
Supplementary Material for Incremental Randomized Sketching for Online Kernel Learning

Xiao Zhang¹ Shizhong Liao¹

Abstract

In this supplementary material, we give the formulation of our incremental randomized sketching, the detailed proofs of the lemmas and theorems in the section ‘‘Theoretical Analysis’’, and more experimental results. Our main theoretical results include:

- The inner product preserving property (Lemma 2).
- The matrix product preserving property (Lemma 3).
- The low-rank approximation property (Theorem 1).
- The regret bound for online kernel learning (Theorem 2).

1. Formulation of Incremental Randomized Sketching

In this section, we construct the incremental randomized sketches of the kernel matrix, formulate the incremental maintenance for the incremental randomized sketches, and build the 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

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

First, we approximate the kernel matrix incrementally. Let $\mathbf{S}_p^{(t+1)} \in \mathbb{R}^{(t+1) \times s_p}$ be an SJLT and $\mathbf{S}_m^{(t+1)} \in \mathbb{R}^{(t+1) \times s_m}$

be a sub-sampling matrix. We use $\mathbf{S}_p^{(t+1)}$ and $\mathbf{S}_m^{(t+1)}$ for reducing the complexity of the problem (2) and the size of the approximate kernel matrix, respectively, and then formulate the *incremental randomized sketches* of $\mathbf{K}^{(t+1)}$ as follows:

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

where

$$\mathbf{C}_p^{(t+1)} = \mathbf{K}^{(t+1)} \mathbf{S}_p^{(t+1)} \quad \text{and} \quad \mathbf{C}_m^{(t+1)} = \mathbf{K}^{(t+1)} \mathbf{S}_m^{(t+1)}.$$

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

$$\mathbf{K}_{sk}^{(t+1)} = \mathbf{C}_m^{(t+1)} \mathbf{F}_{sk}^{(t+1)} \mathbf{C}_m^{(t+1)\top} \approx \mathbf{K}^{(t+1)}, \quad (1)$$

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

$$\begin{aligned} \mathbf{F}_{sk}^{(t+1)} &= \arg \min_{\mathbf{F}} \left\| \mathbf{S}_p^{(t+1)\top} \mathbf{E}_F^{(t+1)} \mathbf{S}_p^{(t+1)} \right\|_{\mathbf{F}}^2 \\ &= \left(\boldsymbol{\Phi}_{pm}^{(t+1)} \right)^\dagger \boldsymbol{\Phi}_{pp}^{(t+1)} \left(\boldsymbol{\Phi}_{pm}^{(t+1)\top} \right)^\dagger, \end{aligned} \quad (2)$$

where $\mathbf{E}_F^{(t+1)}$ is the approximation error at round $t + 1$

$$\mathbf{E}_F^{(t+1)} = \mathbf{C}_m^{(t+1)} \mathbf{F} \mathbf{C}_m^{(t+1)\top} - \mathbf{K}^{(t+1)}.$$

We partition the sketch matrices into block matrices as

$$\mathbf{S}_p^{(t+1)} = \left[\mathbf{S}_p^{(t)\top}, \mathbf{s}_p^{(t+1)} \right]^\top, \quad \mathbf{S}_m^{(t+1)} = \left[\mathbf{S}_m^{(t)\top}, \mathbf{s}_m^{(t+1)} \right]^\top,$$

where $\mathbf{s}_p^{(t+1)} \in \mathbb{R}^{s_p}$ contains d nonzero entries that are determined by d different hash mappings in SJLT. Then the two incremental randomized sketches of the kernel matrix $\mathbf{K}^{(t+1)}$ can be computed incrementally by rank-1 modifications as follows:

- 1) Sketch $\boldsymbol{\Phi}_{pm}^{(t+1)}$

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

The sketch $\Phi_{\text{pm}}^{(t+1)}$ can be maintained as

$$\begin{aligned}
 & \Phi_{\text{pm}}^{(t+1)} \\
 &= \mathbf{S}_{\text{p}}^{(t+1)\top} \mathbf{C}_{\text{m}}^{(t+1)} \\
 &= \mathbf{S}_{\text{p}}^{(t+1)\top} \mathbf{K}^{(t+1)} \mathbf{S}_{\text{m}}^{(t+1)} \\
 &= [\mathbf{S}_{\text{p}}^{(t)\top}, \mathbf{s}_{\text{p}}^{(t+1)}] \begin{bmatrix} \mathbf{K}^{(t)} & \boldsymbol{\psi}^{(t+1)} \\ \boldsymbol{\psi}^{(t+1)\top} & \xi^{(t+1)} \end{bmatrix} \begin{bmatrix} \mathbf{S}_{\text{m}}^{(t)} \\ \mathbf{s}_{\text{m}}^{(t+1)\top} \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{S}_{\text{p}}^{(t)\top} \mathbf{K}^{(t)} + \mathbf{s}_{\text{p}}^{(t+1)} \boldsymbol{\psi}^{(t+1)\top} \\ \mathbf{S}_{\text{p}}^{(t)\top} \boldsymbol{\psi}^{(t+1)} + \xi^{(t+1)} \mathbf{s}_{\text{p}}^{(t+1)} \end{bmatrix}^{\top} \begin{bmatrix} \mathbf{S}_{\text{m}}^{(t)} \\ \mathbf{s}_{\text{m}}^{(t+1)\top} \end{bmatrix} \\
 &= \mathbf{S}_{\text{p}}^{(t)\top} \mathbf{K}^{(t)} \mathbf{S}_{\text{m}}^{(t)} + \mathbf{R}_{\text{pm}}^{(t+1)} + \mathbf{R}_{\text{mp}}^{(t+1)\top} + \mathbf{T}_{\text{pm}}^{(t+1)} \\
 &= \Phi_{\text{pm}}^{(t)} + \mathbf{R}_{\text{pm}}^{(t+1)} + \mathbf{R}_{\text{mp}}^{(t+1)\top} + \mathbf{T}_{\text{pm}}^{(t+1)},
 \end{aligned}$$

where the modifications are performed using the following three rank-1 matrices

$$\begin{aligned}
 \mathbf{R}_{\text{pm}}^{(t+1)} &= \mathbf{s}_{\text{p}}^{(t+1)} \boldsymbol{\psi}^{(t+1)\top} \mathbf{S}_{\text{m}}^{(t)}, \\
 \mathbf{R}_{\text{mp}}^{(t+1)} &= \mathbf{s}_{\text{m}}^{(t+1)} \boldsymbol{\psi}^{(t+1)\top} \mathbf{S}_{\text{p}}^{(t)}, \\
 \mathbf{T}_{\text{pm}}^{(t+1)} &= \xi^{(t+1)} \mathbf{s}_{\text{p}}^{(t+1)} \mathbf{s}_{\text{m}}^{(t+1)\top}.
 \end{aligned}$$

