
Greedy Nyström Approximation Supplementary Material

Ahmed K. Farahat **Ali Ghodsi** **Mohamed S. Kamel**
 University of Waterloo
 Waterloo, Ontario, Canada N2L 3G1
 {afarahat, aghodsib, mkamel}@uwaterloo.ca

1 Projection for rank-1 Nyström approximation

This section provides a proof for the claim that rank-1 Nyström approximation based on the i -th column implicitly projects all data points onto a line which contains data point i and the origin in the high-dimensional feature space defined by the kernel.

The rank-1 Nyström approximation of matrix K based on column i is calculated as:

$$\tilde{K}_{\{i\}} = \frac{1}{K_{ii}} K_{:i} K_{:i}^T, \quad (1)$$

where $K_{:i}$ denotes the i -th column of K , and K_{ii} denotes the i -th diagonal element of K .

Assume that a linear kernel is used. The kernel matrix is calculated as $K = X^T X$, where X is an $m \times n$ data matrix, and m is the number of features. Let $X_{:i}$ be the i -th column of X . $X_{:i}$ also represents a vector that connects data point i and the origin. $K_{:i}$ and K_{ii} can be written as: $K_{:i} = X^T X_{:i}$, and $K_{ii} = \|X_{:i}\|^2$, where $\|\cdot\|$ is the ℓ_2 norm. Based on this, $\tilde{K}_{\{i\}}$ can be expressed as:

$$\tilde{K}_{\{i\}} = X^T \frac{X_{:i}}{\|X_{:i}\|} \frac{X_{:i}^T}{\|X_{:i}\|} X, \quad (2)$$

where $X^T X_{:i} / \|X_{:i}\|$ is a column vector whose j -th element is $X_{:j}^T X_{:i} / \|X_{:i}\|$. This value is the scalar projection of data point j onto $X_{:i}$ (as illustrated in Figure 1). This means that $\tilde{K}_{\{i\}}$ implicitly projects all data points into a 1-dimensional subspace (i.e., a line) which contains data point i and the origin, and then calculates the inner-products between the projected points.

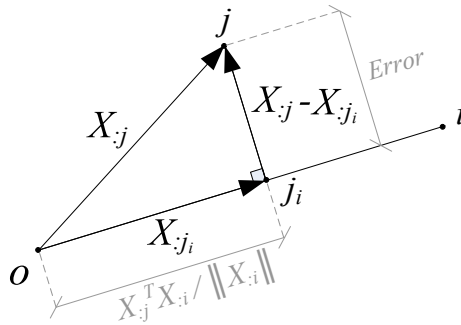


Figure 1: Projection for rank-1 Nyström approximation.

The same proof applies to other kernel types, as any kernel matrix implicitly maps data vectors to a high-dimensional linear space. In this case, the columns of X represent data vectors in the high-dimensional feature space defined by the kernel.

2 Updates formulas for f and g

In this section, the recursive formulas for f and g are derived. The greedy selection criterion at iteration t is:

$$l = \arg \max_i \left\| \frac{1}{\sqrt{E_{ii}}} E_{:i} \right\|^2, \quad (3)$$

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:

$$\begin{aligned} E^{(t)} &= (E - \frac{1}{\alpha} \delta \delta^T)^{(t-1)} = (E - \omega \omega^T)^{(t-1)}, \\ E_{:i}^{(t)} &= (E_{:i} - \frac{\delta_i}{\alpha} \delta)^{(t-1)} = (E_{:i} - \omega_i \omega)^{(t-1)}, \\ E_{ii}^{(t)} &= (E_{ii} - \frac{\delta_i^2}{\alpha})^{(t-1)} = (E_{ii} - \omega_i^2)^{(t-1)}. \end{aligned} \quad (4)$$

Let $f_i = \|E_{:i}\|^2$ and $g_i = E_{ii}$ be the numerator and denominator of the criterion function for data point i respectively. Based on (4), $f_i^{(t)}$ can be calculated as:

$$\begin{aligned} f_i^{(t)} &= (\|E_{:i} - \omega_i \omega\|^2)^{(t-1)} = ((E_{:i} - \omega_i \omega)^T (E_{:i} - \omega_i \omega))^{(t-1)} \\ &= (E_{:i}^T E_{:i} - 2\omega_i E_{:i}^T \omega + \omega_i^2 \|\omega\|^2)^{(t-1)} = (f_i - 2\omega_i E_{:i}^T \omega + \omega_i^2 \|\omega\|^2)^{(t-1)}. \end{aligned} \quad (5)$$

Similarly, $g_i^{(t)}$ can be calculated as:

$$g_i^{(t)} = E_{ii}^{(t)} = (E_{ii} - \omega_i^2)^{(t-1)} = (g_i - \omega_i^2)^{(t-1)}. \quad (6)$$

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

$$\begin{aligned} f^{(t)} &= (f - 2(\omega \circ E\omega) + \|\omega\|^2 (\omega \circ \omega))^{(t-1)}, \\ g^{(t)} &= (g - (\omega \circ \omega))^{(t-1)}, \end{aligned} \quad (7)$$

where \circ represents the Hadamard product operator, and $\|\cdot\|$ is the ℓ_2 norm.

