
Sparse Random Features Algorithm as Coordinate Descent in Hilbert Space

Ian E.H. Yen¹ Ting-Wei Lin² Shou-De Lin² Pradeep Ravikumar¹ Inderjit S. Dhillon¹

Department of Computer Science

1: University of Texas at Austin, 2: National Taiwan University

1: {ianyeh, pradeepr, nderjit}@cs.utexas.edu,

2: {b97083, sdlin}@csie.ntu.edu.tw

Abstract

In this paper, we propose a Sparse Random Features algorithm, which learns a sparse non-linear predictor by minimizing an ℓ_1 -regularized objective function over the Hilbert Space induced from a kernel function. By interpreting the algorithm as Randomized Coordinate Descent in an infinite-dimensional space, we show the proposed approach converges to a solution within ϵ -precision of that using an exact kernel method, by drawing $O(1/\epsilon)$ random features, in contrast to the $O(1/\epsilon^2)$ convergence achieved by current Monte-Carlo analyses of Random Features. In our experiments, the Sparse Random Feature algorithm obtains a sparse solution that requires less memory and prediction time, while maintaining comparable performance on regression and classification tasks. Moreover, as an approximate solver for the infinite-dimensional ℓ_1 -regularized problem, the randomized approach also enjoys better convergence guarantees than a Boosting approach in the setting where the greedy Boosting step cannot be performed exactly.

1 Introduction

Kernel methods have become standard for building non-linear models from simple feature representations, and have proven successful in problems ranging across classification, regression, structured prediction and feature extraction [16, 20]. A caveat however is that they are not scalable as the number of training samples increases. In particular, the size of the models produced by kernel methods scale linearly with the number of training samples, even for sparse kernel methods like support vector machines [17]. This makes the corresponding training and prediction computationally prohibitive for large-scale problems.

A line of research has thus been devoted to kernel approximation methods that aim to preserve predictive performance, while maintaining computational tractability. Among these, Random Features has attracted considerable recent interest due to its simplicity and efficiency [2, 3, 4, 5, 10, 6]. Since first proposed in [2], and extended by several works [3, 4, 5, 10], the Random Features approach is a sampling based approximation to the kernel function, where by drawing D features from the distribution induced from the kernel function, one can guarantee uniform convergence of approximation error to the order of $O(1/\sqrt{D})$. On the flip side, such a rate of convergence suggests that in order to achieve high precision, one might need a large number of random features, which might lead to model sizes even larger than that of the vanilla kernel method.

One approach to remedy this problem would be to employ feature selection techniques to prevent the model size from growing linearly with D . A simple way to do so would be by adding ℓ_1 -regularization to the objective function, so that one can simultaneously increase the number of random features D , while selecting a compact subset of them with non-zero weight. However, the resulting algorithm cannot be justified by existing analyses of Random Features, since the Representer theorem does not hold for the ℓ_1 -regularized problem [15, 16]. In other words, since the prediction

cannot be expressed as a linear combination of kernel evaluations, a small error in approximating the kernel function cannot correspondingly guarantee a small prediction error.

In this paper, we propose a new interpretation of Random Features that justifies its usage with ℓ_1 -regularization — yielding the Sparse Random Features algorithm. In particular, we show that the Sparse Random Feature algorithm can be seen as Randomized Coordinate Descent (RCD) in the Hilbert Space induced from the kernel, and by taking D steps of coordinate descent, one can achieve a solution comparable to exact kernel methods within $O(1/D)$ precision in terms of the objective function. Note that the surprising facet of this analysis is that in the finite-dimensional case, the iteration complexity of RCD increases with number of dimensions [18], which would trivially yield a bound going to infinity for our infinite-dimensional problem. In our experiments, the Sparse Random Features algorithm obtains a sparse solution that requires less memory and prediction time, while maintaining comparable performance on regression and classification tasks with various kernels. Note that our technique is complementary to that proposed in [10], which aims to reduce the cost of evaluating and storing basis functions, while our goal is to reduce the number of basis functions in a model.

Another interesting aspect of our algorithm is that our infinite-dimensional ℓ_1 -regularized objective is also considered in the literature of Boosting [7, 8], which can be interpreted as greedy coordinate descent in the infinite-dimensional space. As an approximate solver for the ℓ_1 -regularized problem, we compare our randomized approach to the boosting approach in theory and also in experiments. As we show, for basis functions that do not allow exact greedy search, a randomized approach enjoys better guarantees.

2 Problem Setup

We are interested in estimating a prediction function $f: \mathcal{X} \rightarrow \mathcal{Y}$ from training data set $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $(\mathbf{x}_n, y_n) \in \mathcal{X} \times \mathcal{Y}$ by solving an optimization problem over some Reproducing Kernel Hilbert Space (RKHS) \mathcal{H} :

$$f^* = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \quad \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 + \frac{1}{N} \sum_{n=1}^N L(f(\mathbf{x}_n), y_n), \quad (1)$$

where $L(z, y)$ is a convex loss function with Lipschitz-continuous derivative satisfying $|L'(z_1, y) - L'(z_2, y)| \leq \beta |z_1 - z_2|$, which includes several standard loss functions such as the *square-loss* $L(z, y) = \frac{1}{2}(z - y)^2$, *square-hinge loss* $L(z, y) = \max(1 - zy, 0)^2$ and *logistic loss* $L(z, y) = \log(1 + \exp(-yz))$.