2) Sketch $\Phi_{\text{pp}}^{(t+1)}$

For sketch $\Phi_{\text{pp}}^{(t+1)}$, we have

$$\begin{aligned}
 & \Phi_{\text{pp}}^{(t+1)} \\
 &= \mathbf{S}_{\text{p}}^{(t+1)\top} \mathbf{C}_{\text{p}}^{(t+1)} \\
 &= \mathbf{S}_{\text{p}}^{(t+1)\top} \mathbf{K}^{(t+1)} \mathbf{S}_{\text{p}}^{(t+1)} \\
 &= [\mathbf{S}_{\text{p}}^{(t)\top}, \mathbf{s}_{\text{p}}^{(t+1)}] \begin{bmatrix} \mathbf{K}^{(t)} & \boldsymbol{\psi}^{(t+1)} \\ \boldsymbol{\psi}^{(t+1)\top} & \xi^{(t+1)} \end{bmatrix} \begin{bmatrix} \mathbf{S}_{\text{p}}^{(t)} \\ \mathbf{s}_{\text{p}}^{(t+1)\top} \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{S}_{\text{p}}^{(t)\top} \mathbf{K}^{(t)} + \mathbf{s}_{\text{p}}^{(t+1)} \boldsymbol{\psi}^{(t+1)\top} \\ \mathbf{S}_{\text{p}}^{(t)\top} \boldsymbol{\psi}^{(t+1)} + \xi^{(t+1)} \mathbf{s}_{\text{p}}^{(t+1)} \end{bmatrix}^{\top} \begin{bmatrix} \mathbf{S}_{\text{p}}^{(t)} \\ \mathbf{s}_{\text{p}}^{(t+1)\top} \end{bmatrix} \\
 &= \mathbf{S}_{\text{p}}^{(t)\top} \mathbf{K}^{(t)} \mathbf{S}_{\text{p}}^{(t)} + \mathbf{R}_{\text{pp}}^{(t+1)} + \mathbf{R}_{\text{pp}}^{(t+1)\top} + \mathbf{T}_{\text{pp}}^{(t+1)}, \\
 &= \Phi_{\text{pp}}^{(t)} + \mathbf{R}_{\text{pp}}^{(t+1)} + \mathbf{R}_{\text{pp}}^{(t+1)\top} + \mathbf{T}_{\text{pp}}^{(t+1)},
 \end{aligned}$$

where the modifications are done by the following two rank-1 matrices

$$\begin{aligned}
 \mathbf{R}_{\text{pp}}^{(t+1)} &= \mathbf{s}_{\text{p}}^{(t+1)} \boldsymbol{\psi}^{(t+1)\top} \mathbf{S}_{\text{p}}^{(t)}, \\
 \mathbf{T}_{\text{pp}}^{(t+1)} &= \xi^{(t+1)} \mathbf{s}_{\text{p}}^{(t+1)} \mathbf{s}_{\text{p}}^{(t+1)\top}.
 \end{aligned}$$

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:

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

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

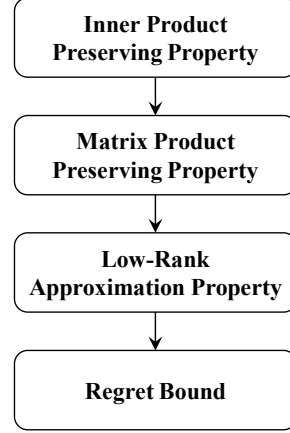


Figure 1. The dependence structure of our theoretical results.

Then $\mathbf{F}_{\text{sk}}^{(t+1)}$ is approximated by

$$\mathbf{F}_{\text{sk}}^{(t+1)} \approx \mathbf{Q}_{t+1} \mathbf{Q}_{t+1}^{\top},$$

where

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

which yields the approximate kernel matrix from (1)

$$\mathbf{K}_{\text{sk}}^{(t+1)} \approx \left(\mathbf{C}_{\text{m}}^{(t+1)} \mathbf{Q}_{t+1} \right) \left(\mathbf{C}_{\text{m}}^{(t+1)} \mathbf{Q}_{t+1} \right)^{\top}.$$

Thus, the kernel function value between the i -th example \mathbf{x}_i and the j -th example \mathbf{x}_j can be approximated by

$$\begin{aligned}
 & \kappa(\mathbf{x}_i, \mathbf{x}_j) \\
 & \approx \left(\left[\mathbf{C}_{\text{m}}^{(t+1)} \right]_{i*} \mathbf{Q}_{t+1} \right) \left(\left[\mathbf{C}_{\text{m}}^{(t+1)} \right]_{j*} \mathbf{Q}_{t+1} \right)^{\top},
 \end{aligned}$$

and the explicit feature mapping can be updated at round $t+1$ by

$$\phi_{t+2}(\cdot) = \left([\kappa(\cdot, \tilde{\mathbf{x}}_1), \dots, \kappa(\cdot, \tilde{\mathbf{x}}_{s_{\text{m}}})] \mathbf{Q}_{t+1} \right)^{\top},$$

where $\{\tilde{\mathbf{x}}_i\}_{i=1}^{s_{\text{m}}}$ are the sampled examples obtained by $\mathbf{S}_{\text{m}}^{(t+1)}$.

