# A novel greedy algorithm for Nyström approximation
## Supplementary material

**Ahmed K. Farahat**      **Ali Ghodsi**      **Mohamed S. Kamel**
University of Waterloo, Waterloo, Ontario, Canada, N2L 3G1

# 1 UPDATE FORMULAS FOR $\boldsymbol{f}$ AND $\boldsymbol{g}$

In this section, the recursive formulas for $\boldsymbol{f}$ and $\boldsymbol{g}$ (Equation (18) in the main paper) are derived. The greedy selection criterion at iteration $t$ is:

$$q = arg\ \max_i \|\frac{1}{\sqrt{E_{ii}}}E_{:i}\|^2, \tag{1}$$

where $E$ is the residual matrix at iteration $t$, $E_{:i}$ denotes the $i$-th column of $E$, and $E_{ii}$ denotes the $i$-th diagonal element of $E$.

As the Nyström approximation is calculated in a recursive manner based on the residual matrix at the previous iteration, $E$, $E_{:i}$, and $E_{ii}$ for a candidate column $i$ can be recursively calculated as follows:

$$E^{(t)} = (E - \frac{1}{\alpha}\boldsymbol{\delta\delta}^T)^{(t-1)} = (E - \boldsymbol{\omega\omega}^T)^{(t-1)},$$

$$E_{:i}^{(t)} = (E_{:i} - \frac{\boldsymbol{\delta}_i}{\alpha}\boldsymbol{\delta})^{(t-1)} = (E_{:i} - \boldsymbol{\omega}_i\boldsymbol{\omega})^{(t-1)}, \tag{2}$$

$$E_{ii}^{(t)} = (E_{ii} - \frac{\boldsymbol{\delta}_i^2}{\alpha})^{(t-1)} = (E_{ii} - \boldsymbol{\omega}_i^2)^{(t-1)}.$$

Let $\boldsymbol{f}_i = \|E_{:i}\|^2$ and $\boldsymbol{g}_i = E_{ii}$ be the numerator and denominator of the criterion function for data point $i$ respectively. Based on (2), $\boldsymbol{f}_i^{(t)}$ can be calculated as:

$$\begin{aligned}
\boldsymbol{f}_i^{(t)} &= \left(\|E_{:i} - \boldsymbol{\omega}_i\boldsymbol{\omega}\|^2\right)^{(t-1)} \\
&= \left((E_{:i} - \boldsymbol{\omega}_i\boldsymbol{\omega})^T(E_{:i} - \boldsymbol{\omega}_i\boldsymbol{\omega})\right)^{(t-1)} \\
&= \left(E_{:i}^T E_{:i} - 2\boldsymbol{\omega}_i E_{:i}^T\boldsymbol{\omega} + \boldsymbol{\omega}_i^2\|\boldsymbol{\omega}\|^2\right)^{(t-1)} \\
&= \left(\boldsymbol{f}_i - 2\boldsymbol{\omega}_i E_{:i}^T\boldsymbol{\omega} + \boldsymbol{\omega}_i^2\|\boldsymbol{\omega}\|^2\right)^{(t-1)}.
\end{aligned} \tag{3}$$

Similarly, $\boldsymbol{g}_i^{(t)}$ can be calculated as:

$$\begin{aligned}
\boldsymbol{g}_i^{(t)} = E_{ii}^{(t)} &= \left(E_{ii} - \boldsymbol{\omega}_i^2\right)^{(t-1)} \\
&= \left(\boldsymbol{g}_i - \boldsymbol{\omega}_i^2\right)^{(t-1)}.
\end{aligned} \tag{4}$$

Let $\boldsymbol{f} = [\boldsymbol{f}_i]_{i=1..n}$ and $\boldsymbol{g} = [\boldsymbol{g}_i]_{i=1..n}$, $\boldsymbol{f}^{(t)}$ and $\boldsymbol{g}^{(t)}$ can be expressed as:

$$\begin{aligned}
\boldsymbol{f}^{(t)} &= \left(\boldsymbol{f} - 2\left(\boldsymbol{\omega} \circ E\boldsymbol{\omega}\right) + \|\boldsymbol{\omega}\|^2\left(\boldsymbol{\omega} \circ \boldsymbol{\omega}\right)\right)^{(t-1)}, \\
\boldsymbol{g}^{(t)} &= \left(\boldsymbol{g} - \left(\boldsymbol{\omega} \circ \boldsymbol{\omega}\right)\right)^{(t-1)},
\end{aligned} \tag{5}$$

where $\circ$ represents the Hadamard product operator, and $\|.\|$ is the $\ell_2$ norm.

