# Random Projections for Linear Support Vector Machines

SAURABH PAUL, Rensselaer Polytechnic Institute
CHRISTOS BOUTSIDIS, IBM T.J. Watson Research Center
MALIK MAGDON-ISMAIL, Rensselaer Polytechnic Institute
PETROS DRINEAS, Rensselaer Polytechnic Institute

Let $\mathbf{X}$ be a data matrix of rank $\rho$, whose rows represent $n$ points in $d$-dimensional space. The linear support vector machine constructs a hyperplane separator that maximizes the 1-norm soft margin. We develop a new oblivious dimension reduction technique which is precomputed and can be applied to any input matrix $\mathbf{X}$. We prove that, with high probability, the margin and minimum enclosing ball in the feature space are preserved to within $\epsilon$-relative error, ensuring comparable generalization as in the original space in the case of classification. For regression, we show that the margin is preserved to $\epsilon$-relative error with high probability. We present extensive experiments with real and synthetic data to support our theory.

## 1. INTRODUCTION

Support Vector Machines (SVM) [Cristianini and Shawe-Taylor 2000] are extremely popular in machine learning today. They have been used in both classification and regression. For classification, the training data set consists of $n$ 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 form of the SVM learning problem is to construct 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" 1-norm margin is maximized. The dual lagrangian formulation of the classification problem leads to the following quadratic program:

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

$$\text{subject to} \sum_{i=1}^{n} y_i \alpha_i = 0, \tag{1}$$

$$0 \le \alpha_i \le C, \quad i = 1 \dots n.$$

In the above formulation, the unknown lagrange multipliers $\{\alpha_i\}_{i=1}^{n}$ are constrained to lie inside the "box constraint" $[0, C]^n$, where $C$ is part of the input. In order to measure the out-of-sample performance of the SVM classifier, we can use the VC-dimension of *fat*-separators. Assuming that the data lie in a ball of radius $B$, and that the hypothesis set consists of hyperplanes of width $\gamma$ (corresponding to the margin), then the $VC$-dimension of this hypothesis set is $O(B^2/\gamma^2)$ [Vapnik 1998]. Now, given

the in-sample error, we can obtain a bound for the out-of-sample error, which is monotonic in the VC-dimension [Vapnik and Chervonenkis 1971].

Analogous to the 1-norm soft margin formulation for SVM classification, we have a similar formulation for regression called the linear $\varepsilon$-insensitive loss SVM [Cristianini and Shawe-Taylor 2000]. The dual problem for $\varepsilon$-insensitive loss SVM regression is formulated as:

$$\max \sum_{i=1}^{n} \alpha_i y_i - \varepsilon \sum_{i=1}^{n} |\alpha_i| - \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{subject to} \sum_{i=1}^{n} \alpha_i = 0,$$

$$- C \leq \alpha_i \leq C, \quad i = 1 \ldots n. \tag{2}$$

Here, $\{\alpha\}_{i=1}^{n}$ are the Lagrange multipliers and they lie in the interval $[-C, C]^n$.

Intuitively, if one can preserve the subspace geometry, then one should be able to preserve the performance of a distance-based algorithm. We construct dimension reduction matrices $\mathbf{R} \in \mathbb{R}^{d \times r}$ which produce $r$-dimensional feature vectors $\tilde{\mathbf{x}}_i = \mathbf{R}^T \mathbf{x}_i$; the matrices $\mathbf{R}$ *do not* depend on the data. We show that for the data in the dimension-reduced space, the margin of separability and the minimum enclosing ball radius are preserved, since the subspace geometry is preserved. So, an SVM with an appropriate structure defined by the margin (width) of the hyperplanes [Vapnik and Chervonenkis 1971] will have comparable VC-dimension and, thus, generalization error. This is true for classification. The $\varepsilon$-insensitive loss SVM regression problem is an unbounded problem and as such, we are not able to infer anything related to the generalization error bounds: we can only infer the preservation of margin.

### 1.1. Notation and SVM Basics

$\mathbf{A}, \mathbf{B}, \ldots$ denote matrices and $\boldsymbol{\alpha}, , \ldots$ 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 $\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{n \times \rho}$ is an orthogonal matrix containing the left singular vectors, $\boldsymbol{\Sigma} \in \mathbb{R}^{\rho \times \rho}$ 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 $\|\mathbf{A}\|_2 = \sigma_1$.

*1.1.1. SVM Classification.* Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the matrix whose rows are the vectors $\mathbf{x}_i^T$, $\mathbf{Y} \in \mathbb{R}^{n \times n}$ be the diagonal matrix with entries $\mathbf{Y}_{ii} = y_i$, and $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \ldots, \alpha_n] \in \mathbb{R}^n$ be the vector of lagrange multipliers to be determined by solving eqn. èqn:svm1. The SVM optimization problem is

$$\max_{\boldsymbol{\alpha}} \mathbf{1}^T \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Y} \mathbf{X} \mathbf{X}^T \mathbf{Y} \boldsymbol{\alpha}$$

$$\text{subject to } \mathbf{1}^T \mathbf{Y} \boldsymbol{\alpha} = 0; \quad \text{and} \quad \mathbf{0} \leq \boldsymbol{\alpha} \leq \mathbf{C}. \tag{3}$$

(In the above, $\mathbf{1}$, $\mathbf{0}$, $\mathbf{C}$ are vectors with the implied constant entry.) Let $\boldsymbol{\alpha}^*$ be an optimal solution of the above problem. The optimal separating hyperplane is given by $\mathbf{w}^* = \mathbf{X}^T \mathbf{Y} \boldsymbol{\alpha}^* = \sum_{i=1}^{n} y_i \alpha_i^* \mathbf{x}_i$, and the points $\mathbf{x}_i$ for which $\alpha_i^* > 0$, i.e., the points which appear in the expansion $\mathbf{w}^*$, are the support vectors. The geometric margin, $\gamma^*$, of this canonical optimal hyperplane is $\gamma^* = 1 / \|\mathbf{w}^*\|_2$, where $\|\mathbf{w}^*\|_2^2 = \sum_{i=1}^{n} \alpha_i^*$. The data radius is $B = \min_{\mathbf{x}^*} \max_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{x}^*\|_2$. It is this $\gamma^*$ and $B$ that factor into the generalization performance of the SVM through the ratio $B/\gamma^*$. It is worth noting, that our results hold for the separable case as well, which amounts to setting $C$ to a large value.

*1.1.2. SVM Regression.* Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the matrix whose rows are the vectors $\mathbf{x}_i^T$, $\mathbf{y}$ be the n-dimensional vector with the target entries, and $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \ldots, \alpha_n] \in \mathbb{R}^n$ be the vector of lagrange multipliers to be determined by solving eqn. èqn:svm4. The SVM optimization problem is

$$\max_{\boldsymbol{\alpha}} \mathbf{y}^T \boldsymbol{\alpha} - \varepsilon \mathbf{1}^T \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{X} \mathbf{X}^T \boldsymbol{\alpha}$$

$$\text{subject to } \mathbf{1}^T \boldsymbol{\alpha} = 0; \quad \text{and} \quad - \mathbf{C} \leq \boldsymbol{\alpha} \leq \mathbf{C}. \tag{4}$$

$\mathbf{1}$, $\mathbf{0}$, $\mathbf{C}$ are vectors with the implied constant entry. Let $\boldsymbol{\alpha}^*$ be an optimal solution of the above problem. The optimal separating hyperplane is given by $\mathbf{w}^* = \boldsymbol{\alpha}^{*T} \mathbf{X} = \sum_{i=1}^{n} \alpha_i^* \mathbf{x}_i$ and the points $\mathbf{x}_i$ for which $\alpha_i^* > 0$, i.e., the points which appear in the expansion $\mathbf{w}^*$, are the support vectors. The geometric margin for regression is defined in the same way as it was done for classification.

## 1.2. Dimension Reduction

Our goal is to study how the SVM performs under (linear) dimensionality reduction transformations in the feature space. Let $\mathbf{R} \in \mathbb{R}^{d \times r}$ be the dimension reduction matrix that reduces the dimensionality of the input from $d$ to $r \ll d$. We will choose $\mathbf{R}$ to be a random projection matrix (see Section 2). The transformed dataset into $r$ dimensions is given by $\tilde{\mathbf{X}} = \mathbf{XR}$, and the SVM optimization problem for classification becomes

$$\max_{\tilde{\alpha}} \mathbf{1}^T \tilde{\alpha} - \frac{1}{2} \tilde{\alpha}^T \mathbf{YXRR}^T \mathbf{X}^T \mathbf{Y} \tilde{\alpha},$$
$$\text{subject to } \mathbf{1}^T \mathbf{Y} \tilde{\alpha} = 0, \qquad \text{and} \qquad \mathbf{0} \leq \tilde{\alpha} \leq C. \tag{5}$$

For regression, the SVM optimization problem becomes

$$\max_{\tilde{\alpha}} \mathbf{y}^T \tilde{\alpha} - \varepsilon \mathbf{1}^T \tilde{\alpha} - \frac{1}{2} \tilde{\alpha}^T \mathbf{XRR}^T \mathbf{X}^T \tilde{\alpha}$$
$$\text{subject to } \mathbf{1}^T \tilde{\alpha} = 0; \qquad \text{and} \qquad -C \leq \tilde{\alpha} \leq C. \tag{6}$$

We will present a construction for $\mathbf{R}$ that leverages the fast Hadamard transform. The running time needed to apply this construction to the original data matrix is $O(nd \log r)$. Notice that while this running time is nearly linear on the size of the original data, it does not take advantage of any sparsity in the input. In order to address this deficiency, we leverage the recent work of Clarkson and Woodruff [2013], Meng and Mahoney [2013] and Nelson and Nguyen [2013], which proposes a construction for $\mathbf{R}$ that can be applied to $\mathbf{X}$ in $O\left(nnz(\mathbf{X}) + poly\left(n\epsilon^{-1}\right)\right)$ time; here $nnz(\mathbf{X})$ denotes the number of non-zero entries of $\mathbf{X}$ and $\rho$ is the rank of $\mathbf{X}$. To the best of our knowledge, this is the first independent implementation and evaluation of this potentially ground-breaking random projection technique (a few experimental results were presented by Clarkson and Woodruff [2013], Meng and Mahoney [2013] and Nelson and Nguyen [2013]). All constructions for $\mathbf{R}$ are oblivious of the data and hence they can be precomputed. Also, the generalization bounds that depend on the final margin and radius of the data will continue to hold for classification, while the bound on the margin holds for regression.

The pratical intent of using linear SVM after random projections is to reduce computational complexity of training SVM and memory. For large-scale datasets, that are too big to fit into memory (see Section 4.1.3 for details), random projections serve as a possible way to estimate out-of-sample error. Random projections reduce the computational complexity of training SVM which is evident from the experiments described in Section 4.