2. Detailed Proofs in Theoretical Analysis

Figure 1 describes the dependence structure of our theoretical results. For convenience, in this section, we denote $\mathbf{S}_{\text{p}}^{(T)}$, $\mathbf{S}_{\text{m}}^{(T)}$, $\mathbf{C}_{\text{m}}^{(T)}$, $\mathbf{F}_{\text{sk}}^{(T)}$ and $\mathbf{K}^{(T)}$ by \mathbf{S}_{p} , \mathbf{S}_{m} , \mathbf{C}_{m} , \mathbf{F}_{sk} and \mathbf{K} , respectively. We first give some extra notations. Let $\mathbf{U}_{\text{m}} \in \mathbb{R}^{T \times s_{\text{m}}}$ be the first s_{m} left singular vectors of $\mathbf{C}_{\text{m}} \in \mathbb{R}^{T \times s_{\text{m}}}$. We denote a matrix with orthonormal columns by $\mathbf{U}_{\text{m}}^{\perp} \in \mathbb{R}^{T \times (T-s_{\text{m}})}$ which satisfies

$$\mathbf{U}_{\text{m}} \mathbf{U}_{\text{m}}^{\top} + \mathbf{U}_{\text{m}}^{\perp} (\mathbf{U}_{\text{m}}^{\perp})^{\top} = \mathbf{I}_T \quad \text{and} \quad \mathbf{U}_{\text{m}}^{\top} \mathbf{U}_{\text{m}}^{\perp} = \mathbf{O}.$$

Then,

$$U_m^\perp (U_m^\perp)^\top K = K - U_m U_m^\top K = K - C_m C_m^\dagger K.$$

The Count Sketch matrix has been used as a specific random projection technique due to the unbiasedness and efficiency while approximating the inner product of two vectors.

Lemma 1 (Lemma 2 from (Pham & Pagh, 2013)). *Given two points $\mathbf{x}, \mathbf{y} \in \mathbb{R}^T$, we denote by $\mathbf{S} \in \mathbb{R}^{T \times s}$ the Count Sketch matrix. Then,*

$$\begin{aligned} \mathbb{E}[\langle \mathbf{S}^\top \mathbf{x}, \mathbf{S}^\top \mathbf{y} \rangle] &= \langle \mathbf{x}, \mathbf{y} \rangle, \\ \text{Var}[\langle \mathbf{S}^\top \mathbf{x}, \mathbf{S}^\top \mathbf{y} \rangle] &\leq \frac{1}{s} (\langle \mathbf{x}, \mathbf{y} \rangle^2 + \|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2). \end{aligned}$$

We first provide the analysis of the expectation and variance while approximating the inner product using \mathbf{S}_p in the proposed incremental randomized sketches.

Lemma 2 (Inner Product Preserving Property). *Given two points $\mathbf{x}, \mathbf{y} \in \mathbb{R}^T$, we denote by $\mathbf{S} \in \mathbb{R}^{T \times s}$ the sketch matrix \mathbf{S}_p in the proposed incremental randomized sketches. Then,*

$$\begin{aligned} \mathbb{E}[\langle \mathbf{S}^\top \mathbf{x}, \mathbf{S}^\top \mathbf{y} \rangle] &= \langle \mathbf{x}, \mathbf{y} \rangle, \\ \text{Var}[\langle \mathbf{S}^\top \mathbf{x}, \mathbf{S}^\top \mathbf{y} \rangle] &\leq \frac{1}{s} (\langle \mathbf{x}, \mathbf{y} \rangle^2 + \|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2). \end{aligned}$$

Proof. Since

$$\mathbf{S}^\top \mathbf{x} = [\mathbf{S}_1^\top \mathbf{x}; \dots; \mathbf{S}_d^\top \mathbf{x}],$$

we have

$$\begin{aligned} &\langle \mathbf{S}^\top \mathbf{x}, \mathbf{S}^\top \mathbf{y} \rangle \\ &= [(\mathbf{S}_1^\top \mathbf{x})^\top, \dots, (\mathbf{S}_d^\top \mathbf{x})^\top] [\mathbf{S}_1^\top \mathbf{y}; \dots; \mathbf{S}_d^\top \mathbf{y}] \\ &= \sum_{i=1}^d \langle \mathbf{S}_i^\top \mathbf{x}, \mathbf{S}_i^\top \mathbf{y} \rangle. \end{aligned} \quad (3)$$

Let $\tilde{\mathbf{S}}_i = \sqrt{d} \mathbf{S}_i$ that is a Count Sketch matrix. By (3) and Lemma 1, we have

$$\begin{aligned} \mathbb{E}[\langle \mathbf{S}^\top \mathbf{x}, \mathbf{S}^\top \mathbf{y} \rangle] &= \sum_{i=1}^d \mathbb{E}[\langle \mathbf{S}_i^\top \mathbf{x}, \mathbf{S}_i^\top \mathbf{y} \rangle] \\ &= \frac{1}{d} \sum_{i=1}^d \mathbb{E}[\langle \tilde{\mathbf{S}}_i^\top \mathbf{x}, \tilde{\mathbf{S}}_i^\top \mathbf{y} \rangle] = \langle \mathbf{x}, \mathbf{y} \rangle. \end{aligned}$$

Similarly, we have

$$\text{Var}[\langle \mathbf{S}^\top \mathbf{x}, \mathbf{S}^\top \mathbf{y} \rangle] = \frac{1}{d^2} \sum_{i=1}^d \text{Var}[\langle \tilde{\mathbf{S}}_i^\top \mathbf{x}, \tilde{\mathbf{S}}_i^\top \mathbf{y} \rangle]$$

and then the result for the variance holds by Lemma 1. \square

Then we demonstrate the unbiasedness of \mathbf{S}_p while approximating matrix products, which shows that our sketched kernel matrix approximation problem is an unbiased estimate of the modified Nyström kernel matrix approximation problem.