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

$$E\omega = \left(K - \sum_{r=1}^{t-2} (\omega \omega^T)^{(r)} \right) \omega = K\omega - \sum_{r=1}^{t-2} (\omega^{(r)T} \omega) \omega^{(r)} \quad (8)$$

Substitute with $E\omega$ in Equation (7), the update formulas for f and g are given as:

$$\begin{aligned} f^{(t)} &= \left(f - 2 \left(\omega \circ \left(K\omega - \sum_{r=1}^{t-2} (\omega^{(r)T} \omega) \omega^{(r)} \right) \right) + \|\omega\|^2 (\omega \circ \omega) \right)^{(t-1)}, \\ g^{(t)} &= (g - (\omega \circ \omega))^{(t-1)}. \end{aligned} \quad (9)$$

3 Partition-based greedy Nyström algorithm

This section describes the partition-based algorithm in details and derives recursive formulas for f and g . The partition-based algorithm first divides data points into $c \ll n$ random groups, and then selects the column of K which best represents the centroids of these groups in the high-dimensional feature space.

Let \mathcal{P}_j be the set of data points that belong to the j -th partition, $P = \{\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_c\}$ be a random partitioning of data points into c groups, and G be an $c \times n$ matrix whose element G_{ji} is the inner-product of the centroid of the j -th group and the i -th data point, weighted with the size of the j -th group. The use of weighted inner-products avoids any bias towards larger groups when calculating the sum of scalar projections. As the scalar projections are implicitly calculated in a high-dimensional linear space defined by the kernel matrix K , G_{ji} can be calculated in this linear space as: $G_{ji} = \sum_{r \in \mathcal{P}_j} K_{ir}$. In general, it requires $\mathcal{O}(n^2)$ to calculate G given K . However, G needs to be calculated only once during the calculation of K . In the case of a linear kernel, this complexity could be significantly reduced by calculating the centroids of each group in the feature space, and then the inner-products between each centroid and all data points. This computational complexity could be reduced further if the data matrix is very sparse. In addition, there is no need to calculate and store the whole kernel matrix in order to calculate G .

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:

$$l = \arg \max_i \left\| \frac{1}{\sqrt{E_{ii}}} H_{:i} \right\|^2, \quad (10)$$

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 $\delta^{(t)}$ be column of E selected at iteration t , $\alpha^{(t)}$ be the corresponding diagonal element of E , and $\gamma^{(t)}$ be the corresponding column of H . Define $\omega^{(t)} = \delta^{(t)}/\sqrt{\alpha^{(t)}}$, and $\mathbf{v}^{(t)} = \gamma^{(t)}/\sqrt{\alpha^{(t)}}$. The rank-1 approximation of $H^{(t)}$ can be calculated as:

$$\tilde{H}_{\{l\}}^{(t)} = \frac{1}{\alpha^{(t)}} \gamma^{(t)} \delta^{T(t)} = \mathbf{v}^{(t)} \omega^{T(t)}, \quad (11)$$

and the new residual matrix H can be calculated as:

$$H^{(t+1)} = H^{(t)} - \mathbf{v}^{(t)} \omega^{T(t)} \quad (12)$$

Based on this recursive formula, the greedy sampling criterion (Equation (10)) can be calculated in a recursive manner as follows. Similar to Equation (4), H , $H_{:i}$, and E_{ii} can be recursively calculated as:

$$\begin{aligned} H^{(t)} &= (H - \mathbf{v} \omega^T)^{(t-1)}, \\ H_{:i}^{(t)} &= (H_{:i} - \omega_i \mathbf{v})^{(t-1)}, \\ E_{ii}^{(t)} &= (E_{ii} - \omega_i^2)^{(t-1)}. \end{aligned} \quad (13)$$