## 1.3. Our Contribution

Our main theoretical results are to show that by solving the SVM optimization problem in the projected space, we get relative-error generalization performance for SVM classification and that, we preserve the margin upto relative error for SVM regression. We briefly discuss the appropriate values of $r$, namely the dimensionality of the dimensionally-reduced problem. If $\mathbf{R}$ is the matrix of the randomized Hadamard transform (see Section 2 for details), then given $\epsilon \in (0, 1/2]$ and $\delta \in (0, 1)$ we set

$$r = O\left(\rho\epsilon^{-2} \cdot \log\left(\rho d\delta^{-1}\right) \cdot \log\left(\rho\epsilon^{-2}\delta^{-1} \log\left(\rho d\delta^{-1}\right)\right)\right). \tag{7}$$

The running time needed to apply the randomized Hadamard transform is $O(nd \log r)$. If $\mathbf{R}$ is constructed as described in Clarkson and Woodruff [2013], Meng and Mahoney [2013] and Nelson and Nguyen [2013], then given $\epsilon \in (0, 1)$ and $\delta \in (0, 1)$ we set

$$r = O\left(\rho\epsilon^{-4} \log\left(\rho/\delta\epsilon\right)\left(\rho + \log\left(1/\delta\epsilon\right)\right)\right). \tag{8}$$

The running time needed to apply this transform is $O\left(nnz(\mathbf{X}) + poly\left(n\epsilon^{-1}\right)\right)$. If $\mathbf{R}$ is a random sign-matrix, then given $\epsilon \in (0, 1/2]$ and we set

$$r = O\left(\rho\epsilon^{-2} \log \rho \log d\right). \tag{9}$$

The running time needed to apply this transform is equal to $O(ndr)$. Finally if $\mathbf{R}$ is a random gaussian matrix, then given $\epsilon \in (0, 1/2]$ and $\delta \in (0, 1)$ we set,

$$r = O\left(\rho\epsilon^{-2} \log\left(\rho/\delta\right)\right). \tag{10}$$

Our main theorem will be stated in terms of the randomized Hadamard Transform, but similar statements can be obtained for the other three transforms.

THEOREM 1.1. *Let $\epsilon \in (0, 1/2]$ be an accuracy paramater and let $\delta \in (0, 1)$ be a failure probability. Let $\mathbf{R} \in \mathbb{R}^{d \times r}$ be the matrix of the randomized Hadamard Transform, with $r$ as in eqn. (7). Let $\gamma^*$ and $\tilde{\gamma}^*$ be the margins obtained by solving the SVM problems using data matrices $\mathbf{X}$ and $\mathbf{XR}$ respectively (eqns. (3) and (5)). Let $B$ be the radius of the minimum ball enclosing all points in the full-dimensional space (rows of $\mathbf{X}$) and let $\tilde{B}$ be the radius of the ball enclosing all points in the dimensionally-reduced space (rows of $\mathbf{XR}$). Then, with probability at least $1 - 2\delta$,*

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

Similar theorems can be stated for the other two constructions of $\mathbf{R}$ by setting the value of $r$ as in eqns. (8), (9) and (10). For the case of SVM regression, we can only show that the margin is preserved up to relative error with probability at least $1 - \delta$, namely

$$\tilde{\gamma}^{*2} \geq (1 - \epsilon) \gamma^{*2}.$$

## 1.4. Prior work

The work most closely related to our results is that of Krishnan et al. [2008], which improved upon Balcazar et al. [2001]. Balcazar et al. [2001] and Balczar et al. [2002] used random sampling techniques for solving the SVM classification and $\varepsilon$-insensitive loss SVM regression problem respectively, but they were not able to implement their algorithms in practice. Krishnan et al. [2008] and Jethava et al. [2009] showed that by using sub-problems based on Gaussian random projections, one can obtain a solution to the SVM classification and regression problem with a margin that is relative-error close to the optimal. Their sampling complexity (the parameter $r$ in our parlance) depends on $B^4$, and, most importantly, on $1/\gamma^{*2}$. This bound is not directly comparable to our result, which only depends on the rank of the data manifold, and holds regardless of the margin of the original problem (which could be arbitrarily small). Our results dramatically improve the running time needed to apply the random projections; our running times are (theoretically) linear in the number of non-zero entries in $\mathbf{X}$, whereas Krishnan et al. [2008] necessitates $O(ndr)$ time to apply $\mathbf{R}$ on $\mathbf{X}$.

Blum [2006] showed relative error margin preservation for linearly separable data by angle preservation between points when using random orthogonal matrix, standard gaussian matrix and the random sign matrix. We show relative-error margin preservation for non-separable data and use methods that improve running time to compute random projections. Shi et al. [2012] establish the conditions under which margins are preserved after random projection and show that error free margins are preserved for both binary and multi-class problems if these conditions are met. They discuss the theory of margin and angle preservation after random projections using Gaussian matrices. They show that margin preservation is closely related to acute angle preservation and inner product preservation. Smaller acute angle leads to better preservation of the angle and the inner product. When the angle is well preserved, the margin is well-preserved too. There are two main differences between their result and ours. They show margin preservation to within additive error, whereas we give margin preservation to within relative error. This is a big difference especially when the margin is small. Moreover, they analyze only the separable case. We analyze the general non-separable dual problem and give a result in terms of the norm of the weight vector. For the separable case, the norm of the weight vector directly relates to the margin. For the non-separable case, one has to analyze the actual quadratic program, and our result essentially claims that the solution in the transformed space will have comparably regularized weights as the solution in the original space.

Shi et al. [2009] used hash kernels which approximately preserved inner product to design a biased approximation of the kernel matrix. The hash kernels can be computed in the number of non-zero terms of a data matrix similar to the method of Clarkson and Woodruff [2013], Meng and Mahoney [2013] and Nelson and Nguyen [2013] that we employed. Shi et al. [2009] used random sign matrices to compute random projections which typically increase the number of non-zero terms of the data matrix. However, the method of Clarkson and Woodruff [2013], Meng and Mahoney [2013] and Nelson and Nguyen [2013] takes advantage of input sparsity. Shi et al. [2009] showed that their generalization bounds on the hash kernel and the original kernel differed by the inverse of the product of the margin and number of datapoints. For smaller margins, this difference will be high. Our generalization bounds are independent of the original margin and hold for arbitrarily small margins.

Zhang et al. [2013] developed algorithms to accurately recover the optimal solution to the original SVM optimization problem using a Gaussian random projection. They compute the dual solution provided that the data matrix has low rank. This is different from our work since we analyze the ratio of radius of the minimum enclosing ball to the margin using random projections and do not try to recover the solution.

Finally, it is worth noting that random projection techniques have been applied extensively in the compressed sensing literature, and our theorems have the same flavor to a number of results in that area. However, to the best of our knowledge, the compressed sensing literature has not investigated the 1-norm soft-margin SVM optimization problem.

## 2. RANDOM PROJECTION MATRICES

Random projections are extremely popular techniques in order to deal with the curse-of-dimensionality. Let the data matrix be $\mathbf{X} \in \mathbb{R}^{n \times d}$ ($n$ data points in $\mathbb{R}^d$) and let $\mathbf{R} \in \mathbb{R}^{d \times r}$ (with $r \ll d$) be a random projection matrix. Then, the projected data matrix is $\tilde{\mathbf{X}} = \mathbf{XR} \in \mathbb{R}^{n \times r}$ ($n$ points in $\mathbb{R}^r$). If $\mathbf{R}$ is carefully chosen, then all pairwise Euclidean distances are preserved with high probability. Thus, the geometry of the set of points in preserved, and it is reasonable to hope that an optimization objective such as the one that appears in SVMs will be only mildly perturbed.

There are many possible constructions for the matrix $\mathbf{R}$ that preserve pairwise distances. The most common one is a matrix $\mathbf{R}$ whose entries are i.i.d. standard Gaussian random variables [Dasgupta and Gupta 2003; Indyk and Motwani 1998] –**RG** for short. Achlioptas [2003] argued that the random sign matrix – **RS** for short – e.g., a matrix whose entries are set to $+1$ or $-1$ with equal probability, also works. Li et al. [2006] used the sparse random projection matrix whose entries were set to $+1$ or $-1$ with probability $1/2\sqrt{d}$ and $0$ with probability $(1 - 1/\sqrt{d})$. These constructions take $O(ndr)$ time to compute $\tilde{\mathbf{X}}$.

More recently, faster methods of constructing random projections have been developed, using, for example, the Fast Hadamard Transform [Ailon and Chazelle 2006] – **FHT** for short. The Hadamard-Walsh matrix for any $d$ that is a power of two is defined as

$$\mathbf{H}_d = \begin{bmatrix} \mathbf{H}_{d/2} & \mathbf{H}_{d/2} \\ \mathbf{H}_{d/2} & -\mathbf{H}_{d/2} \end{bmatrix} \in \mathbb{R}^{d \times d},$$

with $\mathbf{H}_1 = +1$. The normalized Hadamard-Walsh matrix is $\sqrt{\frac{1}{d}}\mathbf{H}_d$, which we simply denote by $\mathbf{H}$. We set:

$$\mathbf{R}_{\text{SRHT}} = \sqrt{\frac{d}{r}}\mathbf{DHS}, \tag{11}$$

a rescaled product of three matrices. $\mathbf{D} \in \mathbb{R}^{d \times d}$ is a random diagonal matrix with $\mathbf{D}_{ii}$ equal to $\pm 1$ with probability $\frac{1}{2}$. $\mathbf{H} \in \mathbb{R}^{d \times d}$ is the normalized Hadamard transform matrix. $\mathbf{S} \in \mathbb{R}^{d \times r}$ is a random *sampling matrix* which randomly samples columns of $\mathbf{DH}$; specifically, each of the $r$ columns of $\mathbf{S}$ is independent and selected uniformly at random (with replacement) from the columns of $\mathbf{I}_d$, the identity matrix. This construction assumes that $d$ is a power of two. If not, we just pad $\mathbf{X}$ with columns of zeros (affecting run times by at most a factor of two). The important property of this transform is that the projected features $\tilde{\mathbf{X}} = \mathbf{XR}$ can be computed efficiently in $O(nd \log r)$ time (see Theorem 2.1 of Ailon and Liberty [2008] for details). An important property of $\mathbf{R}$ (that follows from prior work) is that it preserves orthogonality.

While the randomized Hadamard transform is a major improvement over prior work, it does not take advantage of any sparsity in the input matrix. To fix this, very recent work [Clarkson and Woodruff 2013] shows that carefully constructed random projection matrices can be applied in input sparsity time by making use of generalized sparse embedding matrices. Meng and Mahoney [2013], Nelson and Nguyen [2013] also use a similar construction which runs in input sparsity time. Here we describe the construction of Clarkson and Woodruff [2013]. To understand their construction of $\mathbf{R}$, assume that the rank of $\mathbf{X}$ is $\rho$ and let $r = O\left(\rho\epsilon^{-4}\log(\rho/\delta\epsilon)(\rho + \log(1/\epsilon\delta))\right)$. Then, let $a = \Theta\left(\epsilon^{-1}\log(\rho/\epsilon\delta)\right)$, $v = \Theta\left(\epsilon^{-1}\right)$, and let $q = O\left(\rho\epsilon^{-2}(\rho + \log(1/\epsilon\delta))\right)$ be an integer (by appropriately choosing the constants). The construction starts by letting $h : 1 \ldots d \to 1 \ldots q$ be a random hash function; then, for $i = 1 \ldots q$, let $a_i = |h^{-1}(i)|$ and let $d = \sum_{i=1}^{q} a_i$. The construction proceeds by creating $q$ independent matrices $\mathbf{B}_1 \ldots \mathbf{B}_q$, such that $\mathbf{B}_i \in \mathbb{R}^{va \times a_i}$. Each $\mathbf{B}_i$ is the concatenation (stacking the rows of matrices on top of each other) of the following matrices: $\sqrt{\frac{1}{a}}\mathbf{\Phi}_1\mathbf{D}_1 \ldots \sqrt{\frac{1}{a}}\mathbf{\Phi}_a\mathbf{D}_a$. The matrix $\mathbf{\Phi}_i\mathbf{D}_i \in \mathbb{R}^{v \times a_i}$ is defined as follows: for each $m \in \{1 \ldots v\}$, $h(m) = g'$, where $g'$ is selected from $\{1 \ldots a_i\}$ uniformly at random. $\mathbf{\Phi}_i$ is a $v \times a_i$ binary matrix with $\mathbf{\Phi}_{\mathbf{h(m)},\mathbf{m}} = 1$ and all remaining entries set to zero. $\mathbf{D}$ is an $a_i \times a_i$ random diagonal matrix, with each diagonal entry independently set to be $+1$ or $-1$ with probability $1/2$. Finally, let $\mathbf{S}$ be the block diagonal matrix constructed by stacking the $\mathbf{B}_i$'s across its diagonal and let $\mathbf{P}$ be a $d \times d$ permutation matrix; then, $\mathbf{R} = (\mathbf{SP})^T$. The running time is $O\left(nnz(\mathbf{X}) + poly\left(n\epsilon^{-1}\right)\right)$. We will call the method of Clarkson and Woodruff [2013] to construct a sparse embedding matrix **CW**.

