

THE VARIATIONAL NYSTRÖM METHOD FOR LARGE-SCALE SPECTRAL PROBLEMS

Max Vladymyrov, Google Inc, and Miguel Á. Carreira-Perpiñán, EECS, UC Merced



1 Abstract

Spectral methods for dimensionality reduction and clustering require solving an eigenproblem defined by a sparse affinity matrix. When this matrix is large, one seeks an approximate solution. The standard way to do this is the Nyström method, which first solves a small eigenproblem considering only a subset of landmark points, and then applies an out-of-sample formula to extrapolate the solution to the entire dataset. We show that by constraining the original problem to satisfy the Nyström formula, we obtain an approximation that is computationally simple and efficient, but achieves a lower approximation error using fewer landmarks and less runtime. We also study the role of normalization in the computational cost and quality of the resulting solution.

2 Spectral methods

Consider a spectral problem:

$$\min_{\mathbf{X}} \text{tr}(\mathbf{X}\mathbf{M}\mathbf{X}^T) \quad \text{s.t.} \quad \mathbf{X}\mathbf{X}^T = \mathbf{I} \quad (1)$$

where

- \mathbf{M} is an $N \times N$ symmetric matrix (usually, a graph Laplacian) constructed on a dataset $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$ of $D \times N$,
- $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ are coordinates in \mathbb{R}^d for the N data points (embedding), $d < D$.

This minimization problem occurs in Laplacian Eigenmaps, ISOMAP, Locally Linear Embedding, Spectral Clustering, Kernel PCA etc.

The solution is given by the d trailing eigenvectors \mathbf{U}_M of \mathbf{M} , which is costly to compute when N is large. Our goal is to solve problems of the type (1) approximately.

3 Landmark approximation methods

Select L landmarks $\tilde{\mathbf{Y}} = (\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_L)$ from \mathbf{Y} (e.g. randomly).

W.l.o.g. rearrange \mathbf{M} with landmarks first: $\mathbf{M} = \begin{pmatrix} \mathbf{A} & \mathbf{B}_{21}^T \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{pmatrix}$ and $\mathbf{C} = \begin{pmatrix} \mathbf{A} \\ \mathbf{B}_{21} \end{pmatrix}$ are the columns of \mathbf{M} that correspond to $\tilde{\mathbf{Y}}$.

Ways to approximate the computation of (1):

1. Nyström method

Essentially, an out-of-sample formula:

1. Solve the eigenproblem for $\tilde{\mathbf{Y}}$.
2. Predict the rest of the points with an interpolation formula.

$$\tilde{\mathbf{U}}_M = \begin{pmatrix} \mathbf{U}_A \\ \mathbf{B}_{21} \mathbf{U}_A \mathbf{\Lambda}_A^{-1} \end{pmatrix} = \mathbf{C} \mathbf{U}_A \mathbf{\Lambda}_A^{-1}$$

Problem: ignores non-landmark points for the initial prediction \rightarrow interpolation over the bad solution is bad.

2. Column Sampling

Approximation is given by the left singular vectors of \mathbf{C} :

$$\mathbf{C} = \mathbf{U}_C \mathbf{\Sigma}_C \mathbf{V}_C^T \Rightarrow \tilde{\mathbf{U}}_M = \mathbf{U}_C$$

Rewriting using the SVD of $\mathbf{C}^T \mathbf{C}$ we get $\tilde{\mathbf{U}}_M = \mathbf{C} \mathbf{U}_{C^T \mathbf{C}} \mathbf{\Lambda}_{C^T \mathbf{C}}^{-1/2}$.

Problem: uses more data than Nyström (columns of \mathbf{M}), but still ignores non-landmark/non-landmark interactions \mathbf{B}_{22} for prediction.

3. Locally Linear Landmarks

1. Construct the local linear projection matrix \mathbf{Z} from \mathbf{Y} :

$$\mathbf{y}_n \approx \sum_{l=1}^L z_{nl} \tilde{\mathbf{y}}_l, \quad n = 1, \dots, N \Rightarrow \mathbf{Y} \approx \tilde{\mathbf{Y}} \mathbf{Z} \quad (2)$$

2. Assume that projection is also satisfied in the embedding space: $\mathbf{X} = \tilde{\mathbf{X}} \mathbf{Z}$

3. Adding this constraint to (1) results in a reduced $L \times L$ eigenproblem:

$$\min_{\tilde{\mathbf{X}}} \text{tr}(\tilde{\mathbf{X}} \tilde{\mathbf{M}} \tilde{\mathbf{X}}^T) \quad \text{s.t.} \quad \tilde{\mathbf{X}} \mathbf{Z} \mathbf{Z}^T \tilde{\mathbf{X}}^T = \mathbf{I}$$

with reduced affinities $\tilde{\mathbf{M}} = \mathbf{Z} \mathbf{M} \mathbf{Z}^T$.

