

# Appendices

## A Notation

Table 4: **Notation**

$n$	number of samples
$x_i \in \mathcal{X}, y_i \in \mathbb{R}$	input, target
$Y \in \mathbb{R}^n$	vector of targets $(y_i) \in \mathbb{R}^n$
$X = (x_1, x_2, \dots, x_n)$	tuple of inputs
$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$	positive definite kernel with $\mathbb{H}$ as its RKHS.
$S : \mathbb{H} \rightarrow \mathbb{R}^n$	Sampling operator, defined as $S(f) = f(X) = (f(x_i)) \in \mathbb{R}^n$
$S^* : \mathbb{R}^n \rightarrow \mathbb{H}$	Adjoint of sampling operator, defined as $S^* \alpha = \sum_{i=1}^n K(\cdot, x) \alpha_i$
$\mathcal{K}$	empirical covariance operator using $n$ samples, defined as $\frac{1}{n} S S^* : \mathbb{H} \rightarrow \mathbb{H}$
$K(X, X)$	kernel matrix. Also equals $S S^*$
$B$	mini-batch of indices, subset of $\{1, 2, \dots, n\}$
$m =  B $	batch size
$X[B] \in \mathcal{X}^m$	minibatch of size $m$
$\mathbf{v}[B] \in \mathbb{R}^m$	subvector of $\mathbf{v} \in \mathbb{R}^n$ corresponding to indices $B$
$\mathbf{H}_B$	selector matrix, rows of $I_n$ corresponding to $B$ , so that $\mathbf{v}[B] = \mathbf{H}_B \mathbf{v}$ .
$q$	level of preconditioner
$(\mathbf{E}, \Lambda, \lambda_{q+1})$	top- $q$ eigensystem of $K(X, X)$
$\mathcal{P} : \mathbb{H} \rightarrow \mathbb{H}$	preconditioner in $\mathbb{H}$ defined as $\mathcal{I} - \sum_{i=1}^q \left(1 - \frac{\lambda_{q+1}}{\lambda_i}\right) \psi_i \otimes_{\mathbb{H}} \psi_i$ where $\psi_i = S^* \mathbf{e}_i / \sqrt{\lambda_i} \in \mathbb{H}$ are top- $q$ eigenfunction of $\mathcal{K}$ with eigenvalues $\frac{1}{n} \lambda_i$
$\mathbf{F}$	rescaled eigenvectors defined as $\mathbf{E} \sqrt{I_q - \lambda_{q+1} \Lambda^{-1}}$
$\mathbf{F}[B]$	rows of $\mathbf{F}$ corresponding to the mini-batch $B$
$J$	subset of indices $\{1, 2, \dots, n\}$ to obtain Nyström approximation of $\mathcal{P}$
$s =  J $	size of Nyström subset
$X[J]$	Nyström subsample of size $s$
$(\mathbf{D}, \Delta, \delta_{q+1})$	top- $q$ eigensystem of $\frac{1}{n} K(X[J], X[J])$
$\mathcal{Q} : \mathbb{H} \rightarrow \mathbb{H}$	Nyström approximation of $\mathcal{P}$ defined as $\mathcal{I} - \sum_{i=1}^q \left(1 - \frac{\delta_{q+1}}{\delta_i}\right) \phi_i \otimes_{\mathbb{H}} \phi_i$ where $(\delta_i/s, \phi_i)$ is an eigen-pair of $\frac{1}{s} S_J^* S_J$
$\mathbf{G}$	rescaled eigenvectors defined as $\mathbf{D} \sqrt{(I_q - \delta_{q+1} \Delta^{-1}) \Delta^{-1}}$

## B EigenPro algorithms

**Lemma 9.** Suppose Algorithm 3 is run for  $t$  iterations with hyperparameters given by equation (23), then we have,

$$\|f_t - f^*\|_{\mathbb{H}}^2 \leq \exp\left(-t/\sqrt{\tilde{\kappa}_m \kappa_m}\right) \|f^*\|_{\mathbb{H}}^2. \quad (33)$$

*Proof.* This is an immediate corollary of [20, Theorem 2] and the discussion that ensues therein on the convergence rate of MaSS. Preconditioning only changes the largest eigenvalue from  $\lambda_1$  to  $\lambda_{q+1}$ . This can via the following reduction commonly applied to analyze spectral preconditioning: (16) is actually (13) with the following modified kernel,

$$\tilde{K}(x, z) := K(x, z) - K(x, X) \mathbf{E} \mathbf{Q} \mathbf{E}^\top K(X, z), \quad (34)$$

where  $\mathbf{Q} = \Lambda^{-1}(I_q - \lambda_{q+1} \Lambda^{-1})$ , where  $(\mathbf{E}, \Lambda, \lambda_{q+1})$  is the top- $q$  eigensystem of  $K(X, X)$ . An interested reader can find this reduction argument in [22, 23]. The largest eigenvalue of  $\tilde{K}(X, X)$  can be verified to be  $\lambda_{q+1}$ .

The other key difference is the quantities  $L_1$  and  $\tilde{\kappa}_m$  defined in [20]. Using elementary properties of our Hilbert space  $\mathbb{H}$ , we can easily show that our  $L_1$  and  $\tilde{\kappa}_m$  are upper bounds to these quantities in [20].  $\square$

---

**Algorithm 3** AxlePro

---

1: **Input:** positive definite kernel  $K$ , batch size  $m$ , learning rates  $\eta_1, \eta_2 > 0$ , damp factor  $\gamma \in (0, 1)$ .  
2: **Output:**  $f : \mathcal{X} \rightarrow \mathbb{R}$  solving (9) approximately.  
3: **setup:**  $\alpha \leftarrow \mathbf{0}_n$ , and  $\beta \leftarrow \mathbf{0}_n$ .  
4:  $(\mathbf{E}, \Lambda, \lambda_{q+1}) \leftarrow$  top- $q$  eigensystem of  $K(X, X)$   
5:  $\mathbf{F} \leftarrow \mathbf{E} \sqrt{I_q - \lambda_{q+1} \Lambda^{-1}} \in \mathbb{R}^{n \times q}$   
6: **repeat**  
7:   Fetch batch of indices  $B \subset \{1, \dots, n\}$ ,  $|B| = m$   
8:    $\mathbf{v} \leftarrow K(X[B], X)\beta - Y[B] \in \mathbb{R}^m$   
9:    $\mathbf{w} \leftarrow \mathbf{F} \mathbf{F}[B]^\top \mathbf{v} \in \mathbb{R}^n$   
10:    $\tilde{\alpha} \leftarrow \alpha$  {copy state for momentum}  
11:    $\alpha \leftarrow \beta$   
12:    $\alpha[B]- \leftarrow \eta_1 \mathbf{v}$  {gradient step-1}  
13:    $\alpha+ \leftarrow \eta_1 \mathbf{w}$  {correction-1}  
14:    $\beta \leftarrow \alpha + \gamma(\alpha - \tilde{\alpha})$  {momentum}  
15:    $\beta[B]+ \leftarrow \eta_2 \mathbf{v}$  {gradient step-2}  
16:    $\beta- \leftarrow \eta_2 \mathbf{w}$  {correction-2}  
17: **until** Stopping criterion is reached  
18: **return**  $f(x) = \sum_{i=1}^n \alpha_i K(x, x_i)$

---

Algorithm	FLOPS/batch	memory
EigenPro 2	$nm + sm + 2sq + m + s$	$n + sq$
AxlePro 2	$nm + sm + 2sq + \mathbf{2m} + \mathbf{2s} + \mathbf{n}$	$\mathbf{2n} + sq$
EigenPro 3	$2pm + ps + ms + 2sq + m + p$	$p + sq$
AxlePro 3	$2pm + ps + ms + 2sq + \mathbf{2m} + \mathbf{2p}$	$\mathbf{2p} + sq$

Table 5: The bolded quantities indicate overheads due to preconditioning. Here  $s = |J|$ . We have omitted the lower-order terms of  $n$  and  $2n$  in the per iteration time since they are dominated by  $mn$ . Memory requirement and per iteration complexity of EigenPro and AxlePro. The algorithm AxlePro is based on applying momentum acceleration to EigenPro [22], whereas AxlePro 2 is based on applying momentum acceleration to EigenPro 2 [23]. The overhead of storing the previous parameter is  $n$  as seen the last column.

## C Numerical Experiments (continued)

See Table 6, Table 7, Table 8, Table 9, Table 10 for experiment comparisons.

## D Details on Numerical Experiments

**Hardware.** Experiments were run on the SDSC Expanse GPU cluster with 32GB RAM, with 1 NVIDIA V100 SMX2 with 32GB VRAM, and 1 Xeon Gold 6248 CPU.

**Datasets and Kernels.** Our experiments were conducted on the following datasets: Cifar-10( $n = 50,000, d = 3072$ ), Stellar Classification( $n = 95,000, d = 13$ ), and EMNIST Digits( $n = 240,000, d = 784$ ). We test the performance for both Gaussian kernel and Laplacian kernel. And for Cifar-10 dataset, we also compare the performance with Myrtle-5 kernel [32], which is the state-of-the-art kernel for this dataset.