## 3. GEOMETRY OF SVM IS PRESERVED UNDER RANDOM PROJECTION

We now state and prove our main result, namely that solving the SVM optimization problem in the projected space results in comparable margin and data radius as in the original space. The following lemma will be crucial in our proof.

LEMMA 3.1. *Fix $\epsilon \in (0, \frac{1}{2}]$, $\delta \in (0, 1]$. Let $\mathbf{V} \in \mathbb{R}^{d \times \rho}$ be any matrix with orthonormal columns and let $\mathbf{R} = \mathbf{R}_{\text{SRHT}}$ as in eqn. (11), with $r = O(\rho \epsilon^{-2} \cdot \log(\rho d \delta^{-1}) \cdot \log(\rho \epsilon^{-2} \delta^{-1} \log(\rho d \delta^{-1})))$. Then, with probability at least $1 - \delta$,*
$$\|\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}\|_2 \le \epsilon.$$

PROOF. Consider the matrix $\mathbf{V}^T \mathbf{R} = \mathbf{V}^T \mathbf{DHS}$. Using Lemma 3 of [Drineas et al. 2011],

$$\left\|(\mathbf{HDV})_{(i)}\right\|_2^2 \le \frac{2\rho \ln(40 d\rho)}{d} \Rightarrow (2 \ln(40 d\rho))^{-1} \frac{\left\|(\mathbf{HDV})_{(i)}\right\|_2^2}{\rho} \le \frac{1}{d}$$

holds for all $i = 1, \ldots, d$ with probability at least $1 - \delta$. In the above, the notation $\mathbf{A}_{(i)}$ denotes the $i$-th row of $\mathbf{A}$ as a row vector. Applying Theorem 4 with $\beta = (2 \ln(40 d\rho))^{-1}$ ( [Drineas et al. 2011], Appendix) concludes the lemma. ☐

We now state two similar lemmas that cover two additional constructions for $\mathbf{R}$.

LEMMA 3.2. *Let $\epsilon \in (0, \frac{1}{2}]$ and let $\mathbf{V} \in \mathbb{R}^{d \times \rho}$ be any matrix with orthonormal columns. Let $\mathbf{R} \in \mathbb{R}^{d \times r}$ be a (rescaled) random sign matrix. If $r = O\left(\rho \epsilon^{-2} \log \rho \log d\right)$, then with probability at least $1 - 1/n$,*
$$\|\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}\|_2 \le \epsilon.$$

PROOF. The proof of this result is a simple application of Theorem 3.1(i) of Magen and Zouzias [2011]. ☐

LEMMA 3.3. *Let $\epsilon \in (0, \frac{1}{2}]$, $\delta \in (0, 1)$, and let $\mathbf{V} \in \mathbb{R}^{d \times \rho}$ be any matrix with orthonormal columns. Let $\mathbf{R} \in \mathbb{R}^{d \times r}$ be the CW random projection matrix (see Section 2) with $r = O\left(\rho \epsilon^{-4} \log(\rho/\delta\epsilon)(\rho + \log(1/\epsilon\delta))\right)$. Then, with probability at least $1 - \delta$,*
$$\|\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}\|_2 \le \epsilon.$$

PROOF. The proof of this result follows from Theorem 1 of Meng and Mahoney [2013]. ☐

LEMMA 3.4. *Let $\epsilon \in (0, \frac{1}{2}]$, $\delta \in (0, 1)$, and let $\mathbf{V} \in \mathbb{R}^{d \times \rho}$ be any matrix with orthonormal columns. Let $\mathbf{R} \in \mathbb{R}^{d \times r}$ be the Gaussian random projection matrix with $r = O\left(\rho \epsilon^{-2} \log(\rho/\delta)\right)$. Then with probability at least $1 - \delta$,*
$$\|\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}\|_2 \le \epsilon.$$

PROOF. The proof of this result follows from Corollary 6 of Zhang et al. [2013]. ☐

Lemma 3.4 does not have the log factors as in Lemma 3.1, but Gaussian projections are slower since they require full matrix-matrix multiplications.

THEOREM 3.5. *Let $\epsilon$ be an accuracy parameter and let $\mathbf{R} \in \mathbb{R}^{d \times r}$ be a matrix satisfying $\|\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}\|_2 \le \epsilon$. Let $\gamma^*$ and $\tilde{\gamma}^*$ be the margins obtained by solving the SVM problems using data matrices $\mathbf{X}$ and $\mathbf{XR}$ respectively (eqns. (3) and (5)). Then,*
$$\tilde{\gamma}^{*2} \ge (1 - \epsilon) \cdot \gamma^{*2}.$$

PROOF. Let $\mathbf{E} = \mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}$, and $\boldsymbol{\alpha}^* = [\alpha_1^*, \alpha_2^*, \ldots, \alpha_n^*]^T \in \mathbb{R}^n$ be the vector achieving the optimal solution for the problem of eqn. (3) in Section 1. Then,

$$
\begin{aligned}
Z_{opt} &= \sum_{i=1}^n \alpha_i^* - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y} \mathbf{X} \mathbf{X}^T \mathbf{Y} \boldsymbol{\alpha}^* \\
&= \sum_{i=1}^n \alpha_i^* - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{Y} \boldsymbol{\alpha}^* \\
&= \sum_{i=1}^n \alpha_i^* - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{Y} \boldsymbol{\alpha}^* \\
&\quad - \frac{1}{2} \boldsymbol{\alpha}^{*T} \mathbf{Y} \mathbf{U} \boldsymbol{\Sigma} \mathbf{E} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{Y} \boldsymbol{\alpha}^*.
\end{aligned}
\tag{12}
$$

Let $\tilde{\boldsymbol{\alpha}}^* = [\tilde{\alpha}_1^*, \tilde{\alpha}_2^*, \ldots, \tilde{\alpha}_n^*]^T \in \mathbb{R}^n$ be the vector achieving the optimal solution for the dimensionally-reduced SVM problem of eqn. (5) using $\tilde{\mathbf{X}} = \mathbf{XR}$. Using the SVD of $\mathbf{X}$, we get

$$\tilde{Z}_{opt} = \sum_{i=1}^{n} \tilde{\alpha}_i^* - \frac{1}{2}\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma V}^T\mathbf{RR}^T\mathbf{V\Sigma U}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^*. \tag{13}$$

Since the constraints on $\boldsymbol{\alpha}^*, \tilde{\boldsymbol{\alpha}}^*$ do not depend on the data (see eqns. (3) and (5)), it is clear that $\tilde{\boldsymbol{\alpha}}^*$ is a feasible solution for the problem of eqn. (3). Thus, from the optimality of $\boldsymbol{\alpha}^*$, and using eqn. (13), it follows that

$$\begin{aligned}
Z_{opt} &= \sum_{i=1}^{n} \alpha_i^* - \frac{1}{2}\boldsymbol{\alpha}^{*T}\mathbf{YU\Sigma V}^T\mathbf{RR}^T\mathbf{V\Sigma U}^T\mathbf{Y}\boldsymbol{\alpha}^* \\
&\quad - \frac{1}{2}\boldsymbol{\alpha}^{*T}\mathbf{YU\Sigma E\Sigma U}^T\mathbf{Y}\boldsymbol{\alpha}^* \\
&\geq \sum_{i=1}^{n} \tilde{\alpha}_i^* - \frac{1}{2}\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma V}^T\mathbf{RR}^T\mathbf{V\Sigma U}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^* \\
&\quad - \frac{1}{2}\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma E\Sigma U}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^* \\
&= \tilde{Z}_{opt} - \frac{1}{2}\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma E\Sigma U}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^*. 
\end{aligned} \tag{14}$$

We now analyze the second term using standard sub-multiplicativity properties and $\mathbf{V}^T\mathbf{V} = \mathbf{I}$. Taking $\mathbf{Q} = \tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma}$

$$\begin{aligned}
\frac{1}{2}\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma E\Sigma U}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^* &\leq \frac{1}{2}\|\mathbf{Q}\|_2 \|E\|_2 \|\mathbf{Q}^T\|_2 \\
&= \frac{1}{2}\|\mathbf{E}\|_2 \|\mathbf{Q}\|_2^2 \\
&= \frac{1}{2}\|\mathbf{E}\|_2 \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma V}^T\right\|_2^2 \\
&= \frac{1}{2}\|\mathbf{E}\|_2 \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YX}\right\|_2^2. 
\end{aligned} \tag{15}$$

Combining eqns. (14) and (15), we get

$$Z_{opt} \geq \tilde{Z}_{opt} - \frac{1}{2}\|\mathbf{E}\|_2 \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YX}\right\|_2^2. \tag{16}$$

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

$$\begin{aligned}
&\left|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YXRR}^T\mathbf{X}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^* - \tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YXX}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^*\right| \\
&= \left|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma}\left(\mathbf{V}^T\mathbf{RR}^T\mathbf{V} - \mathbf{V}^T\mathbf{V}\right)\mathbf{\Sigma U}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^*\right| \\
&= \left|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma}\left(-\mathbf{E}\right)\mathbf{\Sigma U}^T\mathbf{Y}\tilde{\boldsymbol{\alpha}}^*\right| \\
&\leq \|\mathbf{E}\|_2 \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma}\right\|_2^2 \\
&= \|\mathbf{E}\|_2 \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YU\Sigma V}^T\right\|_2^2 \\
&= \|\mathbf{E}\|_2 \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YX}\right\|_2^2. 
\end{aligned}$$

We can rewrite the above inequality as $\left|\left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YXR}\right\|_2^2 - \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YX}\right\|_2^2\right| \leq \|\mathbf{E}\|_2 \left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YX}\right\|_2^2$; thus,

$$\left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YX}\right\|_2^2 \leq \frac{1}{1 - \|\mathbf{E}\|_2}\left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YXR}\right\|_2^2.$$

Combining with eqn. (16), we get

$$Z_{opt} \geq \tilde{Z}_{opt} - \frac{1}{2}\left(\frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2}\right)\left\|\tilde{\boldsymbol{\alpha}}^{*T}\mathbf{YXR}\right\|_2^2. \tag{17}$$