2.1 Kernel and Feature Map

There are two ways in practice to specify the space \mathcal{H} . One is via specifying a positive-definite kernel $k(\mathbf{x}, \mathbf{y})$ that encodes similarity between instances, and where \mathcal{H} can be expressed as the completion of the space spanned by $\{k(\mathbf{x}, \cdot)\}_{\mathbf{x} \in \mathcal{X}}$, that is,

$$\mathcal{H} = \left\{ f(\cdot) = \sum_{i=1}^K \alpha_i k(\mathbf{x}_i, \cdot) \mid \alpha_i \in \mathbb{R}, \mathbf{x}_i \in \mathcal{X} \right\}.$$

The other way is to find an explicit feature map $\{\bar{\phi}_h(\mathbf{x})\}_{h \in H}$, where each $h \in H$ defines a basis function $\bar{\phi}_h(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}$. The RKHS \mathcal{H} can then be defined as

$$\mathcal{H} = \left\{ f(\cdot) = \int_{h \in H} w(h) \bar{\phi}_h(\cdot) dh = \langle \mathbf{w}, \bar{\phi}(\cdot) \rangle_{\mathcal{H}} \mid \|f\|_{\mathcal{H}}^2 < \infty \right\}, \quad (2)$$

where $w(h)$ is a weight distribution over the basis $\{\phi_h(\mathbf{x})\}_{h \in H}$. By Mercer's theorem [1], every positive-definite kernel $k(\mathbf{x}, \mathbf{y})$ has a decomposition s.t.

$$k(\mathbf{x}, \mathbf{y}) = \int_{h \in H} p(h) \phi_h(\mathbf{x}) \phi_h(\mathbf{y}) dh = \langle \bar{\phi}(\mathbf{x}), \bar{\phi}(\mathbf{y}) \rangle_{\mathcal{H}}, \quad (3)$$

where $p(h) \geq 0$ and $\bar{\phi}_h(\cdot) = \sqrt{p(h)} \phi_h(\cdot)$, denoted as $\bar{\phi} = \sqrt{p} \circ \phi$. However, the decomposition is not unique. One can derive multiple decompositions from the same kernel $k(\mathbf{x}, \mathbf{y})$ based on

different sets of basis functions $\{\phi_h(\mathbf{x})\}_{h \in H}$. For example, in [2], the Laplacian kernel $k(\mathbf{x}, \mathbf{y}) = \exp(-\gamma\|\mathbf{x} - \mathbf{y}\|_1)$ can be decomposed through both the Fourier basis and the Random Binning basis, while in [7], the Laplacian kernel can be obtained from the integrating of an infinite number of decision trees.

On the other hand, multiple kernels can be derived from the same set of basis functions via different distribution $p(h)$. For example, in [2, 3], a general decomposition method using Fourier basis functions $\{\phi_\omega(\mathbf{x}) = \cos(\omega^T \mathbf{x})\}_{\omega \in \mathbb{R}^d}$ was proposed to find feature map for any *shift-invariant* kernel of the form $k(\mathbf{x} - \mathbf{y})$, where the feature maps (3) of different kernels $k(\Delta)$ differ only in the distribution $p(\omega)$ obtained from the Fourier transform of $k(\Delta)$. Similarly, [5] proposed decomposition based on polynomial basis for any dot-product kernel of the form $k(\langle \mathbf{x}, \mathbf{y} \rangle)$.

2.2 Random Features as Monte-Carlo Approximation

The standard kernel method, often referred to as the “kernel trick,” solves problem (1) through the Representer Theorem [15, 16], which states that the optimal decision function $f^* \in \mathcal{H}$ lies in the span of training samples $\mathcal{H}_{\mathcal{D}} = \left\{ f(\cdot) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \cdot) \mid \alpha_n \in \mathbb{R}, (\mathbf{x}_n, y_n) \in \mathcal{D} \right\}$, which reduces the infinite-dimensional problem (1) to a finite-dimensional problem with N variables $\{\alpha_n\}_{n=1}^N$. However, it is known that even for loss functions with dual-sparsity (e.g. hinge-loss), the number of non-zero α_n increases linearly with data size [17].

Random Features has been proposed as a kernel approximation method [2, 3, 10, 5], where a Monte-Carlo approximation

$$k(\mathbf{x}_i, \mathbf{x}_j) = E_{p(h)}[\phi_h(\mathbf{x}_i)\phi_h(\mathbf{x}_j)] \approx \frac{1}{D} \sum_{k=1}^D \phi_{h_k}(\mathbf{x}_i)\phi_{h_k}(\mathbf{x}_j) = \mathbf{z}(\mathbf{x}_i)^T \mathbf{z}(\mathbf{x}_j) \quad (4)$$

is used to approximate (3), so that the solution to (1) can be obtained by

$$\mathbf{w}_{RF} = \underset{\mathbf{w} \in \mathbb{R}^D}{\operatorname{argmin}} \quad \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{N} \sum_{n=1}^N L(\mathbf{w}^T \mathbf{z}(\mathbf{x}_n), y_n). \quad (5)$$

The corresponding approximation error

$$\left| \mathbf{w}_{RF}^T \mathbf{z}(\mathbf{x}) - f^*(\mathbf{x}) \right| = \left| \sum_{n=1}^N \alpha_n^{RF} \mathbf{z}(\mathbf{x}_n)^T \mathbf{z}(\mathbf{x}) - \sum_{n=1}^N \alpha_n^* k(\mathbf{x}_n, \mathbf{x}) \right|, \quad (6)$$

as proved in [2, Appendix B], can be bounded by ϵ given $D = \Omega(1/\epsilon^2)$ number of random features, which is a direct consequence of the uniform convergence of the sampling approximation (4). Unfortunately, the rate of convergence suggests that to achieve small approximation error ϵ , one needs significant amount of random features, and since the model size of (5) grows linearly with D , such an algorithm might not obtain a sparser model than kernel method. On the other hand, the ℓ_1 -regularized Random-Feature algorithm we are proposing aims to minimize loss with a selected subset of random feature that neither grows linearly with D nor with N . However, (6) does not hold for ℓ_1 -regularization, and thus one cannot transfer guarantee from kernel approximation (4) to the learned decision function.