1. **EigenPro.** We select the precondition level based on the GPU memory so that the batch size at the linear rate regime critical point  $m_1^*$  fully exploits GPU memory. This step involves computing eigenpairs of the kernel matrix and we use Nyström approximation with  $20k$  subsamples. The optimal learning rate can also be computed using the computed eigenpairs.

2. **AxlePro.** For the precondition level, we follow the same way as in Eigenpro. The only hyperparameter needs to be tuned is the smallest eigenvalue of the kernel matrix due to the Nyström approximation and we choose this factor according to the datasets.

3. **FALKON.** According to the GPU memory, we set the number of inducing points to be  $20k$  with the uniform

---

**Algorithm 4** AxlePro 3

---

1: **Input:** positive definite kernel  $K$ , batch size  $m$ , size of approximate preconditioner  $s$ , learning rates  $\eta_1, \eta_2 > 0$ , damping factor  $\gamma \in (0, 1)$ , model centers  $Z = \{z_i\}_{i=1}^p$ .  
2: **Output:**  $f : \mathcal{X} \rightarrow \mathbb{R}$  solving (9) approximately.  
3: **setup:** Sample  $J \subset \{1, 2, \dots, n\}$  with  $|J| = s$ .  
4:  $(\mathbf{D}, \Delta, \delta_{q+1}) \leftarrow$  top- $q$  eigensystem of  $K(X[J], X[J])$   
5:  $\mathbf{G} \leftarrow \mathbf{D} \sqrt{\Delta^{-1}(I_q - \delta_{q+1}\Delta^{-1})} \in \mathbb{R}^{s \times q}$   
6: Initialize  $\boldsymbol{\alpha} = \boldsymbol{\beta} = \mathbf{0}_p$ .  
7: **repeat**  
8:   Fetch batch of indices  $B \subset \{1, \dots, n\}$ ,  $|B| = m$   
9:    $\mathbf{v} \leftarrow K(X[B], Z)\boldsymbol{\beta} + K(X[B], X)\mathbf{b} + K(X[B], X[J])\mathbf{d} - Y[B] \in \mathbb{R}^m$  {gradient}  
10:    $\mathbf{w} \leftarrow \mathbf{G}\mathbf{G}^\top K(X[J], X[B])\mathbf{v} \in \mathbb{R}^s$  {correction}  
11:    $\tilde{\boldsymbol{\alpha}} \leftarrow \boldsymbol{\alpha}$   
12:    $\tilde{\mathbf{c}} \leftarrow \mathbf{c}$   
13:    $\boldsymbol{\alpha} \leftarrow \boldsymbol{\beta}$   
14:    $\boldsymbol{\beta} \leftarrow (1 + \gamma)\boldsymbol{\alpha} - \gamma\tilde{\boldsymbol{\alpha}}$  {momentum-1}  
15:    $\mathbf{A} \leftarrow K(Z, X)\mathbf{a} + K(Z, X[J])\mathbf{c} \in \mathbb{R}^p$   
16:    $\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} + \text{K\_solve}(K, Z, \mathbf{A})$   
17:    $\mathbf{C} \leftarrow K(Z, X)\mathbf{b} + K(Z, X[J])\mathbf{d} \in \mathbb{R}^n$   
18:    $\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} + \text{solve}(K, Z, \mathbf{C})$   
19: **until** Stopping criterion is reached  
20: **return**  $f(x) = \sum_{i=1}^n \alpha_i K(x, z_i)$

---

---

**Algorithm 5** EigenPro 1 [22]

---

**Input:** kernel  $K$ , batch size  $m$ , learning rate  $\eta_0 > 0$ .  
**Output:**  $f : \mathcal{X} \rightarrow \mathbb{R}$  that solves (9) approximately.  
**setup:**  $\boldsymbol{\alpha} \leftarrow \mathbf{0}_n$ .  
 $(\mathbf{E}_q, \Lambda, \lambda_{q+1}) \leftarrow$  top- $q$  eigensystem of  $K(X, X)$   
 $\mathbf{F}_q \leftarrow \mathbf{E}_q \sqrt{I_q - \lambda_{q+1}\Lambda^{-1}}$   
**repeat**  
  Fetch a batch of indices  $B \subset \{1, 2, \dots, n\}$   
   $\mathbf{v} \leftarrow K(X[B], X)\boldsymbol{\alpha} - Y[B] \in \mathbb{R}^m$   
   $\mathbf{w} \leftarrow \mathbf{F}_q \mathbf{F}_q[B]^\top \mathbf{v} \in \mathbb{R}^n$   
   $\boldsymbol{\alpha}[B] \leftarrow \boldsymbol{\alpha}[B] - \eta_0 \mathbf{v}$  {gradient step}  
   $\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} + \eta_0 \mathbf{w}$  {correction}  
**until** Stopping criterion is reached  
**return**  $f_t(x) = \sum_{i=1}^n \alpha_i K(x, x_i)$

---

sampling. We do not use any regularization. We comment here that even in our experiments the MSE does not decrease significantly due to the subsampling, the classification accuracy is indeed increasing during training.

4. PCG. We choose the rank to be 100 for the pivoted Cholesky decomposition.

5. GPYTORCH. We use instances of the class `IndependentMultitaskGPMModel` in order to deal with multiclass classification problems. The optimizer was set to be ‘Adam’ with learning rate 0.05. As in the case of Falkon, while it is hard to notice significant decrease in the logarithmic scale plot of MSE, the classification accuracy is improved during optimization.

Unless otherwise specified, we do not tune  $\lambda_n$  for AxlePro and use the following setup for our experiments.

For Cifar-10 dataset, we scale the data by a factor of 0.05 as our bandwidth selection. With Gaussian kernel, we set  $(m, q, s) = (2000, 500, 20,000)$ , while for Laplacian kernel, we use  $q = 300$  instead. For the Myrtle5 kernel experiment, we follow the original paper [32] to use ZCA preprocessing (without Leave-One-Outtilting and ZCA augmentation) and store the kernel matrix before performing regression. For this kernel, we set  $(m, q, s) = (2000, 300, 10,000)$ .

---

**Algorithm 6** EigenPro 2 [23]

---

**Input:** kernel  $K$ , batch size  $m$ , indices  $J = \{j_\ell\}_{\ell=1}^s \subset \{1, 2, \dots, n\}$ , learning rate  $\eta_0 > 0$ .  
**Output:**  $f : \mathcal{X} \rightarrow \mathbb{R}$  that solves (9) approximately.  
**setup:**  $\alpha \leftarrow \mathbf{0}_n$ .  
 $(D, \Delta, \delta_{q+1}) \leftarrow$  top- $q$  eigensystem of  $K(X[J], X[J])$   
 $G \leftarrow D\sqrt{\Delta^{-1}(I_q - \delta_{q+1}\Delta^{-1})}$   
**repeat**  
  Fetch a batch of indices  $B \subset \{1, 2, \dots, n\}$   
   $v \leftarrow K(X[B], X)\alpha - Y[B] \in \mathbb{R}^m$   
   $w \leftarrow GG^\top K(X[J], X[B])v \in \mathbb{R}^s$   
   $\alpha[B] \leftarrow \alpha[B] - \eta_0 v$  {gradient step}  
   $\alpha[J] \leftarrow \alpha[J] + \eta_0 w$  {correction}  
**until** Stopping criterion is reached  
**return**  $f_t(x) = \sum_{i=1}^n \alpha_i K(x, x_i)$

---

---

**Algorithm 7** EigenPro 3 [1]

---

**Input:** kernel  $K$ , batch size  $m$ , centers  $Z = \{z_i\}$ , indices  $J = \{j_\ell\}_{\ell=1}^s \subset \{1, 2, \dots, n\}$ , learning rate  $\eta_0 > 0$ .  
**Output:**  $f : \mathcal{X} \rightarrow \mathbb{R}$  that solves (9) approximately.  
**setup:**  $\alpha \leftarrow \mathbf{0}_p$ .  
 $(D, \Delta, \delta_{q+1}) \leftarrow$  top- $q$  eigensystem of  $K(X[J], X[J])$   
 $G \leftarrow D\sqrt{\Delta^{-1}(I_q - \delta_{q+1}\Delta^{-1})}$   
**repeat**  
  Fetch a batch of indices  $B \subset \{1, 2, \dots, n\}$   
   $v \leftarrow K(X[B], Z)\alpha - Y[B] \in \mathbb{R}^m$  {gradient}  
   $w \leftarrow GG^\top K(X[J], X[B])v \in \mathbb{R}^s$  {correction}  
   $h \leftarrow K(Z, X[B])v - K(Z, X[J])w$   
   $\theta \leftarrow \text{K\_solve}(K, Z, h)$  {projection}  
   $\alpha \leftarrow \alpha - \eta_0 \cdot \theta$   
**until** Stopping criterion is reached  
**return**  $f_t(x) = \sum_{i=1}^n \alpha_i K(x, z_i)$

---

For Star Classification dataset, we use  $(m, q, s) = (2000, 150, 20,000)$  for Gaussian kernel (bandwidth=2) and  $(m, q, s) = (2000, 600, 20,000)$  for Laplacian kernel (bandwidth=4). We preprocess the data with mean subtraction and standard deviation normalization.