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

$$Z_{opt} = \|\mathbf{w}^*\|_2^2 - \frac{1}{2}\|\mathbf{w}^*\|_2^2 = \frac{1}{2}\|\mathbf{w}^*\|_2^2, \tag{18}$$

$$\tilde{Z}_{opt} = \|\tilde{\mathbf{w}}^*\|_2^2 - \frac{1}{2}\|\tilde{\mathbf{w}}^*\|_2^2 = \frac{1}{2}\|\tilde{\mathbf{w}}^*\|_2^2. \tag{19}$$

Combining eqns. (17), (18), and (19), we get

$$\begin{aligned}
\|\mathbf{w}^*\|_2^2 &\geq \|\tilde{\mathbf{w}}^*\|_2^2 - \left(\frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2}\right)\|\tilde{\mathbf{w}}^*\|_2^2 \\
&= \left(1 - \frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2}\right)\|\tilde{\mathbf{w}}^*\|_2^2. \tag{20}
\end{aligned}$$

Let $\gamma^* = \|\mathbf{w}^*\|_2^{-1}$ be the geometric margin of the problem of eqn. (3) and let $\tilde{\gamma}^* = \|\tilde{\mathbf{w}}^*\|_2^{-1}$ be the geometric margin of the problem of eqn. (5). Then, the above equation implies:

$$\begin{aligned}
\gamma^{*2} &\leq \left(1 - \frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2}\right)^{-1}\tilde{\gamma}^{*2} \\
\Rightarrow \tilde{\gamma}^{*2} &\geq \left(1 - \frac{\|\mathbf{E}\|_2}{1 - \|\mathbf{E}\|_2}\right)\gamma^{*2}. \tag{21}
\end{aligned}$$

$\square$

Our second theorem argues that the radius of the minimum ball enclosing all projected points (the rows of the matrix $\mathbf{X}\mathbf{R}$) is very close to the radius of the minimum ball enclosing all original points (the rows of the matrix $\mathbf{X}$). We will prove this theorem for $\mathbf{R} = \mathbf{R}_{\text{SRHT}}$ as in eqn. (11), but similar results can be proven for the other two constructions for $\mathbf{R}$.

THEOREM 3.6.  *Fix $\epsilon \in (0, \frac{1}{2}]$, $\delta \in (0, 1]$. Let $B$ be the radius of the minimum ball enclosing all points in the full-dimensional space (the rows of the matrix $\mathbf{X}$), and let $\tilde{B}$ be the radius of the ball enclosing all points in the dimensionally reduced space (the rows of the matrix $\mathbf{X}\mathbf{R}$). Then, if $r = O(\rho\epsilon^{-2} \cdot \log(\rho d\delta^{-1}) \cdot \log(\rho\epsilon^{-2}\delta^{-1}\log(\rho d\delta^{-1})))$, with probability at least $1 - \delta$,*

$$\tilde{B}^2 \leq (1 + \epsilon)B^2.$$

PROOF. We consider the matrix $\mathbf{X}_B \in \mathbb{R}^{(n+1)\times d}$ whose first $n$ rows are the rows of $\mathbf{X}$ 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 $\mathbf{X}_B = \mathbf{U}_B\boldsymbol{\Sigma}_B\mathbf{V}_B^T$, where $\mathbf{U}_B \in \mathbb{R}^{(n+1)\times\rho_B}$, $\boldsymbol{\Sigma}_B \in \mathbb{R}^{\rho_B\times\rho_B}$, and $\mathbf{V} \in \mathbb{R}^{d\times\rho_B}$. Here $\rho_B$ is the rank of the matrix $\mathbf{X}_B$ and clearly $\rho_B \leq \rho + 1$. (Recall that $\rho$ is the rank of the matrix $\mathbf{X}$.) Let $B$ be the radius of the minimal radius ball enclosing all $n$ points in the original space. Then, for any $i = 1, \ldots, n$,

$$B^2 \geq \|\mathbf{x}_i - \mathbf{x}_B\|_2^2 = \left\|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\mathbf{X}_B\right\|_2^2. \tag{22}$$

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

$$\begin{aligned}
&\left|\left\|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\mathbf{X}_B\right\|_2^2 - \left\|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\mathbf{X}_B\mathbf{R}\right\|_2^2\right| \\
&= \left|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\left(\mathbf{X}_B\mathbf{X}_B^T - \mathbf{X}_B\mathbf{R}\mathbf{R}^T\mathbf{X}_B^T\right)(\mathbf{e}_i - \mathbf{e}_{n+1})\right| \\
&= \left|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\mathbf{U}_B\boldsymbol{\Sigma}_B\mathbf{E}_B\boldsymbol{\Sigma}_B\mathbf{U}_B^T(\mathbf{e}_i - \mathbf{e}_{n+1})\right| \\
&\leq \|\mathbf{E}_B\|_2\left\|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\mathbf{U}_B\boldsymbol{\Sigma}_B\right\|_2^2 \\
&= \|\mathbf{E}_B\|_2\left\|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\mathbf{U}_B\boldsymbol{\Sigma}_B\mathbf{V}_B^T\right\|_2^2 \\
&= \|\mathbf{E}_B\|_2\left\|(\mathbf{e}_i - \mathbf{e}_{n+1})^T\mathbf{X}_B\right\|_2^2.
\end{aligned}$$

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 $(n+1)$-st row of the matrix $\mathbf{X}_B\mathbf{R}$

(essentially, the projection of the center of the minimal radius enclosing ball for the original points). Let $\tilde{i} = \arg\max_{i=1\ldots n} \left\| (\mathbf{e}_i - \mathbf{e}_{n+1})^T \mathbf{X}_B \mathbf{R} \right\|_2^2$; then, using the above bound and eqn. (22), 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 projected points. Recall that $\tilde{B}$ is defined as the radius of the minimal radius ball that encloses all points in projected subspace; clearly,

$$\tilde{B}^2 \leq (1 + \|\mathbf{E}_B\|_2) B^2.$$

We can now use Lemma 3.1 on $\mathbf{V}_B$ to conclude that (using $\rho_B \leq \rho + 1$) $\|\mathbf{E}_B\|_2 \leq \epsilon$ □

Similar theorems can be proven for the two other constructions of $\mathbf{R}$ by using appropriate values for $r$. We are now ready to conclude the proof of Theorem 1.1.

PROOF. *(of Theorem 1.1)* The proof of Theorem 1.1 follows by combining Theorem 3.5, Lemma 3.1, and Theorem 3.6. The failure probability is at most $2\delta$, by a simple application of the union bound. □

Finally, we state the margin preservation theorem for SVM regression, which is analogous to Theorem 3.5. This theorem holds for all four choices of the random projection matrix $\mathbf{R} \in \mathbb{R}^{d \times r}$ and is identical to the proof of Theorem 3.5.

THEOREM 3.7. *Let $\epsilon$ be an accuracy parameter and let $\mathbf{R} \in \mathbb{R}^{d \times r}$ be a matrix satisfying $\|\mathbf{V}^T \mathbf{V} - \mathbf{V}^T \mathbf{R} \mathbf{R}^T \mathbf{V}\|_2 \leq \epsilon$. Let $\gamma^*$ and $\tilde{\gamma}^*$ be the margins obtained by solving the SVM regression problems using data matrices $\mathbf{X}$ and $\mathbf{XR}$ respectively (eqns. (4) and (6)). Then,*

$$\tilde{\gamma}^{*2} \geq (1 - \epsilon) \cdot \gamma^{*2}.$$

## 4. EXPERIMENTS

In our experimental evaluations, we implemented random projections using four different methods: RG, RS, FHT, and CW (see Section 2 for definitions) in MATLAB version 7.13.0.564 (R2011b). We ran the algorithms using the same values of $r$ (the dimension of the projected feature space) for all algorithms, but we varied $r$ across different datasets. We used LIBLINEAR [Fan et al. 2008] and LIBSVM [Chang and Lin 2011] as our linear SVM solver with default settings. In all cases, we ran our experiments on the original full data (referred to as "full" in the results), as well as on the projected data. For large-scale datasets, we use LIBLINEAR which is a faster SVM solver than LIBSVM, while for medium-scale datasets we use LIBSVM. We partitioned the data randomly for ten-fold cross-validation in order to estimate out-of-sample error. We repeated this partitioning ten times to get ten ten-fold cross-validation experiments. In the case where the dataset is already available in the form of a training and test-set, we do not perform ten-fold cross validation and use the given training and test set instead. In order to estimate the effect of the randomness in the construction of the random projection matrices, we repeated our cross-validation experiments ten times using ten different random projection matrices for all datasets. For classification experiments, we report in-sample error ($\epsilon_{in}$), out-of-sample error ($\epsilon_{out}$), the time to compute random projections ($t_{rp}$), the total time needed to both compute random projections *and* run SVMs on the lower-dimensional problem ($t_{run}$), and the margin ($\gamma$). For regression experiments, we report the margin, the combined running-time of random projections and SVM, mean-squared error ($mse$) and the squared correlation-coefficient ($\beta$) of $\epsilon_{in}$. All results are averaged over the ten cross-validation experiments and the ten choices of random projection matrices. For each of the aforementioned quantities, we report both its mean value $\mu$ and its standard deviation $\sigma$.

### 4.1. Experiments on SVM Classification

We describe experimental evaluations on three real-world datasets, namely a collection of document-term matrices (the TechTC-300 dataset [Davidov et al. 2004]), a subset of the Reuters Corpus dataset (RCV1 Dataset [Lewis et al. 2004]) and a population genetics dataset (the joint Human Genome Diversity Panel or HGDP [Li et al. 2008] and the HapMap Phase 3 data [Paschou et al. 2010]) and also on three synthetic datasets. The synthetic datasets, a subset of the RCV1 dataset and the TechTC-300 dataset correspond to binary classification tasks while the joint HapMap-HGDP dataset and a subset of the RCV1 dataset correspond to multi-class classification tasks; our algorithms perform well in

multi-class classification as well. For the multi-class experiments of Section 4.1.3, we do not report a margin. We use LIBLINEAR as our SVM solver for Hapmap-HGDP [1] and the RCV1 datasets, while for the remaining datasets we use LIBSVM as our solver. For multi-class experiments, we use the method of Crammer and Singer [Crammer and Singer 2000] implemented in LIBLINEAR.

