

Supplementary Material for “Large-scale Online Kernel Learning with Random Feature Reparameterization”

Tu Dinh Nguyen[†], Trung Le[‡], Hung Bui[‡], Dinh Phung[†]

[‡]Adobe Research, Adobe Systems, Inc.

[†]Center for Pattern Recognition and Data Analytics, Deakin University, Australia
{tu.nguyen, trung.l, dinh.phung}@deakin.edu.au, hubui@adobe.com

This document presents supplementary material to complement the manuscript entitled “*Reparameterized Random Feature for Large-scale Online Kernel Learning*”, accepted at the 26th International Joint Conference on Artificial Intelligence (IJCAI), 2017. We first provide definitions of loss functions and their gradients w.r.t parameters. We then present full derivations of some equations in theoretical analysis, and add more details on gradient learning of our proposed model.

1 LOSS FUNCTIONS

We use 5 loss functions in our proposed model, namely Hinge, logistic, ℓ_2 , ℓ_1 and ε -insensitive. Let $\mathbb{I}\{s\}$ denote the indicator function that returns 1 if the logical statement s is true, and 0 otherwise. These losses and their gradients w.r.t parameters \mathbf{w} are listed in Table 1.

Loss	$\ell(y, \mathbf{w}^\top \phi(\mathbf{x}))$	$\nabla_{\mathbf{w}} \ell(y, \mathbf{w}^\top \phi(\mathbf{x}))$
Hinge	$\max(0, 1 - y\mathbf{w}^\top \phi(\mathbf{x}))$	$-\mathbb{I}\{y\mathbf{w}^\top \phi(\mathbf{x}) \leq 1\} y \phi(\mathbf{x})$
Logistic	$\log[1 + \exp(-y\mathbf{w}^\top \phi(\mathbf{x}))]$	$-y \phi(\mathbf{x}) \frac{\exp(-y\mathbf{w}^\top \phi(\mathbf{x}))}{\exp(-y\mathbf{w}^\top \phi(\mathbf{x})) + 1}$
ℓ_2	$\frac{1}{2} (y - \mathbf{w}^\top \phi(\mathbf{x}))^2$	$(\mathbf{w}^\top \phi(\mathbf{x}) - y) \phi(\mathbf{x})$
ℓ_1	$ y - \mathbf{w}^\top \phi(\mathbf{x}) $	$\text{sign}(\mathbf{w}^\top \phi(\mathbf{x}) - y) \phi(\mathbf{x})$
ε -insensitive	$\max\{0, y - \mathbf{w}^\top \phi(\mathbf{x}) - \varepsilon\}$	$\mathbb{I}\{ y - \mathbf{w}^\top \phi(\mathbf{x}) > \varepsilon\} \times \text{sign}(\mathbf{w}^\top \phi(\mathbf{x}) - y) \phi(\mathbf{x})$

Table 1: Typical loss functions and their gradients.

For *multiclass classification*, we follow the work of Crammer and Singer (2001) in which the objective function reads:

$$\mathcal{J}(\mathbf{w}) \triangleq \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \mathbb{E}_{\mathbb{P}} [\ell(\mathbf{w}_y^\top \phi(\mathbf{x}) - \mathbf{w}_z^\top \phi(\mathbf{x}))]$$

where $z = \text{argmax}_{k \neq y} \mathbf{w}_k^\top \phi(\mathbf{x})$ and the loss function can be either Hinge loss or logistic loss as listed in Table 2.

Loss	Hinge	Logistic
$\ell(y, \mathbf{w}^\top \phi(\mathbf{x}))$	$\max(0, 1 - \mathbf{w}_y^\top \phi(\mathbf{x}) + \mathbf{w}_z^\top \phi(\mathbf{x}))$	$\log[1 + \exp(\mathbf{w}_z^\top \phi(\mathbf{x}) - \mathbf{w}_y^\top \phi(\mathbf{x}))]$
$\nabla_{\mathbf{w}_y} \ell(y, \mathbf{w}^\top \phi(\mathbf{x}))$	$-\mathbb{I}\{\mathbf{w}_y^\top \phi(\mathbf{x}) - \mathbf{w}_z^\top \phi(\mathbf{x}) \leq 1\} \phi(\mathbf{x})$	$-\phi(\mathbf{x}) \frac{\exp(\mathbf{w}_z^\top \phi(\mathbf{x}) - \mathbf{w}_y^\top \phi(\mathbf{x}))}{\exp(\mathbf{w}_z^\top \phi(\mathbf{x}) - \mathbf{w}_y^\top \phi(\mathbf{x})) + 1}$
$\nabla_{\mathbf{w}_z} \ell(y, \mathbf{w}^\top \phi(\mathbf{x}))$	$\mathbb{I}\{\mathbf{w}_y^\top \phi(\mathbf{x}) - \mathbf{w}_z^\top \phi(\mathbf{x}) \leq 1\} \phi(\mathbf{x})$	$\phi(\mathbf{x}) \frac{\exp(\mathbf{w}_z^\top \phi(\mathbf{x}) - \mathbf{w}_y^\top \phi(\mathbf{x}))}{\exp(\mathbf{w}_z^\top \phi(\mathbf{x}) - \mathbf{w}_y^\top \phi(\mathbf{x})) + 1}$

Table 2: Hinge loss, logistic loss and their gradients for multiclass classification.