3 Sparse Random Feature as Coordinate Descent

In this section, we present the Sparse Random Features algorithm and analyze its convergence by interpreting it as a fully-corrective randomized coordinate descent in a Hilbert space. Given a feature map of orthogonal basic functions $\{\bar{\phi}_h(\mathbf{x}) = \sqrt{p(h)}\phi_h(\mathbf{x})\}_{h \in H}$, the optimization program (1) can be written as the infinite-dimensional optimization problem

$$\min_{\mathbf{w} \in \mathcal{H}} \quad \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N L(\langle \mathbf{w}, \bar{\phi}(\mathbf{x}_n) \rangle_{\mathcal{H}}, y_n). \quad (7)$$

Instead of directly minimizing (7), the *Sparse Random Features* algorithm optimizes the related ℓ_1 -regularized problem defined as

$$\min_{\bar{\mathbf{w}} \in \mathcal{H}} F(\bar{\mathbf{w}}) = \lambda \|\bar{\mathbf{w}}\|_1 + \frac{1}{N} \sum_{n=1}^N L(\langle \bar{\mathbf{w}}, \phi(\mathbf{x}_n) \rangle_{\mathcal{H}}, y_n), \quad (8)$$

where $\bar{\phi}(\mathbf{x}) = \sqrt{\bar{\mathbf{p}}} \circ \phi(\mathbf{x})$ is replaced by $\phi(\mathbf{x})$ and $\|\bar{\mathbf{w}}\|_1$ is defined as the ℓ_1 -norm in function space $\|\bar{\mathbf{w}}\|_1 = \int_{h \in H} |\bar{w}(h)| dh$. The whole procedure is depicted in Algorithm 1. At each iteration, we draw R coordinates h_1, h_2, \dots, h_R from distribution $p(h)$, add them into a working set A^t , and minimize (8) w.r.t. the working set A^t as

$$\min_{\bar{w}(h), h \in A^t} \lambda \sum_{h \in A^t} |\bar{w}(h)| + \frac{1}{N} \sum_{n=1}^N L\left(\sum_{h \in A^t} \bar{w}(h) \phi_h(\mathbf{x}_n), y_n\right). \quad (9)$$

At the end of each iteration, the algorithm removes features with zero weight to maintain a compact working set.

Algorithm 1 Sparse Random-Feature Algorithm

Initialize $\bar{\mathbf{w}}^0 = \mathbf{0}$, working set $A^{(0)} = \{\}$, and $t = 0$.

repeat

1. Sample h_1, h_2, \dots, h_R i.i.d. from distribution $p(h)$.
2. Add h_1, h_2, \dots, h_R to the set $A^{(t)}$.
3. Obtain $\bar{\mathbf{w}}^{t+1}$ by solving (9).
4. $A^{(t+1)} = A^{(t)} \setminus \{h \mid \bar{w}^{t+1}(h) = 0\}$.
5. $t \leftarrow t + 1$.

until $t = T$

3.1 Convergence Analysis

In this section, we analyze the convergence behavior of Algorithm 1. The analysis comprises of two parts. First, we estimate the number of iterations Algorithm 1 takes to produce a solution \mathbf{w}^t that is at most ϵ away from some arbitrary reference solution \mathbf{w}^{ref} on the ℓ_1 -regularized program (8). Then, by taking \mathbf{w}^{ref} as the optimal solution \mathbf{w}^* of (7), we obtain an approximation guarantee for \mathbf{w}^t with respect to \mathbf{w}^* . The proofs for most lemmas and corollaries will be in the appendix.

Lemma 1. *Suppose loss function $L(z, y)$ has β -Lipschitz-continuous derivative and $|\phi_h(\mathbf{x})| \leq B, \forall h \in \mathcal{H}, \forall \mathbf{x} \in \mathcal{X}$. The loss term $Loss(\bar{\mathbf{w}}; \phi) = \frac{1}{N} \sum_{n=1}^N L(\langle \bar{\mathbf{w}}, \phi(\mathbf{x}_n) \rangle, y_n)$ in (8) has*

$$Loss(\bar{\mathbf{w}} + \eta \delta_h; \phi) - Loss(\bar{\mathbf{w}}; \phi) \leq g_h \eta + \frac{\gamma}{2} \eta^2,$$

where $\delta_h = \delta(\|x - h\|)$ is a Dirac function centered at h , and $g_h = \nabla_{\bar{\mathbf{w}}} Loss(\bar{\mathbf{w}}; \phi)(h)$ is the Frechet derivative of the loss term evaluated at h , and $\gamma = \beta B^2$.

The above lemma states smoothness of the loss term, which is essential to guarantee descent amount obtained by taking a coordinate descent step. In particular, we aim to express the expected progress made by Algorithm 1 as the *proximal-gradient* magnitude of $\bar{F}(\mathbf{w}) = F(\sqrt{\bar{\mathbf{p}}} \circ \mathbf{w})$ defined as

$$\bar{F}(\mathbf{w}) = \lambda \|\sqrt{\bar{\mathbf{p}}} \circ \mathbf{w}\|_1 + \frac{1}{N} \sum_{n=1}^N L(\langle \mathbf{w}, \bar{\phi}(\mathbf{x}_n) \rangle, y_n). \quad (10)$$

. Let $\mathbf{g} = \nabla_{\bar{\mathbf{w}}} Loss(\bar{\mathbf{w}}; \phi)$, $\bar{\mathbf{g}} = \nabla_{\mathbf{w}} Loss(\mathbf{w}, \bar{\phi})$ be the gradients of loss terms in (8), (10) respectively, and let $\rho \in \partial(\lambda \|\bar{\mathbf{w}}\|_1)$. We have following relations between (8) and (10):

$$\bar{\rho} := \sqrt{\bar{\mathbf{p}}} \circ \rho \in \partial(\lambda \|\sqrt{\bar{\mathbf{p}}} \circ \mathbf{w}\|_1), \quad \bar{\mathbf{g}} = \sqrt{\bar{\mathbf{p}}} \circ \mathbf{g}, \quad (11)$$

by simple applications of the chain rule. We then analyze the progress made by each iteration of Algorithm 1. Recalling that we used R to denote the number of samples drawn in step 1 of our algorithm, we will first assume $R = 1$, and then show that same result holds also for $R > 1$.