Table I. $\epsilon_{out}$ and $\gamma$ of Synthetic Data

| $\epsilon_{out}$ | | Projected Dimension $r$ | | | | $\gamma$ | | Projected Dimension $r$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 256 | 512 | 1024 | **full** | | | 256 | 512 | 1024 | **full** |
| D1 | CW ($\mu$) | 24.08 | 19.45 | 16.66 | **15.10** | D1 | CW ($\mu$) | 5.72 | 6.67 | 7.16 | **7.74** |
| | ($\sigma$) | 4.52 | 4.15 | 3.52 | **2.60** | | ($\sigma$) | 0.58 | 0.58 | 0.59 | **0.59** |
| | RS ($\mu$) | 24.1.0 | 19.46 | 16.36 | **15.10** | | RS ($\mu$) | 5.73 | 6.66 | 7.18 | **7.74** |
| | ($\sigma$) | 4.45 | 3.79 | 3.22 | **2.60** | | ($\sigma$) | 0.57 | 0.55 | 0.55 | **0.59** |
| | FHT ($\mu$) | 23.52 | 19.59 | 16.67 | **15.10** | | FHT ($\mu$) | 5.76 | 6.64 | 7.15 | **7.74** |
| | ($\sigma$) | 4.21 | 4.05 | 3.37 | **2.60** | | ($\sigma$) | 0.56 | 0.58 | 0.56 | **0.59** |
| | RG ($\mu$) | 24.34 | 19.73 | 16.69 | **15.10** | | RG ($\mu$) | 5.67 | 6.60 | 7.13 | **7.74** |
| | ($\sigma$) | 4.44 | 3.86 | 3.28 | **2.60** | | ($\sigma$) | 0.57 | 0.51 | 0.54 | **0.59** |
| D2 | CW ($\mu$) | 25.94 | 21.07 | 17.33 | **15.44** | D2 | CW ($\mu$) | 6.62 | 8.09 | 8.88 | **9.78** |
| | ($\sigma$) | 4.13 | 4.16 | 3.45 | **2.54** | | ($\sigma$) | 0.64 | 0.62 | 0.59 | **0.66** |
| | RS ($\mu$) | 25.80 | 20.80 | 17.47 | **15.44** | | RS ($\mu$) | 6.65 | 8.10 | 8.88 | **9.78** |
| | ($\sigma$) | 4.40 | 3.93 | 3.42 | **2.54** | | ($\sigma$) | 0.64 | 0.60 | 0.63 | **0.66** |
| | FHT ($\mu$) | 25.33 | 21.23 | 17.58 | **15.44** | | FHT ($\mu$) | 6.66 | 8.06 | 8.84 | **9.78** |
| | ($\sigma$) | 3.69 | 4.24 | 3.53 | **2.54** | | ($\sigma$) | 0.63 | 0.65 | 0.63 | **0.66** |
| | RG ($\mu$) | 25.43 | 20.54 | 17.25 | **15.44** | | RG ($\mu$) | 6.66 | 8.13 | 8.90 | **9.78** |
| | ($\sigma$) | 4.03 | 3.65 | 3.38 | **2.54** | | ($\sigma$) | 0.65 | 0.60 | 0.63 | **0.66** |
| D3 | CW ($\mu$) | 27.62 | 22.97 | 18.93 | **15.83** | D3 | CW ($\mu$) | 7.69 | 9.84 | 11.07 | **12.46** |
| | ($\sigma$) | 3.46 | 3.22 | 3.32 | **2.00** | | ($\sigma$) | 0.67 | 0.60 | 0.71 | **0.69** |
| | RS ($\mu$) | 28.15 | 23.00 | 18.72 | **15.83** | | RS ($\mu$) | 7.61 | 9.85 | 11.05 | **12.46** |
| | ($\sigma$) | 3.02 | 3.48 | 2.78 | **2.00** | | ($\sigma$) | 0.59 | 0.6212 | 0.62 | **0.69** |
| | FHT ($\mu$) | 27.92 | 23.41 | 18.73 | **15.83** | | FHT ($\mu$) | 7.63 | 9.83 | 11.11 | **12.46** |
| | ($\sigma$) | 3.46 | 3.60 | 3.02 | **2.00** | | ($\sigma$) | 0.67 | 0.64 | 0.64 | **0.69** |
| | RG ($\mu$) | 27.71 | 22.85 | 18.96 | **15.83** | | RG ($\mu$) | 7.69 | 9.85 | 11.04 | **12.46** |
| | ($\sigma$) | 3.38 | 3.29 | 3.33 | **2.00** | | ($\sigma$) | 0.67 | 0.61 | 0.7 | **0.69** |

*Synthetic data:* $\epsilon_{out}$ decreases and $\gamma$ increases as a function of $r$ in all three families of matrices, using any of the four random projection methods. $\mu$ and $\sigma$ indicate the mean and the standard deviation of $\epsilon_{out}$ over ten matrices in each family $D1$, $D2$, and $D3$, ten ten-fold cross-validation experiments, and ten choices of random projection matrices for the four methods that we investigated (a total of 1,000 experiments for each family of matrices).

*4.1.1. Synthetic datasets.* The synthetic datasets are separable by construction. More specifically, we first constructed a weight vector $\mathbf{w} \in \mathbb{R}^d$, whose entries were selected in i.i.d. trials from a Gaussian distribution $\mathcal{N}(\mu, \sigma)$ of mean $\mu$ and standard-deviation $\sigma$. We experimented with the following three distributions: $\mathcal{N}(0,1)$, $\mathcal{N}(1, 1.5)$, and $\mathcal{N}(2, 2)$. Then, we normalized $\mathbf{w}$ to create $\hat{\mathbf{w}} = \mathbf{w} / \|\mathbf{w}\|_2$. Let $\mathbf{X}_{ij} = \mathcal{N}(0,1)$; then, we set $\mathbf{x}_i$ to be equal to the $i$-th row of $\mathbf{X}$, while $\mathbf{y}_i = sign\left(\hat{\mathbf{w}}^T \mathbf{x}_i\right)$. We generated families of matrices of different dimensions. More specifically, family $D1$ contained matrices in $\mathbb{R}^{200 \times 5,000}$; family $D2$ contained matrices in $\mathbb{R}^{250 \times 10,000}$; and family $D3$ contained matrices in $\mathbb{R}^{300 \times 20,000}$. We generated ten datasets for each of the families $D1$, $D2$, and $D3$, and we report average results over the ten datasets. We set $r$ to 256, 512, and 1024 and set $C$ to 1,000 in LIBSVM for all the experiments. Tables I shows $\epsilon_{out}$ and $\gamma$ for the three datasets $D1$, $D2$, and $D3$. $\epsilon_{in}$ is zero for all three data families. As expected, $\epsilon_{out}$ and $\gamma$ improve as $r$ grows for all four random projection methods. Also, the time needed to compute random projections is very small compared to the time needed to run SVMs on the projected data. Figure 1 shows the combined running time of random projections and SVMs, which is nearly the same for all four random projection methods. It is obvious that this combined running time is much smaller that the time needed to run SVMs on the full dataset (without any dimensionality reduction). For instance, for $r = 1024$, $t_{run}$ for $D1$, $D2$, and $D3$ is (respectively) 6, 9, and 25 times smaller than $t_{run}$ on the full-data.

*4.1.2. The TechTC-300 dataset.* For our first real dataset, we use the TechTC-300 data, consisting of a family of 295 document-term data matrices. The TechTC-300 dataset comes from the Open Directory Project (ODP), which is a large, comprehensive directory of the web, maintained by volunteer editors. Each matrix in the TechTC-300 dataset contains a pair of categories from the ODP. Each category corresponds to a label, and thus the resulting classification task is binary. The documents that are

---

[1] In Paul et al. [2013], the experiments on Hapmap-HGDP dataset were done using LIBSVM's one-against-one multi-class classification method. LIBLINEAR does not have the one-against-one method implemented in the package. So we use the Crammer and Singer method [Crammer and Singer 2000].

Fig. 1. Total (average) running times, in seconds, of random projections *and* SVMs on the lower-dimensional data for each of the three families of synthetic data. Vertical bars indicate the, relatively small, standard deviation (see the caption of Table I).

collected from the union of all the subcategories within each category are represented in the bag-of-words model, with the words constituting the features of the data [Davidov et al. 2004]. Each data matrix consists of 150-280 documents (the rows of the data matrix **X**), and each document is described with respect to 10,000-40,000 words (features, columns of the matrix **X**). Thus, TechTC-300 provides a diverse collection of data sets for a systematic study of the performance of the SVM on the projected versus full data. We set the parameter $C$ to 500 in LIBSVM for all 295 document-term matrices and

Table II. TechTC-300 Dataset

|  |  | Projected Dimension $r$ | | | |
|  |  | 128 | 256 | 512 | **full** |
|---|---|---|---|---|---|
| $\epsilon_{out}$ | CW$(\mu)$ | 24.63 | 22.84 | 21.26 | **17.35** |
|  | $(\sigma)$ | 10.57 | 10.37 | 10.17 | **9.45** |
|  | RS$(\mu)$ | 24.58 | 22.90 | 21.38 | **17.35** |
|  | $(\sigma)$ | 10.57 | 10.39 | 10.23 | **9.45** |
|  | FHT $(\mu)$ | 24.63 | 22.93 | 21.35 | **17.35** |
|  | $(\sigma)$ | 10.66 | 10.39 | 10.2 | **9.45** |
|  | RG $(\mu)$ | 24.59 | 22.96 | 21.36 | **17.35** |
|  | $(\sigma)$ | 10.54 | 10.5 | 10.18 | **9.45** |
| $\gamma$ | CW $(\mu)$ | 1.66 | 1.88 | 1.99 | **2.09** |
|  | $(\sigma)$ | 3.68 | 3.79 | 3.92 | **4.00** |
|  | RS $(\mu)$ | 1.66 | 1.88 | 1.99 | **2.09** |
|  | $(\sigma)$ | 3.65 | 3.80 | 3.91 | **4.00** |
|  | FHT $(\mu)$ | 1.66 | 1.88 | 1.98 | **2.09** |
|  | $(\sigma)$ | 3.65 | 3.81 | 3.88 | **4.00** |
|  | RG $(\mu)$ | 1.66 | 1.88 | 1.99 | **2.09** |
|  | $(\sigma)$ | 3.70 | 3.83 | 3.91 | **4.00** |
| $t_{rp}$ | CW $(\mu)$ | 0.0046 | 0.0059 | 0.0075 | $--$ |
|  | $(\sigma)$ | 0.0019 | 0.0026 | 0.0033 | $--$ |
|  | RS $(\mu)$ | 0.0429 | 0.0855 | 0.1719 | $--$ |
|  | $(\sigma)$ | 0.0178 | 0.0356 | 0.072 | $--$ |
|  | FHT $(\mu)$ | 0.0443 | 0.0882 | 0.1764 | $--$ |
|  | $(\sigma)$ | 0.0206 | 0.0413 | 0.0825 | $--$ |
|  | RG $(\mu)$ | 0.039 | 0.078 | 0.1567 | $--$ |
|  | $(\sigma)$ | 0.0159 | 0.0318 | 0.0642 | $--$ |
| $t_{run}$ | CW $(\mu)$ | 1.23 | 2.22 | 4.63 | **4.85** |
|  | $(\sigma)$ | 0.87 | 0.93 | 1.93 | **2.12** |
|  | RS $(\mu)$ | 0.99 | 1.53 | 3.02 | **4.85** |
|  | $(\sigma)$ | 0.97 | 0.59 | 1.12 | **2.12** |
|  | FHT $(\mu)$ | 0.95 | 1.46 | 2.83 | **4.85** |
|  | $(\sigma)$ | 0.96 | 0.55 | 1.02 | **2.12** |
|  | RG $(\mu)$ | 0.82 | 1.23 | 2.48 | **4.85** |
|  | $(\sigma)$ | 0.83 | 0.45 | 0.84 | **2.12** |

*TechTC300:* Results on the TechTC300 dataset, averaged over 295 data matrices using four different random projection methods. The table shows how $\epsilon_{out}$, $\gamma$, $t_{rp}$ (in seconds), and $t_{run}$ (in seconds) depend on $r$. $\mu$ and $\sigma$ indicate the mean and the standard deviation of each quantity over 295 matrices, ten ten-fold cross-validation experiments, and ten choices of random projection matrices for the four methods that we investigated.

set $r$ to 128, 256, and 512. We use a lower value of C than for the other data sets for computational reasons: larger $C$ is less efficient. We note that our classification accuracy is slightly worse (on the full data) than the accuracy presented in Section 4.4 of Davidov et al. [2004], because we did not fine-tune the SVM parameters as they did, since that is not the focus of this study. For every dataset and every value of $r$ we tried, the in-sample error on the projected data matched the in-sample error on the full data. We thus focus on $\epsilon_{out}$, the margin $\gamma$, the time needed to compute random projections $t_{rp}$, and the total running time $t_{run}$. We report our results averaged over 295 data matrices. Table II shows the behavior of these parameters for different choices of $r$. As expected, $\epsilon_{out}$ and the margin $\gamma$ improve as $r$ increases, and they are nearly identical for all four random projection methods. The time needed to compute random projections is smallest for CW, followed by RG, RS and FHT. As a matter of fact, $t_{rp}$ for CW is ten to 20 times faster than RG, RS and FHT for different values of $r$. This is predicted by the theory in Clarkson and Woodruff [2013], since CW is optimized to take advantage of input sparsity. However, this advantage is lost when SVMs are applied on the dimensionally-reduced data. Indeed, the combined running time $t_{run}$ is fastest for RG, followed by FHT, RS and CW. In all cases, the total running time is smaller than the SVM running time on full dataset. For example, in the case of RG or

FHT, setting $r = 512$ achieves a running time $t_{run}$ which is about twice as fast as running SVMs on the full dataset; $\epsilon_{out}$ increases by less than $4\%$.