2 REPARAMETERIZED RANDOM FEATURES

In this section, we first roughly demonstrate the fast convergence of the upper bound of approximation error probability proven in Lemma 2 in the manuscript. Fig. 1 illustrates that the bound exponentially converges to zero when the number of random features passes a certain point for each precision ε .

We now present full derivations of some equations in Section 3.1 in the manuscript.

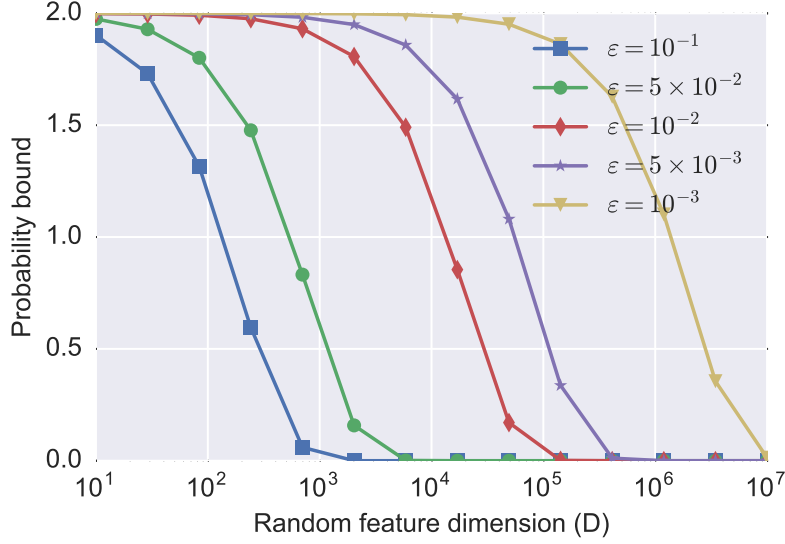


Figure 1: The convergence of the upper bound of approximation error probability.

- Eq (6) in the manuscript:

$$\begin{aligned}
\mathbb{E} [L_g^2] &= \mathbb{E} [\|\nabla g(\mathbf{u}^*)\|^2] = \mathbb{E} [\|\nabla \tilde{k}_\epsilon(\mathbf{u}^*) - \nabla k(\mathbf{u}^*)\|^2] \\
&= \mathbb{E} [\|\nabla \tilde{k}_\epsilon(\mathbf{u}^*)\|^2] - \|\nabla k(\mathbf{u}^*)\|^2 \leq \mathbb{E} [\|\nabla \tilde{k}_\epsilon(\mathbf{u}^*)\|^2] \\
&= \mathbb{E} \left[\left\| \frac{1}{D} \sum_{d=1}^D \sin \left((\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}_d)^\top \mathbf{u}^* \right) (\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}_d)^\top \right\|^2 \right] \\
&\leq \frac{1}{D^2} \mathbb{E} \left[\left(\sum_{d=1}^D \left\| \sin \left((\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}_d)^\top \mathbf{u}^* \right) (\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}_d)^\top \right\| \right)^2 \right] \\
&\leq \frac{1}{D^2} \mathbb{E} \left[\left(\sum_{d=1}^D \|\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}_d\| \right)^2 \right] \tag{1} \\
&\leq \frac{1}{D^2} \mathbb{E} \left[D \sum_{d=1}^D \|\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}_d\|^2 \right] \tag{2} \\
&= \frac{1}{D} \sum_{d=1}^D \mathbb{E} [\|\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}_d\|^2] \\
&= \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}\|^2] = \|\boldsymbol{\sigma}\|^2 \tag{3}
\end{aligned}$$

Here we note that to go from Eq. (1) to Eq. (2), we have used Schwartz-Cauchy inequality $\left(\sum_{d=1}^D a_d \right)^2 \leq D \sum_{d=1}^D a_d^2$ and to obtain Eq. (3) we derive as:

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\text{diag}(\boldsymbol{\sigma}) \boldsymbol{\epsilon}\|^2] &= \sum_{n=1}^N \sigma_n^2 \mathbb{E}_{\epsilon_n \sim \mathcal{N}(0,1)} [\epsilon_n^2] \\
&= \sum_{n=1}^N \sigma_n^2 = \|\boldsymbol{\sigma}\|^2
\end{aligned}$$

where ϵ_n^2 follows Chi square distribution with 1 degree of freedom, and hence $\mathbb{E} [\epsilon_n^2] = 1$.