Lemma 3 (Matrix Product Preserving Property). *Let $\mathbf{A} \in \mathbb{R}^{T \times m}$, $\mathbf{B} \in \mathbb{R}^{p \times T}$. If $\mathbf{S} \in \mathbb{R}^{T \times s}$ is the sketch matrix \mathbf{S}_p in the proposed incremental randomized sketches, then,*

- 1) $\mathbb{E}[\|\mathbf{S}^\top \mathbf{A}\|_F^2] = \|\mathbf{A}\|_F^2, \quad \text{Var}[\|\mathbf{S}^\top \mathbf{A}\|_F^2] \leq \frac{2}{s} \|\mathbf{A}\|_F^4.$
- 2) $\mathbb{E}[\mathbf{B} \mathbf{S} \mathbf{S}^\top \mathbf{A}] = \mathbf{B} \mathbf{A}.$
- 3) $\mathbb{E}[\|\mathbf{B} \mathbf{A} - \mathbf{B} \mathbf{S} \mathbf{S}^\top \mathbf{A}\|_F^2] \leq \frac{2}{s} \|\mathbf{B}\|_F^2 \|\mathbf{A}\|_F^2.$

Proof. 1) It is obvious that

$$\mathbb{E}[\|\mathbf{S}^\top \mathbf{A}\|_F^2] = \mathbb{E}\left[\sum_{i=1}^m X_i\right] = \sum_{i=1}^m \mathbb{E}[X_i],$$

where $X_i = \|\mathbf{S}^\top [\mathbf{A}]_{*i}\|_2^2$.

Using Lemma 2, we have

$$\mathbb{E}[X_i] = \|\mathbf{A}\|_{*i,2}^2,$$

and then

$$\mathbb{E}[\|\mathbf{S}^\top \mathbf{A}\|_F^2] = \sum_{i=1}^m \|\mathbf{A}\|_{*i,2}^2 = \|\mathbf{A}\|_F^2.$$

We can observe that

$$\begin{aligned} \mathbb{E}[\|\mathbf{S}^\top \mathbf{A}\|_F^4] &= \mathbb{E}\left[\left(\sum_{i=1}^m X_i\right)^2\right] \\ &= \mathbb{E}\left[\sum_{i,j=1}^m X_i X_j\right] \\ &= \sum_{i,j=1}^m \mathbb{E}[X_i X_j]. \end{aligned} \quad (4)$$

By Lemma 2, we obtain

$$\text{Var}[X_i] \leq \frac{2}{s} \|\mathbf{A}\|_{*i,2}^4,$$

and thus

$$\begin{aligned} \mathbb{E}[X_i X_j] &= \text{Cov}[X_i, X_j] + \mathbb{E}[X_i] \mathbb{E}[X_j] \\ &\leq \sqrt{\text{Var}[X_i] \text{Var}[X_j]} + \mathbb{E}[X_i] \mathbb{E}[X_j] \\ &= \frac{s+2}{s} \|\mathbf{A}\|_{*i,2}^2 \|\mathbf{A}\|_{*j,2}^2. \end{aligned} \quad (5)$$

From (4) and (5), we have

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{S}^\top \mathbf{A}\|_{\mathbb{F}}^4 \right] &\leq \frac{s+2}{s} \sum_{i,j=1}^m \left(\|\mathbf{A}_{*i}\|_2^2 \|\mathbf{A}_{*j}\|_2^2 \right) \\ &= \frac{s+2}{s} \|\mathbf{A}\|_{\mathbb{F}}^4. \end{aligned}$$

Consequently,

$$\begin{aligned} \text{Var} \left[\|\mathbf{S}^\top \mathbf{A}\|_{\mathbb{F}}^2 \right] &= \mathbb{E} \left[\|\mathbf{S}^\top \mathbf{A}\|_{\mathbb{F}}^4 \right] - \mathbb{E}^2 \left[\|\mathbf{S}^\top \mathbf{A}\|_{\mathbb{F}}^2 \right] \\ &\leq \frac{s+2}{s} \|\mathbf{A}\|_{\mathbb{F}}^4 - \|\mathbf{A}\|_{\mathbb{F}}^4 \\ &= \frac{2}{s} \|\mathbf{A}\|_{\mathbb{F}}^4. \end{aligned}$$

2) Since

$$[\mathbf{B}\mathbf{S}\mathbf{S}^\top \mathbf{A}]_{ij} = [\mathbf{B}]_{i*}^\top \mathbf{S}\mathbf{S}^\top [\mathbf{A}]_{*j} = \langle \mathbf{S}^\top [\mathbf{B}]_{i*}, \mathbf{S}^\top [\mathbf{A}]_{*j} \rangle,$$

by Lemma 2 we can obtain

$$\mathbb{E} [[\mathbf{B}\mathbf{S}\mathbf{S}^\top \mathbf{A}]_{ij}] = \langle [\mathbf{B}]_{i*}, [\mathbf{A}]_{*j} \rangle = [\mathbf{B}\mathbf{A}]_{ij}. \quad (6)$$

3) Let $\mathbf{Y}_{ij} = [\mathbf{B}\mathbf{A} - \mathbf{B}\mathbf{S}\mathbf{S}^\top \mathbf{A}]_{ij}$. By (6) we have $\mathbb{E} [\mathbf{Y}_{ij}] = 0$. Then,

$$\mathbb{E} [\mathbf{Y}_{ij}^2] = \text{Var} [\mathbf{Y}_{ij}] + \mathbb{E}^2 [\mathbf{Y}_{ij}] = \text{Var} [\mathbf{Y}_{ij}].$$

By Lemma 2 and the Cauchy-Schwarz inequality we obtain

$$\begin{aligned} \mathbb{E} [\mathbf{Y}_{ij}^2] &= \text{Var} [\mathbf{Y}_{ij}] \\ &= \text{Var} [[\mathbf{B}\mathbf{S}\mathbf{S}^\top \mathbf{A}]_{ij}] \\ &= \text{Var} [\langle \mathbf{S}^\top [\mathbf{B}]_{i*}, \mathbf{S}^\top [\mathbf{A}]_{*j} \rangle] \\ &\leq \frac{1}{s} \left(\langle [\mathbf{B}]_{i*}, [\mathbf{A}]_{*j} \rangle^2 + \|\mathbf{B}\|_{\mathbb{F}}^2 \|\mathbf{A}\|_{\mathbb{F}}^2 \right) \\ &\leq \frac{2}{s} \|\mathbf{B}\|_{\mathbb{F}}^2 \|\mathbf{A}\|_{\mathbb{F}}^2. \end{aligned}$$

Thus

$$\mathbb{E} \left[\sum_{ij} \mathbf{Y}_{ij}^2 \right] = \sum_{ij} \mathbb{E} [\mathbf{Y}_{ij}^2] \leq \frac{2}{s} \|\mathbf{B}\|_{\mathbb{F}}^2 \|\mathbf{A}\|_{\mathbb{F}}^2.$$

□

The following lemma provides that all singular values of $\mathbf{S}^\top \mathbf{U}$ lie in $[1 - \epsilon_0, 1 + \epsilon_0]$ with high probability.