Based on the recursive formula of $E$, the term $E\boldsymbol{\omega}$ at iteration $(t-1)$ can be expressed as:

$$\begin{aligned}
E\boldsymbol{\omega} &= \left(K - \Sigma_{r=1}^{t-2}\left(\boldsymbol{\omega\omega}^T\right)^{(r)}\right)\boldsymbol{\omega} \\
&= K\boldsymbol{\omega} - \Sigma_{r=1}^{t-2}\left(\boldsymbol{\omega}^{(r)^T}\boldsymbol{\omega}\right)\boldsymbol{\omega}^{(r)}
\end{aligned} \tag{6}$$

Substitute with $E\boldsymbol{\omega}$ in Equation (5), the update formulas for $\boldsymbol{f}$ and $\boldsymbol{g}$ are given as:

$$\begin{aligned}
\boldsymbol{f}^{(t)} &= \Big(\boldsymbol{f} - 2\Big(\boldsymbol{\omega} \circ \Big(K\boldsymbol{\omega} - \Sigma_{r=1}^{t-2}\left(\boldsymbol{\omega}^{(r)^T}\boldsymbol{\omega}\right)\boldsymbol{\omega}^{(r)}\Big)\Big) \\
&\quad + \|\boldsymbol{\omega}\|^2\left(\boldsymbol{\omega} \circ \boldsymbol{\omega}\right)\Big)^{(t-1)}, \\
\boldsymbol{g}^{(t)} &= \left(\boldsymbol{g} - \left(\boldsymbol{\omega} \circ \boldsymbol{\omega}\right)\right)^{(t-1)}.
\end{aligned} \tag{7}$$

In the case of partition-based greedy Nyström algorithm, the update formulas for $\boldsymbol{f}$ and $\boldsymbol{g}$ (Equation (19) in the main paper) can be derived as follows.

Let $E^{(t)}$ and $H^{(t)}$ be the residual matrices of $K$ and $G$ at iteration $t$ respectively. The efficient sampling criterion based on centroids can be expressed as follows:

$$q = arg\ \max_i \|\frac{1}{\sqrt{E_{ii}}}H_{:i}\|^2, \tag{8}$$

where $H_{:i}$ denotes the $i$-th column of $H$, and $E_{ii}$ denotes the $i$-th diagonal element of $E$. The term $H_{ji}/\sqrt{E_{ii}}$ is the scalar projection of the $j$-th centroid onto $X_{:i}$. Let $\boldsymbol{\delta}^{(t)}$ be column of $E$ selected at iteration $t$, $\alpha^{(t)}$ be the corresponding diagonal element of $E$, and $\boldsymbol{\gamma}^{(t)}$ be the corresponding column of $H$. Define $\boldsymbol{\omega}^{(t)} = \boldsymbol{\delta}^{(t)}/\sqrt{\alpha^{(t)}}$, and $\boldsymbol{v}^{(t)} = \boldsymbol{\gamma}^{(t)}/\sqrt{\alpha^{(t)}}$. The rank-1 approximation of $H^{(t)}$ can be calculated as:

$$\tilde{H}_{\{q\}}^{(t)} = \frac{1}{\alpha^{(t)}}\boldsymbol{\gamma}^{(t)}\boldsymbol{\delta}^{T(t)} = \boldsymbol{v}^{(t)}\boldsymbol{\omega}^{T(t)}, \tag{9}$$

and the new residual matrix $H$ can be calculated as:

$$H^{(t+1)} = H^{(t)} - \boldsymbol{v}^{(t)}\boldsymbol{\omega}^{T(t)} \tag{10}$$

Based on this recursive formula, the greedy sampling criterion (Equation 8) can be calculated in a recursive manner as follows. Similar to Equation (2), $H$, $H_{:i}$, and $E_{ii}$ can be recursively calculated as:

$$H^{(t)} = (H - \boldsymbol{v}\boldsymbol{\omega}^T)^{(t-1)},$$
$$H_{:i}^{(t)} = (H_{:i} - \boldsymbol{\omega}_i\boldsymbol{v})^{(t-1)}, \quad (11)$$
$$E_{ii}^{(t)} = (E_{ii} - \boldsymbol{\omega}_i^2)^{(t-1)}.$$