Fig. 2.   $\epsilon_{out}$ as a function of $r$ in the Hapmap-HGDP dataset for four different random projection methods and two different classification tasks. Vertical bars indicate the standard-deviation over the ten ten-fold cross-validation experiments and the ten choices of the random projection matrices for each of the four methods.

*4.1.3. The HapMap-HGDP dataset.* Predicting ancestry of individuals using a set of genetic markers is a well-studied classification problem. We use a population genetics dataset from the Human Genome Diversity Panel (HGDP) and the HapMap Phase 3 dataset (see Paschou et al. [2010] for details), in order to classify individuals into broad geographic regions, as well as into (finer-scale) populations. We study a total of 2,250 individuals from approximately 50 populations and five broad geographic regions (Asia, Africa, Europe, the Americas, and Oceania). The features in this dataset correspond to $492,516$ Single Nucleotide Polymorphisms (SNPs), which are well-known biallelic loci of genetic variation across the human genome. Each entry in the resulting $2,250 \times 492,516$ matrix is set to $+1$ (homozygotic in one allele), $-1$ (homozygotic in the other allele), or $0$ (heterozygotic), depending on the genotype of the respective SNP for a particular sample. Missing entries were filled in with $-1$, $+1$, or $0$, with probability 1/3. Each sample has a known population and region of origin, which constitute its label. We set $r$ to $256$, $512$, $1024$, and $2048$ in our experiments. Since this task is a multi-class classification problem, we used LIBLINEAR's Crammer and Singer technique for classification. We ran two sets of experiments: in the first set, the classification problem is to assign samples to broad regions of origin, while in the second experiment, our goal is to classify samples into (fine-scale) populations. We set $C$ to 1,000 in LIBLINEAR for all the experiments. The in-sample error is zero in all cases. Figure 2 shows the out-of-sample error for regions and populations classification, which are nearly identical for all four random projection methods. For regional classification, we estimated $\epsilon_{out}$ to be close to $2\%$, and for population-level classification, $\epsilon_{out}$ is close to $20\%$. This experiment strongly supports the **computational benefits** of our methods in terms of main memory. $\mathbf{X}$ is a $2,250 \times 492,516$ matrix, which is too large to fit into memory in order to run SVMs. Figure 3 shows that the combined running time for four different random projection methods are nearly identical for both regions and population classification tasks. However, the time needed to compute the random projections is different from one method to

Total running time: regional classification



Total running time: population-level classification



Time needed to compute random projections

Fig. 3.  Total running time in seconds (random projections *and* SVM classification on the dimensionally-reduced data) for Hapmap-HGDP dataset for four different projection methods using both regional and population-level labels. Notice that the time needed to compute random projection is independent of the classification labels. Vertical bars indicate standard-deviation, as in Figure 2.

the next. FHT is fastest, followed by RS, RG and CW. In this particular case, the input matrix is dense, and CW seems to be outperformed by the other methods as the running time of CW depends on the number of non-zeros of the matrix.

  *4.1.4. The RCV1 dataset.* The RCV1 dataset [Lewis et al. 2004][2] is a benchmark dataset on text categorization. We use the RCV1 dataset for both binary and multi-class classification tasks. The RCV1

_____

[2] The RCV1 dataset is available publicly at http://www.csie.ntu.edu.tw/ cjlin/libsvmtools/datasets and contains a training-set and a test-set of predesignated size.

Fig. 4. Ratio of number of non-zero entries of projected data and full-data for RCV1 dataset.

binary classification dataset had one training set containing 20,242 data-points and 47,236 features. We generate ten different test-sets each containing 20,000 points and 47,236 features from the available test-set for the binary classification task. The RCV1 binary classification dataset contains CCAT and ECAT as the positive classes and GCAT and MCAT as the negative classes. Instances in both positive and negative classes are absent. The RCV1 multi-class dataset had one training set 15,564 training points with 47,236 features and ten different test-sets each containing 16,000 data points and 47,236 features were generated from the available test-set. There are 53 classes in the RCV1 multi-class dataset. We set $r$ to $2048$, $4096$ and $8192$. We use LIBLINEAR as our SVM solver. We use the L2-regularized L2-loss support vector classification in the primal mode with default settings for binary classification and the method of Crammer and Singer [Crammer and Singer 2000] for multi-class classification. We set $C = 10$ for both multi-class classification and binary classification. The RCV1 dataset is very sparse with 0.16% non-zero entries in the binary-class dataset and 0.14% non-zero entries in the multi-class dataset.

Tables III and IV show the results for RCV1 binary and multi class datasets. $t_{svm}$ denotes the SVM-training time on the projected data. For both binary and multi-class tasks, we observe that $\epsilon_{out}$ is close to that of full dataset and $\epsilon_{out}$ decreases with increase in number of projections. The SVM training time is smaller than that of full dataset for CW method only. For the other methods like FHT, RS and RG, the number of non-zeros in the projected data increases which increases the SVM training time. This is evident from Figure 4 which shows the ratio of number of non-zeros of projected data and full data. For all methods except CW, the number of non-zeros increases with increase in value of $r$. This follows from the theory predicted in Clarkson and Woodruff [2013], since CW method takes advantage of input sparsity. The combined running time is smaller than that of full-dataset for CW method. The margin of the projected data is close to that of full data for RCV1 binary class dataset.

### 4.2. PCA vs Random Projections

Principal Components Analysis (PCA) constructs a small number of linear features that summarize the input data. PCA is computed by first mean-centering the features of the original data and then computing a low-rank approximation of the data matrix using SVD. Thus the PCA feature matrix is given by $\mathbf{Z} = \mathbf{X}_c \mathbf{V}_k$, where $\mathbf{X}_c$ represents the centered data matrix $\mathbf{X}$ and $\mathbf{V}_k$ represents the top $k$ right singular vectors of $\mathbf{X}_c$. To the best of our knowledge, there is no known theoretical framework connecting PCA with margin or generalization error of SVM, so we only provide empirical evidence of the comparison.

Our goal is to evaluate if the data matrix represented by a small number of principal components can give the same or better performance than random projections when combined with SVMs, in terms of both running time and out-of-sample error. Note that the number of random projections is always greater than the rank of the matrix. For PCA we retain a number of principal components that is less than or equal to the rank of the matrix in order to compare its performance to random projections.

We used the TechTC300 dataset for experimental evaluation and used MATLAB's SVD solver in "econ" mode to compute PCA. We kept $k$ equal to 32, 64, and $\rho$ (where $\rho$ is the rank of the matrix) principal components. The results corresponding to a number of principal components equal to the rank of the data matrices are referred to as "full-rank", while "full" refers to the results on the full-dimensional dataset. We ran PCA experiments on 294 datasets with $k = 64$, since one of them had rank less than 64. For $k = 32$, we used the entire set of 295 TechTC300 matrices. $t_{pca}$ denotes the time to compute PCA on the dataset and we set $C = 500$ and $C = 1$ in our experiments. The out-of-sample error for PCA is equal to or sometimes slightly better compared to the error on the full-dimensional datasets. Even though a smaller number of principal components achieve better out-of-sample error when compared to random projections, the combined running time of SVMs and PCA is typically higher than that of random projections. The combined running time of SVMs and PCA is sensitive

Table III. RCV1 Dataset (Multi-class)

| | | Projected Dimension $r$ | | | |
|---|---|---|---|---|---|
| | | 2048 | 4096 | 8192 | **full** |
| $\epsilon_{\mathbf{out}}$ | CW($\mu$) | 18.30 | 15.15 | 13.50 | **11.83** |
| | ($\sigma$) | 0.157 | 0.117 | 0.079 | **0.18** |
| | RS($\mu$) | 17.57 | 14.66 | 13.21 | **11.83** |
| | ($\sigma$) | 0.172 | 0.134 | 0.1094 | **0.18** |
| | FHT ($\mu$) | 17.50 | 14.58 | 13.2 | **11.83** |
| | ($\sigma$) | 0.2033 | 0.0905 | 0.0597 | **0.18** |
| | RG ($\mu$) | 17.47 | 14.61 | 13.21 | **11.83** |
| | ($\sigma$) | 0.101 | 0.146 | 0.036 | **0.18** |
| $\mathbf{t_{rp}}$ | CW ($\mu$) | 0.0093 | 0.0189 | 0.0357 | —— |
| | ($\sigma$) | 0.0005 | 0.0007 | 0.0010 | —— |
| | RS ($\mu$) | 2.412 | 4.782 | 9.49 | —— |
| | ($\sigma$) | 0.0323 | 0.0847 | 0.0636 | —— |
| | FHT ($\mu$) | 2.484 | 4.85 | 9.966 | —— |
| | ($\sigma$) | 0.1621 | 0.023 | 0.23 | —— |
| | RG ($\mu$) | 13.559 | 28.875 | 55.05 | —— |
| | ($\sigma$) | 0.0715 | 1.9801 | 1.4315 | —— |
| $\mathbf{t_{svm}}$ | CW ($\mu$) | 1.99 | 2.23 | 2.80 | **2.96** |
| | ($\sigma$) | 0.2309 | 0.4030 | 0.2929 | —— |
| | RS ($\mu$) | 45.103 | 94.652 | 207.278 | **2.96** |
| | ($\sigma$) | 2.82 | 5.705 | 15.084 | —— |
| | FHT ($\mu$) | 45.468 | 98.287 | 201.558 | **2.96** |
| | ($\sigma$) | 3.143 | 6.689 | 10.765 | —— |
| | RG ($\mu$) | 47.458 | 121.208 | 263.392 | **2.96** |
| | ($\sigma$) | 7.222 | 23.393 | 45.814 | —— |

*RCV1 Multi-class: $\mu$ and $\sigma$ represents the mean and standard deviation of the results which have been averaged over ten different random projection matrices. Since there was one training set, there is no standard deviation for $t_{svm}$ of "full" data.*

to the value of C, while the running time for random projections and SVM do not vary greatly, like PCA, by change of $C$.[3] Random projections are therefore a faster and more stable method than PCA. However, PCA appears to perform better than random projections when the number of components is equal to the rank of the matrix. Table V shows the results of PCA experiments; note that the standard deviation of $t_{run}$ for $C = 500$ is quite high because of the varied running times of SVMs on the various TechTC300 matrices.

