
Incremental Randomized Sketching for Online Kernel Learning

Xiao Zhang¹ Shizhong Liao¹

Abstract

Randomized sketching has been used in offline kernel learning, but it cannot be applied directly to online kernel learning due to the lack of incremental maintenances for randomized sketches with regret guarantees. To address these issues, we propose a novel incremental randomized sketching approach for online kernel learning, which has efficient incremental maintenances with theoretical guarantees. We construct two incremental randomized sketches using the sparse transform matrix and the sampling matrix for kernel matrix approximation, update the incremental randomized sketches using rank-1 modifications, and construct a time-varying explicit feature mapping for online kernel learning. We prove that the proposed incremental randomized sketching is statistically unbiased for the matrix product approximation, obtains a $1 + \epsilon$ relative-error bound for the kernel matrix approximation, enjoys a sublinear regret bound for online kernel learning, and has constant time and space complexities at each round for incremental maintenances. Experimental results demonstrate that the incremental randomized sketching achieves a better learning performance in terms of accuracy and efficiency even in adversarial environments.

1. Introduction

Online kernel learning must rise to the challenge of incremental maintenance in a low computational complexity with accuracy and convergence guarantees at each round. Online gradient descent is a judicious choice for online kernel learning due to its efficiency and effectiveness (Bottou, 2010; Shalev-Shwartz et al., 2011). However, the kernelization of the online gradient descent may cause a linear growth with respect to the model size and increase the hypothesis

updating time significantly (Wang et al., 2012). To solve the problem, budgeted online kernel learning has been proposed to maintain a buffer of support vectors (SVs) with a fixed budget, which limits the model size to reduce the computational complexities (Crammer et al., 2003). Several budget maintenance strategies were introduced to bound the buffer size of SVs. Dekel et al. (2005) maintained the budget by discarding the oldest support vector. Orabona et al. (2008) projected the new example onto the linear span of SVs in the feature space to reduce the buffer size. But the buffer of SVs cannot be used for kernel matrix approximation directly. Besides, it is difficult to provide a lower bound of the buffer size that is key to the theoretical guarantees for budgeted online kernel learning.

Unlike the budget maintenance strategies of SVs, sketches of the kernel matrix have been introduced into online kernel learning to approximate the kernel incrementally. Engel et al. (2004) used the approximate linear dependency (ALD) rule to construct the sketches for kernel approximation in kernel regression. But the kernel matrix approximation error is not compared with the best rank- k approximation. Sun et al. (2012) proved that the size of sketches constructed by ALD rule grows sublinearly with respect to the data size when the eigenvalues decay fast enough. Recently, randomized sketches of kernel matrix were used as surrogates to perform the offline kernel approximation with finite sample bounds (Sarlós, 2006; Woodruff, 2014; Raskutti & Mahoney, 2015), but few efforts have been made to formulate suitable randomized sketches for online kernel learning. Existing online kernel learning approaches using randomized sketches usually adopted a random sampling strategy. Lu et al. (2016) applied the Nyström approach for kernel function approximation and constructed the explicit feature mapping using the singular value decomposition, but it fixed the randomized sketches after initialization, which is equivalent to the offline Nyström approach using uniform sampling that may suffer from a low approximation accuracy in adversarial environments due to the non-uniformity structure of the data (Kumar et al., 2009). Calandriello et al. (2017b) assumed that the kernel function can be expressed as an inner product in the explicit feature space, and proposed a second order online kernel learning approach using online sampling strategy, but the construction of the explicit feature mapping needs to be further explored. Calandriello et al. (2017a)

¹College of Intelligence and Computing, Tianjin University, Tianjin 300350, China. Correspondence to: Shizhong Liao <szliao@tju.edu.cn>.

formulated an explicit feature mapping using a Nyström approximation, and performed effective second-order updates for hypothesis updating, which enjoys a logarithmic regret and has a per-step cost that grows with the effective dimension.

In this paper, we propose an incremental randomized sketching approach for online kernel learning. The proposed incremental randomized sketching inherits the efficient incremental maintenance for matrix approximations and explicit feature mappings, has a constant time complexity at each round, achieves a $1 + \epsilon$ relative-error for kernel matrix approximation, and enjoys a sublinear regret for online kernel learning, therefore presenting an effective and efficient sketching approach to online kernel learning.

2. Notations and Preliminaries

Let $[\mathbf{A}]_{*i}$ and $[\mathbf{A}]_{i*}$ denote the i -th column of the matrix \mathbf{A} and the column vector consisting of the entries from the i -th row of \mathbf{A} respectively, $\|\mathbf{A}\|_F$ and $\|\mathbf{A}\|_2$ the Frobenius and spectral norm of \mathbf{A} respectively, and \mathbf{A}^\dagger the Moore-Penrose pseudoinverse of \mathbf{A} . Let $[M] = \{1, 2, \dots, M\}$, $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^T \subseteq (\mathcal{X} \times \mathcal{Y})^T$ be the sequence of T instances, where $\mathcal{X} \subseteq \mathbb{R}^D$ and $\mathcal{Y} = \{-1, 1\}$. We denote the loss function by $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+ \cup \{0\}$, the kernel function by $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and its corresponding kernel matrix by $\mathbf{K} = (\kappa(\mathbf{x}_i, \mathbf{x}_j)) \in \mathbb{R}^{T \times T}$, the i -th eigenvalue of \mathbf{K} by $\lambda_i(\mathbf{K})$, $\lambda_1(\mathbf{K}) \geq \lambda_2(\mathbf{K}) \geq \dots \geq \lambda_T(\mathbf{K})$, and the reproducing kernel Hilbert space (RKHS) associated with κ by $\mathcal{H}_\kappa = \text{span}\{\kappa(\cdot, \mathbf{x}) : \mathbf{x} \in \mathcal{X}\}$.

2.1. Kernel Matrix Approximation

Given a kernel matrix $\mathbf{K} \in \mathbb{R}^{T \times T}$, the sketched kernel matrix approximation problem in the offline setting is described as follows (Wang et al., 2016):

$$\mathbf{F}_{\text{fast}} = \arg \min_{\mathbf{F} \in \mathbb{R}^{B \times B}} \|\mathbf{S}^\top (\mathbf{C} \mathbf{F} \mathbf{C}^\top - \mathbf{K}) \mathbf{S}\|_F^2,$$

from which we can obtain the approximate kernel matrix $\mathbf{K}_{\text{fast}} = \mathbf{C} \mathbf{F}_{\text{fast}} \mathbf{C}^\top \approx \mathbf{K}$, where $\mathbf{S} \in \mathbb{R}^{T \times s}$ is a sketch matrix, $\mathbf{C} \in \mathbb{R}^{T \times B}$ and $s, B > 0$ are the sketch sizes. The Nyström approach uses the column selection matrix \mathbf{P} as the sketch matrix (Williams & Seeger, 2001), i.e., $\mathbf{S} = \mathbf{P} \in \mathbb{R}^{T \times B}$, and formulate \mathbf{C} as $\mathbf{C} = \mathbf{K} \mathbf{P}$, which yields the approximate kernel matrix $\mathbf{K}_{\text{ny}} = \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^\top$, where $\mathbf{W} = \mathbf{P}^\top \mathbf{K} \mathbf{P} \in \mathbb{R}^{B \times B}$. While $\mathbf{S} = \mathbf{I}_T \in \mathbb{R}^{T \times T}$, the modified Nyström approach (also called the prototype model) (Wang & Zhang, 2013) is obtained as

$$\mathbf{K}_{\text{mod}} = \mathbf{C} \mathbf{F}_{\text{mod}} \mathbf{C}^\top \approx \mathbf{K},$$

where $\mathbf{F}_{\text{mod}} = \mathbf{C}^\dagger \mathbf{K} (\mathbf{C}^\top)^\dagger$, which is more accurate but more computationally expensive than the standard Nyström. But in the context of online kernel learning, the size of

the kernel matrix increases over time, a scenario for which the offline kernel matrix approximation approaches do not readily apply due to the lack of efficient incremental maintenances with regret guarantees.