Theorem 1 (Descent Amount). *The expected descent of the iterates of Algorithm 1 satisfies*

$$E[F(\bar{\mathbf{w}}^{t+1})] - F(\bar{\mathbf{w}}^t) \leq -\frac{\gamma\|\bar{\boldsymbol{\eta}}^t\|^2}{2}, \quad (12)$$

where $\bar{\boldsymbol{\eta}}$ is the proximal gradient of (10), that is,

$$\bar{\boldsymbol{\eta}} = \underset{\boldsymbol{\eta}}{\operatorname{argmin}} \quad \lambda\|\sqrt{\bar{\mathbf{p}}} \circ (\mathbf{w}^t + \boldsymbol{\eta})\|_1 - \lambda\|\sqrt{\bar{\mathbf{p}}} \circ \mathbf{w}^t\|_1 + \langle \bar{\mathbf{g}}, \boldsymbol{\eta} \rangle + \frac{\gamma}{2}\|\boldsymbol{\eta}\|^2 \quad (13)$$

and $\bar{\mathbf{g}} = \nabla_{\mathbf{w}} \operatorname{Loss}(\mathbf{w}^t, \bar{\boldsymbol{\phi}})$ is the derivative of loss term w.r.t. \mathbf{w} .

Proof. Let $g_h = \nabla_{\bar{\mathbf{w}}} \operatorname{Loss}(\bar{\mathbf{w}}^t, \boldsymbol{\phi})(h)$. By Corollary 1, we have

$$F(\bar{\mathbf{w}}^t + \eta \boldsymbol{\delta}_h) - F(\bar{\mathbf{w}}^t) \leq \lambda|\bar{w}^t(h) + \eta| - \lambda|\bar{w}^t(h)| + g_h \eta + \frac{\gamma}{2}\eta^2. \quad (14)$$

Minimizing RHS w.r.t. η , the minimizer η_h should satisfy

$$g_h + \rho_h + \gamma\eta_h = 0 \quad (15)$$

for some sub-gradient $\rho_h \in \partial(\lambda|\bar{w}^t(h) + \eta_h|)$. Then by definition of sub-gradient and (15) we have

$$\lambda|\bar{w}^t(h) + \eta| - \lambda|\bar{w}^t(h)| + g_h \eta + \frac{\gamma}{2}\eta^2 \leq \rho_h \eta_h + g_h \eta_h + \frac{\gamma}{2}\eta_h^2 \quad (16)$$

$$= -\gamma\eta_h^2 + \frac{\gamma}{2}\eta_h^2 = -\frac{\gamma}{2}\eta_h^2. \quad (17)$$

Note the equality in (16) holds if $\bar{w}^t(h) = 0$ or the optimal $\eta_h = 0$, which is true for Algorithm 1. Since $\bar{\mathbf{w}}^{t+1}$ minimizes (9) over a block A^t containing h , we have $F(\bar{\mathbf{w}}^{t+1}) \leq F(\bar{\mathbf{w}}^t + \eta_h \boldsymbol{\delta}_h)$. Combining (14) and (16), taking expectation over h on both sides, and then we have

$$E[F(\bar{\mathbf{w}}^{t+1})] - F(\bar{\mathbf{w}}^t) \leq -\frac{\gamma}{2}E[\eta_h^2] = -\frac{\gamma}{2}\|\sqrt{\bar{\mathbf{p}}} \circ \boldsymbol{\eta}\|^2 = -\frac{\gamma}{2}\|\bar{\boldsymbol{\eta}}\|^2$$

Then it remains to verify that $\bar{\boldsymbol{\eta}} = \sqrt{\bar{\mathbf{p}}} \circ \boldsymbol{\eta}$ is the proximal gradient (13) of $\bar{F}(\mathbf{w}^t)$, which is true since $\bar{\boldsymbol{\eta}}$ satisfies the optimality condition of (13)

$$\bar{\mathbf{g}} + \bar{\boldsymbol{\rho}} + \gamma\bar{\boldsymbol{\eta}} = \sqrt{\bar{\mathbf{p}}} \circ (\mathbf{g} + \boldsymbol{\rho} + \gamma\boldsymbol{\eta}) = \mathbf{0},$$

where first equality is from (11) and the second is from (15). \square

Theorem 2 (Convergence Rate). *Given any reference solution \mathbf{w}^{ref} , the sequence $\{\mathbf{w}^t\}_{t=1}^{\infty}$ satisfies*

$$E[\bar{F}(\mathbf{w}^t)] \leq \bar{F}(\mathbf{w}^{ref}) + \frac{2\gamma\|\mathbf{w}^{ref}\|^2}{k}, \quad (18)$$

where $k = \max\{t - c, 0\}$ and $c = \frac{2(\bar{F}(\mathbf{0}) - \bar{F}(\mathbf{w}^{ref}))}{\gamma\|\mathbf{w}^{ref}\|^2}$ is a constant.