For a comparison of PCA to random projections, we consider the case of $r = 512$ and the random gaussian matrix and randomized Hadamard Transform (see Table II for $C = 500$), which has the best combined running time for random projections and SVMs. The combined running time of SVMs and "full-rank" PCA is smaller than that of RG (FHT) and SVMs by 0.19 (0.16) seconds, while the out-of-sample error of the former is only 4% better. However, the time needed to compute PCA is 1.73 seconds, while random projections take negligible time to be computed; applying SVMs on the dimensionally-reduced matrix is the bottleneck of the computation. If PCA retains only 32 or 64 principal components, the running time of our method is smaller by factors of 30 and 5 respectively.

For $C = 1$ and $r = 512$, SVMs and "full-rank" PCA is smaller than that of RG (FHT) and SVMs by 0.43 (0.15) seconds, while the out-of-sample error of the former is again 4% better. For $C = 1$, if PCA retains only 32 or 64 principal components, the running time of our method is smaller by a few seconds.

These clearly show the advantage of using random projections over PCA, especially when a small number of principal components is desired. The PCA feature matrix is sensitive to the value of $C$. For a higher value of $C$, SVM takes a longer time to train the inseparable data of the PCA feature matrix.

### 4.3. Experiments on SVM regression

We describe experimental evaluations on real world datasets, namely Yalefaces dataset [Cai et al. 2006] and a gene expression dataset (NCI60 [Ross et al. 2000]). We convert the multi-label classification tasks into a regression problem. We use LIBSVM with default settings and use $C = 1$ in all our experiments. The general observations from our experiments are as follows: (i) the combined runtime of random projections and SVM is smaller than the runtime of SVM on full dataset, (ii) the margin increases with an increase in the number of projections.

---

[3] We repeated all experiments on TechTC-300 using $C = 1$ and noticed the same pattern in the results as we did for $C = 500$. RG and FHT are faster than the remaining two methods.

Table IV. RCV1 Dataset (binary-class)

|  |  | Projected Dimension $r$ | | | |
|---|---|---|---|---|---|
|  |  | 2048 | 4096 | 8192 | **full** |
| $\epsilon_{\text{out}}$ | CW($\mu$) | 12.90 | 9.57 | 6.54 | 4.51 |
|  | ($\sigma$) | 1.234 | 0.619 | 0.190 | **0.1161** |
|  | RS($\mu$) | 9.26 | 7.82 | 6.22 | 4.51 |
|  | ($\sigma$) | 0.134 | 0.116 | 0.077 | **0.1161** |
|  | FHT($\mu$) | 9.23 | 7.94 | 6.29 | 4.51 |
|  | ($\sigma$) | 0.166 | 0.145 | 0.150 | **0.1161** |
|  | RG($\mu$) | 9.24 | 7.75 | 6.20 | 4.51 |
|  | ($\sigma$) | 0.255 | 0.113 | 0.109 | **0.1161** |
| $\gamma$ | CW($\mu$) | 0.0123 | 0.0094 | 0.0125 | **0.0152** |
|  | ($\sigma$) | 0.0002 | 0.0002 | 0.0003 | $--$ |
|  | RS($\mu$) | 0.0158 | 0.0105 | 0.0127 | **0.0152** |
|  | ($\sigma$) | 0.0005 | 0.00009 | 0.0001 | $--$ |
|  | FHT($\mu$) | 0.0160 | 0.0105 | 0.0127 | **0.0152** |
|  | ($\sigma$) | 0.0003 | 0.0001 | 0.0001 | $--$ |
|  | RG($\mu$) | 0.0158 | 0.0106 | 0.0127 | **0.0152** |
|  | ($\sigma$) | 0.0004 | 0.0001 | 0.00009 | $--$ |
| $t_{\text{rp}}$ | CW($\mu$) | 0.0112 | 0.0231 | 0.0456 | $--$ |
|  | ($\sigma$) | 0.0008 | 0.0014 | 0.0008 | $--$ |
|  | RS($\mu$) | 4.34 | 8.94 | 18.89 | $--$ |
|  | ($\sigma$) | 0.48 | 1.33 | 3.48 | $--$ |
|  | FHT($\mu$) | 4.29 | 8.21 | 16.37 | $--$ |
|  | ($\sigma$) | 0.577 | 1.003 | 1.484 | $--$ |
|  | RG($\mu$) | 20.52 | 48.18 | 86.384 | $--$ |
|  | ($\sigma$) | 0.248 | 3.32 | 6.15 | $--$ |
| $t_{\text{svm}}$ | CW($\mu$) | 0.0726 | 0.1512 | 0.2909 | **0.368** |
|  | ($\sigma$) | 0.0082 | 0.0113 | 0.009 | $--$ |
|  | RS($\mu$) | 9.22 | 20.11 | 41.29 | **0.368** |
|  | ($\sigma$) | 0.9396 | 3.245 | 7.3966 | $--$ |
|  | FHT($\mu$) | 8.78 | 18.95 | 37.18 | **0.368** |
|  | ($\sigma$) | 1.29 | 2.79 | 4.76 | $--$ |
|  | RG($\mu$) | 11.17 | 30.17 | 44.76 | **0.368** |
|  | ($\sigma$) | 0.364 | 4.204 | 2.079 | $--$ |

*RCV1 Binary-class: $\mu$ and $\sigma$ represents the mean and standard deviation of the results which have been averaged over ten different random projection matrices. Since there was one training set, there is no standard deviation for $t_{svm}$ and $\gamma$ of "full" data.*

Table V. TechTC300 PCA Experiments

| C=500 |  | Projected Dimension $k$ | | | | | C=1 |  | Projected Dimension $k$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | 32 | 64 | **full-rank** | **full** |  |  |  | 32 | 64 | **full-rank** | **full** |
| $\epsilon_{\text{out}}$ | PCA ($\mu$) | 15.02 | 17.35 | 17.35 | **17.35** |  | $\epsilon_{\text{out}}$ | PCA ($\mu$) | 13.33 | 15.72 | 17.21 | **17.21** |
|  | ($\sigma$) | 9.32 | 9.53 | 9.45 | **9.45** |  |  | ($\sigma$) | 8.10 | 9.20 | 9.44 | **9.44** |
| $t_{\text{pca}}$ | PCA ($\mu$) | 1.73 | 1.73 | 1.73 | $--$ |  | $t_{\text{pca}}$ | PCA ($\mu$) | 1.73 | 1.73 | 1.73 | $--$ |
|  | ($\sigma$) | 1.15 | 1.15 | 1.15 | $--$ |  |  | ($\sigma$) | 1.15 | 1.15 | 1.15 | $--$ |
| $t_{\text{run}}$ | PCA ($\mu$) | 86.73 | 14.04 | 2.67 | **4.85** |  | $t_{\text{run}}$ | PCA ($\mu$) | 5.70 | 2.99 | 2.80 | **4.92** |
|  | ($\sigma$) | 174.35 | 81.12 | 1.56 | **2.12** |  |  | ($\sigma$) | 10.68 | 5.83 | 1.62 | **2.16** |

*PCA Experiments:* Results on the TechTC300 dataset, averaged over all data matrices using PCA. The table shows how $\epsilon_{out}$, $t_{pca}$ (in seconds), and $t_{run}$ (in seconds) depend on $r$. $\mu$ and $\sigma$ indicate the mean and the standard deviation of each quantity over the data matrices and ten ten-fold cross-validation experiments.

*4.3.1. Yalefaces Dataset.* The Yalefaces dataset [Cai et al. 2006] consists of 165 grayscale images of 15 individuals. There were eleven images per subject, one per different facial expression (happy, sad, etc) or configuration (center-light, left-light, etc). The dataset has 165 datapoints and 4,096 features with 15 classes. The classes were used as the labels for regression. We set the value of $r$ to 256, 512, and 1024. $t_{run}$ for SVM and random projections is approximately 9, 7 and 4 times smaller than that of full-dataset. The margin increases as the number of random projections increases. The mean-squared in-sample error decreases with an increase in the number of random projections.

*4.3.2. NCI60 Dataset.* The NCI60 dataset [Ross et al. 2000] consists of 1375 gene expression profiles of 60 human cancer cell lines. The dataset contains 1375 features and 60 datapoints with ten classes. The features contain the log-ratio of the expression levels. The classes were used as labels for regression. We set the value of $r$ to 128, 256, and 512. The running time of the four methods are nearly the same. The squared correlation-coefficient is very close to one and is not influenced by the number of projections, $r$. The mean squared $\epsilon_{in}$ remains the same for all values of $r$.

Table VI. Yalefaces Dataset

| | | Projected Dimension $r$ | | | |
|---|---|---|---|---|---|
| | | 256 | 512 | 1024 | **full** |
| $\epsilon_{\mathbf{in}}$ | CW($mse$) | 0.2697 | 0.0257 | 0.0103 | **0.0098** |
| | ($\beta$) | 0.9859 | 0.9987 | 0.9995 | **0.9995** |
| | RS($mse$) | 0.2201 | 0.0187 | 0.0107 | **0.0098** |
| | ($\beta$) | 0.9886 | 0.9991 | 0.9995 | **0.9995** |
| | FHT ($mse$) | 0.2533 | 0.0233 | 0.0102 | **0.0098** |
| | ($\beta$) | 0.9868 | 0.9988 | 0.9995 | **0.9995** |
| | RG ($mse$) | 0.281 | 0.0174 | 0.0105 | **0.0098** |
| | ($\beta$) | 0.9853 | 0.9991 | 0.9995 | **0.9995** |
| $\gamma$ | CW ($\mu$) | 0.11 | 0.12 | 0.13 | **0.14** |
| | ($\sigma$) | 0.0032 | 0.0028 | 0.0020 | **0.0031** |
| | RS ($\mu$) | 0.11 | 0.12 | 0.13 | **0.14** |
| | ($\sigma$) | 0.0043 | 0.0024 | 0.0018 | **0.0031** |
| | FHT ($\mu$) | 0.11 | 0.12 | 0.13 | **0.14** |
| | ($\sigma$) | 0.0026 | 0.0022 | 0.0031 | **0.0031** |
| | RG ($\mu$) | 0.11 | 0.12 | 0.13 | **0.14** |
| | ($\sigma$) | 0.0034 | 0.0030 | 0.0025 | **0.0031** |
| $\mathbf{t_{run}}$ | CW ($\mu$) | 4.14 | 5.06 | 8.62 | **34.93** |
| | ($\sigma$) | 0.17 | 0.12 | 0.24 | **0.02** |
| | RS ($\mu$) | 3.92 | 4.81 | 8.45 | **34.93** |
| | ($\sigma$) | 0.17 | 0.10 | 0.19 | **0.02** |
| | FHT ($\mu$) | 4.12 | 5.09 | 8.69 | **34.93** |
| | ($\sigma$) | 0.23 | 0.17 | 0.25 | **0.02** |
| | RG ($\mu$) | 4.37 | 5.49 | 9.43 | **34.93** |
| | ($\sigma$) | 0.15 | 0.09 | 0.06 | **0.02** |

*Yalefaces Dataset:* Results on the Yalefaces dataset using four different random projection methods. The table shows how $\epsilon_{in}$, $\gamma$ and $t_{run}$ (in seconds) depend on $r$. $mse$ and $\beta$ indicate the mean-squared error and the squared correlation coefficient, while $\mu$ and $\sigma$ represent the mean and standard deviation over ten ten-fold cross-validation experiments, and ten choices of random projection matrices for the four methods that we investigated.