2.2. Randomized Sketches

Randomized sketches via hash functions can be described in general using hash-based sketch matrices. Without loss of generality, we assume that the sketch size s is divisible by d . Let h_1, \dots, h_d be hash functions from $\{1, \dots, n\}$ to $\{1, \dots, s/d\}$, and $\sigma_1, \dots, \sigma_d$ the other d hash functions. We denote the hash-based sketch submatrix by $\mathbf{S}_k \in \mathbb{R}^{n \times (s/d)}$ and the hash-based sketch matrix by $\mathbf{S} = [\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_d] \in \mathbb{R}^{n \times s}$, where $[\mathbf{S}_k]_{i,j} = \sigma_k(i)$ for $j = h_k(i)$ and $[\mathbf{S}_k]_{i,j} = 0$ for $j \neq h_k(i)$, $k = 1, \dots, d$. The randomized sketch of $\mathbf{A} \in \mathbb{R}^{T \times n}$ is denoted by $\mathbf{A} \mathbf{S} \in \mathbb{R}^{T \times s}$. Sketch matrices of different randomized sketches can be specified by different parameters s, d and functions $\{\sigma_k\}_{k=1}^d$ as follows.

Count Sketch Matrix (Charikar et al., 2004): $d = 1$ and $\sigma_1(i) \in \{-1, +1\}$ is 2-wise independent hash function.

Sparsier Johnson-Lindenstrauss Transform (SJLT) (Kane & Nelson, 2014): For the block construction of SJLT, $\sigma_k(i) \in \{-1/\sqrt{d}, +1/\sqrt{d}\}$ is $O(\log T)$ -wise independent hash function for $k = 1, \dots, d$, where d is the number of blocks. SJLT is a generalization of the Count Sketch matrix.

3. Novel Incremental Randomized Sketching

In this section, we propose a novel incremental randomized sketching approach. Specifically, we construct the incremental randomized sketches for the kernel matrix approximation, formulate the incremental maintenance for the incremental randomized sketches, and build a time-varying explicit feature mapping. In the online setting, at round $t + 1$, a new example \mathbf{x}_{t+1} arrives and the kernel matrix $\mathbf{K}^{(t+1)}$ can be represented as a bordered matrix as follows:

$$\mathbf{K}^{(t+1)} = \begin{bmatrix} \mathbf{K}^{(t)} & \boldsymbol{\psi}^{(t+1)} \\ \boldsymbol{\psi}^{(t+1)\top} & \xi^{(t+1)} \end{bmatrix} \in \mathbb{R}^{(t+1) \times (t+1)},$$

where $\xi^{(t+1)} = \kappa(\mathbf{x}_{t+1}, \mathbf{x}_{t+1})$ and

$$\boldsymbol{\psi}^{(t+1)} = [\kappa(\mathbf{x}_{t+1}, \mathbf{x}_1), \kappa(\mathbf{x}_{t+1}, \mathbf{x}_2), \dots, \kappa(\mathbf{x}_{t+1}, \mathbf{x}_t)]^\top.$$

To approximate the kernel matrix incrementally, we introduce a sketch matrix $\mathbf{S}_p^{(t+1)}$ to reduce the complexity of the problem (2) and another sketch matrix $\mathbf{S}_m^{(t+1)}$ to reduce the size of the approximate kernel matrix, where $\mathbf{S}_p^{(t+1)} \in \mathbb{R}^{(t+1) \times s_p}$ is an SJLT (defined in Section 2.2) and $\mathbf{S}_m^{(t+1)} \in \mathbb{R}^{(t+1) \times s_m}$ is a sub-sampling matrix. Using the two sketch matrices, we first formulate the *incremental randomized sketches* of $\mathbf{K}^{(t+1)}$ as follows:

$$\boldsymbol{\Phi}_{\text{pm}}^{(t+1)} = \mathbf{S}_p^{(t+1)\top} \mathbf{C}_m^{(t+1)} \quad \text{and} \quad \boldsymbol{\Phi}_{\text{pp}}^{(t+1)} = \mathbf{S}_p^{(t+1)\top} \mathbf{C}_p^{(t+1)},$$

where

$$C_p^{(t+1)} = K^{(t+1)} S_p^{(t+1)} \quad \text{and} \quad C_m^{(t+1)} = K^{(t+1)} S_m^{(t+1)}.$$

Then $K^{(t+1)}$ can be approximated by

$$K_{\text{sk}}^{(t+1)} = C_m^{(t+1)} F_{\text{sk}}^{(t+1)} C_m^{(t+1)\top} \approx K^{(t+1)}, \quad (1)$$

where $F_{\text{sk}}^{(t+1)} \in \mathbb{R}^{s_m \times s_m}$ is obtained by solving the following sketched matrix approximation problem at round $t+1$

$$\begin{aligned} F_{\text{sk}}^{(t+1)} &= \arg \min_F \left\| S_p^{(t+1)\top} E_F^{(t+1)} S_p^{(t+1)} \right\|_F^2 \\ &= \left(\Phi_{\text{pm}}^{(t+1)} \right)^\dagger \Phi_{\text{pp}}^{(t+1)} \left(\Phi_{\text{pm}}^{(t+1)\top} \right)^\dagger, \end{aligned} \quad (2)$$

where $E_F^{(t+1)} = C_m^{(t+1)} F C_m^{(t+1)\top} - K^{(t+1)}$.

Next, we design the incremental maintenance for the incremental randomized sketches. We partition the sketch matrices into block matrices as

$$S_p^{(t+1)} = \left[S_p^{(t)\top}, s_p^{(t+1)} \right]^\top, \quad S_m^{(t+1)} = \left[S_m^{(t)\top}, s_m^{(t+1)} \right]^\top,$$

where $s_m^{(t+1)} \in \mathbb{R}^{s_m}$ is a sampling vector, $s_p^{(t+1)} \in \mathbb{R}^{s_p}$ is a vector that contains d nonzero entries determined by d different hash mappings in SJLT. Then we can update the sketch $\Phi_{\text{pm}}^{(t+1)} \in \mathbb{R}^{s_p \times s_m}$ incrementally by rank-1 modifications

$$\Phi_{\text{pm}}^{(t+1)} = \Phi_{\text{pm}}^{(t)} + \Delta_{\text{pm}}^{(t+1)},$$

where $\Delta_{\text{pm}}^{(t+1)} = R_{\text{pm}}^{(t+1)} + R_{\text{mp}}^{(t+1)\top} + T_{\text{pm}}^{(t+1)}$ consists of the following three rank-1 matrices

$$\begin{aligned} R_{\text{pm}}^{(t+1)} &= s_p^{(t+1)} \psi^{(t+1)\top} S_m^{(t)}, \\ R_{\text{mp}}^{(t+1)} &= s_m^{(t+1)} \psi^{(t+1)\top} S_p^{(t)}, \\ T_{\text{pm}}^{(t+1)} &= \xi^{(t+1)} s_p^{(t+1)} s_m^{(t+1)\top}. \end{aligned}$$

Similarly, we update the sketch $\Phi_{\text{pp}}^{(t+1)} \in \mathbb{R}^{s_p \times s_p}$ by rank-1 modifications as follows:

$$\Phi_{\text{pp}}^{(t+1)} = \Phi_{\text{pp}}^{(t)} + \Delta_{\text{pp}}^{(t+1)},$$

where $\Delta_{\text{pp}}^{(t+1)} = R_{\text{pp}}^{(t+1)} + R_{\text{pp}}^{(t+1)\top} + T_{\text{pp}}^{(t+1)}$ and the modifications are performed using two rank-1 matrices

$$\begin{aligned} R_{\text{pp}}^{(t+1)} &= s_p^{(t+1)} \psi^{(t+1)\top} S_p^{(t)}, \\ T_{\text{pp}}^{(t+1)} &= \xi^{(t+1)} s_p^{(t+1)} s_p^{(t+1)\top}. \end{aligned}$$

Figure 1 illustrates the kernel matrix approximation using the incremental randomized sketching.