- Eq (7) in the manuscript:

$$\begin{aligned}
& 2\sqrt{\tau_1 \tau_2 \left(\frac{\tau_1}{\tau_2}\right)^{\frac{-N+2}{N+2}}} = 2\tau_1^{\frac{2}{N+2}} \tau_2^{\frac{N}{N+2}} \\
& = 2 \times 4^{\frac{2N+1}{N+2}} \left(\frac{\text{diam}(\mathcal{X}) \|\boldsymbol{\sigma}\|^2}{\varepsilon^2}\right)^{\frac{N}{N+2}} \exp\left(\frac{-D\varepsilon^2}{N+2}\right) \\
& \leq 2^5 \exp\left(-\frac{D\varepsilon^2}{N+2} + \frac{N}{N+2} \log\left(\frac{\text{diam}(\mathcal{X}) \|\boldsymbol{\sigma}\|^2}{\varepsilon^2}\right)\right) \\
& \leq 2^5 \exp\left(-\frac{D\varepsilon^2}{N+2} + \log\left(\frac{\text{diam}(\mathcal{X}) \|\boldsymbol{\sigma}\|^2}{\varepsilon^2}\right)\right) \\
& = 2^5 \left(\frac{\text{diam}(\mathcal{X}) \|\boldsymbol{\sigma}\|^2}{\varepsilon^2}\right) \exp\left(-\frac{D\varepsilon^2}{N+2}\right)
\end{aligned}$$

Next we provide details on derivations of gradients to learn the kernel parameters in Section 3.2 in the manuscript. Recall that the gradient of loss function w.r.t to kernel width parameters $\boldsymbol{\sigma}$ in Eq. (10) in the manuscript is computed by:

$$\begin{aligned}
\nabla_{\boldsymbol{\sigma}} \ell(y, \mathbf{w}^\top \tilde{\boldsymbol{\phi}}(\mathbf{x})) &= \sum_{d=1}^D \left(\frac{\partial [\hat{\mathbf{z}}_{\boldsymbol{\sigma}}]_d(\mathbf{x})}{\partial \boldsymbol{\omega}_d} \right)^\top \frac{\partial \ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))}{\partial [\hat{\mathbf{z}}_{\boldsymbol{\sigma}}]_d(\mathbf{x})} \odot \frac{\partial \boldsymbol{\omega}_d}{\partial \boldsymbol{\sigma}} \\
&= \sum_{d=1}^D [-\mathbf{x} \sin(\boldsymbol{\omega}_d^\top \mathbf{x}), \mathbf{x} \cos(\boldsymbol{\omega}_d^\top \mathbf{x})]^\top \\
&\quad \times \left(\frac{\partial \ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))}{\partial [\hat{\mathbf{z}}_{\boldsymbol{\sigma}}]_d(\mathbf{x})} \right)^\top \odot \boldsymbol{\epsilon}_d
\end{aligned}$$

Similarity to those w.r.t \mathbf{w} , the derivative of loss function ℓ w.r.t $[\hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})]_d$ can be computed as shown in Table 3.

Loss	$\ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))$	$\nabla_{[\hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})]_d} \ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))$
Hinge	$\max(0, 1 - y \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))$	$-\mathbb{I}\{y \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}) \leq 1\} y \mathbf{w}$
Logistic	$\log[1 + \exp(-y \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))]$	$-y \mathbf{w} \frac{\exp(-y \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))}{\exp(-y \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})) + 1}$
ℓ_2	$\frac{1}{2} (y - \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}))^2$	$(\mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}) - y) \mathbf{w}$
ℓ_1	$ y - \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}) $	$\text{sign}(\mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}) - y) \mathbf{w}$
ε -insensitive	$\max\{0, y - \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}) - \varepsilon\}$	$\mathbb{I}\{ y - \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}) > \varepsilon\} \times \text{sign}(\mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x}) - y) \mathbf{w}$

Table 3: Typical loss functions and their gradients w.r.t $[\hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})]_d$.

As the variance $\boldsymbol{\sigma}$ is constrained to be positive, we can learn the log-variance $\gamma = \log(\boldsymbol{\sigma})$ instead, hence naturally induce a positive variance. The gradient of loss function w.r.t to γ can be computed as:

$$\begin{aligned}
\nabla_{\gamma} \ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})) &= \nabla_{\boldsymbol{\sigma}} \ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})) \odot \nabla_{\gamma} \boldsymbol{\sigma} \\
&= \nabla_{\boldsymbol{\sigma}} \ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})) \odot \nabla_{\gamma} e^{\gamma} \\
&= \nabla_{\boldsymbol{\sigma}} \ell(y, \mathbf{w}^\top \hat{\mathbf{z}}_{\boldsymbol{\sigma}}(\mathbf{x})) \odot e^{\gamma}
\end{aligned}$$

References

Koby Crammer and Yoram Singer. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research (JMLR)*, 2(Dec):265–292, 2001.