4. Reconstruct final embedding using $\mathbf{X} = \tilde{\mathbf{X}} \mathbf{Z}$

Final eigenvector approximation: $\tilde{\mathbf{U}}_M = \mathbf{Z} \tilde{\mathbf{X}}$.

Problem: Local linearly assumption may not be always true. Strong dependence on features \mathbf{Y} .

4. Random Projection method

1. Use a random matrix $\mathbf{S}_{N \times L}$ to form a low-dimensional sample matrix $\mathbf{M}_S = \mathbf{M} \mathbf{S}$,
2. Compute SVD of the projection of \mathbf{M} onto QR decomposition of \mathbf{M}_S ,
3. Project the results back to the original space.

4 Generalizing approximations

All the methods above can be generalized as follows:

1. Define an out-of-sample matrix \mathbf{Z} .
2. Compute a reduced eigenproblem and a matrix $\mathbf{Q}_{L \times d}$ that depends on it.
3. Final approximation is equal to $\mathbf{U}_M = \mathbf{Z} \mathbf{Q}$

Algorithm	$\mathbf{Z}_{N \times L}$	$\mathbf{Q}_{L \times d}$	Eigenproblem $\mathbf{A} \mathbf{U} = \mathbf{B} \mathbf{U} \mathbf{\Lambda}$
Nyström	\mathbf{C}	$\mathbf{U} \mathbf{\Lambda}^{-1}$	\mathbf{A}, \mathbf{I}
Column sampling	\mathbf{C}	$\mathbf{U} \mathbf{\Lambda}^{-1/2}$	$\mathbf{Z}^T \mathbf{Z}, \mathbf{I}$
Random Projection	$\text{qr}(\mathbf{M}^q \mathbf{S})$	\mathbf{U}	$\mathbf{Z}^T \mathbf{M} \mathbf{Z}, \mathbf{I}$
LLL	eq. (2)	\mathbf{U}	$\mathbf{Z}^T \mathbf{M} \mathbf{Z}, \mathbf{Z}^T \mathbf{Z}$
Variational Nyström	\mathbf{C}	\mathbf{U}	$\mathbf{Z}^T \mathbf{M} \mathbf{Z}, \mathbf{Z}^T \mathbf{Z}$

5 Variational Nyström

Since the embedding is defined by the Nyström out-of-sample formula $\mathbf{X} = \tilde{\mathbf{X}} \mathbf{C}^T$, we add this as a constraint to the spectral problem (1):

$$\min_{\tilde{\mathbf{X}}} \text{tr}(\tilde{\mathbf{X}} \mathbf{C}^T \mathbf{M} \mathbf{C} \tilde{\mathbf{X}}^T) \quad \text{s.t.} \quad \mathbf{X} \mathbf{X}^T = \mathbf{I}, \quad \mathbf{X} = \tilde{\mathbf{X}} \mathbf{C}^T$$

Which result in a reduced eigenproblem:

$$\min_{\tilde{\mathbf{X}}} \text{tr}(\tilde{\mathbf{X}} \mathbf{C}^T \mathbf{M} \mathbf{C} \tilde{\mathbf{X}}^T) \quad \text{s.t.} \quad \tilde{\mathbf{X}} \mathbf{C}^T \mathbf{C} \tilde{\mathbf{X}}^T = \mathbf{I}$$

From LLL perspective:

- replace customary built out-of-sample matrix \mathbf{Z} with a readily available column matrix \mathbf{C} ,
- abandon local linearity assumption of the weights \mathbf{Z} and dependence on the features \mathbf{Y} ,
- save computation of \mathbf{Z} .