Finally, we construct the time-varying explicit feature mapping using the incremental randomized sketches. We decompose $\Phi_{\text{pp}}^{(t+1)}$ via the rank- k singular value decomposition (SVD) as follows:

$$\Phi_{\text{pp}}^{(t+1)} \approx V^{(t+1)} \Sigma^{(t+1)} V^{(t+1)\top}, \quad (3)$$

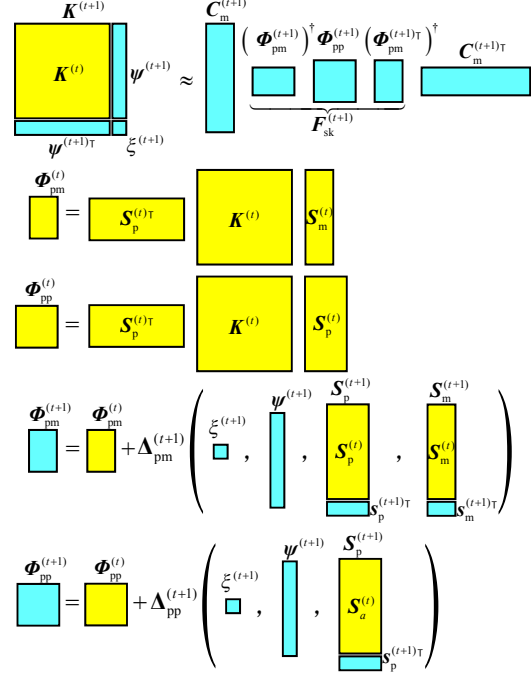


Figure 1. The proposed incremental randomized sketching for kernel matrix approximation at round $t+1$, where $K_{\text{sk}}^{(t+1)} = C_m^{(t+1)} F_{\text{sk}}^{(t+1)} C_m^{(t+1)\top}$ is the approximate kernel matrix, $\Phi_{\text{pm}}^{(t+1)}$ and $\Phi_{\text{pp}}^{(t+1)}$ are two incremental randomized sketches for constructing $F_{\text{sk}}^{(t+1)}$, $\Delta_{\text{pp}}^{(t+1)}$ can be seen as a modification function with input variables $\xi^{(t+1)}$, $\psi^{(t+1)}$, $S_p^{(t+1)}$, and $\Delta_{\text{pm}}^{(t+1)}$ uses an additional input variable $S_m^{(t+1)}$. Since the incremental randomized sketches retain the key information with the efficient rank-1 modifications, our incremental randomized sketching can be computed in constant time and space complexities and enjoys a sublinear regret while applied to online kernel learning (see Section 5).

where $V^{(t+1)} \in \mathbb{R}^{s_p \times k}$, $\Sigma^{(t+1)} \in \mathbb{R}^{k \times k}$ and $\text{rank } k \leq s_p$.

Then $F_{\text{sk}}^{(t+1)}$ in (2) is approximated by

$$F_{\text{sk}}^{(t+1)} \approx Q_{t+1} Q_{t+1}^\top,$$

where

$$Q_{t+1} = \left(\Phi_{\text{pm}}^{(t+1)} \right)^\dagger V^{(t+1)} \left(\Sigma^{(t+1)} \right)^{\frac{1}{2}},$$

which yields the approximate kernel matrix from (1)

$$K_{\text{sk}}^{(t+1)} \approx \left(C_m^{(t+1)} Q_{t+1} \right) \left(C_m^{(t+1)} Q_{t+1} \right)^\top.$$

Further, the time-varying explicit feature mapping can be updated at round $t+1$ by

$$\phi_{t+2}(\cdot) = \left([\kappa(\cdot, \tilde{x}_1), \dots, \kappa(\cdot, \tilde{x}_{s_m})] \right) Q_{t+1}^\top, \quad (4)$$

where $\{\tilde{x}_i\}_{i=1}^{s_m}$ are the sampled examples obtained by $S_m^{(t+1)}$, by which the kernel function value between the examples x_i and x_j can be approximated by $\kappa(x_i, x_j) \approx \langle \phi_{t+2}(x_i), \phi_{t+2}(x_j) \rangle$.

4. Application to Online Kernel Learning

Using the time-varying explicit feature mapping $\phi_t(\cdot)$, we formulate the hypothesis $f_t(\mathbf{x}_t)$ for prediction at round t , which is an approximation of the optimal hypothesis on the first $t - 1$ instances¹:

$$f_t(\mathbf{x}_t) = \langle \mathbf{w}_t, \phi_t(\mathbf{x}_t) \rangle \approx \sum_{i \in [t-1]} \alpha_i^* \kappa(\mathbf{x}_t, \mathbf{x}_i),$$

where \mathbf{w}_t is the weight vector and $\{\alpha_i^*\}_{i \in [t-1]}$ is the set of the optimal coefficients of the optimal hypothesis over the first $t - 1$ instances. When applying the hypotheses $\{f_t\}_{t \in [T]}$ to online kernel learning, the learning procedure consists of the following two stages. At the first stage, we use a buffer \mathcal{V}_t to store the examples with nonzero losses, restrict the size of the buffer under a fixed budget B ($s_p, s_m < B \ll T$), and perform the Kernelized Online Gradient Descent (KOGD) (Kivinen et al., 2001) until the size of \mathcal{V}_t reaches the budget B . Assuming that T_0 rounds have passed while the budget has been reached, let $\mathbf{S}_m^{(T_0)}$ be the uniform sampling matrix, we sample $\{\tilde{\mathbf{x}}_i\}_{i=1}^{s_m}$ using $\mathbf{S}_m^{(T_0)}$ and initialize the sketches $\Phi_{pp}^{(T_0)}$, $\Phi_{pm}^{(T_0)}$ and \mathbf{Q}_{T_0} . Then we obtain $\phi_{T_0+1}(\cdot)$ and initialize the hypothesis \mathbf{w}_{T_0+1} such that

$$\mathbf{w}_{T_0+1}^\top \phi_{T_0+1}(\mathbf{x}_{T_0}) = \sum_{\mathbf{x}_i \in \mathcal{V}_{T_0}} \alpha_i \kappa(\mathbf{x}_{T_0}, \mathbf{x}_i),$$

which yields

$$\mathbf{w}_{T_0+1}^\top = \sum_{\mathbf{x}_i \in \mathcal{V}_{T_0}} \alpha_i \kappa(\mathbf{x}_{T_0}, \mathbf{x}_i) \frac{\phi_{T_0+1}(\mathbf{x}_{T_0})^\top}{\|\phi_{T_0+1}(\mathbf{x}_{T_0})\|_2^2},$$

where $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_B]^\top$ is obtained by KOGD. At the second stage, we update the hypothesis using Online Gradient Descent (OGD) (Zinkevich, 2003). Denote the hypothesis at round t using the updated explicit feature mapping $\phi_{t+1}(\cdot)$ by

$$\bar{f}_t(\mathbf{x}_t) = \langle \bar{\mathbf{w}}_t, \phi_{t+1}(\mathbf{x}_t) \rangle. \quad (5)$$

We can obtain $\bar{\mathbf{w}}_t$ by setting $\bar{f}_t(\mathbf{x}_t) = f_t(\mathbf{x}_t)$, which yields

$$\bar{\mathbf{w}}_t^\top = f_t(\mathbf{x}_t) \phi_{t+1}(\mathbf{x}_t)^\dagger = \frac{f_t(\mathbf{x}_t) \cdot \phi_{t+1}(\mathbf{x}_t)^\top}{\|\phi_{t+1}(\mathbf{x}_t)\|_2^2}.$$

Then we update the hypothesis $\bar{f}_t(\cdot) = \langle \bar{\mathbf{w}}_t, \phi_{t+1}(\cdot) \rangle$ in (5) as follows:

$$\mathbf{w}_{t+1} = \bar{\mathbf{w}}_t - \eta \nabla \mathcal{L}_t(\bar{\mathbf{w}}_t), \quad (6)$$

where $\ell_t(\bar{\mathbf{w}}_t) = \ell_t(\bar{f}_t) := \ell(\bar{f}_t(\mathbf{x}_t), y_t)$ and $\mathcal{L}_t(\bar{\mathbf{w}}_t) := \ell_t(\bar{\mathbf{w}}_t) + \frac{\lambda}{2} \|\bar{\mathbf{w}}_t\|_2^2$.