Proof. First, the equality actually holds in inequality (16), since for $h \notin A^{(t-1)}$, we have $w^t(h) = 0$, which implies $\lambda|w^t(h) + \eta| - \lambda|w^t(h)| = \rho\eta$, $\rho \in \partial(\lambda|w^t(h) + \eta|)$, and for $h \in A^{(t-1)}$ we have $\bar{\eta}_h = 0$, which gives 0 to both LHS and RHS. Therefore, we have

$$-\frac{\gamma}{2}\|\bar{\boldsymbol{\eta}}\|^2 = \min_{\boldsymbol{\eta}} \quad \lambda\|\sqrt{\bar{\mathbf{p}}} \circ (\mathbf{w}^t + \boldsymbol{\eta})\|_1 - \lambda\|\sqrt{\bar{\mathbf{p}}} \circ \mathbf{w}^t\|_1 + \bar{\mathbf{g}}^T \boldsymbol{\eta} + \frac{\gamma}{2}\|\boldsymbol{\eta}\|^2. \quad (19)$$

Note the minimization in (19) is separable for different coordinates. For $h \in A^{(t-1)}$, the weight $w^t(h)$ is already optimal in the beginning of iteration t , so we have $\bar{\rho}_h + \bar{g}_h = 0$ for some $\bar{\rho}_h \in \partial(|\sqrt{p(h)}w(h)|)$. Therefore, $\eta_h = 0$, $h \in A^{(t-1)}$ is optimal both to $(|\sqrt{p(h)}(w(h) + \eta_h)| + \bar{g}_h \eta_h)$ and to $\frac{\gamma}{2}\eta_h^2$. Set $\eta_h = 0$ for the latter, we have

$$\begin{aligned} -\frac{\gamma}{2}\|\bar{\boldsymbol{\eta}}\|^2 &= \min_{\boldsymbol{\eta}} \left\{ \lambda\|\sqrt{\bar{\mathbf{p}}} \circ (\mathbf{w}^t + \boldsymbol{\eta})\|_1 - \lambda\|\sqrt{\bar{\mathbf{p}}} \circ \mathbf{w}^t\|_1 + \langle \bar{\mathbf{g}}, \boldsymbol{\eta} \rangle + \frac{\gamma}{2} \int_{h \notin A^{(t-1)}} \eta_h^2 dh \right\} \\ &\leq \min_{\boldsymbol{\eta}} \left\{ \bar{F}(\mathbf{w}^t + \boldsymbol{\eta}) - \bar{F}(\mathbf{w}^t) + \frac{\gamma}{2} \int_{h \notin A^{(t-1)}} \eta_h^2 dh \right\} \end{aligned}$$

from convexity of $\bar{F}(\mathbf{w})$. Consider solution of the form $\boldsymbol{\eta} = \alpha(\mathbf{w}^{ref} - \mathbf{w}^t)$, we have

$$\begin{aligned} -\frac{\gamma}{2}\|\bar{\boldsymbol{\eta}}\|^2 &\leq \min_{\alpha \in [0,1]} \left\{ \bar{F}(\mathbf{w}^t + \alpha(\mathbf{w}^{ref} - \mathbf{w}^t)) - \bar{F}(\mathbf{w}^t) + \frac{\gamma\alpha^2}{2} \int_{h \notin A^{(t-1)}} (w^{ref}(h) - w^t(h))^2 dh \right\} \\ &\leq \min_{\alpha \in [0,1]} \left\{ \bar{F}(\mathbf{w}^t) + \alpha(\bar{F}(\mathbf{w}^{ref}) - \bar{F}(\mathbf{w}^t)) - \bar{F}(\mathbf{w}^t) + \frac{\gamma\alpha^2}{2} \int_{h \notin A^{(t-1)}} w^{ref}(h)^2 dh \right\} \\ &\leq \min_{\alpha \in [0,1]} \left\{ -\alpha(\bar{F}(\mathbf{w}^t) - \bar{F}(\mathbf{w}^{ref})) + \frac{\gamma\alpha^2}{2} \|\mathbf{w}^{ref}\|^2 \right\}, \end{aligned}$$

where the second inequality results from $w^t(h) = 0, h \notin A^{(t-1)}$. Minimizing last expression w.r.t. α , we have $\alpha^* = \min\left(\frac{\bar{F}(\mathbf{w}^t) - \bar{F}(\mathbf{w}^{ref})}{\gamma\|\mathbf{w}^{ref}\|^2}, 1\right)$ and

$$-\frac{\gamma}{2}\|\bar{\boldsymbol{\eta}}\|^2 \leq \begin{cases} -(\bar{F}(\mathbf{w}^t) - \bar{F}(\mathbf{w}^{ref}))^2 / (2\gamma\|\mathbf{w}^{ref}\|^2) & , \text{ if } \bar{F}(\mathbf{w}^t) - \bar{F}(\mathbf{w}^{ref}) < \gamma\|\mathbf{w}^{ref}\|^2 \\ -\frac{\gamma}{2}\|\mathbf{w}^{ref}\|^2 & , \text{ o.w.} \end{cases} \quad (20)$$

Note, since the function value $\{\bar{F}(\mathbf{w}^t)\}_{t=1}^\infty$ is non-increasing, only iterations in the beginning fall in second case of (20), and the number of such iterations is at most $c = \lceil \frac{2(\bar{F}(\mathbf{0}) - \bar{F}(\mathbf{w}^{ref}))}{\gamma\|\mathbf{w}^{ref}\|^2} \rceil$. For $t > c$, we have

$$E[\bar{F}(\mathbf{w}^{t+1})] - \bar{F}(\mathbf{w}^t) \leq -\frac{\gamma\|\bar{\boldsymbol{\eta}}^t\|_2^2}{2} \leq -\frac{(\bar{F}(\mathbf{w}^t) - \bar{F}(\mathbf{w}^{ref}))^2}{2\gamma\|\mathbf{w}^{ref}\|^2}. \quad (21)$$

The recursion then leads to the result. \square

Note the above bound does not yield useful result if $\|\mathbf{w}^{ref}\|^2 \rightarrow \infty$. Fortunately, the optimal solution of our target problem (7) has finite $\|\mathbf{w}^*\|^2$ as long as in (7) $\lambda > 0$, so it always give a useful bound when plugged into (18), as following corollary shows.