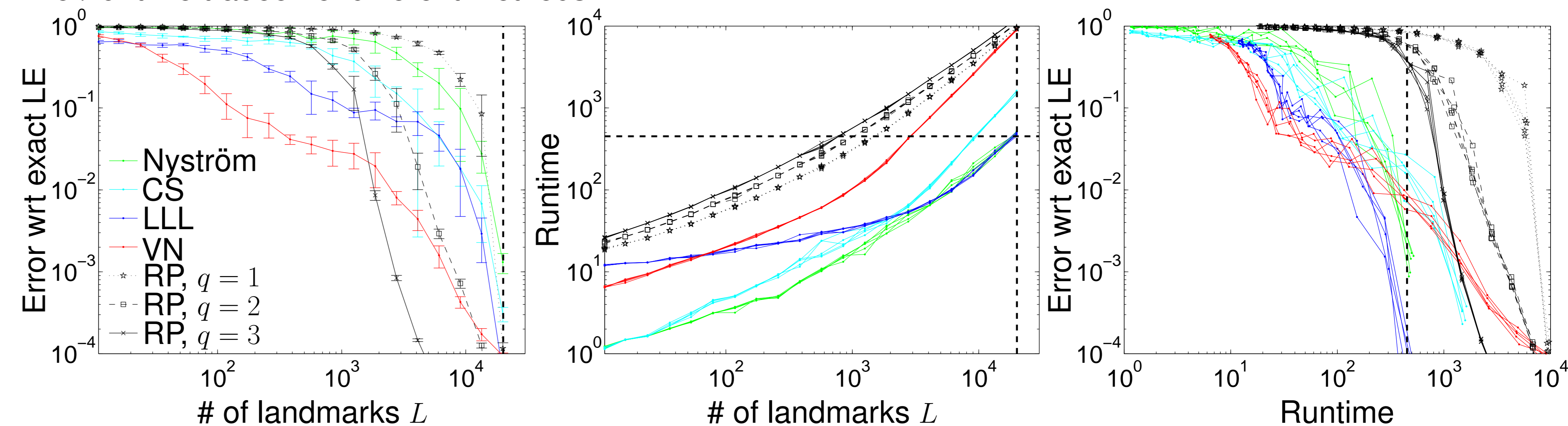
From Nyström perspective:

- for fixed $\tilde{\mathbf{Y}}$ gives better approx. than Nyström or Column Sampling (optimal for the out-of-sample kernel \mathbf{C}).
- use the same out-of-sample matrix \mathbf{C} , but optimize the choice of the reduced eigenproblem,
- use all the elements from \mathbf{M} to construct the reduced eigenproblem,
- forgo the interpolating property of Nyström.

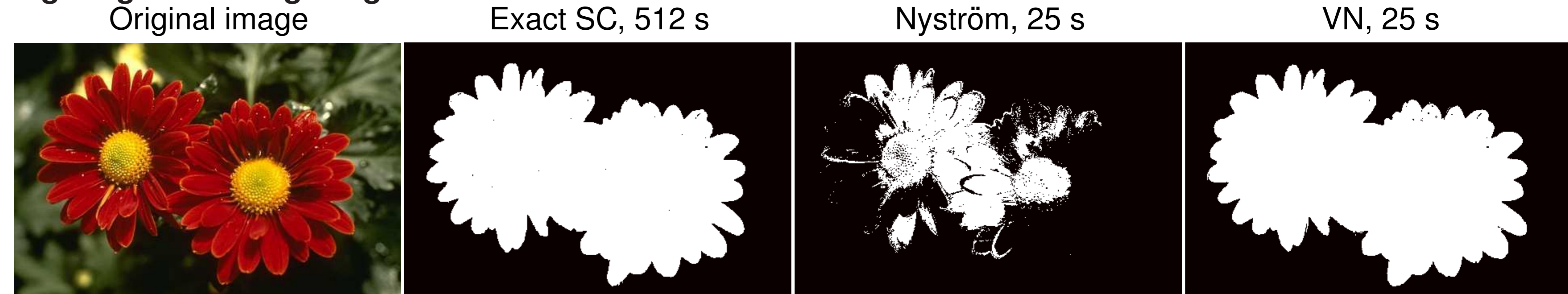
7 Experiments

Medium size experiment. Reduce dimensionality of 20 000 random points from MNIST to $d = 10$.