¹It follows from the representer theorem that it is sufficient to approximate the optimal hypothesis of the linear combination form.

To avoid the unbounded sizes of $\psi^{(t+1)}$, $\mathbf{S}_p^{(t+1)}$ and $\mathbf{S}_m^{(t+1)}$, we fix the sampled examples $\{\tilde{\mathbf{x}}_i\}_{i=1}^{s_m}$ in $\phi_t(\cdot)$, $t = T_0 + 1, \dots, T$, and adopt a periodic updating strategy for \mathbf{Q}_t . Specifically, we update \mathbf{Q}_t in (4) once for every ρ examples at the second stage, where $\rho \in [T - B]$ is the update cycle. To reduce the time complexity of the matrix decomposition, we adopt the randomized SVD (Halko et al., 2011) for (3). Besides, we also use the randomized SVD for computing the Moore-Penrose pseudoinverse of $\Phi_{pm}^{(t+1)}$. Finally, we summarize the above stages into Algorithm 1, called sketched online gradient descent (SkeGD).

Algorithm 1 SkeGD Algorithm

Require: sketch sizes s_p and s_m , stepsize η , rank k , budget B , update cycle ρ , number of blocks d , regularization parameter λ

- 1: Initialize $f_1 = 0$ and $\mathcal{V}_1 = \emptyset$
- 2: **for** $t = 1, \dots, T$ **do**
- 3: Receive new example \mathbf{x}_t
- 4: Predict $\hat{y}_t = \text{sgn}(f_t(\mathbf{x}_t))$
- 5: **if** $|\mathcal{V}_t| < B$ **then**
- 6: $\mathcal{V}_{t+1} = \mathcal{V}_t \cup \{\mathbf{x}_t\}$ whenever the loss is nonzero
- 7: Update f_t by kernelized online gradient descent
- 8: **else**
- 9: **if** $|\mathcal{V}_t| = B$ **then**
- 10: Initialize $T_0 = t$
- 11: Initialize $\Phi_{pp}^{(T_0)}$, $\Phi_{pm}^{(T_0)}$ and \mathbf{Q}_{T_0}
- 12: Initialize $\phi_{T_0+1}(\cdot)$ and \mathbf{w}_{T_0+1}
- 13: **else**
- 14: **if** $t \bmod \rho = 1$ **then**
- 15: Update $\Phi_{pp}^{(t)}$ and $\Phi_{pm}^{(t)}$ using rank-1 modifications
- 16: Compute \mathbf{Q}_t using $\Phi_{pp}^{(t)}$ and $\Phi_{pm}^{(t)}$
- 17: $\phi_{t+1}(\mathbf{x}_t) = ([\kappa(\mathbf{x}_t, \tilde{\mathbf{x}}_1), \dots, \kappa(\mathbf{x}_t, \tilde{\mathbf{x}}_{s_m})] \mathbf{Q}_t)^\top$
- 18: $\bar{\mathbf{w}}_t^\top = f_t(\mathbf{x}_t) \phi_{t+1}(\mathbf{x}_t)^\top / \|\phi_{t+1}(\mathbf{x}_t)\|_2^2$
- 19: **else**
- 20: $\mathbf{Q}_t = \mathbf{Q}_{t-1}$, $\Phi_{pp}^{(t)} = \Phi_{pp}^{(t-1)}$, $\Phi_{pm}^{(t)} = \Phi_{pm}^{(t-1)}$
- 21: $\bar{\mathbf{w}}_t = \mathbf{w}_t$
- 22: **end if**
- 23: $\mathbf{w}_{t+1} = \bar{\mathbf{w}}_t - \eta \nabla \mathcal{L}_t(\bar{\mathbf{w}}_t)$
- 24: **end if**
- 25: **end if**
- 26: **end for**

5. Theoretical Analysis

In this section, we demonstrate that the proposed incremental randomized sketching satisfies the low-rank approximation property, analyze the regret for online kernel learning and give the complexity analysis for SkeGD. The detailed proofs can be found in the supplementary material. For convenience, we denote $\mathbf{S}_p^{(T)}$, $\mathbf{S}_m^{(T)}$, $\mathbf{C}_m^{(T)}$, $\mathbf{F}_{sk}^{(T)}$ and $\mathbf{K}^{(T)}$ by \mathbf{S}_p , \mathbf{S}_m , \mathbf{C}_m , \mathbf{F}_{sk} and \mathbf{K} , respectively.

5.1. Low-Rank Approximation Property

We first present the bias and variance analysis for approximating the matrix product using \mathbf{S}_p in the proposed incremental randomized sketches, which shows that our sketched

kernel matrix approximation problem (2) is an unbiased estimate of the modified Nyström kernel matrix approximation problem.

Lemma 1 (Matrix Product Preserving Property). *Let $\mathbf{A} \in \mathbb{R}^{T \times m}$ and $\mathbf{S}_p \in \mathbb{R}^{T \times s_p}$. Then,*

$$\mathbb{E} [\|\mathbf{S}_p^\top \mathbf{A}\|_{\mathbb{F}}^2] = \|\mathbf{A}\|_{\mathbb{F}}^2, \quad \text{Var} [\|\mathbf{S}_p^\top \mathbf{A}\|_{\mathbb{F}}^2] \leq \frac{2}{s_p} \|\mathbf{A}\|_{\mathbb{F}}^4.$$

Proof Sketch. We first partition \mathbf{S}_p into d hash-based sketch submatrices, and demonstrate the expectation and variance while approximating the inner product by \mathbf{S}_p , which yields the unbiasedness for the matrix product approximation. Then we derive the expectation of $\|\mathbf{S}_p^\top \mathbf{A}\|_{\mathbb{F}}^2$ and obtain the variance of $\|\mathbf{S}_p^\top \mathbf{A}\|_{\mathbb{F}}^2$. \square

By Lemma 1, we can prove that the proposed incremental randomized sketching is nearly as accurate as the modified Nyström for kernel matrix approximation, which is more accurate than the standard Nyström approach.