Similar to Section 2, let $\mathbf{f}_i = \|H_{:i}\|^2$ and $\mathbf{g}_i = E_{ii}$ be the numerator and denominator of the criterion function for data point i . $\mathbf{f}_i^{(t)}$ and $\mathbf{g}_i^{(t)}$ can be calculated as follows:

$$\begin{aligned} \mathbf{f}_i^{(t)} &= (\|H_{:i} - \omega_i \mathbf{v}\|^2)^{(t-1)} = (\mathbf{f}_i - 2\omega_i H_{:i}^T \mathbf{v} + \omega_i^2 \|\mathbf{v}\|^2)^{(t-1)}, \\ \mathbf{g}_i^{(t)} &= E_{ii}^{(t)} = (E_{ii} - \omega_i^2)^{(t-1)} = (\mathbf{g}_i - \omega_i^2)^{(t-1)}. \end{aligned} \quad (14)$$

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

$$\begin{aligned} \mathbf{f}^{(t)} &= (\mathbf{f} - 2(\omega \circ H^T \mathbf{v}) + \|\mathbf{v}\|^2 (\omega \circ \omega))^{(t-1)}, \\ \mathbf{g}^{(t)} &= (\mathbf{g} - (\omega \circ \omega))^{(t-1)}, \end{aligned} \quad (15)$$

where \circ represents the Hadamard product operator, and $\|\cdot\|$ is the ℓ_2 norm.

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

$$H^T \mathbf{v} = \left(G^T - \sum_{r=1}^{t-2} (\omega \mathbf{v}^T)^{(r)} \right) \mathbf{v} = G^T \mathbf{v} - \sum_{r=1}^{t-2} (\mathbf{v}^{(r)T} \mathbf{v}) \omega^{(r)}$$

Substitute with $H^T \mathbf{v}$ in Equation (15), the update formulas for \mathbf{f} and \mathbf{g} are given as:

$$\begin{aligned} \mathbf{f}^{(t)} &= \left(\mathbf{f} - 2 \left(\omega \circ \left(G^T \mathbf{v} - \sum_{r=1}^{t-2} (\mathbf{v}^{(r)T} \mathbf{v}) \omega^{(r)} \right) \right) + \|\mathbf{v}\|^2 (\omega \circ \omega) \right)^{(t-1)}, \\ \mathbf{g}^{(t)} &= (\mathbf{g} - (\omega \circ \omega))^{(t-1)}. \end{aligned} \quad (16)$$

Table 1: Properties of data sets used to evaluate different Nyström methods. n and m are the number of instances and features respectively.

Data set	Type	n	m
Reuters-21578	Documents	5946	18933
Reviews	Documents	4069	36746
LAI	Documents	3204	29714
MNIST-4K	Digit Images	4000	784
PIE-20	Face Images	3400	1024
Yale-B-38	Face Images	2414	1024

4 The complete partition-based greedy Nyström algorithm

Inputs: $K, k, d,$

Outputs: $\mathcal{S}, \tilde{K}_{\mathcal{S}}, \tilde{K}_{\mathcal{S},d}, Y$

Steps:

1. Initialize $\mathcal{S} = \{ \}$
2. Generate a random partitioning P , Calculate G : $G_{ji} = \sum_{r \in \mathcal{P}_j} K_{ir}$
3. Initialize $\mathbf{f}_i^{(0)} = \|G_{:i}\|^2, \mathbf{g}_i^{(0)} = K_{ii}, \mathbf{f}^{(0)} = [\mathbf{f}_i^{(0)}]_{i=1..n}$ and $\mathbf{g}^{(0)} = [\mathbf{g}_i^{(0)}]_{i=1..n}$
4. Repeat $t = 1 \rightarrow k$:
 - (a) $l = \arg \max_i \mathbf{f}_i^{(t)} / \mathbf{g}_i^{(t)}$
 - (b) $\mathcal{S} = \mathcal{S} \cup \{l\}$
 - (c) Update $\delta^{(t)} = K_{:l} - \sum_{r=1}^{t-1} \omega_l^{(r)} \omega^{(r)}, \gamma^{(t)} = G_{:l} - \sum_{r=1}^{t-1} \omega_l^{(r)} \mathbf{v}^{(r)}, \alpha^{(t)} = \delta_l^{(r)}$
 - (d) Calculate $\omega^{(t)} = \delta^{(t)} / \sqrt{\alpha^{(t)}}, \mathbf{v}^{(t)} = \gamma^{(t)} / \sqrt{\alpha^{(t)}}$
 - (e) Update \mathbf{f}, \mathbf{g} :
$$\mathbf{f}^{(t)} = \left(\mathbf{f} - 2 \left(\omega \circ \left(G^T \mathbf{v} - \sum_{r=1}^{t-2} \left(\mathbf{v}^{(r)T} \mathbf{v} \right) \omega^{(r)} \right) \right) + \|\mathbf{v}\|^2 (\omega \circ \omega) \right)^{(t-1)},$$