Table VII. NCI60 Dataset

| | | Projected Dimension $r$ | | | |
|---|---|---|---|---|---|
| | | 128 | 256 | 512 | **full** |
| $\epsilon_{\mathbf{in}}$ | CW($mse$) | 0.0098 | 0.0097 | 0.0097 | **0.0097** |
| | ($\beta$) | 0.9987 | 0.9989 | 0.9989 | **0.9990** |
| | RS($mse$) | 0.0098 | 0.0097 | 0.0097 | **0.0097** |
| | ($\beta$) | 0.9987 | 0.9988 | 0.9990 | **0.9990** |
| | FHT ($mse$) | 0.0097 | 0.0098 | 0.0097 | **0.0097** |
| | ($\beta$) | 0.9987 | 0.9988 | 0.9989 | **0.9990** |
| | RG ($mse$) | 0.0098 | 0.0098 | 0.0097 | **0.0097** |
| | ($\beta$) | 0.9987 | 0.9988 | 0.9989 | **0.9990** |
| $\gamma$ | CW ($\mu$) | 1.89 | 2.12 | 2.22 | **2.33** |
| | ($\sigma$) | 0.05 | 0.07 | 0.07 | **0.12** |
| | RS ($\mu$) | 1.90 | 2.09 | 2.25 | **2.33** |
| | ($\sigma$) | 0.10 | 0.06 | 0.07 | **0.12** |
| | FHT ($\mu$) | 1.88 | 2.10 | 2.22 | **2.33** |
| | ($\sigma$) | 0.10 | 0.09 | 0.05 | **0.12** |
| | RG ($\mu$) | 1.87 | 2.12 | 2.20 | **2.33** |
| | ($\sigma$) | 0.10 | 0.10 | 0.06 | **0.12** |
| $\mathbf{t_{run}}$ | CW ($\mu$) | 0.26 | 0.39 | 0.76 | **2.19** |
| | ($\sigma$) | 0.01 | 0.01 | 0.01 | **0.003** |
| | RS ($\mu$) | 0.27 | 0.40 | 0.78 | **2.19** |
| | ($\sigma$) | 0.02 | 0.01 | 0.01 | **0.003** |
| | FHT ($\mu$) | 0.29 | 0.43 | 0.85 | **2.19** |
| | ($\sigma$) | 0.04 | 0.05 | 0.10 | **0.003** |
| | RG ($\mu$) | 0.26 | 0.38 | 0.74 | **2.19** |
| | ($\sigma$) | 0.02 | 0.04 | 0.05 | **0.003** |

*NCI60 Dataset:* Results on the NCI60 dataset using four different random projection methods. The table shows how $\epsilon_{in}$, $\gamma$ and $t_{run}$ (in seconds) depend on $r$. See caption of Table VI for an explanation of $mse$, $\beta$, $\mu$ and $\sigma$.

## 5. CONCLUSIONS AND OPEN PROBLEMS

We present theoretical and empirical results indicating that random projections are a useful dimensionality reduction technique for SVM classification and regression problems that handle sparse or

dense data in high-dimensional feature spaces. Our theory predicts that the dimensionality of the projected space (denoted by $r$) has to grow essentially *linearly* (up to logarithmic factors) in $\rho$ (the rank of the data matrix) in order to achieve relative error approximations to the margin and the radius of the minimum ball enclosing the data. Such relative-error approximations imply excellent generalization performance. However, our experiments show that considerably smaller values for $r$ results in classification that is essentially as accurate as running SVMs on all available features, despite the fact that the matrices have full numerical rank. This seems to imply that our theoretical results can be improved. We implemented and tested random projection methods that work well on dense matrices (the RS and FHT methods of Section 2), as well as a very recent random projection method that works well with sparse matrices (the CW method of Section 2). We also experimented with different SVM solvers for lage and medium-scale datasets. As expected, FHT, RG and RS work well on dense data while CW is an excellent choice for sparse data, as indicated by the SVM classification experiments. For large-scale sparse data, CW is the method of choice as the other methods outweigh the benefits of performing random projections. For SVM regression experiments, the combined running times using the four methods are the same for dense datasets. The mean squared error and the squared correlation-coefficient of $\epsilon_{in}$ of the projected data are non-zero as opposed to the SVM classification experiments. Finally, we compare random projections with a popular method of dimensionality reduction, namely PCA and see that the combined running time of random projection and SVM is faster than that of SVM and PCA, with a slightly worse out-of-sample error. All our experiments are on matrices of approximately low-rank, while the theory holds for matrices of exactly low-rank. It is not known if the theory extends to matrices of approximately low rank. This is an open problem and needs further investigation.

**REFERENCES**

D. Achlioptas. 2003. Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *J. Comput. System Sci.* 66, 4 (2003), 671–687.

N. Ailon and B. Chazelle. 2006. Approximate nearest neighbors and the fast Johnson-Lindenstrauss transform. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing*. 557–563.

N. Ailon and E. Liberty. 2008. Fast dimension reduction using Rademacher series on dual BCH codes. In *Proceedings of the 19th Annual ACM-SIAM Symposium on Discrete Algorithms*. 1–9.

J.L. Balcazar, Y. Dai, and O. Watanabe. 2001. A Random Sampling Technique for Training Support Vector Machines. In *Proceedings of the 12th International Conference on Algorithmic Learning Theory*. 119–134.

J.L. Balcazar, Y. Dai, and O. Watanabe. 2002. Provably Fast Support Vector Regression using Random Sampling. In *In Proceedings of SIAM Workshop in Discrete Mathematics and Data Mining*.

A. Blum. 2006. Random Projection, Margins, Kernels, and Feature-Selection. In *In Proceedings of the International Conference on Subspace, Latent Structure and Feature Selection*. 52–68.

D. Cai, X. He, J. Han, and H-J. Zhang. 2006. Orthogonal Laplacianfaces for Face Recognition. *IEEE Transactions on Image Processing* 15, 11 (2006), 3608–3614.

C-C. Chang and C-J. Lin. 2011. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* 2 (2011), 27:1–27:27. Issue 3. Software available at http://www.csie.ntu.edu.tw/~cjlin/libsvm.

K.L. Clarkson and D.W. Woodruff. 2013. Low Rank Approximation and Regression in Input Sparsity Time. In *Proceedings of the 45th ACM Symposium on the Theory of Computing*.

K. Crammer and Y. Singer. 2000. On the learnability and design of output codes for multi-class problems.. In *In Computational Learning Theory*. 35–46.

N. Cristianini and J. Shawe-Taylor. 2000. *Support Vector Machines and other kernel-based learning methods*. Cambridge University Press.

S. Dasgupta and A. Gupta. 2003. An elementary proof of a theorem of Johnson and Lindenstrauss. *Random Structures and Algorithms* 22, 1 (2003), 60–65.

D. Davidov, E. Gabrilovich, and S. Markovitch. 2004. Parameterized generation of labeled datasets for text categorization based on a hierarchical directory. In *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. 250–257. http://techtc.cs.technion.ac.il/techtc300/techtc300.html.

P. Drineas, M.W. Mahoney, S. Muthukrishnan, and T. Sarlos. 2011. Faster least squares approximation. *Numerische. Math.* 117, 2 (2011), 219–249.

R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. 2008. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research* (2008), 1871 –1874.

P. Indyk and R. Motwani. 1998. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the 30th Annual ACM Symposium on Theory of Computing*. 604–613.

V. Jethava, K. Suresh, C. Bhattacharyya, and R. Hariharan. 2009. Randomized Algorithms for Large scale SVMs. *CoRR* abs/0909.3609 (2009). http://arxiv.org/abs/0909.3609.

S. Krishnan, C. Bhattacharyya, and R. Hariharan. 2008. A Randomized Algorithm for Large Scale Support Vector Learning. In *Advances in 20th Neural Information Processing Systems*. 793–800.

D. D. Lewis, Y. Yang, T. G. Rose, and F. Li. 2004. RCV1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research* (2004), 361–397.

J.Z. Li, D.M. Absher, H. Tang, A.M. Southwick, A.M. Casto, S. Ramachandran, H.M. Cann, G.S. Barsh, M. Feldman, L.L. Cavalli-Sforza, and R.M. Myers. 2008. Worldwide human relationships inferred from genome-wide patterns of variation. *Science* 319, 5866 (2008), 1100–1104.

P. Li, T.J. Hastie, and W.K. Church. 2006. Very sparse random projections. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 287–296.

A. Magen and A. Zouzias. 2011. Low rank matrix-valued Chernoff bounds and approximate matrix multiplication. In *Proceedings of the 22nd Annual ACM-SIAM Symposium on Discrete Algorithms*. 1422–1436.

X. Meng and M.W. Mahoney. 2013. Low-distortion Subspace Embeddings in Input-Sparsity Time and Applications to Robust Linear Regression. In *Proceedings of the 45th ACM Symposium on the Theory of Computing*.

J. Nelson and H.L. Nguyen. 2013. OSNAP: Faster numerical linear algebra algorithms via sparser subspace embeddings. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*.

P. Paschou, J. Lewis, A. Javed, and P. Drineas. 2010. Ancestry informative markers for fine-scale individual assignment to worldwide populations. *Journal of Medical Genetics* 47, 12 (2010), 835–47.

S. Paul, C. Boutsidis, M. Magdon-Ismail, and P. Drineas. 2013. Random Projections for Support Vector Machines. In *Proceedings of the 16th International Conference on Artificial Intelligence & Statistics, JMLR W& CP*. 498 –506.

D.T. Ross, U. Scherf, M.B. Eisen, C.M. Perou, P. Spellman, V. Iyer, S.S. Jeffrey, M. Van de Rijn, M. Waltham, A. Pergamenschikov, J.C.F Lee, D. Lashkari, D. Shalon, T.G. Myers, J.N. Weinstein, D. Botstein, and P.O. Brown. 2000. Systematic Variation in Gene Expression Patterns in Human Cancer Cell Lines. *Nature Genetics* 24, 3 (2000), 227–234.

Q. Shi, J. Petterson, G. Dror, J. Langford, A. Smola, and S.V.N. Vishwanathan. 2009. Hash Kernels for Structured Data. *Journal of Machine Learning Research* 10 (2009), 2615–2637.

Q. Shi, C. Shen, R. Hill, and A.V.D Hengel. 2012. Is margin preserved after random projection ?. In *Proceedings of 29th International Conference on Machine Learning*. 591–598.

V.N. Vapnik and A. Chervonenkis. 1971. On the Uniform Convergence of Relative Frequencies of Events to their Probabilities. *Theory of Probability and its Applications* 16 (1971), 264–280.

V. N. Vapnik. 1998. Statistical Learning Theory. *Theory of Probability and its Applications* 16 (1998), 264–280.

L. Zhang, M. Mahdavi, R. Jin, and T. Yang. 2013. Recovering Optimal Solution by Dual Random Projection. In *Conference on Learning Theory (COLT) JMLR W & CP*, Vol. 30. 135–157. http://arxiv.org/abs/1211.3046.