Theorem 1 (Low-Rank Approximation Property). *Let $\mathbf{K} \in \mathbb{R}^{T \times T}$ be a symmetric matrix, $\epsilon_0 \in (0, 1)$. $\mathbf{F}_{\text{sk}} \in \mathbb{R}^{s_m \times s_m}$, $\mathbf{C}_m \in \mathbb{R}^{T \times s_m}$ are matrices defined in (1). Set $\tau = s_m/s_p$ and $d = \Theta(\log^3(s_m))$ for $\mathbf{S}_p \in \mathbb{R}^{T \times s_p}$. Let $\mathbf{U}_m \in \mathbb{R}^{T \times s_m}$ be a matrix with orthonormal columns, $\mathbf{U}_m^\perp \in \mathbb{R}^{T \times (T-s_m)}$ be a matrix satisfying $\mathbf{U}_m \mathbf{U}_m^\top + \mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top = \mathbf{I}_T$ and $\mathbf{U}_m^\top \mathbf{U}_m^\perp = \mathbf{O}$, and*

$$s_p = \Omega(s_m \text{polylog}(s_m \delta_0^{-1}) / \epsilon_0^2),$$

then with probability at least $1 - \delta_0$ all singular values of $\mathbf{S}_p^\top \mathbf{U}_m$ are $1 \pm \epsilon_0$, and with probability at least $1 - \delta$

$$\|\mathbf{C}_m \mathbf{F}_{\text{sk}} \mathbf{C}_m^\top - \mathbf{K}\|_{\mathbb{F}}^2 \leq (1 + \epsilon) \|\mathbf{C}_m \mathbf{F}_{\text{mod}} \mathbf{C}_m^\top - \mathbf{K}\|_{\mathbb{F}}^2,$$

where

$$\sqrt{\epsilon} = 2\tau \sqrt{\frac{T}{\delta_1 \delta_2}} + \sqrt{\frac{2\tau}{\delta_2}} (\epsilon_0^2 + 2\epsilon_0 + 2), \quad (7)$$

δ_i is the failure probability of matrix product preserving as

$$\Pr \left\{ \frac{\|\mathbf{B}_i \mathbf{A}_i - \mathbf{B}_i \mathbf{S}_p \mathbf{S}_p^\top \mathbf{A}_i\|_{\mathbb{F}}^2}{\|\mathbf{B}_i\|_{\mathbb{F}}^2 \|\mathbf{A}_i\|_{\mathbb{F}}^2} > \frac{2}{\delta_i s_p} \right\} \leq \delta_i, \quad i = 1, 2,$$

$\mathbf{A}_1 = \mathbf{U}_m$, $\mathbf{B}_1 = \mathbf{I}_T$, $\mathbf{A}_2 = \mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}$, $\mathbf{B}_2 = \mathbf{U}_m^\top$, and $\delta = \delta_0 + \delta_1 + \delta_2$.

Proof Sketch. Let $\mathbf{A} \in \mathbb{R}^{T \times m}$, $\mathbf{B} \in \mathbb{R}^{p \times T}$. By Lemma 1 and the Markov's inequality, we derive the error bound for approximating $\mathbf{B}\mathbf{A}$ by $\mathbf{B}\mathbf{S}_p \mathbf{S}_p^\top \mathbf{A}$, which yields the bound of the kernel matrix approximation error combined with the singular values of $\mathbf{S}_p \mathbf{U}_m$. Then we bound the singular values of $\mathbf{S}_p \mathbf{U}_m$ and obtain the final upper bound of the approximation error. \square

Remark 1. *Theorem 1 shows that the proposed incremental randomized sketching achieves a relative-error bound for kernel matrix approximation, for the modified Nyström is a $1 + \epsilon'$ relative-error approximation with respect to the best rank- k approximation (Wang & Zhang, 2013). The kernel approximation error of the modified Nyström approach is determined by s_m . Thus, larger s_m leads to a tighter kernel approximation error bound. Besides, from (7) we can obtain $\sqrt{\tau} = \Theta(T^{-\frac{1}{4}} \epsilon^{\frac{1}{4}})$. Clearly, small τ can guarantee the approximation error ϵ . Finally, the choice of s_p depends on τ and s_m , i.e., $s_p = s_m/\tau$.*

5.2. Regret Analysis

Let $\mathbf{K}_{B,\rho} \in \mathbb{R}^{(B+\lfloor(T-B)/\rho\rfloor) \times (B+\lfloor(T-B)/\rho\rfloor)}$ be the intersection matrix of \mathbf{K} which is constructed by $B + \lfloor(T-B)/\rho\rfloor$ examples, $\mu(\mathbf{K}_{B,\rho})$ be the coherence of $\mathbf{K}_{B,\rho}$ as

$$\mu(\mathbf{K}_{B,\rho}) = \frac{B + \lfloor(T-B)/\rho\rfloor}{\text{rank}(\mathbf{K}_{B,\rho})} \max_i \|(U_{B,\rho})_{i,:}\|_2^2,$$

where $U_{B,\rho}$ is the singular vector matrix of $\mathbf{K}_{B,\rho}$. Obviously, $\mu(\mathbf{K}_{B,\rho})$ is independent of T when $\rho = \Theta(T-B)$. We demonstrate the regret bound for online kernel learning using our incremental randomized sketching as follows².

Theorem 2 (Regret Bound). *Let $\mathbf{K} \in \mathbb{R}^{T \times T}$ be a kernel matrix with $\kappa(\mathbf{x}_i, \mathbf{x}_j) \leq 1$, k ($k \leq s_p$) be the rank in the incremental randomized sketches, $\epsilon_0 \in (0, 1)$ and δ_i ($i = 0, 1, 2$) be the failure probabilities defined in Theorem 1. Set the update cycle $\rho = \lfloor\theta(T-B)\rfloor$, $\theta \in (0, 1)$, $d = \Theta(\log^3(s_m))$ and*

$$s_p = \Omega(s_m \text{polylog}(s_m \delta_0^{-1}) / \epsilon_0^2),$$

$$s_m = \Omega(\mu(\mathbf{K}_{B,\rho}) k \log k),$$

for $\mathbf{S}_p \in \mathbb{R}^{T \times s_p}$ and $\mathbf{S}_m \in \mathbb{R}^{T \times s_m}$. Assume $\ell_t(\cdot)$ is a convex loss function that is Lipschitz continuous with Lipschitz constant L , and the eigenvalues of \mathbf{K} decay polynomially with decay rate $\beta > 1$. Let \mathbf{w}_t , $t \in [T]$ be the sequence of hypotheses generated by (6), satisfying $|f_t(\mathbf{x}_t)| = |\langle \mathbf{w}_t, \phi_t(\mathbf{x}_t) \rangle| \leq C_f$, $t \in [T]$. Then, for the optimal hypothesis f^* that minimize $\frac{\lambda}{2} \|f\|_{\mathcal{H}_\kappa}^2 + \frac{1}{T} \sum_{t=1}^T \ell_t(f)$ in the original reproducing kernel Hilbert space \mathcal{H}_κ , with probability at least $1 - \delta$

$$\begin{aligned} \mathcal{R}_{\text{reg}}(T, f^*) &\leq \frac{\|\mathbf{w}_T^*\|_2^2}{2\theta\eta} + \frac{\eta L^2 T}{2} + \frac{\sqrt{1+\epsilon}}{\lambda} O(\sqrt{B}) + \\ &\frac{C_f^2}{2\theta\eta} + \frac{1}{\lambda(\beta-1)} \left(\frac{3}{2} - \frac{B+1/\theta}{T} \right), \end{aligned}$$

where $\mathcal{R}_{\text{reg}}(T, f^*) = \sum_{t=1}^T (\mathcal{L}_t(\mathbf{w}_t) - \mathcal{L}_t(f^*))$, $\delta = \delta_0 + \delta_1 + \delta_2$, ϵ is defined in (7), \mathbf{w}_t^* is the optimal hypothesis

²In the theoretical analysis, we assume $T_0 = B$ and omit KOGD used at the first stage that enjoys a $O(\sqrt{B})$ regret bound.

on the incremental randomized sketches over the first t instances and $Z = \arg \max_{i \in [T]} \|\mathbf{w}_i^*\|_2$.