Corollary 1 (Approximation Guarantee). *The output of Algorithm 1 satisfies*

$$E \left[\lambda \|\bar{\mathbf{w}}^{(D)}\|_1 + \text{Loss}(\bar{\mathbf{w}}^{(D)}; \boldsymbol{\phi}) \right] \leq \{ \lambda \|\mathbf{w}^*\|_2 + \text{Loss}(\mathbf{w}^*; \bar{\boldsymbol{\phi}}) \} + \frac{2\gamma\|\mathbf{w}^*\|_2^2}{D'} \quad (22)$$

with $D' = \max\{D - c, 0\}$, where \mathbf{w}^* is the optimal solution of problem (7), c is a constant defined in Theorem 2.

Then the following two corollaries extend the guarantee (22) to any $R \geq 1$, and a bound holds with high probability. The latter is a direct result of [18, Theorem 1] applied to the recursion (21).

Corollary 2. *The bound (22) holds for any $R \geq 1$ in Algorithm 1, where if there are T iterations then $D = TR$.*

Corollary 3. *For $D \geq \frac{2\gamma\|\mathbf{w}^*\|_2^2}{\epsilon}(1 + \log \frac{1}{\rho}) + 2 - \frac{4}{c} + c$, the output of Algorithm 1 has*

$$\lambda \|\bar{\mathbf{w}}^{(D)}\|_1 + \text{Loss}(\bar{\mathbf{w}}^{(D)}; \boldsymbol{\phi}) \leq \{ \lambda \|\mathbf{w}^*\|_2 + \text{Loss}(\mathbf{w}^*; \bar{\boldsymbol{\phi}}) \} + \epsilon \quad (23)$$

with probability $1 - \rho$, where c is as defined in Theorem 2 and \mathbf{w}^* is the optimal solution of (7).

3.2 Relation to the Kernel Method

Our result (23) states that, for D large enough, the Sparse Random Features algorithm achieves either a comparable loss to that of the vanilla kernel method, or a model complexity (measured in ℓ_1 -norm) less than that of kernel method (measured in ℓ_2 -norm). Furthermore, since \mathbf{w}^* is not the optimal solution of the ℓ_1 -regularized program (8), it is possible for the LHS of (23) to be much smaller than the RHS. On the other hand, since any \mathbf{w}^* of finite ℓ_2 -norm can be the reference solution \mathbf{w}^{ref} , the λ used in solving the ℓ_1 -regularized problem (8) can be different from the λ used in the kernel method. The tightest bound is achieved by minimizing the RHS of (23), which is equivalent to minimizing (7) with some unknown $\tilde{\lambda}(\lambda)$ due to the difference of $\|\mathbf{w}\|_1$ and $\|\mathbf{w}\|_2^2$. In practice, we can follow a regularization path to find small enough λ that yields comparable predictive performance while maintains model as compact as possible. Note, when using different sampling distribution $p(h)$ from the decomposition (3), our analysis provides different bounds (23) for the Randomized Coordinate Descent in Hilbert Space. This is in contrast to the analysis in the finite-dimensional case, where RCD with different sampling distribution converges to the same solution [18].

3.3 Relation to the Boosting Method

Boosting is a well-known approach to minimize infinite-dimensional problems with ℓ_1 -regularization [8, 9], and which in this setting, performs greedy coordinate descent on (8). For each iteration t , the algorithm finds the coordinate $h^{(t)}$ yielding steepest descent in the loss term

$$h^{(t)} = \underset{h \in H}{\operatorname{argmin}} \quad \frac{1}{N} \sum_{n=1}^N L'_n \phi_h(\mathbf{x}_n) \quad (24)$$

to add into a working set A^t and minimize (8) w.r.t. A^t . When the greedy step (24) can be solved exactly, Boosting has fast convergence to the optimal solution of (8) [13, 14]. On the contrary, randomized coordinate descent can only converge to a sub-optimal solution in finite time when there are infinite number of dimensions. However, in practice, only a very limited class of basis functions allow the greedy step in (24) to be performed exactly. For most basis functions (weak learners) such as perceptrons and decision trees, the greedy step (24) can only be solved approximately. In such cases, Boosting might have no convergence guarantee, while the randomized approach is still guaranteed to find a comparable solution to that of the kernel method. In our experiments, we found that the randomized coordinate descent performs considerably better than approximate Boosting with the perceptron basis functions (weak learners), where as adopted in the Boosting literature [19, 8], a convex surrogate loss is used to solve (24) approximately.

4 Experiments

In this section, we compare Sparse Random Features (Sparse-RF) to the existing Random Features algorithm (RF) and the kernel method (Kernel) on regression and classification problems with kernels set to Gaussian RBF, Laplacian RBF [2], and Perceptron kernel [7]¹. For Gaussian and Laplacian RBF kernel, we use Fourier basis function with corresponding distribution $p(h)$ derived in [2]; for Perceptron kernel, we use perceptron basis function with distribution $p(h)$ being uniform over unit-sphere as shown in [7]. For regression, we solve kernel ridge regression (1) and RF regression (6) in closed-form as in [10] using Eigen, a standard C++ library of numerical linear algebra. For Sparse-RF, we solve the LASSO sub-problem (9) by standard RCD algorithm. In classification, we use LIBSVM² as solver of kernel method, and use Newton-CG method and Coordinate Descent method in LIBLINEAR [12] to solve the RF approximation (6) and Sparse-RF sub-problem (9) respectively. We set $\lambda_N = N\lambda = 1$ for the kernel and RF methods, and for Sparse-RF, we choose $\lambda_N \in \{1, 10, 100, 1000\}$ that gives RMSE (accuracy) closest to the RF method to compare sparsity and efficiency. The results are in Tables 1 and 2, where the cost of kernel method grows at least quadratically in the number of training samples. For *YearPred*, we use $D = 5000$ to maintain tractability of the RF method. Note for *Covtype* dataset, the ℓ_2 -norm $\|\mathbf{w}^*\|_2$ from kernel machine is significantly larger than that of others, so according to (22), a larger number of random features D are required to obtain similar performance, as shown in Figure 1.