Let $\boldsymbol{f}_i = \|H_{:i}\|^2$ and $\boldsymbol{g}_i = E_{ii}$ be the numerator and denominator of the criterion function for data point $i$ respectively. $\boldsymbol{f}_i^{(t)}$ and $\boldsymbol{g}_i^{(t)}$ can be calculated as follows:

$$\boldsymbol{f}_i^{(t)} = \left(\|H_{:i} - \boldsymbol{\omega}_i\boldsymbol{v}\|^2\right)^{(t-1)}$$
$$= \left(\boldsymbol{f}_i - 2\boldsymbol{\omega}_i H_{:i}^T\boldsymbol{v} + \boldsymbol{\omega}_i^2\|\boldsymbol{v}\|^2\right)^{(t-1)},$$
$$\boldsymbol{g}_i^{(t)} = E_{ii}^{(t)} = \left(E_{ii} - \boldsymbol{\omega}_i^2\right)^{(t-1)} \quad (12)$$
$$= \left(\boldsymbol{g}_i - \boldsymbol{\omega}_i^2\right)^{(t-1)}.$$

Let $\boldsymbol{f} = [\boldsymbol{f}_i]_{i=1..n}$ and $\boldsymbol{g} = [\boldsymbol{g}_i]_{i=1..n}$, $\boldsymbol{f}^{(t)}$ and $\boldsymbol{g}^{(t)}$ can be expressed as:

$$\boldsymbol{f}^{(t)} = \left(\boldsymbol{f} - 2\left(\boldsymbol{\omega}\circ H^T\boldsymbol{v}\right) + \|\boldsymbol{v}\|^2\left(\boldsymbol{\omega}\circ\boldsymbol{\omega}\right)\right)^{(t-1)},$$
$$\boldsymbol{g}^{(t)} = \left(\boldsymbol{g} - \left(\boldsymbol{\omega}\circ\boldsymbol{\omega}\right)\right)^{(t-1)}, \quad (13)$$

where $\circ$ represents the Hadamard product operator, and $\|.\|$ is the $\ell_2$ norm.

The term $H^T\boldsymbol{v}$ at iteration $(t-1)$ can be calculated recursively as:

$$H^T\boldsymbol{v} = \left(G^T - \Sigma_{r=1}^{t-2}\left(\boldsymbol{\omega}\boldsymbol{v}^T\right)^{(r)}\right)\boldsymbol{v}$$
$$= G^T\boldsymbol{v} - \Sigma_{r=1}^{t-2}\left(\boldsymbol{v}^{(r)T}\boldsymbol{v}\right)\boldsymbol{\omega}^{(r)} \quad (14)$$

Substitute with $H^T\boldsymbol{v}$ in Equation (13), the update formulas for $\boldsymbol{f}$ and $\boldsymbol{g}$ are given as:

$$\boldsymbol{f}^{(t)} = \left(\boldsymbol{f} - 2\left(\boldsymbol{\omega}\circ\left(G^T\boldsymbol{v} - \Sigma_{r=1}^{t-2}\left(\boldsymbol{v}^{(r)T}\boldsymbol{v}\right)\boldsymbol{\omega}^{(r)}\right)\right)\right.$$
$$\left. + \|\boldsymbol{v}\|^2\left(\boldsymbol{\omega}\circ\boldsymbol{\omega}\right)\right)^{(t-1)},$$
$$\boldsymbol{g}^{(t)} = \left(\boldsymbol{g} - \left(\boldsymbol{\omega}\circ\boldsymbol{\omega}\right)\right)^{(t-1)}. \quad (15)$$

## 2 COMPARISON WITH ENSEMBLE NYSTRÖM

In this section, the proposed greedy Nyström methods (**GreedyNyström** and **PartGreedyNys**) are compared to the ensemble Nytröm algorithm (**EnsembleNyström**) proposed by Kumar et al. (2009). The ensemble Nyström method constructs a low-rank approximation of a kernel matrix using an ensemble of $p$ Nyström approximations. As suggested by Kumar et al. (2009), the ridge regression algorithm can be used to learn the mixture weights of different approximations using a validation set of columns sampled from the original kernel matrix. In this experiment, an ensemble of $p = 10$ Nyström approximations is used, and $l$ columns are sampled to calculate each low-rank approximation. A validation set of $s = 20$ columns is used for estimating the mixture weights of the ensemble, and a hold-out set of $s' = 20$ columns is used to estimate the ridge parameter.

Tables 1 and 2 show the relative accuracies and run time of different methods. Two values are used for $l$: $l = 3\%n$ and $l = 5\%n$, with $k = 1\%n$. It can be observed that both **PartGreedyNys** and **GreedyNyström** outperforms the ensemble method (**EnsembleNyström**) in term of approximation accuracy for most data sets.

## References

S. Kumar, M. Mohri, and A. Talwalkar. Ensemble Nyström Method. In *Advances in Neural Information Processing Systems 22*, pages 1060–1068, 2009.