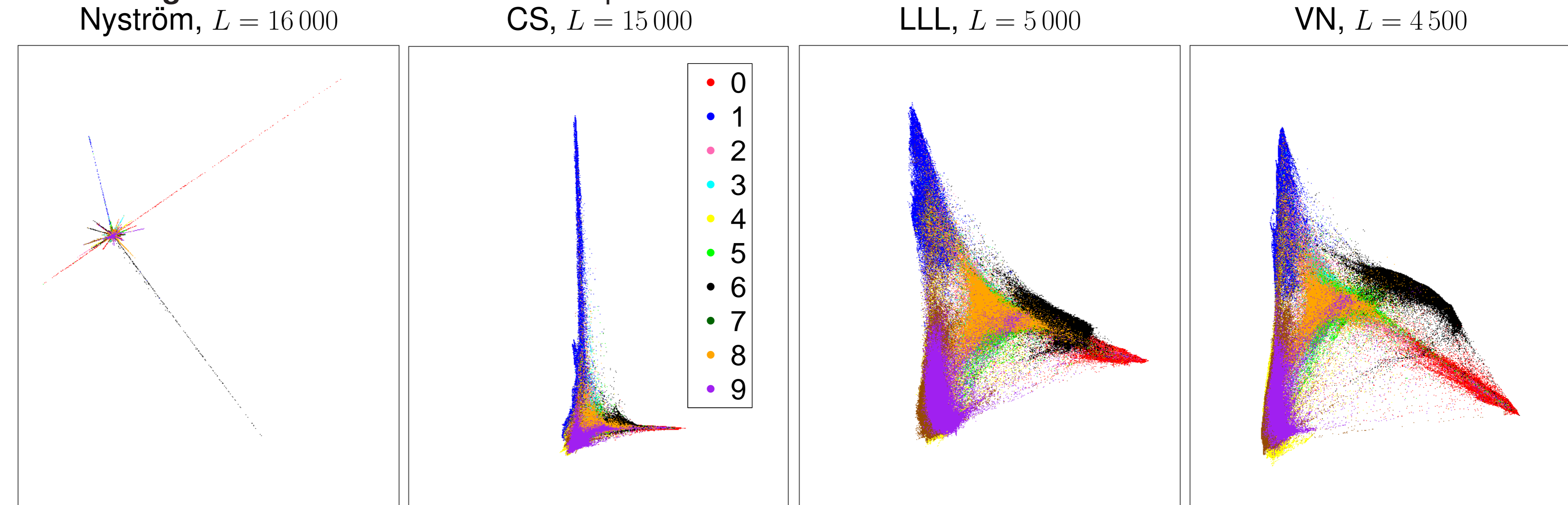
Error/runtime tradeoff of different methods.



Figure/ground image segmentation.



Embedding of the infinite MNIST. 1 020 000 points. Fix the runtime to 10 min.



6 Subsampling Graph Laplacians

Consider \mathbf{M} given by graph Laplacian:

$$\mathbf{L} = \mathbf{D}^{-1/2} (\mathbf{D} - \mathbf{W}) \mathbf{D}^{-1/2},$$

(as in Laplacian Eigenmaps, spectral clustering etc.) where

$$w_{nm} = \exp(-\|\mathbf{y}_n - \mathbf{y}_m\|/2\sigma^2) \quad \mathbf{D} = \text{diag} \left(\sum_{m=1}^N w_{nm} \right)$$

- \mathbf{L} is a data dependent kernel: Graph Laplacian computed for a subset of L input points is not equal to the $L \times L$ subset of graph Laplacian constructed for N points.

- This can be problematic for methods that depend on subsampling, such as Nyström and Variational Nyström. Not a problem for LLL, since there is no subsampling involved.

- Our solution: normalize out-of-sample kernel separately, but in a way that (1) interpolates over the landmarks and (2) gives exact solution when $L = N$:

- We show that for the Variational Nyström final normalization is more general and has much simpler form than for Nyström method.

8 Conclusions

1. The Variational Nyström method is the optimal way to use the out-of-sample Nyström formula to solve an eigenproblem approximately. It is able to achieve a low-to-medium accuracy solution faster than Nyström and other methods.
2. We present a simple unified model of spectral approximations, combining many existing algorithms such as Nyström, Column Sampling, LLL.
3. We study the role of normalization in subsampling of the graph Laplacian kernel.

ICML@NYC

International Conference on Machine Learning

Partially supported by NSF award IIS-1423515.