Proof Sketch. We first decompose $\mathcal{L}_t(\mathbf{w}_t) - \mathcal{L}_t(\mathbf{w}^*)$ into two terms as follows:

$$\underbrace{\mathcal{L}_t(\mathbf{w}_t) - \mathcal{L}_t(\mathbf{w}_t^*)}_{\text{Optimization error}} + \underbrace{\mathcal{L}_t(\mathbf{w}_t^*) - \mathcal{L}_t(\mathbf{w}^*)}_{\text{Estimation error}},$$

where \mathbf{w}^* is the optimal hypothesis on the incremental randomized sketches in hindsight. Then we bound the errors and complete the proof. \square

Remark 2. *The assumption of polynomial decay for eigenvalues of the kernel matrix is a general assumption, met for shift invariant kernels, finite rank kernels and convolution kernels (Liu & Liao, 2015; Belkin, 2018), that is $\exists \beta > 1$ such that $\lambda_i(\mathbf{K}) = O(i^{-\beta})$. Since the regret bound contains the term $1/\theta$ and the term $-1/\theta$ while SkeGD requires an additional $O(\theta^{-2})$ runtime for updating (see Section 5.3), a suitable update cycle $\rho = \lfloor \theta(T - B) \rfloor$ needs to be tuned, which is analyzed in the experiments in adversarial environments. From Remark 1, smaller τ leads to a tighter kernel matrix approximation error bound, which yields a tighter regret bound. Besides, the lower bounds of the sketch size s_p and s_m are independent of the number of rounds T .*

In Theorem 2, we can obtain a $O(\sqrt{T})$ regularized regret bound by setting $\eta = 1/\sqrt{T}$, and achieve a $O(\sqrt{T})$ unregularized regret bound by setting $\eta, \lambda = 1/\sqrt{T}$. The existing regret bound of first-order online kernel learning with Nyström is $O(\sqrt{T})$ while requiring the budget $B = O(T)$ (Lu et al., 2016), but the proposed SkeGD does not need this requirement. Under the same assumption that the eigenvalues decay polynomially with decay rate $\beta > 1$, the regret bound of existing second-order online kernel learning via adaptive sketching is $O(T^{1/\beta} \log^2 T)$ (Calandriello et al., 2017a), which is tighter than our bound when $\beta \gg 2$. But this second-order approach has a $O(T^{2/\beta})$ time complexity at each round, while our SkeGD has a constant time complexity at each round with respect to T (see Section 5.3).

5.3. Complexity Analysis for SkeGD

The running time of SkeGD is mainly consumed by the matrix decomposition and multiplication while updating the incremental randomized sketches. Table 1 summarizes the computational complexities of SkeGD at each round³. At the first stage ($t \leq T_0$), SkeGD has constant time and space complexities per round. At the second stage ($t > T_0$), the time complexity of SkeGD at each updating round (\mathbf{Q}_t is updated) is $O(\nu + s_p^2 + s_p s_m)$, and the space complexity of SkeGD is $O(s_p^2 + s_p s_m + \nu)$, where $\nu = B + \lfloor (T -$

³Since we set $k, d < s_m < s_p$ in the experiments, we omit k and d in the complexity analysis for SkeGD.

$B)/\rho] = B + \lfloor \theta^{-1} \rfloor$ when $\rho = \lfloor \theta(T - B) \rfloor$ ($0 < \theta < 1$). Since \mathbf{Q}_t is updated once for every ρ examples, SkeGD has a constant space complexity, and a $O(s_m)$ time complexity at each round except for $\lfloor \theta^{-1} \rfloor$ updating rounds which are in $O(B + s_p^2 + s_p s_m)$ time complexity. Thus, the overall time complexity of SkeGD at the second stage is

$$O\left(T s_m + \left\lfloor \frac{T - B}{\rho} \right\rfloor (\nu + s_p^2 + s_p s_m)\right).$$

Table 1. The computational complexities of SkeGD at each round, where $\nu = B + \lfloor \theta^{-1} \rfloor$ when $\rho = \lfloor \theta(T - B) \rfloor$ ($0 < \theta < 1$).

Item	Round	Operation	Time	Space
KOGD	$t < T_0$	-	$O(B)$	$O(B)$
$\Phi_{pp}^{(T_0)}$	$t = T_0$	Initialization	$O(B s_p)$	$O(B s_p)$
$\Phi_{pm}^{(T_0)}$	$t = T_0$	Initialization	$O(B s_m)$	$O(B s_m)$
$\Phi_{pp}^{(t)}$	$t > T_0$	Updating	$O(\nu)$	$O(s_p^2 + \nu)$
$\Phi_{pm}^{(t)}$	$t > T_0$	Updating	$O(s_m)$	$O(s_p s_m)$
$\Phi_{pp}^{(t)}$	$t > T_0$	Decomposition	$O(s_p^2)$	$O(s_p^2)$
$\Phi_{pm}^{(t)}$	$t > T_0$	Decomposition	$O(s_p s_m)$	$O(s_p s_m)$
\mathbf{Q}_t	$t > T_0$	Multiplication	$O(s_p s_m)$	$O(s_p s_m)$
$\phi_{t+1}(\mathbf{x})$	$t > T_0$	Multiplication	$O(s_m)$	$O(s_m)$

6. Experiments

In this section, we evaluate the effectiveness and efficiency of the proposed incremental randomized sketching for kernel approximation and online kernel learning.

6.1. Experimental Setups

All experiments are performed on a machine with 4-core Intel Core i7 3.60 GHz CPU and 16GB memory. We compare our SkeGD with the state-of-the-art online kernel learning algorithms on the well-known classification benchmarks⁴. We merge the training and testing data into a single dataset for each benchmark, which covers the number of instances ranging from 1,000 to 581,012. All the experiments are performed over 20 different random permutations of the datasets. The stepsizes η of all the gradient descent based algorithms are tuned in $10^{[-5:+1:0]}$, and the regularization parameters λ are tuned in $10^{[-4:+1:1]}$. Besides, we use the Gaussian kernel $\kappa(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / 2\sigma^2)$, where the set $\sigma \in \{2^{[-5:+0.5:7]}\}$ are adopted as the candidate kernel set. Since kernel alignment (KA) is computationally efficient and independent of learning algorithms (Cortes et al., 2012; Ding et al., 2018), which is suitable for kernel selection of online kernel learning (Zhang & Liao, 2018), we select the optimal kernels using KA for each dataset.

6.2. Kernel Approximation

First, we evaluate the performance of the proposed incremental randomized sketching (IRS) in terms of kernel ap-

⁴<https://www.csie.ntu.edu.tw/~cjlin/libsvm>

proximation. The kernel values are approximated using the explicit feature mapping (4) in our IRS. We compare with the incremental Nyström (INys) approach (Lu et al., 2016), which constructs the fixed intersection matrix using the first B examples in the sequence of T instances. For our IRS, we

Table 2. Results of kernel approximation using IRS and INys.

Approach	german ($B = 100$)		german ($B = 200$)	
	Relative error	Time (s)	Relative error	Time (s)
INys	0.078 ± 0.012	0.054	0.047 ± 0.006	0.081
IRS	0.059 ± 0.008	0.047	0.031 ± 0.003	0.066
Approach	svmguide3 ($B = 100$)		svmguide3 ($B = 200$)	
	Relative error	Time (s)	Relative error	Time (s)
INys	0.037 ± 0.001	0.083	0.031 ± 0.001	0.092
IRS	0.018 ± 0.001	0.065	0.012 ± 0.002	0.077

set $s_p = 3B/4$, $s_m = B/2$, $k = 0.2B$ and $d = 4$. Table 2 shows the relative error $\|\widetilde{\mathbf{K}}^{(T)} - \mathbf{K}\|_F^2 / \|\mathbf{K}\|_F^2$ and running time on *german* and *svmguide3* with different budget B , where $\widetilde{\mathbf{K}}^{(T)}$ is the approximate kernel matrix after T rounds. We can observe that IRS outperforms the INys in terms of the relative approximation error on both datasets. Besides, our IRS performs a randomized SVD and only needs to sample $s_m = B/2$ examples while updating the incremental randomized sketches, hence is more efficient.