For EMNIST Digits dataset, we preprocess the data by mean subtraction and scale by a factor of 0.001. We set  $(m, q, s) = (700, 600, 20,000)$  for Gaussian kernel (bandwidth=1) and  $(m, q, s) = (800, 600, 20,000)$  for Laplacian kernel (bandwidth=1).

## E Exact preconditioner

Calculating the exact preconditioner requires finding the top- $q$  eigensystem of  $K(X, X)$  which can cost  $n^2q$  FLOPs. Instead, we can use a Nyström extension. Let  $J = \{j_1, j_2, \dots, j_s\} \subseteq \{1, 2, \dots, n\}$ . Next, let  $(\tilde{E}, \Lambda, \lambda_{q+1})$  be the top- $q$  eigensystem of  $K(X[J], X[J])$ , where the columns of  $\tilde{E} \in \mathbb{R}^{s \times q}$  are the eigenvectors of  $K(X[J], X[J])$ . Finally, obtain approximate eigenvectors  $E \in \mathbb{R}^{n \times q}$  of  $K(X, X)$  as

$$E \leftarrow K(X, X[J])\tilde{E}$$

Due to the matrix multiplication, the vectors  $E$  obtained as such are not orthonormal. Hence we run a thin QR decomposition to refine  $E$  as,

$$E_{,-} \leftarrow \text{thinQR}(E).$$

This only has a total complexity  $s^2q + qsn + nq^2$ , where  $s^2q$  is the cost of getting eigenvectors  $\tilde{E}$ , and  $sqn$  is the cost of Nyström extension, and  $nq^2$  is the cost of thin QR decomposition. Since  $q \leq s \leq n$ , the term  $nsq$  dominates.

---

**Algorithm 8** EigenPro 4 [2]

---

**Input:** kernel  $K$ , batch size  $m$ , centers  $Z = \{z_i\}$ , indices  $J = \{j_\ell\}_{\ell=1}^s \subset \{1, 2, \dots, n\}$ , learning rate  $\eta_0 > 0$ .

**Output:**  $f : \mathcal{X} \rightarrow \mathbb{R}$  that solves (9) approximately.

**setup:**  $\alpha \leftarrow \mathbf{0}_p, \mathbf{a} \leftarrow \mathbf{0}_n, \mathbf{c} \leftarrow \mathbf{0}_{|J|}, \mathbf{h} \leftarrow \mathbf{0}_0$

$(\mathbf{D}, \Delta, \delta_{q+1}) \leftarrow \text{top-}q \text{ eigensystem of } K(X[J], X[J])$

$\mathbf{G} \leftarrow \mathbf{D} \sqrt{\Delta^{-1}} (\mathbf{I}_q - \delta_{q+1} \Delta^{-1})$

**repeat**

Fetch a batch of indices  $B \subset \{1, 2, \dots, n\}$

$\mathbf{v} \leftarrow K(X[B], Z)\alpha + K(X[B], X)\mathbf{a} + K(X[B], X[J])\mathbf{c} - Y[B] \in \mathbb{R}^m$

{gradient}

$\mathbf{w} \leftarrow \mathbf{G}\mathbf{G}^\top K(X[J], X[B])\mathbf{v} \in \mathbb{R}^s$

{correction}

$\mathbf{a}[B] \leftarrow -\eta_0 \mathbf{v}$

{update temporary weights-1}

$\mathbf{c} \leftarrow \mathbf{c} + \eta_0 \mathbf{w}$

{update temporary weights-2}

$\mathbf{h} \leftarrow \mathbf{h} - \eta_0 (K(Z, X[B])\mathbf{v} + K(Z, X[J])\mathbf{w})$

{accumulate gradients}

**if** projection condition holds **then**

$\boldsymbol{\theta} \leftarrow \text{K\_solve}(K, Z, \mathbf{h})$

{delayed projection}

$\alpha \leftarrow \alpha - \eta_0 \cdot \boldsymbol{\theta}$

Reset  $\mathbf{a} \leftarrow \mathbf{0}_n, \mathbf{c} \leftarrow \mathbf{0}_s, \mathbf{h} \leftarrow \mathbf{0}_p$ ,

**end if**

**until** Stopping criterion is reached

**return**  $f_t(x) = \sum_{i=1}^n \alpha_i K(x, z_i)$

---

## F Details on convergence analysis

*Proof of Lemma 7(d).* We consider the norms restricted to the subspace  $\text{span}(\{K(\cdot, x_i)\}_{i=1}^n)$  since we only care about  $\|f\|_{\mathbb{H}}$  and  $\|\hat{f}\|_{\mathbb{H}_q}$  for  $f = \sum_{i=1}^n \alpha_i K(\cdot, x_i)$ .

Now, suppose  $f = \hat{f} = \sum_{i=1}^n \alpha_i \hat{K}_q(\cdot, x_i)$ .

Claim:  $\text{span}(\{K(\cdot, x_i)\}_{i=1}^n) = \text{span}(\{\hat{K}_q(\cdot, x_i)\}_{i=1}^n)$ .

W.L.O.G, we assume  $J = (1, 2, \dots, s)$ . Denote  $\mathbf{v}_i := \mathbf{G}\mathbf{G}^\top K(X[J], x_i) \in \mathbb{R}^s$

$$\begin{aligned} \hat{K}_q(\cdot, x_i) &= K(\cdot, x_i) - K(\cdot, X[J])\mathbf{G}\mathbf{G}^\top K(X[J], x_i) \\ &= K(\cdot, X)(e_i - [\mathbf{v}_i \ \mathbf{0}_{n-s}]^T) \end{aligned}$$

So  $\sum_{i=1}^n \alpha_i \hat{K}_q(\cdot, x_i) = \hat{K}_q(\cdot, X)\alpha = K(\cdot, X)A\alpha$ , where  $A := \mathbf{I}_n - \begin{pmatrix} \mathbf{G}\mathbf{G}^\top K(J, X) \\ \mathbf{0}_{(n-s) \times n} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_s - \mathbf{G}\mathbf{G}^\top K(J, J) & -\mathbf{G}\mathbf{G}^\top K(J, J^C) \\ \mathbf{0}_{(n-s) \times s} & \mathbf{I}_{(n-s) \times (n-s)} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_s - \mathbf{D}(\mathbf{I}_q - \delta_{q+1}\Delta^{-1})\mathbf{D}^\top & -\mathbf{G}\mathbf{G}^\top K(J, J^C) \\ \mathbf{0}_{(n-s) \times s} & \mathbf{I}_{(n-s) \times (n-s)} \end{pmatrix}$ . Since  $A$  is an upper triangular block matrix, we can easily get  $\frac{\delta_{q+1}}{\delta_1} = \min(\frac{\delta_{q+1}}{\delta_1}, 1) \leq \lambda(A) \leq \max(\frac{\delta_{q+1}}{\delta_q}, 1) = 1$ . This also implies  $A$  is invertible and hence our Claim holds.

As a consequence,

$$\begin{aligned} \|f\|_{\mathbb{H}}^2 &= (A\alpha)^\top K(X, X)(A\alpha) \\ &\leq \lambda_{\max}(K(X, X)) \|A\alpha\|_2^2 \\ &\leq \lambda_{\max}(K(X, X)) \lambda_{\max}^2(A) \|\alpha\|_2^2 \\ &\leq \lambda_{\max}(K(X, X)) \|\alpha\|_2^2. \end{aligned}$$

Similarly,  $\|f\|_{\mathbb{H}}^2 \geq \lambda_{\min}(K(X, X)) \lambda_{\min}^2(A) \|\alpha\|_2^2 \geq \lambda_{\min}(K(X, X)) \frac{\delta_{q+1}^2}{\delta_1^2} \|\alpha\|_2^2$

On the other hand,

$$\|\hat{f}\|_{\mathbb{H}_q}^2 = \alpha^\top \hat{K}_q(X, X)\alpha \leq \lambda_{\max}(\hat{K}_q(X, X)) \|\alpha\|_2^2.$$

Similarly,  $\|\hat{f}\|_{\mathbb{H}_q}^2 \geq \lambda_{\min}(\widehat{K}_q(X, X)) \|\alpha\|_2^2$ .

So by proposition 8, with high probability, there exist constants  $a, a'$  depending (polynomially) on the eigenvalues of  $K(X, X)$  such that  $a \|f\|_{\mathbb{H}} \leq \|\hat{f}\|_{\mathbb{H}_q} \leq a' \|f\|_{\mathbb{H}}$ .  $\square$

### F.1 Condition numbers for analysis of MaSS

Let  $\lambda_1$  and  $\lambda_n$  be the largest and smallest non-zero eigenvalues of  $\mathcal{K}$ .

Let a mini-batch of size  $m$  be  $\{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^m$ . Defining mini-batch covariance operator as follows

$$\tilde{\mathcal{K}}^{(m)} := \frac{1}{m} \sum_{i=1}^m K(\tilde{x}_i, \cdot) \otimes K(\tilde{x}_i, \cdot) \quad (35)$$

$$\tilde{\mathcal{K}}_{\mathcal{P}}^{(m)} := \frac{1}{m} \sum_{i=1}^m k_{\mathcal{P}}(\tilde{x}_i, \cdot) \otimes k_{\mathcal{P}}(\tilde{x}_i, \cdot) \quad (36)$$

where  $\tilde{\mathcal{K}}_{\mathcal{P}}$  is the covariance operator after preconditioning. The kernel function  $k_{\mathcal{P}}$  can be defined in the same way as in Lemma 7 with the exact preconditioner.