$$\mathbf{g}^{(t)} = (\mathbf{g} - (\omega \circ \omega))^{(t-1)}.$$
5. $W = [w^{(1)} \quad w^{(2)} \quad \dots \quad w^{(k)}]^T, \tilde{K}_{\mathcal{S}} = W^T W$
6. $\Omega = \text{eigvec}(WW^T)$
7. $Y = \Omega_d^T W, \tilde{K}_{\mathcal{S},d} = Y^T Y$

5 Properties of data sets

Experiments have been conducted on six benchmark data sets, whose properties are summarized in Table 1¹. The *Reuters-21578* is the training set of the Reuters-21578 collection [1], which has been used in the evaluation of many document clustering and classification techniques. The *Reviews* and *LAI* are documents data sets from TREC collections². The pre-processed versions of *Reviews* and *LAI* that are distributed with the CLUTO Toolkit³ [2] were used. The *MNIST-4K* is a subset of the MNIST data set of handwritten digits, which has been used to evaluate different digit recognition

¹The data sets *Reuters-21578*, *MNIST-4K*, *PIE-20* and *YaleB-38* are available in MAT format at: <http://www.zjucadcg.cn/dengcai/Data/data.html>. *PIE-20* is a subset of *PIE-32x32* with the images of the first 20 persons.

²<http://trec.nist.gov>

³<http://www.cs.umn.edu/~karypis/cluto>

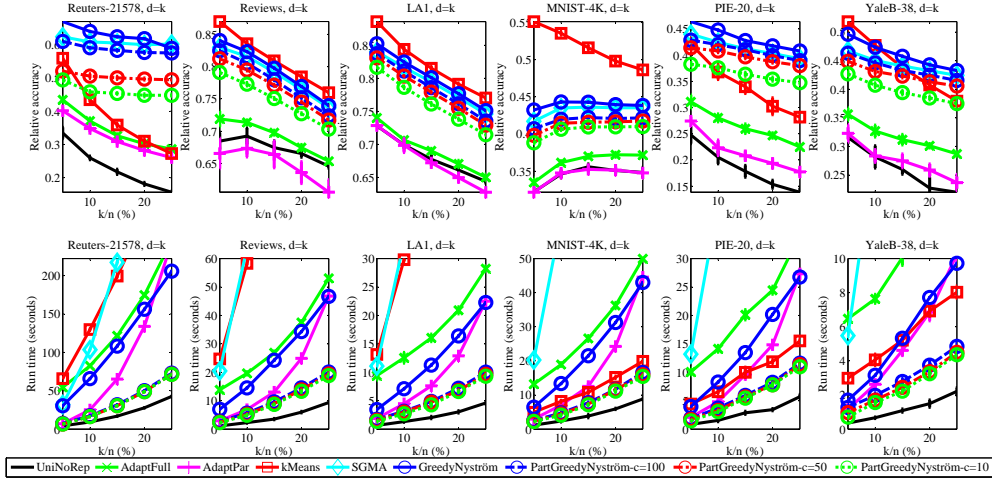


Figure 2: The relative accuracy and run time for the calculation of rank- k approximations \tilde{K}_S .

algorithms⁴. The *PIE-20* and *YaleB-38* are subsets of the CMU PIE [3] and Extended Yale Face [4] data sets respectively. These two data sets contain face images and have been used in the evaluation of many face recognition algorithms.

6 Results for rank- k Nyström approximation

Figure 2 shows the relative accuracy and run times for different methods when calculating the rank- k Nyström approximations \tilde{K}_S . It can be observed from these results that the proposed algorithms obtain very accurate low-rank approximations, with minimum overhead in run time.

References

- [1] D.D. Lewis. Reuters-21578 text categorization test collection distribution 1.0, 1999.
- [2] G. Karypis. CLUTO - a clustering toolkit. Technical Report #02-017, University of Minnesota, Department of Computer Science, 2003.
- [3] T. Sim, S. Baker, and M. Bsat. The CMU pose, illumination, and expression database. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 25(12):1615–1618, 2003.
- [4] K.C. Lee, J. Ho, and D. Kriegman. Acquiring linear subspaces for face recognition under variable lighting. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 27(5):684–698, 2005.

⁴<http://yann.lecun.com/exdb/mnist>