2048	$11.13\ (0.66)$	$10.63 \ (0.62)$	10.04 (0.66)	8.56(0.54)	9.97(0.55)

49 TechTC-300 documents did not change much, so we report the results of C=1. We set the number of features to 300, 400 and 500. Figs 1 and 3 show the out-of-sample error for the 49 datasets for r1=300, 400 and 500. For the supervised feature selection, BSS is comparable to RFE and leverage-score sampling and better than RRQR, LPSVM, full-data and uniform sampling in terms of out-of-sample error. For LPSVM, the number

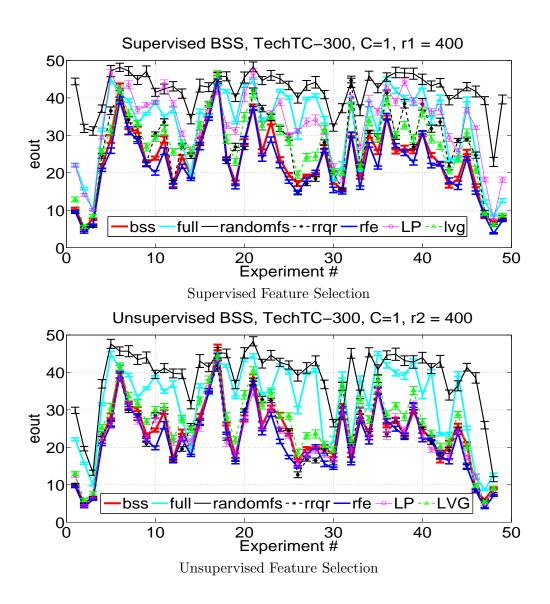


Figure 3: Plots of out-of-sample error of Supervised and Unsupervised BSS and leverage-score compared with other methods for 49 TechTC-300 documents averaged over ten ten-fold cross validation experiments. Vertical bars represent standard deviation.

of selected features averaged over 49 datasets was greater than 500, but it performed worse than BSS and leverage-score sampling. Leverage-score sampling is comparable to BSS and better than RRQR, LPSVM, full-data and uniform sampling and slightly worse than RFE. We list the most frequently occurring words selected by supervised BSS and leverage-score for the $r_1 = 300$ case for five TechTC-300 datasets over 100 training sets. Table 2 shows the names of the five TechTC-300 document-term matrices. The words shown in Table 3 were selected in all cross-validation experiments for these five datasets. The words are closely related to the categories to which the documents belong, which shows that BSS and Leverage-score sampling select important features from the support-vector matrix. For example, for the document-pair (ii), where the documents belong to the category of "Arts:Music:Styles:Opera" and "US:Navy: Decommisioned Attack Submarines", the BSS

algorithm selects submarine, hullnumber, opera, tickets and Leverage-score sampling selects hullnumber, opera, music, tickets which are closely related to the two classes. Thus, we see that using only 2%-4% of all features we are able to obtain good out-of-sample error.

5.4 Experiments on Unsupervised Feature Selection

For the unsupervised feature selection case, we performed experiments on the same 49 TechTC-300 datasets and set r_2 to 300, 400 and 500. We include the results for $r_2 = 300$ and $r_2 = 500$ in Figs 2 and 3. For LPSVM, the number of selected features averaged over 49 datasets was close to 300. In the unsupervised case, BSS and leverage-score sampling are comparable to each other and also comparable to the other methods RRQR, LPSVM and RFE. These methods are better than random feature selection and full-data without feature selection. This shows that unsupervised BSS and leverage-score sampling are competitive feature selection algorithms.

Supervised feature selection is comparable to unsupervised feature selection for BSS, Leverage-score sampling and RFE, while unsupervised RRQR and LPSVM are better than their supervised versions. Running BSS (or leverage-score sampling) on the support-vector set is equivalent to running BSS (or leverage-score sampling) on the training data. However, RRQR and LPSVM are primarily used as unsupervised feature selection techniques and so they perform well in that setting. RFE is a heuristic based on SVM and running RFE on the support-vectors is equivalent to running RFE on the training data.