### F.2 Identifying quantities defined by [20]

Suppose we run MaSS to solve  $\min_{f \in \mathbb{H}} \frac{1}{n} \sum_{i=1}^n (S_{\{i\}} f - Y_i)^2$ , with Hessian  $H = \mathcal{K}$ . Note that  $S_{\{i\}} f = \langle K(x_i, \cdot), f \rangle_{\mathbb{H}}$ . The authors define the quantities  $L, \mu$  to be the largest and smallest eigenvalues of the Hessian  $H : \mathbb{H} \mapsto \mathbb{H}$ , and  $\kappa = L/\mu$  to be the condition number. In our case,  $L = \lambda_1/n$ ,  $\mu = \lambda_n/n$  and  $\kappa = \frac{\lambda_1}{\lambda_n}$ .

They also define the quantity  $L_1$  to be the smallest number such that

$$\mathbb{E}[\|K(\tilde{x}, \cdot)\|_{\mathbb{H}}^2 K(\tilde{x}, \cdot) \otimes K(\tilde{x}, \cdot)] \lesssim L_1 H \quad (37)$$

where  $\mathbb{E}$  is over the random variable  $\tilde{x}$  from the empirical distribution. Note that in our case  $\|K(\tilde{x}, \cdot)\|_{\mathbb{H}}^2 = K(\tilde{x}, \tilde{x})$  since,

$$\|K(x, \cdot)\|_{\mathbb{H}}^2 = \langle K(x, \cdot), K(x, \cdot) \rangle_{\mathbb{H}} = K(x, x) \leq \beta := \max_i K(x_i, x_i) \quad (38)$$

Then we have  $L_1 \leq \beta$ , since

$$\mathbb{E}[\|K(\tilde{x}, \cdot)\|^2 K(\tilde{x}, \cdot) \otimes K(\tilde{x}, \cdot)] \preceq \beta \mathbb{E}[K(\tilde{x}, \cdot) \otimes K(\tilde{x}, \cdot)] = \beta \mathcal{K}$$

Thus

$$L_m := \frac{L_1}{m} + \frac{(m-1)L}{m} \leq \frac{\beta + (m-1)\frac{\lambda_1}{n}}{m} \quad (39)$$

Defining  $\tilde{\kappa}$  as the smallest positive real number such that

$$\mathbb{E}[\|K(\tilde{x}, \cdot)\|_{\mathcal{K}^{-1}}^2 K(\tilde{x}, \cdot) \otimes K(\tilde{x}, \cdot)] \preceq \tilde{\kappa} \mathcal{K} \quad (40)$$

Deriving  $\|K(x_i, \cdot)\|_{\mathcal{K}^{-1}}^2$ . Recall that  $\psi_i = S^* \mathbf{e}_i / \sqrt{\lambda_i}$ . Also note that  $\mathcal{K}^{-1} = \sum_{i=1}^n \frac{n}{\lambda_i} \psi_i \otimes \psi_i$ .

$$\begin{aligned} \|K(x_i, \cdot)\|_{\mathcal{K}^{-1}}^2 &= \langle K(x_i, \cdot), \mathcal{K}^{-1} K(x_i, \cdot) \rangle_{\mathbb{H}} \\ &= \left\langle K(x_i, \cdot), \left( \sum_{j=1}^n \frac{n}{\lambda_j} \psi_j \otimes \psi_j \right) K(x_i, \cdot) \right\rangle_{\mathbb{H}} \\ &= \sum_{j=1}^n \frac{n}{\lambda_j} \langle \psi_j, K(x_i, \cdot) \rangle_{\mathbb{H}}^2 = \sum_{j=1}^n \frac{n}{\lambda_j} \psi_j^2(x_i) \stackrel{(a)}{=} \sum_{j=1}^n \frac{1}{\lambda_j} n \lambda_j e_{ji}^2 \end{aligned} \quad (41a)$$

$$= n \sum_{j=1}^n e_{ji}^2 = n \|\mathbf{e}_j\|^2 = n \quad (41b)$$

where (a) follows from the fact that  $\psi_j(x_i) = S_{\{i\}}\psi_j = S_{\{i\}}S^*e_j/\sqrt{\lambda_j} = \mathbf{H}_{\{i\}}SS^*e_j/\sqrt{\lambda_j} = \mathbf{H}_{\{i\}}e_j\sqrt{\lambda_j}$ .

Deriving  $\tilde{\kappa}$

$$\mathbb{E} \left[ \|K(\tilde{x}, \cdot)\|_{\mathcal{K}^{-1}}^2 K(\tilde{x}, \cdot) \otimes K(\tilde{x}, \cdot) \right] \stackrel{(a)}{=} n \mathbb{E} [K(\tilde{x}, \cdot) \otimes K(\tilde{x}, \cdot)] = n \mathbb{E} [\tilde{\mathcal{K}}^{(1)}] = n \mathcal{K} \stackrel{(b)}{\Rightarrow} \tilde{\kappa} = n$$

where (a) follows from equation (41b) and (b) from definition of  $\tilde{\kappa}$  in equation (40). This also implies  $\tilde{c}_m = \tilde{\kappa}_m$  in lemma 7 since (b) holds for any kernel function.

### F.3 Formulae for hyperparameters

$$L_1 = \beta, L = \frac{\lambda_1}{n}, \tilde{\kappa} = n$$

$$L_m = \frac{L_1}{m} + \frac{(m-1)L}{m} = \frac{\beta + (m-1)\frac{\lambda_1}{n}}{m} \quad (42a)$$

$$\kappa_m = \frac{nL_m}{\lambda_n} \quad (42b)$$

$$\tilde{\kappa}_m = \frac{\tilde{\kappa}}{m} + \frac{m-1}{m} = 1 + \frac{n-1}{m} \quad (42c)$$

$$\eta_1(m) = \frac{1}{L_m} = \frac{m}{\beta + (m-1)\frac{\lambda_1}{n}} \quad (42d)$$

$$\eta_2(m) = \eta_1 \frac{\sqrt{\kappa_m \tilde{\kappa}_m}}{1 + \sqrt{\kappa_m \tilde{\kappa}_m}} \left( 1 - \frac{1}{\tilde{\kappa}_m} \right) \quad (42e)$$

$$\gamma(m) = \frac{\sqrt{\kappa_m \tilde{\kappa}_m} - 1}{\sqrt{\kappa_m \tilde{\kappa}_m} + 1} \quad (42f)$$

### F.4 Condition numbers after preconditioning

Note that after preconditioning we are operating in Hilbert space  $\mathbb{H}_{\mathcal{P}}$ . Also, the largest eigenvalue of  $\mathcal{K}_{\mathcal{P}}$  is  $\lambda_{q+1}$  i.e.,  $L = \lambda_{q+1}$

Similar to equation (37), defining  $L_1$  as the smallest positive number such that

$$\mathbb{E} \left[ \|k_{\mathcal{P}}(\tilde{x}, \cdot)\|^2 k_{\mathcal{P}}(\tilde{x}, \cdot) \otimes k_{\mathcal{P}}(\tilde{x}, \cdot) \right] \preceq L_1 \mathcal{K}_{\mathcal{P}} \quad (43)$$

Deriving  $\|k_{\mathcal{P}}(\mathbf{x}_i, \cdot)\|^2$

$$\begin{aligned} \|k_{\mathcal{P}}(\mathbf{x}_i, \cdot)\|^2 &= \langle k_{\mathcal{P}}(\mathbf{x}_i, \cdot), k_{\mathcal{P}}(\mathbf{x}_i, \cdot) \rangle_{\mathbb{H}_{\mathcal{P}}} \\ &= k_{\mathcal{P}}(\mathbf{x}_i, \mathbf{x}_i) \\ &\stackrel{(a)}{=} K(\mathbf{x}_i, \mathbf{x}_i) - \sum_{j=1}^q \left( 1 - \frac{\lambda_{q+1}}{\lambda_j} \right) \lambda_j e_{ji}^2 \\ &= K(\mathbf{x}_i, \mathbf{x}_i) - \sum_{j=1}^q (\lambda_j - \lambda_{q+1}) e_{ji}^2 \end{aligned} \quad (44)$$

where (a) is from the definition of  $k_{\mathcal{P}}$ .

Define  $\beta_{\mathcal{P}}$  as the maximum norm of  $k_{\mathcal{P}}(\mathbf{x}_i, \cdot)$  among the samples

$$\beta_{\mathcal{P}} := \max_i \|k_{\mathcal{P}}(\mathbf{x}_i, \cdot)\|^2 = \max_i \left\{ K(\mathbf{x}_i, \mathbf{x}_i) - n \sum_{j=1}^q (\lambda_j - \lambda_{q+1}) e_{ji}^2 \right\} \quad (45)$$

Deriving  $L_1$

$$\mathbb{E} \left[ \|k_{\mathcal{P}}(\tilde{x}, \cdot)\|^2 k_{\mathcal{P}}(\tilde{x}, \cdot) \otimes k_{\mathcal{P}}(\tilde{x}, \cdot) \right] \preceq \beta_{\mathcal{P}} \mathbb{E} [k_{\mathcal{P}}(\tilde{x}, \cdot) \otimes k_{\mathcal{P}}(\tilde{x}, \cdot)] = \beta_{\mathcal{P}} \mathbb{E} [\tilde{\mathcal{K}}_{\mathcal{P}}^{(1)}] = \beta_{\mathcal{P}} \mathcal{K}_{\mathcal{P}}$$

which implies  $L_1 \leq \beta_{\mathcal{P}}$  due to the definition of  $L_1$  in equation (43).