In Figure 1, we compare Sparse-RF (randomized coordinate descent) to Boosting (greedy coordinate descent) and the bound (23) obtained from SVM with Perceptron kernel and basis function (weak learner). The figure shows that Sparse-RF always converges to a solution comparable to that of the kernel method, while Boosting with approximate greedy steps (using convex surrogate loss) converges to a higher objective value, due to bias from the approximation.

Acknowledgement

S.-D.Lin acknowledges the support of Telecommunication Lab., Chunghwa Telecom Co., Ltd via TL-103-8201, AOARD via No. FA2386-13-1-4045, Ministry of Science and Technology, National Taiwan University and Intel Co. via MOST102-2911-I-002-001, NTU103R7501, 102-2923-E-002-007-MY2, 102-2221-E-002-170, 103-2221-E-002-104-MY2. P.R. acknowledges the support of ARO via W911NF-12-1-0390 and NSF via IIS-1149803, IIS-1320894, IIS-1447574, and DMS-1264033. This research was also supported by NSF grants CCF-1320746 and CCF-1117055.

²Data set for classification can be downloaded from *LIBSVM data set* web page, and data set for regression can be found at UCI Machine Learning Repository and Ali Rahimi's page for the paper [2].

²We follow the FAQ page of LIBSVM to replace hinge-loss by square-hinge-loss for comparison.

Table 1: Results for Kernel Ridge Regression. Fields from top to bottom are model size (# of support vectors or # of random features or # of non-zero weights respectively), testing RMSE, training time, testing prediction time, and memory usage during training.

Data set	Gaussian RBF			Laplacian RBF			Perceptron Kernel		
	Kernel	RF	Sparse-RF	Kernel	RF	Sparse-RF	Kernel	RF	Sparse-RF
CPU $N_{tr} = 6554$ $N_t = 819$ $d = 21$	SV=6554 RMSE=0.038 $T_{tr}=154$ s $T_t=2.59$ s Mem=1.36 G	D=10000 0.037 875 s 6 s 4.71 G	NZ=57 0.032 22 s 0.04 s 0.069 G	SV=6554 0.034 157 s 3.13 s 1.35 G	D=10000 . 0.035 803 s 6.99 s 4.71 G	NZ=289 0.027 43 s 0.18 s 0.095 G	SV=6554 0.026 151 s 2.48 s 1.36 G	D=10000 0.038 776 s 6.37 s 4.71 G	NZ=251 0.027 27 s 0.13 s 0.090 G
Census $N_{tr} = 18186$ $N_t = 2273$ $d = 119$	SV=18186 RMSE= 0.029 $T_{tr}=2719$ s $T_t=74$ s Mem=10 G	D=10000 0.032 1615 s 80 s 8.2 G	NZ=1174 0.030 229 s 8.6 s 0.55 G	SV=18186 0.146 3268 s 68 s 10 G	D=10000 0.168 1633 s 88 s 8.2 G	NZ=5269 0.179 225 s 38s 1.7 G	SV=18186 0.010 2674 s 67.45 s 10 G	D=10000 0.016 1587 s 76 s 8.2 G	NZ=976 0.016 185 s 6.7 s 0.49 G
YearPred $N_{tr} = 463715$ $N_t = 51630$ $d = 90$	SV=# RMSE=# $T_{tr}=#$ $T_t=#$ Mem=#	D=5000 0.103 7697 s 697 s 76.7G	NZ=1865 0.104 1618 s 97 s 45.6G	SV=# # # # #	D=5000 0.286 9417 s 715 s 76.6 G	NZ=3739 0.273 1453 s 209 s 54.3 G	SV=# # # # #	D=5000 0.105 8636 s 688 s 76.7 G	NZ=896 0.105 680 s 51 s 38.1 G

Table 2: Results for Kernel Support Vector Machine. Fields from top to bottom are model size (# of support vectors or # of random features or # of non-zero weights respectively), testing accuracy, training time, testing prediction time, and memory usage during training.

Data set	Gaussian RBF			Laplacian RBF			Perceptron Kernel		
	Kernel	RF	Sparse-RF	Kernel	RF	Sparse-RF	Kernel	RF	Sparse-RF
Cod-RNA $N_{tr} = 59535$ $N_t = 10000$ $d = 8$	SV=14762 Acc= 0.966 $T_{tr}=95$ s $T_t=15$ s Mem=3.8 G	D=10000 0.964 214 s 56 s 9.5 G	NZ=180 0.964 180 s 0.61 s 0.66 G	SV=13769 0.971 290 s 15 s 3.6 G	D=10000 . 0.969 137 s 46 s 9.6 G	NZ=1195 0.970 137 s 6.41 s 1.8 G	SV=15201 0.967 57.34 s 7.01 s 3.6 G	D=10000 0.964 197 s 71.9 s 9.6 G	NZ=1148 0.963 131 s 3.81 s 1.4 G
IJCNN $N_{tr} = 127591$ $N_t = 14100$ $d = 22$	SV=16888 Acc= 0.991 $T_{tr}=636$ s $T_t=34$ s Mem=12 G	D=10000 0.989 601 s 88 s 20 G	NZ=1392 0.989 292 s 11 s 7.5 G	SV=16761 0.995 988 s 34 s 12 G	D=10000 0.992 379 s 86 s 20 G	NZ=2508 0.992 566 s 25 s 9.9 G	SV=26563 0.991 634 s 16 s 11 G	D=10000 0.987 381 s 77 s 20 G	NZ=1530 0.988 490 s 11 s 7.8 G
Covtype $N_{tr} = 464810$ $N_t = 116202$ $d = 54$	SV=335606 Acc= 0.849 $T_{tr}=74891$ s $T_t=3012$ s Mem=78.5 G	D=10000 0.829 9909 s 735 s 74.7 G	NZ=3421 0.836 6273 s 132 s 28.1 G	SV=224373 0.954 64172 s 2004 s 80.8 G	D=10000 0.888 10170 s 635 s 74.6 G	NZ=3141 0.869 2788 s 175 s 56.5 G	SV=358174 0.905 79010 s 1774 s 80.5 G	D=10000 0.835 6969 s 664 s 74.7 G	NZ=1401 0.836 1706 s 70 s 44.4 G