6.3. Online Learning under a Fixed Budget

To demonstrate the performance of SkeGD for online kernel learning, we compare SkeGD with the existing sketching and budgeted online kernel learning algorithms under a fixed budget B , including RBP (Cavallanti et al., 2007), Forgetron (Dekel et al., 2005), Projectron and Projectron++ (Orabona et al., 2008), BPA-S (Wang & Vucetic, 2010), BOGD (Zhao et al., 2012), and NOGD (Lu et al., 2016). The compared algorithms are obtained from the LIBOL v0.3.0 toolbox and LSOKL toolbox⁵. For all the algorithms, we adopt the hinge loss function, set $B = 100$ for small datasets ($T \leq 10,000$) and $B = 200$ for other datasets. Besides, we set $\tau = 0.2$, $s_p = 3B/4$, $s_m = \tau s_p$, $d = 4$ and $\rho = 0.3T$ in our SkeGD if not specially specified, and the rank $k = 0.1B$ for NOGD and SkeGD. The mistake rate is used to evaluate the accuracy of online kernel learning, which is computed by $\sum_{t=1}^T I(y_t f_t(\mathbf{x}_t) < 0) / T \times 100$.

Experimental results from Table 3 show that SkeGD consistently performs better than the other algorithms in terms of the mistake rate of online kernel learning. For most datasets, SkeGD is more efficient than the other algorithms under the same budget B , because the time complexity per round of SkeGD is $O(s_m)$ ($s_m < B$) for most rounds as shown in complexity analysis, while the compared algo-

gorithms are in $\Omega(B)$ time complexity at each round. Then,

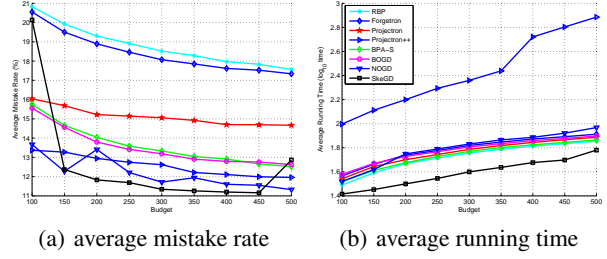


Figure 2. The average mistake rates and average running time w.r.t. the budget B on *cod-rna*.

we set $s_p = 3B/4$, $s_m = B/2$ and study the influence of the budget B . From Figure 2, we observe that the increasing budget B yields lower mistake rate but leads to higher computation time cost. This is because larger budget B means that more examples are preserved to construct the explicit feature mapping and a higher computation cost is needed to obtain a better approximation quality. For fixed budget

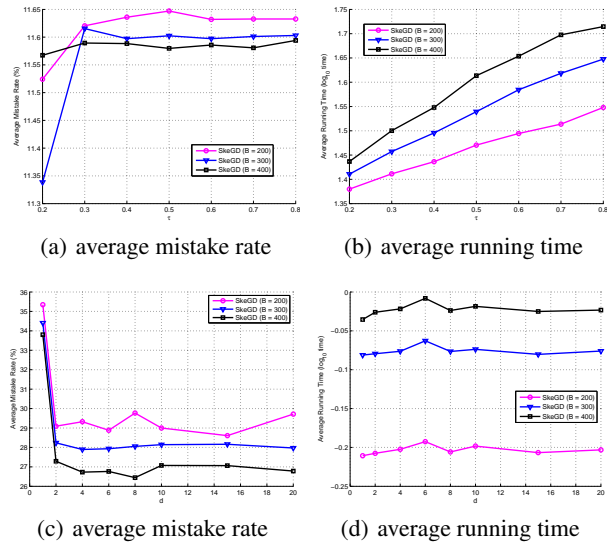


Figure 3. (a), (b): average mistake rate, average running time w.r.t. $\tau = s_m/s_p$ on *cod-rna*; (c), (d): average mistake rate and average running time w.r.t. the number of blocks d on *spambase*.

($B = 200, 300, 400$), we set $s_p = 3B/4$, $s_m = \tau s_p$ and vary τ from 0.2 to 0.8. From Figure 3 (a), (b) we can see that SkeGD with smaller τ yields better empirical performance in terms of the average mistake rate and the average running time, which conforms to the theoretical results. For different B , we set $s_p = 3B/4$, $s_m = B/2$, and explore the influence of d that is the number of blocks in the sketch matrix S_p . The results in Figure 3 (c), (d) verify that large d can significantly improve the accuracy of online kernel

⁵<http://libol.stevenhoi.org/>
<http://lsokl.stevenhoi.org/>

Table 3. Comparisons among RBP, Forgetron, Projectron, Projectron++, BPA-S, BOGD, NOGD and our SkeGD w.r.t. the mistake rates (%) and the running time (s).

Algorithm	german		svmguid3		spambase		a9a	
	Mistake rate	Time	Mistake rate	Time	Mistake rate	Time	Mistake rate	Time
RBP	38.924 ± 1.754	0.071	34.091 ± 1.230	0.080	35.533 ± 0.769	0.391	28.297 ± 0.115	15.338
Forgetron	38.387 ± 1.345	0.089	35.023 ± 1.161	0.119	36.082 ± 0.692	0.502	28.977 ± 0.380	19.503
Projectron	37.008 ± 1.232	0.069	34.172 ± 0.989	0.081	33.627 ± 1.531	0.381	21.376 ± 0.315	15.310
Projectron++	35.396 ± 1.262	0.162	28.902 ± 0.650	0.161	32.424 ± 0.932	0.928	18.702 ± 0.389	23.525
BPA-S	35.345 ± 1.119	0.072	28.948 ± 0.765	0.081	33.359 ± 1.351	0.392	21.840 ± 0.391	15.304
BOGD	33.985 ± 2.013	0.085	30.013 ± 1.109	0.088	33.112 ± 0.692	0.447	29.614 ± 0.244	15.569
NOGD	29.881 ± 0.510	0.075	22.891 ± 0.043	0.085	34.332 ± 1.671	0.396	18.843 ± 0.366	15.416
SkeGD	27.932 ± 0.481	0.076	21.388 ± 0.122	0.082	31.301 ± 0.531	0.322	18.524 ± 0.283	11.736

Algorithm	w7a		ijcnn1		cod-rna		covtype.binary	
	Mistake rate	Time	Mistake rate	Time	Mistake rate	Time	Mistake rate	Time
RBP	5.246 ± 0.034	17.753	16.702 ± 0.120	11.887	19.359 ± 0.114	46.182	41.321 ± 0.024	217.216
Forgetron	5.552 ± 0.045	18.579	17.092 ± 0.203	13.648	18.874 ± 0.048	54.248	41.334 ± 0.076	299.276
Projectron	4.749 ± 0.455	17.639	12.231 ± 0.020	11.682	16.138 ± 0.452	50.473	32.748 ± 0.424	226.050
Projectron++	4.236 ± 1.042	17.101	9.588 ± 0.012	17.930	12.743 ± 0.560	162.823	30.749 ± 0.311	466.637
BPA-S	3.204 ± 0.010	17.937	11.368 ± 0.035	11.610	13.967 ± 0.235	47.952	36.538 ± 0.409	215.262
BOGD	3.482 ± 0.133	17.969	12.064 ± 0.079	11.993	13.818 ± 0.049	53.349	41.858 ± 0.014	223.180
NOGD	2.903 ± 0.009	17.923	9.540 ± 0.003	11.984	12.393 ± 0.583	56.845	34.021 ± 0.617	216.327
SkeGD	2.782 ± 0.001	7.452	9.348 ± 0.001	8.034	11.638 ± 0.050	25.918	30.493 ± 1.063	142.739

learning with little sacrifice of the time efficiency, where d is equivalent to the number of nonzero entries in each row of S_p .

6.4. Online Learning in Adversarial Environments

In this experiment, we compare SkeGD with NOGD, FOGD (Lu et al., 2016) and PROS-N-KONS (Calandriello et al., 2017a) in adversarial environments. We set the budget $B = 100$ for SkeGD and NOGD, and use identical parameters to those used in Section 6.3. Besides, we use the same experimental settings for FOGD (feature dimension = $4B$) and PROS-N-KONS as provided in the original papers. Inspired by the adversarial settings in (Calandriello et al., 2017a; Wang et al., 2018), we build adversarial datasets using the benchmark `german`. We decompose the online learning game into k_b blocks, where each block includes k_r rounds. At each round of the same block, we receive the same instance sampled from `german`. Moreover, the revealed labels are flipped at each even block by multiplying -1 . Thus, the sample size of the adversarial dataset is $k_b \times k_r$. We set $k_b = 500$, $k_r = 10$ in `german-1`, and $k_b = 500$, $k_r = 20$ in `german-2`.