Deriving  $L_m$

$$L_m := \frac{L_1}{m} + \frac{(m-1)L}{m} \leq \frac{\beta_{\mathcal{P}} + (m-1)\lambda_{q+1}}{m} \quad (46)$$

Defining  $\tilde{\kappa}$  as the smallest positive real number such that

$$\mathbb{E} \left[ \|k_{\mathcal{P}}(\tilde{x}, \cdot)\|_{\mathcal{K}_{\mathcal{P}}^{-1}}^2 k_{\mathcal{P}}(\tilde{x}, \cdot) \otimes k_{\mathcal{P}}(\tilde{x}, \cdot) \right] \preceq \tilde{\kappa} \mathcal{K}_{\mathcal{P}} \quad (47)$$

Deriving  $\|k_{\mathcal{P}}(\mathbf{x}_i, \cdot)\|_{\mathcal{K}_{\mathcal{P}}^{-1}}^2$

$$\begin{aligned} \|k_{\mathcal{P}}(\mathbf{x}_i, \cdot)\|_{\mathcal{K}_{\mathcal{P}}^{-1}}^2 &= \langle k_{\mathcal{P}}(\mathbf{x}_i, \cdot), \mathcal{K}_{\mathcal{P}}^{-1} k_{\mathcal{P}}(\mathbf{x}_i, \cdot) \rangle_{\mathbb{H}_{\mathcal{P}}} \\ &\stackrel{(a)}{=} \left\langle k_{\mathcal{P}}(\mathbf{x}_i, \cdot), \left( \sum_{j=1}^q \frac{1}{\lambda_{q+1}} \psi'_j \otimes \psi'_j + \sum_{j=q+1}^n \frac{1}{\lambda_j} \psi'_j \otimes \psi'_j \right) k_{\mathcal{P}}(\mathbf{x}_i, \cdot) \right\rangle_{\mathbb{H}_{\mathcal{P}}} \\ &= \sum_{j=1}^q \frac{1}{\lambda_{q+1}} \langle \psi'_j, k_{\mathcal{P}}(\mathbf{x}_i, \cdot) \rangle_{\mathbb{H}_{\mathcal{P}}}^2 + \sum_{j=q+1}^n \frac{1}{\lambda_j} \langle \psi'_j, k_{\mathcal{P}}(\mathbf{x}_i, \cdot) \rangle_{\mathbb{H}_{\mathcal{P}}}^2 \\ &\stackrel{(b)}{=} n \left\{ \sum_{j=1}^q e_{ji}^2 + \sum_{j=q+1}^n e_{ji}^2 \right\} \\ &\stackrel{(c)}{=} n \end{aligned} \quad (48)$$

where (a) is from the eigendecomposition of  $\mathcal{K}_{\mathcal{P}}$ , (b) follows from using eigenfunction evaluation with  $\psi'$  in  $\mathbb{H}_{\mathcal{P}}$  and (c) follows from the fact that  $\mathbf{E} = [\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_n]$  is an orthonormal matrix which implies  $\mathbf{E}^\top \mathbf{E} = \mathbf{I} = \mathbf{E} \mathbf{E}^\top$ .

Deriving  $\tilde{\kappa}$

$$\mathbb{E} \left[ \|k_{\mathcal{P}}(\tilde{x}, \cdot)\|_{\mathcal{K}_{\mathcal{P}}^{-1}}^2 k_{\mathcal{P}}(\tilde{x}, \cdot) \otimes k_{\mathcal{P}}(\tilde{x}, \cdot) \right] \quad (49)$$

$$\stackrel{(a)}{=} n \mathbb{E} [k_{\mathcal{P}}(\tilde{x}, \cdot) \otimes k_{\mathcal{P}}(\tilde{x}, \cdot)] \quad (50)$$

$$\begin{aligned} &= n \mathbb{E} \left[ \widetilde{\mathcal{K}_{\mathcal{P}}}^{(1)} \right] \\ &= n \mathcal{K}_{\mathcal{P}} \end{aligned} \quad (51)$$

$$\stackrel{(b)}{\implies} \tilde{\kappa} = n \quad (52)$$

where (a) follows from equation (48) and (b) from definition of  $\tilde{\kappa}$  in equation (47).

## F.5 MaSS parameters after preconditioning

Using (46),  $\eta_1$  defined in (23a) can be written as

$$\eta_1^*(m) = \frac{m}{\beta_{\mathcal{P}} + (m-1)\lambda_{q+1}}$$

Using (52), (42), we can write  $\eta_2, \gamma$  defined in (23b), (23c) as

$$\begin{aligned} \eta_2^*(m) &= \frac{\eta_1^*(m) \sqrt{(\beta_{\mathcal{P}} + (m-1)\lambda_{q+1})(n+m-1)}}{\sqrt{(\beta_{\mathcal{P}} + (m-1)\lambda_{q+1})(n+m-1)} + m\sqrt{\lambda_n}} \\ \gamma^*(m) &= \frac{\sqrt{(\beta_{\mathcal{P}} + (m-1)\lambda_{q+1})(n+m-1)} - m\sqrt{\lambda_n}}{\sqrt{(\beta_{\mathcal{P}} + (m-1)\lambda_{q+1})(n+m-1)} + m\sqrt{\lambda_n}} \end{aligned}$$

Since the kernel matrix is positive definite we see that  $\lambda_{q+1} \geq \beta/n$  and using (52) we get regime critical points



$m_1^*, m_2^*$  defined in [20] as

$$m_1^* = \min \left( \frac{\beta_{\mathcal{P}}}{\lambda_{q+1}}, n \right) = \frac{\beta_{\mathcal{P}}}{\lambda_{q+1}}$$

$$m_2^* = \max \left( \frac{\beta_{\mathcal{P}}}{\lambda_{q+1}}, n \right) = n$$

Note that there is no saturation regime for kernel methods

If optimal mini-batch size for linear regime  $m_* := m_1^*$  is used, then optimal MaSS parameters are

$$\eta_1^* = \frac{m_*^2}{\beta_{\mathcal{P}}(2m_* - 1)}$$

$$\eta_2^* = \frac{m_*^2 \sqrt{n + m_* - 1}}{\beta_{\mathcal{P}}(2m_* - 1) \sqrt{n + m_* - 1} + m_* \sqrt{m_* \lambda_n \beta_{\mathcal{P}}(2m_* - 1)}}$$

$$\gamma^* = \frac{\sqrt{\beta_{\mathcal{P}}(2m_* - 1)} \sqrt{n + m_* - 1} - m_* \sqrt{m_* \lambda_n}}{\sqrt{\beta_{\mathcal{P}}(2m_* - 1)} \sqrt{n + m_* - 1} + m_* \sqrt{m_* \lambda_n}}$$

**Informal:** Assume large natural image dataset with gaussian/laplacian kernel. Due to the eigenvalue decay,  $m_*$  is reasonably large and  $\beta_{\mathcal{P}}$  is usually very close to 1. The MaSS parameters can be approximated as follows

$$\eta_1^* \approx \frac{m_*}{2}$$

$$\eta_2^* \approx \frac{m_* \sqrt{n + m_* - 1}}{2 \sqrt{n + m_* - 1} + m_* \sqrt{2 \lambda_n}}$$

$$\gamma^* \approx \frac{\sqrt{2} \sqrt{n + m_* - 1} - m_* \sqrt{\lambda_n}}{\sqrt{2} \sqrt{n + m_* - 1} + m_* \sqrt{\lambda_n}}$$

**Proposition 10.** Suppose  $m \leq \frac{n}{2}$ , we have

$$m \sqrt{\tilde{\kappa}_m \kappa_m} \leq n \sqrt{\kappa}$$

*Proof.* According to the definitions of  $\tilde{\kappa}_m$  and  $\kappa_m$ , the statement is equivalent to

$$m^2 \left( \frac{n}{m} + \frac{m-1}{m} \right) \left( \frac{L_1 + (m-1)\lambda_{q+1}}{\lambda_n m} \right) \leq n^2 \frac{\lambda_1}{\lambda_n},$$

i.e.

$$(n + m - 1)(L_1 + (m - 1)\lambda_{q+1}) \leq n^2 \lambda_1.$$

Since  $L_1 = \max_i K(x_i, x_i) \leq \lambda_1$ , we have

$$\begin{aligned} (n + m - 1)(L_1 + (m - 1)\lambda_{q+1}) &\leq 2n(\lambda_1 + (m - 1)\lambda_1) \\ &\leq 2mn\lambda_1 \\ &\leq n^2 \lambda_1. \end{aligned}$$

□

## G Acceleration of EigenPro 4

Accelerated iteration of EigenPro 3 in RKHS:

$$f_{t+1} \leftarrow \text{proj}_{\mathcal{Z}}(g_t - \eta_1 \mathcal{P}_s \tilde{\nabla}_f \mathcal{L}(g_t))$$

$$g_{t+1} \leftarrow \text{proj}_{\mathcal{Z}}((1 + \gamma)f_{t+1} - \gamma f_t + \eta_2 \mathcal{P}_s \tilde{\nabla}_f \mathcal{L}(g_t))$$

Suppose you project after every  $T$  iterations.