Lemma 4 ((Nelson & Nguyen, 2013)). *Let $\mathbf{U} \in \mathbb{R}^{T \times c}$ be a matrix with orthonormal columns, $\mathbf{S} \in \mathbb{R}^{T \times s}$ the sketch matrix \mathbf{S}_p in the proposed incremental randomized sketches. Set $d = \Theta(\epsilon_0^{-1} \log^3(c\delta_0^{-1}))$ for \mathbf{S} . For $\epsilon_0 \in (0, 1)$, with probability at least $1 - \delta_0$ all singular values of $\mathbf{S}^\top \mathbf{U}$*

$$\sigma_i(\mathbf{S}^\top \mathbf{U}) = 1 \pm \epsilon_0$$

as long as

$$s \geq \frac{c \log^8(c\delta_0^{-1})}{\epsilon_0^2}.$$

Further, this holds if the hash function h and σ defining the \mathbf{S} is $\Omega(\log(c\delta_0^{-1}))$ -wise independent.

Then we demonstrate that the proposed incremental randomized sketching is nearly as accurate as the modified Nyström for kernel matrix approximation, which 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.

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). If $\mathbf{S}_p \in \mathbb{R}^{T \times s_p}$ is the sketch matrix in the proposed incremental randomized sketches with $d = \Theta(\log^3(s_m))$, let $\tau = s_m/s_p$ and assume*

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

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

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

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. Let $\mathbf{A} \in \mathbb{R}^{T \times m}$, $\mathbf{B} \in \mathbb{R}^{p \times T}$. By 3) in Lemma 3 and the Markov's inequality, with probability at least $1 - \hat{\delta}$

$$\|\mathbf{B}\mathbf{A} - \mathbf{B}\mathbf{S}_p \mathbf{S}_p^\top \mathbf{A}\|_{\mathbb{F}}^2 \leq \frac{2}{\hat{\delta} s_p} \|\mathbf{B}\|_{\mathbb{F}}^2 \|\mathbf{A}\|_{\mathbb{F}}^2. \quad (7)$$

From (7) with probability at least $1 - \delta_1$

$$\begin{aligned} \|\mathbf{I}_T \mathbf{S}_p \mathbf{S}_p^\top \mathbf{U}_m\|_2 &\leq \|\mathbf{I}_T \mathbf{S}_p \mathbf{S}_p^\top \mathbf{U}_m - \mathbf{I}_T \mathbf{U}_m\|_2 + \|\mathbf{I}_T \mathbf{U}_m\|_2 \\ &\leq \|\mathbf{I}_T \mathbf{S}_p \mathbf{S}_p^\top \mathbf{U}_m - \mathbf{I}_T \mathbf{U}_m\|_{\mathbb{F}} + 1 \\ &\leq \sqrt{\frac{2}{\delta_1 s_p}} \|\mathbf{I}_T\|_{\mathbb{F}} \|\mathbf{U}_m\|_{\mathbb{F}} + 1 \\ &= \sqrt{\frac{2T s_m}{\delta_1 s_p}} + 1, \end{aligned}$$

and with probability at least $1 - \delta_2$

$$\begin{aligned}
 & \|(\mathbf{S}_p^\top \mathbf{U}_m)^\top \mathbf{S}_p^\top \mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}\|_F \\
 &= \|(\mathbf{S}_p^\top \mathbf{U}_m)^\top \mathbf{S}_p^\top \mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K} - \mathbf{U}_m^\top \mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}\|_F \\
 &\leq \sqrt{\frac{2}{\delta_2 s_p}} \|\mathbf{U}_m\|_F \|\mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}\|_F \\
 &= \sqrt{\frac{2s_m}{\delta_2 s_p}} \|\mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}\|_F.
 \end{aligned}$$

Assume $s_p \geq s_m \log^8(s_m \delta_0^{-1})/\epsilon_0^2$. By Lemma 4, with probability at least $1 - \delta_0$

$$\sigma_{\max}^2(\mathbf{S}_p \mathbf{U}_m) \leq (1 + \epsilon_0)^2.$$

Let $\mathbf{E} = \mathbf{C}_m \mathbf{F}_{\text{mod}} \mathbf{C}_m^\top - \mathbf{K}$. According to the above discussion, by Lemma 17 from (Wang et al., 2015), with probability at least $1 - (\delta_0 + \delta_1 + \delta_2)$

$$\begin{aligned}
 & \|(\mathbf{S}_p^\top \mathbf{U}_m)^\top \mathbf{S}_p^\top \mathbf{E} \mathbf{S}_p (\mathbf{S}_p^\top \mathbf{U}_m)\|_F \\
 &\leq (\|\mathbf{I}_T \mathbf{S}_p \mathbf{S}_p^\top \mathbf{U}_m\|_2 + \sigma_{\max}^2(\mathbf{S}_p \mathbf{U}_m)) \\
 &\quad \|(\mathbf{S}_p^\top \mathbf{U}_m)^\top \mathbf{S}_p^\top \mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}\|_F \\
 &\leq \sqrt{\epsilon} \|\mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}\|_F,
 \end{aligned}$$

where

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

It finally follows from Lemma 17 in (Wang et al., 2015) that with probability at least $1 - (\delta_0 + \delta_1 + \delta_2)$

$$\begin{aligned}
 & \|\mathbf{C}_m \mathbf{F}_{\text{sk}} \mathbf{C}_m^\top - \mathbf{K}\|_F^2 \\
 &\leq \|\mathbf{E}\|_F^2 + \sigma_{\min}^{-8}(\mathbf{S}_p \mathbf{U}_m) \|(\mathbf{S}_p^\top \mathbf{U}_m)^\top \mathbf{S}_p^\top \mathbf{E} \mathbf{S}_p (\mathbf{S}_p^\top \mathbf{U}_m)\|_F^2 \\
 &\leq \|\mathbf{E}\|_F^2 + \epsilon \|\mathbf{U}_m^\perp (\mathbf{U}_m^\perp)^\top \mathbf{K}\|_F^2 \\
 &= \|\mathbf{E}\|_F^2 + \epsilon \|\mathbf{K} - \mathbf{C}_m \mathbf{C}_m^\top \mathbf{K}\|_F^2 \\
 &\leq (1 + \epsilon) \|\mathbf{E}\|_F^2.
 \end{aligned}$$