5.5 Approximate BSS

We describe a heuristic to make supervised BSS scalable to large-scale datasets. For datasets with large number of support vectors, we premultiply the support vector matrix X with a random gaussian matrix $\mathbf{G} \in \mathbb{R}^{t \times p}$ to obtain $\hat{\mathbf{X}} = \mathbf{G}\mathbf{X}$ and then use BSS to select features from the right singular vectors of $\hat{\mathbf{X}}$. The right singular vectors of $\hat{\mathbf{X}}$ closely approximates the right singular vectors of X. Hence the columns selected from X will be approximately same as the columns selected from X. We include the algorithm as Algorithm 2. We performed experiments on a subset of Reuters Corpus dataset, namely reuters-CCAT, which contains binary classification task. We used the L2-regularized L2-loss SVM formulation in the dual form in LIBLINEAR and set the value of C to 10. We experimented with different values of C on the full-dataset, and since there was small change in classification accuracy among the different values of C, we chose C=10 for our experiments. We pre-multiplied the support vector matrix with a random gaussian matrix of size $t \times p$, where p is the number of support vectors and t was set to 128 and 256. We repeated our experiments five times using five different random gaussian matrices to get around the randomness. We set the value of r_1 in BSS to 1024 and 2048. LPSVM selects 1898 features for CCAT. Table 4 shows the results of our experiments. We observe that the out-of-sample error using approx-BSS is close to that of RRQR and comparable to RFE, LPSVM and full-data. The out-of-sample error of approx-BSS decreases with an increase in the value of t. This shows that we get a good approximation of the right singular vectors of the support vector matrix with an increase in number of projections.

Input: Support vector matrix $\mathbf{X} \in \mathbb{R}^{p \times d}$, t, r. Output: Matrices $\mathbf{S} \in \mathbb{R}^{d \times r}$, $\mathbf{D} \in \mathbb{R}^{r \times r}$.

- 1. Generate a random Gaussian matrix, $\mathbf{G} \in \mathbb{R}^{t \times p}$.
- 2. Compute $\hat{\mathbf{X}} = \mathbf{G}\mathbf{X}$.
- 3. Compute right singular vectors \mathbf{V} of $\hat{\mathbf{X}}$ using SVD.
- 4. Run Algorithm 1 using V and r as inputs and get matrices S and D as outputs.
- 5. Return **S** and **D**.

Algorithm 2: Approximate BSS

6 Conclusions

Our simple method of extending an unsupervised feature selection method into a supervised one for SVM not only has a provable guarantee, but also works well empirically: BSS and leverage-score sampling are comparable and often better than prior state-of-the-art feature selection methods for SVM, and those methods don't come with guarantees.

Our supervised sparsification algorithms only preserve the margin for the support vectors in the feature space. We do not make any claims about the margin of the full data in the feature space constructed from the support vectors. This appears challenging and it would be interesting to see progress made in this direction: can one choose O(#support vectors) features for the full data set and obtain provable guarantees on the margin and data radius? There have been recent advances in approximate leverage-scores for large-scale datasets. A possible future work in this direction would be to see if those algorithms indeed work well with SVMs.

7 Acknowledgements

PD and SP are supported by NSF IIS-1447283 and IIS-1319280 respectively.

References

- [1] A. Dasgupta, P. Drineas, B. Harb, V. Josifovski, and M.W. Mahoney. Feature selection methods for text classification. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 230–239, 2007.
- [2] N. Cristianini and J. Shawe-Taylor. Support Vector Machines and other kernel-based learning methods. Cambridge University Press, 2000.
- [3] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik. Gene selection for cancer classification using support vector machines. *Machine Learning*, 46(1-3):389–422, 2002.
- [4] M. Glenn Fung and O.L. Mangasarian. A feature selection newton method for support vector machine classification. *Comput. Optim. Appl.*, 28(2):185–202, 2004.