Consider the AxlePro2 iterations

$$f_{t+1} = g_t - \eta \mathcal{P}_s \tilde{\nabla}_f \mathbf{L}(g_t) \quad (53)$$

$$g_{t+1} = (1 + \gamma)f_{t+1} - \gamma f_t + \eta_2 \mathcal{P}_s \tilde{\nabla}_f \mathbf{L}(g_t) \quad (54)$$

for  $t = 1, 2, \dots, T$  and  $f_T = \text{proj}_{\mathcal{Z}}(f_T)$  which requires  $f_T(Z)$  and  $g_T(Z)$  which will be computed iteratively.

Suppose  $B_t$  is the batch at step  $t$ , define  $C_t := C_{t-1} \cup B_t$  and  $C_0 = \emptyset$  and let  $f_t = S_Z^* \alpha_t^Z + S_{C_t}^* \alpha_t^X + S_J^* \alpha_t^J$  and  $g_t = S_Z^* \beta_t^Z + S_{C_t}^* \beta_t^X + S_J^* \beta_t^J$ .

Then  $\tilde{\nabla}_f \mathbf{L}(g_t) \in \text{range}(S_{B_{t+1}}^*)$  given by

$$\begin{aligned} \tilde{\nabla}_f \mathbf{L}(g_t) &= S_{B_{t+1}}^* \mathbf{v}_t \\ \mathbf{v}_t &:= g_t(X[B_{t+1}]) - y_t \\ g_t(X[B_{t+1}]) &= K(X[B_{t+1}], Z) \beta_t^Z + K(X[B_{t+1}], X[C_t]) \beta_t^X + K(X[B_{t+1}], X[J]) \beta_t^J \\ \mathbf{w}_t &:= \mathbf{G} \mathbf{G}^\top K(X[J], X[B_{t+1}]) \mathbf{v}_t \\ \mathcal{P}_s \tilde{\nabla}_f \mathbf{L}(g_t) &= S_{B_{t+1}}^* \mathbf{v}_t - S_J^* \mathbf{w}_t \end{aligned}$$

$$\begin{aligned} S_Z^* \alpha_{t+1}^Z + S_{C_{t+1}}^* \alpha_{t+1}^X + S_J^* \alpha_{t+1}^J &\leftarrow S_Z^* \beta_t^Z + S_{C_t}^* \beta_t^X + S_J^* \beta_t^J - \eta_1 (S_{B_{t+1}}^* \mathbf{v}_t - S_J^* \mathbf{w}_t) \\ S_Z^* \beta_{t+1}^Z + S_{C_{t+1}}^* \beta_{t+1}^X + S_J^* \beta_{t+1}^J &\leftarrow (1 + \gamma)(S_Z^* \alpha_{t+1}^Z + S_{C_t}^* \alpha_{t+1}^X + S_J^* \alpha_{t+1}^J) - \gamma(S_Z^* \alpha_t^Z + S_{C_t}^* \alpha_t^X + S_J^* \alpha_t^J) \\ &\quad + \eta_2 (S_{B_{t+1}}^* \mathbf{v}_t - S_J^* \mathbf{w}_t) \end{aligned}$$

This gives us

$$\alpha_{t+1}^Z \leftarrow \beta_t^Z \quad (55a)$$

$$\beta_{t+1}^Z \leftarrow (1 + \gamma) \alpha_{t+1}^Z - \gamma \alpha_t^Z \quad (55b)$$

$$\alpha_{t+1}^J \leftarrow \beta_t^J + \eta_1 \mathbf{w}_t \quad (55c)$$

$$\beta_{t+1}^J \leftarrow (1 + \gamma) \alpha_{t+1}^J - \gamma \alpha_t^J - \eta_2 \mathbf{w}_t \quad (55d)$$

$$\alpha_{t+1}^X[C_t] \leftarrow \beta_t^X[C_t] \quad (55e)$$

$$\alpha_{t+1}^X[B_{t+1}] \leftarrow -\eta_1 \mathbf{v}_t \quad (55f)$$

$$\beta_{t+1}^X[C_t] \leftarrow (1 + \gamma) \alpha_{t+1}^X[C_t] - \gamma \alpha_t^X[C_t] \quad (55g)$$

$$\beta_{t+1}^X[B_{t+1}] \leftarrow (1 + \gamma) \alpha_{t+1}^X[B_{t+1}] + \eta_2 \mathbf{v}_t = (\eta_2 - (1 + \gamma) \eta_1) \mathbf{v}_t \quad (55h)$$

*Claim 1.* If we run the extrapolation iteration  $a_{t+1} = (1 + \gamma)a_t - \gamma a_{t-1}$  initialized at  $a_1 \neq a_0$ , then

$$a_t = a_0 + \frac{1 - \gamma^t}{1 - \gamma} (a_1 - a_0). \quad (56)$$

*Proof.* We will prove this by induction. Observe that the  $t = 1$  case hold trivially. Assume that the formula is true for  $t \leq \tau$  for some  $\tau \geq 1$ . If we prove that the formula holds for  $t = \tau + 1$ , we are done.

$$a_{\tau+1} = (1 + \gamma)a_\tau - \gamma a_{\tau-1} = (1 + \gamma) \left( a_0 + \frac{1 - \gamma^\tau}{1 - \gamma} (a_1 - a_0) \right) - \gamma \left( a_0 + \frac{1 - \gamma^{\tau-1}}{1 - \gamma} (a_1 - a_0) \right) \quad (57)$$

$$= a_0 + \frac{a_1 - a_0}{1 - \gamma} ((1 + \gamma)(1 - \gamma^\tau) - \gamma(1 - \gamma^{\tau-1})) \quad (58)$$

$$= a_0 + \frac{a_1 - a_0}{1 - \gamma} (1 + \gamma - \gamma^\tau - \gamma^{\tau+1} - \gamma + \gamma^\tau) = a_0 + \frac{1 - \gamma^{\tau+1}}{1 - \gamma} (a_1 - a_0). \quad (59)$$

which proves the claim.  $\square$

Let  $t_B$  be the time step when batch  $B$  is chosen for the first time, i.e.,  $B_t = B$ .

$$\beta_t^X[B] = c_t \mathbf{v}_{t_B} \quad \text{where} \quad c_t = \begin{cases} 0 & t < t_B \\ c_1 & t = t_B \\ (c_1 + \frac{1-\gamma^{(t-t_B+1)}}{1-\gamma} c_2) & t > t_B \end{cases} \quad (60)$$

where  $c_1 = -((1+\gamma)\eta_1 - \eta_2)$ , and  $c_2 = \eta_2 - \eta_1$

Finally, after  $T$  iterations, we need to project  $f_T$  and  $g_T$  back onto  $\text{span}(S_Z^*)$ . To that end, we must calculate  $g_T(Z)$ . To that end, consider

$$f_T(Z) = K(Z, Z)\alpha_T^Z + K(Z, X[C_T])\alpha_T^X + K(Z, X[J])\alpha_T^J \quad (61)$$

$$g_T(Z) = K(Z, Z)\beta_T^Z + K(Z, X[C_T])\beta_T^X + K(Z, X[J])\beta_T^J \quad (62)$$

We can compute

$$f_{T+1} = \text{proj}_Z(f_T) \quad (63)$$

$$g_{T+1} = \text{proj}_Z(g_T) \quad (64)$$

by

$$f_{T+1} = S_Z^* K(Z, Z)^{-1} f_T(Z) \quad (65)$$

$$g_{T+1} = S_Z^* K(Z, Z)^{-1} g_T(Z) \quad (66)$$

which can be implemented as

$$\alpha_{T+1}^Z = K(Z, Z)^{-1} f_T(Z) = \alpha_T^Z + K(Z, Z)^{-1} (K(Z, X[C_T])\alpha_T^X + K(Z, X[J])\alpha_T^J) \quad (67)$$

$$\beta_{T+1}^Z = K(Z, Z)^{-1} g_T(Z) = \beta_T^Z + K(Z, Z)^{-1} (K(Z, X[C_T])\beta_T^X + K(Z, X[J])\beta_T^J) \quad (68)$$

$$\mathbf{v}_t \leftarrow K(X[B], Z)\beta_t - Y[B] \in \mathbb{R}^m \quad (69a)$$

$$\mathbf{w}_t \leftarrow K(Z, X[B])\mathbf{v}_t - K(Z, X[J])\mathbf{G}\mathbf{G}^\top K(X[J], X[B])\mathbf{v}_t \in \mathbb{R}^{|Z|} \quad (69b)$$

$$\mathbf{u}_t \leftarrow \text{solve } K(Z, Z)\mathbf{u}_t = \mathbf{w}_t \quad (69c)$$

$$\alpha_{t+1} \leftarrow \beta_t - \eta_1 \mathbf{u}_t \quad (69d)$$

$$\beta_{t+1} \leftarrow (1+\gamma)\alpha_{t+1} - \gamma\alpha_t + \eta_2 \mathbf{u}_t \quad (69e)$$

emulate the updates Acceleration algorithm of EigenPro 3.