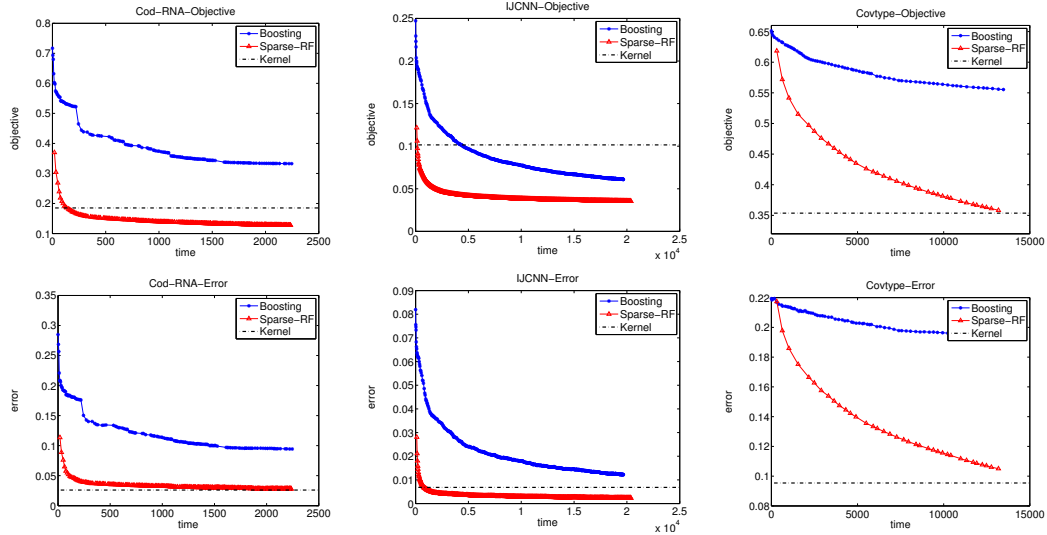


Figure 1: The ℓ_1 -regularized objective (8) (top) and error rate (bottom) achieved by *Sparse Random Feature* (randomized coordinate descent) and *Boosting* (greedy coordinate descent) using *perceptron* basis function (weak learner). The dashed line shows the ℓ_2 -norm plus loss achieved by kernel method (RHS of (22)) and the corresponding error rate using *perceptron kernel* [7].

References

- [1] Mercer, J. Functions of positive and negative type and their connection with the theory of integral equations. Royal Society London, A 209:415 446, 1909.
- [2] Rahimi, A. and Recht, B. Random features for large-scale kernel machines. NIPS 20, 2007.
- [3] Rahimi, A. and Recht, B. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. NIPS 21, 2008.
- [4] Vedaldi, A., Zisserman, A.: Efficient additive kernels via explicit feature maps. In CVPR. (2010)
- [5] P. Kar and H. Karnick. Random feature maps for dot product kernels. In Proceedings of AIS-TATS'12, pages 583 591, 2012.
- [6] T. Yang, Y.-F. Li, M. Mahdavi, R. Jin, and Z.-H. Zhou. Nystrom method vs. random Fourier features: A theoretical and empirical comparison. In Adv. NIPS, 2012.
- [7] Hsuan-Tien Lin, Ling Li, Support Vector Machinery for Infinite Ensemble Learnings. JMLR 2008.
- [8] Saharon Rosset, Ji Zhu, and Trevor Hastie. Boosting as a Regularized Path to a Maximum Margin Classifier. JMLR, 2004.
- [9] Saharon Rosset, Grzegorz Swirszcz, Nathan Srebro, and Ji Zhu. ℓ_1 -regularization in infinite dimensional feature spaces. In Learning Theory: 20th Annual Conference on Learning Theory, 2007.
- [10] Q. Le, T. Sarlos, and A. J. Smola. Fastfood - approximating kernel expansions in loglinear time. In The 30th International Conference on Machine Learning, 2013.
- [11] C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. ACM Transactions on Intelligent Systems and Technology, 2011.
- [12] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification. Journal of Machine Learning Research, 9:1871 1874, 2008.
- [13] Gunnar Ratsch, Sebastian Mika, and Manfred K. Warmuth. On the convergence of leveraging. In NIPS, 2001.
- [14] Matus Telgarsky. The Fast Convergence of Boosting. In NIPS, 2011.
- [15] Kimeldorf, G. S. and Wahba, G. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. Annals of Mathematical Statistics, 41:495502, 1970.
- [16] Scholkopf, Bernhard and Smola, A. J. Learning with Kernels. MIT Press, Cambridge, MA, 2002.
- [17] Steinwart, Ingo and Christmann, Andreas. Support Vector Machines. Springer, 2008.
- [18] P. Ricktarik and M. Takac, Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function, School of Mathematics, University of Edinburgh, Tech. Rep., 2011.
- [19] Chen, S.-T., Lin, H.-T. and Lu, C.-J. An online boosting algorithm with theoretical justifications. ICML 2012.
- [20] Taskar, B., Guestrin, C., and Koller, D. Max-margin Markov networks. NIPS 16, 2004.
- [21] G. Song et.al. Reproducing kernel Banach spaces with the l_1 norm. Journal of Applied and Computational Harmonic Analysis, 2011.