- [5] V.N. Vapnik. Statistical Learning Theory. Theory of Probability and its Applications, 16:264–280, 1998.
- [6] V.N. Vapnik and A. Chervonenkis. On the Uniform Convergence of Relative Frequencies of Events to their Probabilities. *Theory of Probability and its Applications*, 16:264–280, 1971.
- [7] J.D. Batson, D.A. Spielman, and N. Srivastava. Twice-ramanujan sparsifiers. In *Proceedings* of the 41st annual ACM STOC, pages 255–262, 2009.
- [8] A. Rakotomamonjy. Variable selection using svm based criteria. JMLR, 3:1357–1370, 2003.
- [9] J. Weston, S. Mukherjee, O. Chapelle, M. Pontil, T. Poggio, and V. Vapnik. Feature selection for svms. In *NIPS*, volume 12, pages 668–674, 2000.
- [10] J. Weston, A. Elisseeff, B. Schölkopf, and M. Tipping. Use of the zero norm with linear models and kernel methods. *JMLR*, 3:1439–1461, 2003.
- [11] M. Tan, L. Wang, and I.W. Tsang. Learning sparse sym for feature selection on very high dimensional datasets. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pages 1047–1054, 2010.
- [12] H. Do, A. Kalousis, and M. Hilario. Feature weighting using margin and radius based error bound optimization in syms. In *European Conference on Machine Learning (ECML)*, pages 315–329, 2009b.
- [13] A. Kalousis and H.T. Do. Convex formulations of radius-margin based support vector machines. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, pages 169–177, 2013.
- [14] L. Wang, J. Zhu, and H. Zou. The doubly regularized support vector machine. Statistica Sinica, 16(2):589, 2006.
- [15] L. Wang, J. Zhu, and H. Zou. Hybrid huberized support vector machines for microarray classification and gene selection. *Bioinformatics*, 24(3):412–419, 2008.
- [16] G. Ye, Y. Chen, and X. Xie. Efficient variable selection in support vector machines via the alternating direction method of multipliers. In AISTATS, pages 832–840, 2011.
- [17] R. Gilad-Bachrach, A. Navot, and N. Tishby. Margin based feature selection-theory and algorithms. In *Proceedings of the twenty-first international conference on Machine learning (ICML)*, page 43, 2004.
- [18] C. Park, K-R. Kim, R. Myung, and J-Y. Koo. Oracle properties of scad-penalized support vector machine. *Journal of Statistical Planning and Inference*, 142(8):2257–2270, 2012.
- [19] S. Paul, C. Boutsidis, M. Magdon-Ismail, and P. Drineas. Random projections for support vector machines. In Sixteenth International Conference on Artificial Intelligence and Statistics (AISTATS), pages 498–506. JMLR W&CP 31, 2013.
- [20] S. Paul, C. Boutsidis, M. Magdon-Ismail, and P. Drineas. Random projections for linear support vector machines. *ACM Trans. Knowl. Discov. Data*, 8(4):22:1–22:25, 2014.
- [21] C. Boutsidis and M. Magdon-Ismail. Deterministic feature selection for k-means clustering. *IEEE Transactions on Information Theory*, 59(9):6099–6110, 2013.
- [22] C. Boutsidis, M. W. Mahoney, and P. Drineas. Unsupervised feature selection for the k-means clustering problem. *Advances in Neural Information Processing Systems (NIPS)*, 2009.
- [23] S. Paul and P. Drineas. Deterministic feature selection for regularized least squares classification. In *Machine Learning and Knowledge Discovery in Databases (ECML-PKDD)*, volume 8725 of *LNCS*, pages 533–548, 2014.
- [24] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li. Rcv1: A new benchmark collection for text categorization research. *JMLR*, pages 361–397, 2004.

- [25] C-C. Chang and C-J. Lin. Libsvm: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1-27:27, 2011. Software available at http://www.csie.ntu.edu.tw/~cjlin/libsvm.
- [26] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. Liblinear: A library for large linear classification. *JMLR*, pages 1871 –1874, 2008.
- [27] C. Bhattacharyya. Second order cone programming formulations for feature selection. *JMLR*, 5:1417–1433, 2004.
- [28] D. Davidov, E. Gabrilovich, and S. Markovitch. Parameterized generation of labeled datasets for text categorization based on a hierarchical directory. In *Proceedings of the 27th Annual International ACM SIGIR Conference*, pages 250–257, 2004. http://techtc.cs.technion.ac.il/techtc300/techtc300.html.