## H ADDITIONAL PROOFS

*Proof of Proposition 1.* For  $f_t = S^* \alpha_t$ , a stochastic gradient with respect to the mini-batch  $(X[B], Y[B])$  is given by

$$\tilde{\nabla}_f \mathcal{L}(f_t) = \frac{1}{|B|} S_B^* (S_B f_t - Y[B]) \quad (70a)$$

$$= \frac{1}{|B|} S_B^* (S_B S^* \alpha_t - Y[B]) = S_B^* \mathbf{v}_t \quad (70b)$$

Observe that  $S_B = \mathbf{H}_B S$ , whereby  $S_B^* = S^* \mathbf{H}_B^\top$ . The claim follows immediately.  $\square$

*Proof of Proposition 3.* We start by showing that

$$\mathcal{P} S_B^* = S_B^* - S^* \mathbf{F} \mathbf{F}^\top \mathbf{H}_B^\top. \quad (71)$$

For a vector  $\mathbf{u} \in \mathbb{R}^m$ , observe that

$$\mathcal{P}S_B^* \mathbf{u} = S_B^* \mathbf{u} - \sum_{i=1}^q \left(1 - \frac{\lambda_{q+1}}{\lambda_i}\right) \psi_i \otimes_{\mathbb{H}} \psi_i S_B^* \mathbf{u} \quad (72)$$

Now the term  $\psi_i \otimes_{\mathbb{H}} \psi_i S_B^* \mathbf{u}$  simplifies as

$$\begin{aligned} \psi_i \langle \psi_i, S_B^* \mathbf{u} \rangle_{\mathbb{H}} &= \frac{1}{\lambda_i} S^* \mathbf{e}_i \langle S^* \mathbf{e}_i, S_B^* \mathbf{u} \rangle_{\mathbb{H}} \\ &= \frac{1}{\lambda_i} S^* \mathbf{e}_i \langle \mathbf{e}_i, S S^* \mathbf{H}_B^{\top} \mathbf{u} \rangle_{\mathbb{R}^n} \\ &= \frac{1}{\lambda_i} S^* \mathbf{e}_i \mathbf{e}_i^{\top} K(X, X) \mathbf{H}_B^{\top} \mathbf{u} \\ &= S^* \mathbf{e}_i \mathbf{e}_i^{\top} \mathbf{H}_B^{\top} \mathbf{u} \end{aligned}$$

where we have used the definition of adjoint and the fact that  $S_B^* = (\mathbf{H}_B S)^* = S^* \mathbf{H}_B^{\top}$ . Summing the  $q$  terms, and observing that  $\mathbf{F} = \sum_{i=1}^q (1 - \frac{\lambda_{q+1}}{\lambda_i}) \mathbf{e}_i \mathbf{e}_i^{\top}$ , we have  $\mathcal{P}S_B^* \mathbf{u} = S_B^* \mathbf{u} - S^* \mathbf{F} \mathbf{F}^{\top} \mathbf{H}_B^{\top} \mathbf{u}$  for all  $\mathbf{u} \in \mathbb{R}^m$ . Since  $\mathbf{u}$  is arbitrary, this proves the claim in equation (71).

Next, observe that  $\mathcal{P}\tilde{\nabla}_f \mathbf{L}(g_t) = \mathcal{P}S_B^* (S_B S^* \boldsymbol{\beta}_t - Y[B]) = \mathcal{P}S_B^* \mathbf{v}_t$ , and rewriting the updates in equation (27a) using the relation  $f_t = S^* \boldsymbol{\alpha}_t$  and  $g_t = S^* \boldsymbol{\beta}_t$ , we get

$$\begin{aligned} S^* \boldsymbol{\alpha}_{t+1} &\leftarrow S^* \boldsymbol{\beta}_t - \eta_1 \mathcal{P}S_B^* \mathbf{v}_t \\ &= S^* \boldsymbol{\beta}_t - \eta_1 S_B^* \mathbf{v}_t + \eta_1 S^* \mathbf{F} \mathbf{F}^{\top} \mathbf{H}_B^{\top} \mathbf{v} \\ &= S^* (\boldsymbol{\beta}_t - \eta_1 \mathbf{H}_B^{\top} \mathbf{v}_t + \eta_1 \mathbf{w}_t) \end{aligned}$$

which is indeed equation (17c). A similar calculation can also show equation (17d) emulates equation (27b).  $\square$

*Proof of Proposition 4.* We start by showing that

$$\mathcal{Q}S_B^* = S_B^* - S_J^* \mathbf{G} \mathbf{G}^{\top} K(X[J], X[B]). \quad (73)$$

For a vector  $\mathbf{u} \in \mathbb{R}^m$ , observe that

$$\mathcal{Q}S_B^* \mathbf{u} = S_B^* \mathbf{u} - \sum_{i=1}^q \left(1 - \frac{\delta_{q+1}}{\delta_i}\right) \phi_i \otimes_{\mathbb{H}} \phi_i S_B^* \mathbf{u} \quad (74)$$

Now the term  $\phi_i \otimes_{\mathbb{H}} \phi_i S_B^* \mathbf{u}$  simplifies as

$$\begin{aligned} \phi_i \langle \phi_i, S_B^* \mathbf{u} \rangle_{\mathbb{H}} &= \frac{1}{\delta_i} S_J^* \mathbf{d}_i \langle S_J^* \mathbf{d}_i, S_B^* \mathbf{u} \rangle_{\mathbb{H}} \\ &= \frac{1}{\delta_i} S_J^* \mathbf{d}_i \langle \mathbf{d}_i, S_J S_J^* \mathbf{u} \rangle_{\mathbb{R}^m} \\ &= \frac{1}{\delta_i} S_J^* \mathbf{d}_i \mathbf{d}_i^{\top} K(X[J], X[B]) \mathbf{u}. \end{aligned}$$

Summing the  $q$  terms, and observing that  $\mathbf{G} = \sum_{i=1}^q \frac{1}{\delta_i} (1 - \frac{\delta_{q+1}}{\delta_i}) \mathbf{d}_i \mathbf{d}_i^{\top}$  we have  $\mathcal{Q}S_B^* \mathbf{u} = S_B^* \mathbf{u} - S_J^* \mathbf{G} \mathbf{G}^{\top} K(X[J], X[B]) \mathbf{u}$  for all  $\mathbf{u} \in \mathbb{R}^m$ . Since  $\mathbf{u}$  is arbitrary, this proves the claim in equation (73). The rest of the proof proceeds similar to the proof of Proposition 3.  $\square$

Table 6: Laplacian Kernel

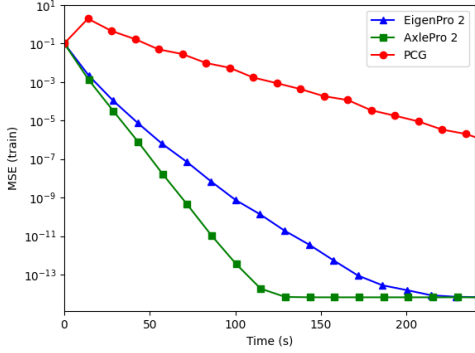
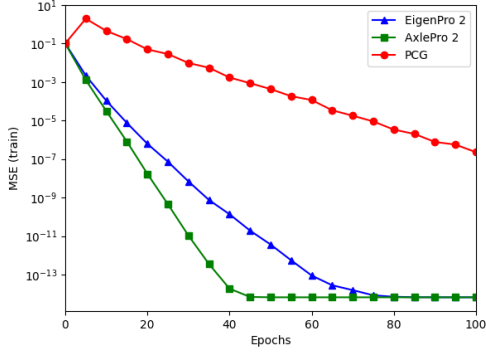
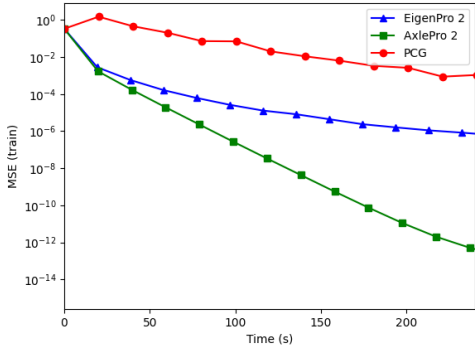
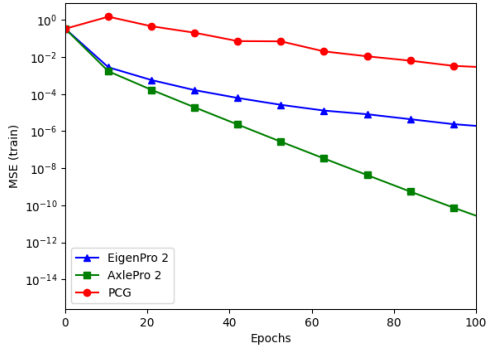
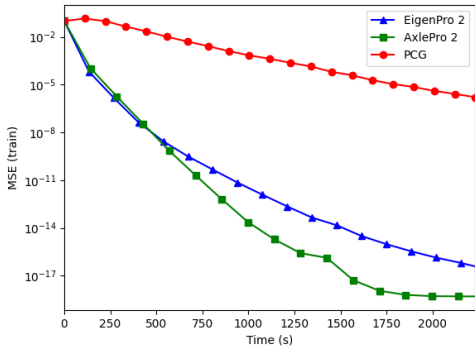
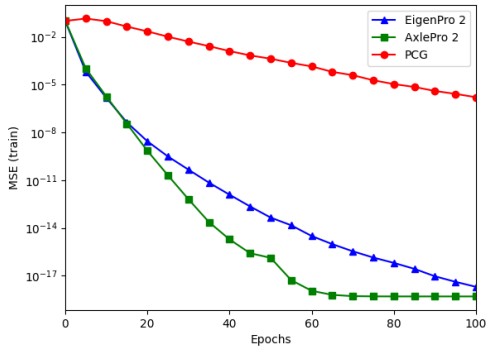
Dataset	MSE v/s Time	MSE v/s Epochs
CIFAR10		
Stellar Classification		
EMNIST Digits		