From Table 4, we can observe that SkeGD is significantly more accurate than other algorithms in adversarial environments, which demonstrates the effectiveness of our incremental randomized sketches. We also notice that SkeGD is much more efficient than the second-order algorithm PROS-N-KONS, and has the time costs comparable to the other first-order algorithms. Additionally, too large or too small values of the update cycle may lead to an increased loss that conforms to Remark 2, and $\rho \approx 0.005T$ performs well in terms of accuracy while having a comparable efficiency.

 Table 4. Comparison of online kernel learning algorithms in adversarial environments w.r.t. the mistake rates (%) and the running time (s), where $\rho = \lfloor \theta(T - B) \rfloor$ is the update cycle of SkeGD.

Algorithm	german-1		german-2	
	Mistake rate	Time	Mistake rate	Time
FOGD	37.493 ± 0.724	0.140	32.433 ± 0.196	0.265
NOGD	30.918 ± 0.003	0.405	26.737 ± 0.002	0.778
PROS-N-KONS	27.633 ± 0.416	33.984	17.737 ± 0.900	98.873
SkeGD ($\theta = 0.1$)	17.320 ± 0.136	0.329	7.865 ± 0.059	0.597
SkeGD ($\theta = 0.01$)	17.272 ± 0.112	0.402	7.407 ± 0.086	0.633
SkeGD ($\theta = 0.005$)	16.578 ± 0.360	0.484	7.266 ± 0.065	0.672
SkeGD ($\theta = 0.001$)	16.687 ± 0.155	1.183	6.835 ± 0.136	1.856

7. Conclusion

In this paper, we have proposed the novel incremental randomized sketching for online kernel learning, which maintains the incremental randomized sketches using efficient rank-1 modifications and constructs a time-varying explicit feature mapping that is suitable for adversarial environments. The proposed sketching approach preserves the kernel matrix approximation accuracy with a $1 + \epsilon$ relative-error, has the lower bounds of the sketch sizes independently of the number of rounds, and enjoys a sublinear regret bound for online kernel learning. Under the same assumptions about kernel eigenvalue decay, in contrast to the second-order online kernel learning via sketching that has per-round computational complexities depending on the number of rounds, our sketching has constant computational complexities at each round, independently of the number of rounds. The theoretical formulation and algorithmic implementation may provide a sketch scheme for both online and offline large-scale kernel learning.

Acknowledgements

This work was supported in part by the National Natural Science Foundation of China (No. 61673293).

References

- Belkin, M. Approximation beats concentration? An approximation view on inference with smooth radial kernels. In *Proceedings of the 31st Conference on Learning Theory*, pp. 1348–1361, 2018.
- Bottou, L. Large-scale machine learning with stochastic gradient descent. In *Proceedings of the 19th International Conference on Computational Statistics*, pp. 177–186, 2010.
- Calandriello, D., Lazaric, A., and Valko, M. Efficient second-order online kernel learning with adaptive embedding. In *Advances in Neural Information Processing Systems 30*, pp. 6140–6150, 2017a.
- Calandriello, D., Lazaric, A., and Valko, M. Second-order kernel online convex optimization with adaptive sketching. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 645–653, 2017b.
- Cavallanti, G., Cesa-Bianchi, N., and Gentile, C. Tracking the best hyperplane with a simple budget perceptron. *Machine Learning*, 69(2):143–167, 2007.
- Charikar, M., Chen, K., and Farach-Colton, M. Finding frequent items in data streams. *Theoretical Computer Science*, 312(1):3–15, 2004.
- Cortes, C., Mohri, M., and Rostamizadeh, A. Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research*, 13:795–828, 2012.
- Crammer, K., Kandola, J., and Singer, Y. Online classification on a budget. In *Advances in Neural Information Processing Systems 16*, pp. 225–232, 2003.
- Dekel, O., Shalev-Shwartz, S., and Singer, Y. The forgetron: A kernel-based perceptron on a fixed budget. In *Advances in Neural Information Processing Systems 18*, pp. 259–266, 2005.
- Ding, L., Liao, S., Liu, Y., Yang, P., and Gao, X. Randomized kernel selection with spectra of multilevel circulant matrices. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, pp. 2910–2917, 2018.
- Engel, Y., Mannor, S., and Meir, R. The kernel recursive least-squares algorithm. *IEEE Transactions on Signal Processing*, 52(8):2275–2285, 2004.
- Halko, N., Martinsson, P.-G., and Tropp, J. A. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288, 2011.
- Kane, D. M. and Nelson, J. Sparser Johnson-Lindenstrauss transforms. *Journal of the ACM*, 61(1):4:1–23, 2014.
- Kivinen, J., Smola, A. J., and Williamson, R. C. Online learning with kernels. In *Advances in Neural Information Processing Systems 14*, pp. 785–792, 2001.
- Kumar, S., Mohri, M., and Talwalkar, A. On sampling-based approximate spectral decomposition. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 553–560, 2009.
- Liu, Y. and Liao, S. Eigenvalues ratio for kernel selection of kernel methods. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence*, pp. 2814–2820, 2015.
- Lu, J., Hoi, S. C., Wang, J., Zhao, P., and Liu, Z. Large scale online kernel learning. *Journal of Machine Learning Research*, 17:1613–1655, 2016.
- Orabona, F., Keshet, J., and Caputo, B. The projectron: A bounded kernel-based perceptron. In *Proceedings of the 25th International Conference on Machine Learning*, pp. 720–727, 2008.
- Raskutti, G. and Mahoney, M. Statistical and algorithmic perspectives on randomized sketching for ordinary least-squares. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 617–625, 2015.
- Sarlós, T. Improved approximation algorithms for large matrices via random projections. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*, pp. 143–152, 2006.
- Shalev-Shwartz, S., Singer, Y., Srebro, N., and Cotter, A. Pegasos: Primal estimated sub-gradient solver for SVM. *Mathematical Programming*, 127(1):3–30, 2011.
- Sun, Y., Schmidhuber, J., and Gomez, F. J. On the size of the online kernel sparsification dictionary. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 329–336, 2012.
- Wang, G., Zhao, D., and Zhang, L. Minimizing adaptive regret with one gradient per iteration. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 2762–2768, 2018.
- Wang, S. and Zhang, Z. Improving CUR matrix decomposition and the Nyström approximation via adaptive sampling. *Journal of Machine Learning Research*, 14:2729–2769, 2013.

- Wang, S., Zhang, Z., and Zhang, T. Towards more efficient SPSD matrix approximation and CUR matrix decomposition. *Journal of Machine Learning Research*, 17:1–49, 2016.
- Wang, Z. and Vucetic, S. Online passive-aggressive algorithms on a budget. In *Proceedings of the 13rd International Conference on Artificial Intelligence and Statistics*, pp. 908–915, 2010.
- Wang, Z., Crammer, K., and Vucetic, S. Breaking the curse of kernelization: Budgeted stochastic gradient descent for large-scale SVM training. *Journal of Machine Learning Research*, 13:3103–3131, 2012.
- Williams, C. and Seeger, M. Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems 14*, pp. 682–688, 2001.
- Woodruff, D. P. Sketching as a tool for numerical linear algebra. *Foundations and Trends[®] in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- Zhang, X. and Liao, S. Online kernel selection via incremental sketched kernel alignment. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 3118–3124, 2018.
- Zhao, P., Wang, J., Wu, P., Jin, R., and Hoi, S. C. Fast bounded online gradient descent algorithms for scalable kernel-based online learning. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 169–176, 2012.
- Zinkevich, M. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning*, pp. 928–936, 2003.