□

Denote the hypothesis at round t using the updated explicit feature mapping $\phi_{t+1}(\cdot)$ by

$$\bar{f}_t(\mathbf{x}_t) = \langle \bar{\mathbf{w}}_t, \phi_{t+1}(\mathbf{x}_t) \rangle. \quad (8)$$

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

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

Then we update the hypothesis $\bar{f}_t(\cdot) = \langle \bar{\mathbf{w}}_t, \phi_{t+1}(\cdot) \rangle$ in (8) as follows:

$$\mathbf{w}_{t+1} = \bar{\mathbf{w}}_t - \eta \nabla \mathcal{L}_t(\bar{\mathbf{w}}_t), \quad (9)$$

where

$$\ell_t(\bar{\mathbf{w}}_t) = \ell_t(\bar{f}_t) = \ell(\bar{f}_t(\mathbf{x}_t), y_t)$$

and

$$\mathcal{L}_t(\bar{\mathbf{w}}_t) = \ell_t(\bar{\mathbf{w}}_t) + \frac{\lambda}{2} \|\bar{\mathbf{w}}_t\|_2^2.$$

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

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

where $\mathbf{U}_{B,\rho}$ is the singular vector matrix of $\mathbf{K}_{B,\rho}$. We finally obtain the following regret bound for online kernel learning¹.

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

$$d = \Theta(\log^3(s_m)) \quad \text{and} \quad \tau = s_m/s_p,$$

for the sketch matrix $\mathbf{S}_p \in \mathbb{R}^{T \times s_p}$ in the proposed incremental randomized sketches. Assume ℓ_t is a convex loss function that is Lipschitz continuous with the 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 (9), satisfying

$$|f_t(\mathbf{x}_t)| = |\langle \mathbf{w}_t, \phi_t(\mathbf{x}_t) \rangle| \leq C_f, \quad t \in [T].$$

For the optimal hypothesis f^* that minimizes

$$\mathcal{L}(f) = \frac{\lambda}{2} \|f\|_{\mathcal{H}_\kappa}^2 + \frac{1}{T} \sum_{t=1}^T \ell_t(f),$$

in the original reproducing kernel Hilbert space \mathcal{H}_κ , if

$$\begin{aligned}
 s_p &= \Omega(s_m \text{polylog}(s_m \delta_0^{-1})/\epsilon_0^2), \\
 s_m &= \Omega(\mu(\mathbf{K}_{B,\rho}) k \log k),
 \end{aligned}$$

then with probability at least $1 - \delta$

$$\begin{aligned}
 & \sum_{t=1}^T (\mathcal{L}_t(\mathbf{w}_t) - \mathcal{L}_t(f^*)) \\
 &\leq \frac{C_f^2 + \|\mathbf{w}_Z^*\|_2^2}{2\theta\eta} + \frac{\eta L^2 T}{2} \\
 &\quad + \frac{1}{\lambda(\beta-1)} \left(\frac{3}{2} - \frac{B+1/\theta}{T} \right) + \frac{\sqrt{1+\epsilon}}{\lambda} O(\sqrt{B}),
 \end{aligned}$$

¹In the theoretical analysis, we assume $T_0 = B$ and omit KOGD used at the first stage that enjoys a $O(\sqrt{B})$ regret bound.

where $\delta = \delta_0 + \delta_1 + \delta_2$, \mathbf{w}_i^* is the optimal hypothesis on the incremental randomized sketches over the first t instances and $Z = \arg \max_{i \in [T]} \|\mathbf{w}_i^*\|_2$,

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

Proof. By Theorem 2 in (Yang et al., 2012), we have

$$\mathcal{L}(\mathbf{w}^*) - \mathcal{L}(f^*) \leq \frac{1}{2T\lambda} \|\mathbf{K}_{\text{sk}}^{(T)} - \mathbf{K}\|_2,$$

where \mathbf{w}^* is the optimal hypothesis on the incremental randomized sketches in hindsight, which yields

$$\begin{aligned} & \sum_{t=1}^T (\mathcal{L}_t(\mathbf{w}^*) - \mathcal{L}_t(f^*)) \\ & \leq \frac{1}{2\lambda} \left\| \widehat{[\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho}} - \mathbf{K} \right\|_2 \\ & \leq \frac{1}{2\lambda} \left(\left\| \widehat{[\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho}} - \widehat{\mathbf{K}}_{B,\rho} \right\|_2 + \left\| \widehat{\mathbf{K}}_{B,\rho} - \mathbf{K} \right\|_2 \right) \\ & = \frac{1}{2\lambda} \left(\left\| [\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho} - \mathbf{K}_{B,\rho} \right\|_2 + \left\| \widehat{\mathbf{K}}_{B,\rho} - \mathbf{K} \right\|_2 \right) \end{aligned} \quad (10)$$

where $\mathbf{K}_{B,\rho} \in \mathbb{R}^{(B+\lfloor(T-B)/\rho\rfloor) \times (B+\lfloor(T-B)/\rho\rfloor)}$ is the intersection matrix of \mathbf{K} which is constructed by $B + \lfloor(T-B)/\rho\rfloor$ examples, $[\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho}$ is the approximate matrix for $\mathbf{K}_{B,\rho}$ using the proposed incremental randomized sketching with rank parameter k , \mathbf{O} is a zero matrix of size $(T-B - \lfloor(T-B)/\rho\rfloor) \times (T-B - \lfloor(T-B)/\rho\rfloor)$ and

$$\begin{aligned} \widehat{\mathbf{K}}_{B,\rho} &= \text{diag} \{ \mathbf{K}_{B,\rho}, \mathbf{O} \} \in \mathbb{R}^{T \times T}, \\ \widehat{[\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho}} &= \text{diag} \{ [\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho}, \mathbf{O} \} \in \mathbb{R}^{T \times T}. \end{aligned}$$