Table 7: Gaussian Kernel

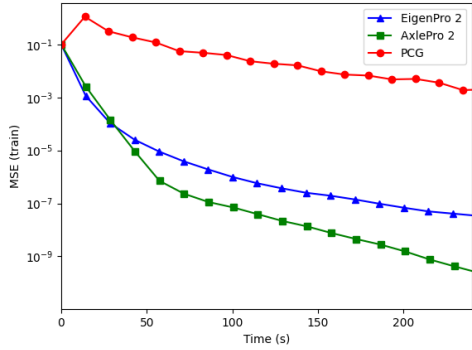
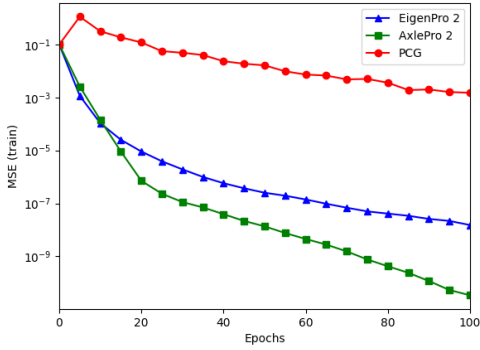
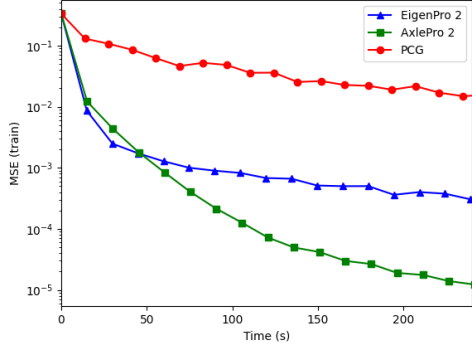
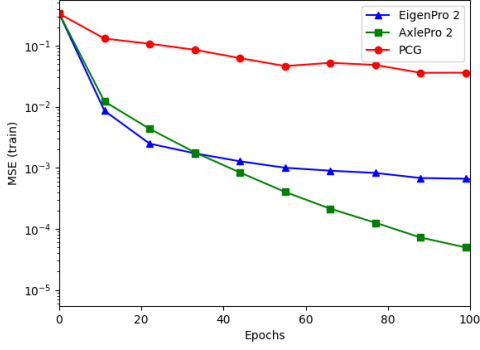
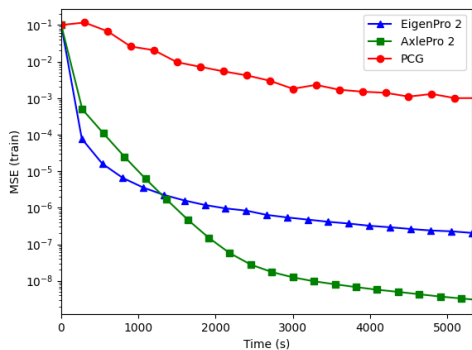
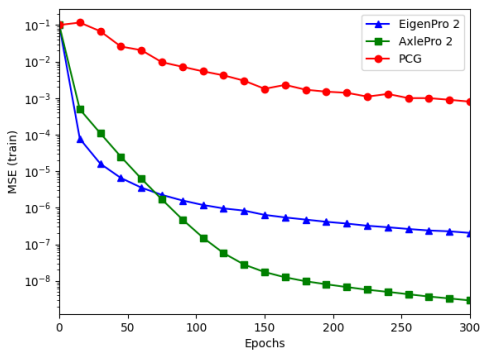
Dataset	MSE v/s Time	MSE v/s Epochs
CIFAR10		
Stellar Classification		
EMNIST Digits		

Table 8: Laplacian Kernel with single precision

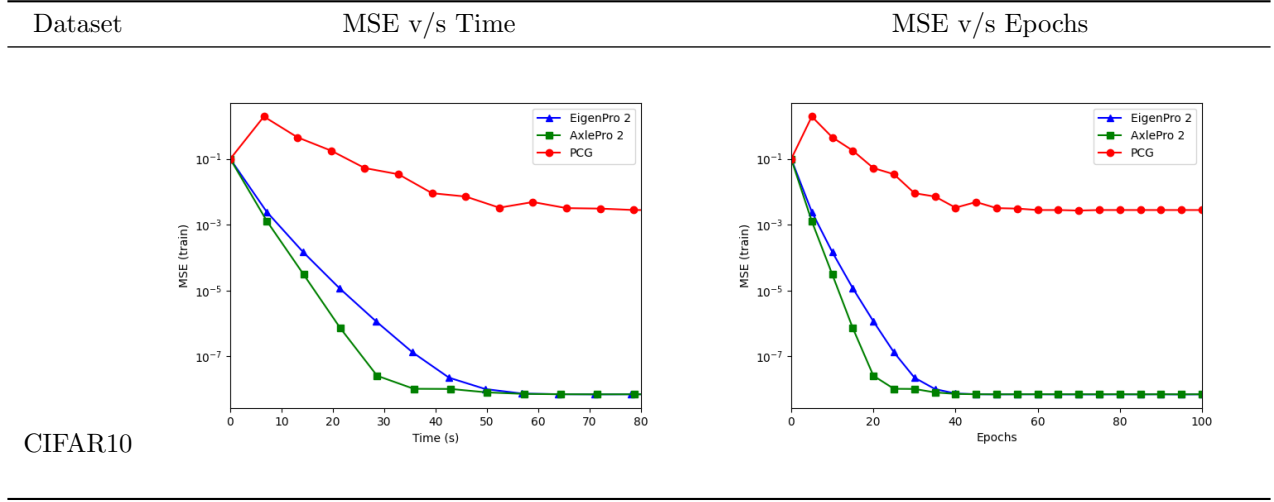


Table 9: Myrtle5 Kernel with stored kernel matrix

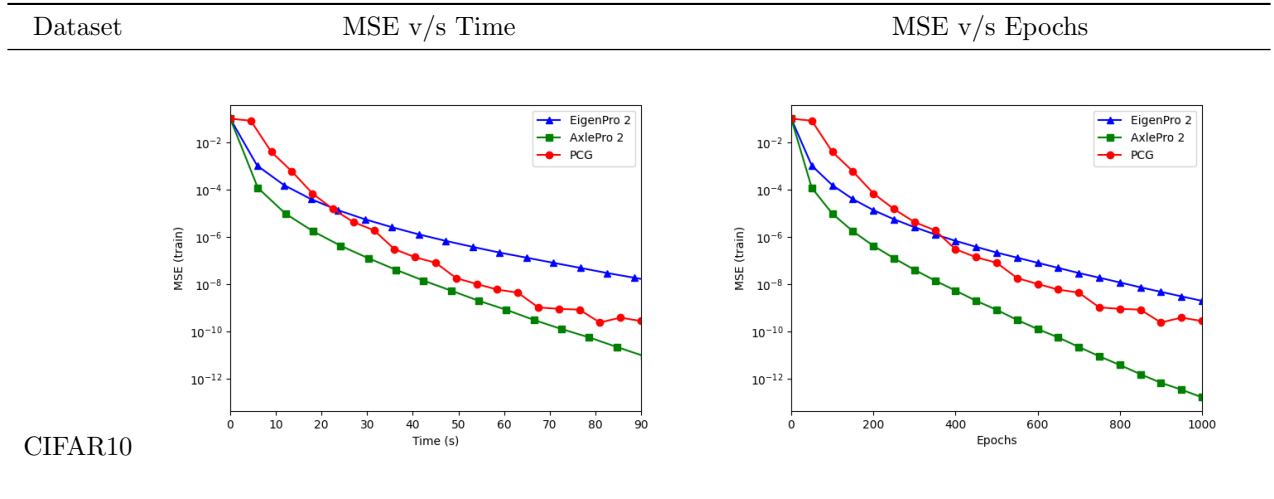


Table 10: Acceleration by Nyström approximation

Dataset	MSE v/s Time	MSE v/s Epochs
$n = 5k,$ $s = 800$		