Since the eigenvalues of the kernel matrix decay polynomially with decay rate $\beta > 1$, the following bound holds

$$\begin{aligned} & \left\| \widehat{\mathbf{K}}_{B,\rho} - \mathbf{K} \right\|_2 \\ & \leq \frac{T-B - \lfloor(T-B)/\rho\rfloor}{T} \sum_{i=1}^T i^{-\beta} \\ & \leq \frac{T-B - \lfloor(T-B)/\rho\rfloor}{T} \int_1^T i^{-\beta} di \\ & = \frac{T-B - \lfloor(T-B)/\rho\rfloor}{T} \frac{1}{\beta-1} \left(1 - \frac{1}{T^{\beta-1}} \right) \\ & \leq \frac{1}{\beta-1} \left(1 - \frac{B + \lfloor(T-B)/\rho\rfloor}{T} \right). \end{aligned} \quad (11)$$

From Theorem 1, with probability at least $1 - \delta$,

$$\begin{aligned} & \left\| [\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho} - \mathbf{K}_{B,\rho} \right\|_2 \\ & \leq \sqrt{1 + \epsilon} \left\| [\mathbf{C}_m \mathbf{F}_{\text{mod}} \mathbf{C}_m^T]_{B,\rho} - \mathbf{K}_{B,\rho} \right\|_F, \end{aligned} \quad (12)$$

where $[\mathbf{C}_m \mathbf{F}_{\text{mod}} \mathbf{C}_m^T]_{B,\rho}$ is the approximate matrix for $\mathbf{K}_{B,\rho}$ using the modified Nyström approach with rank parameter k .

Denote the best rank- k approximation of \mathbf{A} by $(\mathbf{A})_k$. Since the eigenvalues of \mathbf{K} decay polynomially with decay rate $\beta > 1$, there exists $\beta > 1$ such that $\lambda_i(\mathbf{K}) = O(i^{-\beta})$, which yields

$$\begin{aligned} & \left\| \mathbf{K}_{B,\rho} - (\mathbf{K}_{B,\rho})_k \right\|_F \\ & = \sqrt{B + \lfloor(T-B)/\rho\rfloor - k} \cdot (k+1)^{-\beta} \\ & = O(\sqrt{B}). \end{aligned} \quad (13)$$

Given $\epsilon' \in (0, 1)$, when $s_m = \Omega(\mu(\mathbf{K}_{B,\rho})k \log k)$, the following bound holds (Wang et al., 2016)

$$\begin{aligned} & \left\| [\mathbf{C}_m \mathbf{F}_{\text{mod}} \mathbf{C}_m^T]_{B,\rho} - \mathbf{K}_{B,\rho} \right\|_F \\ & \leq \sqrt{1 + \epsilon'} \left\| \mathbf{K}_{B,\rho} - (\mathbf{K}_{B,\rho})_k \right\|_F, \end{aligned} \quad (14)$$

where $\mu(\mathbf{K}_{B,\rho})$ is the coherence of $\mathbf{K}_{B,\rho}$. Combining (12), (13) with (14), we obtain

$$\left\| [\mathbf{K}_{\text{sk}}^{(T)}]_{B,\rho} - \mathbf{K}_{B,\rho} \right\|_2 \leq \sqrt{1 + \epsilon} O(\sqrt{B}). \quad (15)$$

Substituting (11) and (15) into (10), we have

$$\begin{aligned} & \sum_{t=1}^T (\mathcal{L}_t(\mathbf{w}^*) - \mathcal{L}_t(f^*)) \\ & \leq \frac{1}{2\lambda(\beta-1)} \left(1 - \frac{B + \lfloor(T-B)/\rho\rfloor}{T} \right) + \\ & \quad \frac{\sqrt{1 + \epsilon}}{2\lambda} O(\sqrt{B}). \end{aligned} \quad (16)$$

Then we analyze the regret caused by hypothesis updating on the incremental randomized sketches. We first decompose $\mathcal{L}_t(\mathbf{w}_t) - \mathcal{L}_t(\mathbf{w}^*)$ into two terms as follows:

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

where $f_t^*(\cdot) = \langle \mathbf{w}_t^*, \phi_t(\cdot) \rangle$ is the optimal hypothesis on the incremental randomized sketches over the first t instances, and \mathbf{w}^* is the optimal hypothesis on the incremental randomized sketches in hindsight. The optimization error measures the discrepancy between the hypothesis generated by our sketched online gradient descent and the optimal hypothesis on the incremental randomized sketches at each round, and the estimation error measures the difference between the optimal hypotheses on the incremental randomized sketches over the first t instances and all the T instances respectively. For the optimization error, by the convexity of loss function,

we have

$$\begin{aligned}
 & \sum_{t=1}^T (\mathcal{L}_t(\mathbf{w}_t) - \mathcal{L}_t(\mathbf{w}_t^*)) \\
 & \leq \sum_{i=1}^{\lfloor (T-B)/\rho \rfloor} \frac{\|\bar{\mathbf{w}}_{B+(i-1)\rho} - \mathbf{w}_{B+i\rho}^*\|_2^2}{2\eta} + \frac{\eta L^2 T}{2} \\
 & \leq \sum_{i=1}^{\lfloor (T-B)/\rho \rfloor} \frac{\|\bar{\mathbf{w}}_{B+(i-1)\rho}\|_2^2 + \|\mathbf{w}_{B+i\rho}^*\|_2^2}{2\eta} + \frac{\eta L^2 T}{2} \\
 & \leq \sum_{i=1}^{\lfloor (T-B)/\rho \rfloor} \frac{|f_{B+(i-1)\rho}(\mathbf{x}_{B+(i-1)\rho})|^2 + \|\mathbf{w}_{B+i\rho}^*\|_2^2}{2\eta} + \\
 & \quad \frac{\eta L^2 T}{2} \\
 & \leq \left\lfloor \frac{T-B}{\rho} \right\rfloor \frac{C_f^2 + \|\mathbf{w}_Z^*\|_2^2}{2\eta} + \frac{\eta L^2 T}{2}, \tag{17}
 \end{aligned}$$

where $\bar{\mathbf{w}}_B = \mathbf{w}_{T_0+1}$, $Z = \arg \max_{i \in [T]} \|\mathbf{w}_i^*\|_2$ and $|f_t(\mathbf{x}_t)| \leq C_f$, $t \in [T]$. For the estimation error, we obtain the following upper bound

$$\begin{aligned}
 & \sum_{t=1}^T (\mathcal{L}_t(\mathbf{w}_t^*) - \mathcal{L}_t(\mathbf{w}^*)) \\
 & \leq \frac{1}{2\lambda} \left\| \mathbf{K}_{\text{sk}}^{(T_0)} - \mathbf{K}_{\text{sk}}^{(T)} \right\|_2 \\
 & \leq \frac{1}{2\lambda} \left(\left\| \mathbf{K}_{\text{sk}}^{(T_0)} - \mathbf{K}^{(T_0)} \right\|_2 + \left\| \mathbf{K}^{(T_0)} - \mathbf{K} \right\|_2 + \left\| \mathbf{K}_{\text{sk}}^{(T)} - \mathbf{K} \right\|_2 \right) \tag{18} \\
 & \leq \frac{1}{2\lambda} \left[\sqrt{1+\epsilon} O(\sqrt{B}) + \frac{1}{\beta-1} \left(1 - \frac{B}{T} \right) + \left\| \mathbf{K}_{\text{sk}}^{(T)} - \mathbf{K} \right\|_2 \right].
 \end{aligned}$$

Finally, the three inequalities (16), (17) and (18) combined give the following bound

$$\begin{aligned}
 & \sum_{t=1}^T (\mathcal{L}_t(\mathbf{w}_t) - \mathcal{L}_t(f^*)) \\
 & \leq \left\lfloor \frac{T-B}{\rho} \right\rfloor \frac{C_f^2 + \|\mathbf{w}_Z^*\|_2^2}{2\eta} + \frac{\eta L^2 T}{2} + \\
 & \quad \frac{1}{\lambda(\beta-1)} \left(\frac{3}{2} - \frac{B + \lfloor (T-B)/\rho \rfloor}{T} \right) + \\
 & \quad \frac{\sqrt{1+\epsilon}}{\lambda} O(\sqrt{B}).
 \end{aligned}$$

□

3. More Experimental Results

To further analyze the convergence of the compared algorithms and SkeGD, we give the convergence curves in terms

of the mistake rates on `a9a` and `cod-rna`. As the results shown in Figure 2, the mistake rates of our SkeGD converge much faster than the other online kernel learning algorithms, which demonstrates the efficiency and effectiveness of SkeGD. In Figure 2 (b), the growing mistake rates of the compared algorithms indicate that the buffer of support vectors cannot retain the key information with a small budget. Whereas SkeGD has a decreasing mistake rate, which demonstrates the effectiveness of the proposed incremental randomized sketches.

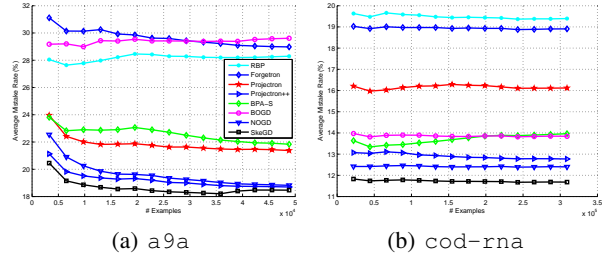


Figure 2. The average mistake rates of RBP, Forgetron, Projection, Projection++, BPA-S, BOGD, NOGD and our SkeGD.

Table 1 shows the experimental results on the adversarial datasets `spambase-1` and `spambase-2`. As the same adversarial settings in `german-1` and `german-1`, we construct these two datasets using the benchmark `spambase`, and set $k_b = 500$, $k_r = 10$ in `spambase-1` and $k_b = 500$, $k_r = 20$ in `spambase-2`. From the results, we can observe that our SkeGD achieves the highest accuracy in adversarial environments, and is much more efficient than the second-order algorithm PROS-N-KONS, while having a comparable efficiency to the other first-order online kernel learning algorithms. Besides, an appropriate update cycle results in better performances with respect to the accuracy, which conforms to Remark 2.

Table 1. 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	spambase-1		spambase-2	
	Mistake rate	Time	Mistake rate	Time
FOGD	39.793 ± 0.140	0.235	29.257 ± 0.186	0.478
NOGD	44.714 ± 0.001	0.468	40.586 ± 0.002	0.938
PROS-N-KONS	35.953 ± 0.423	265.26	26.800 ± 0.802	696.39
SkeGD ($\theta = 0.1$)	29.034 ± 1.584	0.344	16.879 ± 3.779	0.646
SkeGD ($\theta = 0.01$)	25.367 ± 1.134	0.494	14.700 ± 0.338	0.743
SkeGD ($\theta = 0.005$)	25.170 ± 0.847	0.682	14.733 ± 2.562	0.848
SkeGD ($\theta = 0.001$)	25.626 ± 0.966	2.337	14.845 ± 1.766	2.890

References

- Nelson, J. and Nguyễn, H. L. OSNAP: Faster numerical linear algebra algorithms via sparser subspace embeddings. In *Proceedings of the 54th Annual Symposium on Foundations of Computer Science*, pp. 117–126, 2013.
- Pham, N. and Pagh, R. Fast and scalable polynomial kernels via explicit feature maps. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 239–247, 2013.
- Wang, S., Zhang, Z., and Zhang, T. Towards more efficient symmetric matrix sketching and the CUR matrix decomposition. *arXiv:1503.08395v4*, 2015.
- Wang, S., Luo, L., and Zhang, Z. SPSD matrix approximation via column selection: Theories, algorithms, and extensions. *Journal of Machine Learning Research*, 17: 1–49, 2016.
- Yang, T., Li, Y.-F., Mahdavi, M., Jin, R., and Zhou, Z.-H. Nyström method vs random fourier features: A theoretical and empirical comparison. In *Advances in Neural Information Processing Systems 25*, pp. 476–484, 2012.