

10-605/10-805: Machine Learning with Large Datasets

FALL 2022

Kernel Approximation

Outline

1. Background: Kernel Methods
2. Kernel Ridge Regression
3. SVD
4. Large-Scale Kernel Methods (via Approximate SVD)

RECALL

Empirical risk minimization

$$\hat{f}_n = \arg \min_{f \in F} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

- A popular approach in supervised ML
- Given a loss ℓ and data $(x_1, y_1), \dots (x_n, y_n)$, we estimate a predictor f by minimizing the *empirical risk*
- We typically restrict this predictor to lie in some class, F
 - Could reflect our prior knowledge about the task
 - Or may be for computational convenience

Question: how should we select our function class, F , and our loss function, ℓ ??

Empirical risk minimization

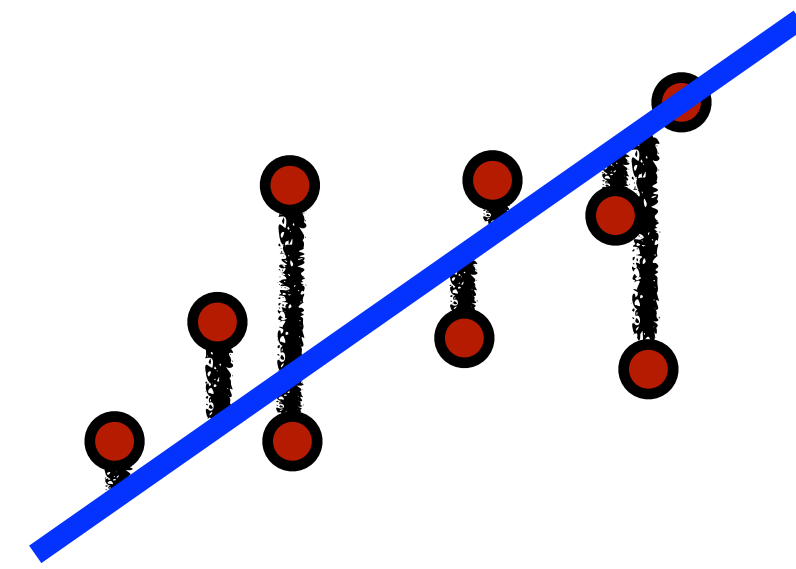
$$\hat{f}_n = \arg \min_{f \in F} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

linear (actually, affine) functions

square loss

- A popular approach in supervised ML
- Given a loss ℓ and data $(x_1, y_1), \dots (x_n, y_n)$, we estimate a predictor f by minimizing the *empirical risk*
- We typically restrict this predictor to lie in some class, F
 - Could reflect our prior knowledge about the task
 - Or may be for computational convenience

Question: how should we select our function class, F , and our loss function, ℓ ??



Given n training points with k features, we define:

- $\mathbf{X} \in \mathbb{R}^{n \times k}$: matrix storing points
- $\mathbf{y} \in \mathbb{R}^n$: real-valued labels
- $\hat{\mathbf{y}} \in \mathbb{R}^n$: predicted labels, where $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- $\mathbf{w} \in \mathbb{R}^k$: regression parameters / model to learn

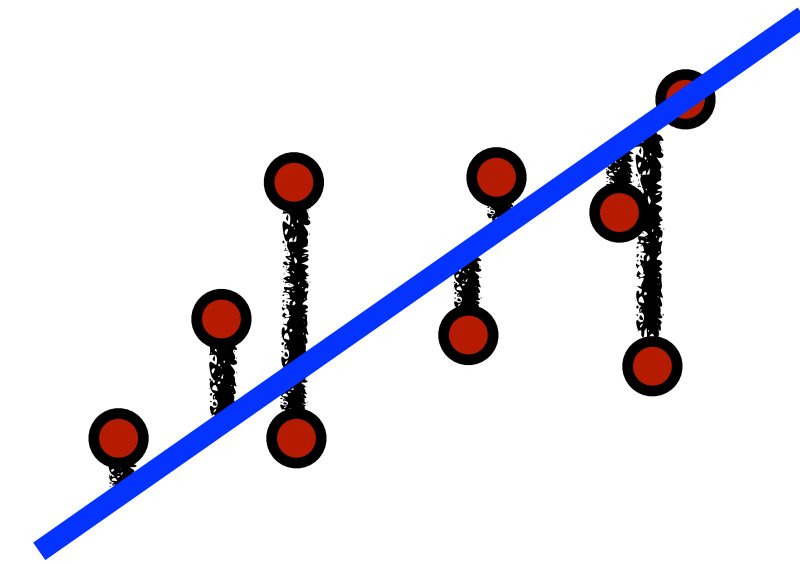
Least Squares Regression: Learn mapping (\mathbf{w}) from features to labels that minimizes residual sum of squares:

$$\min_{\mathbf{w}} ||\mathbf{X}\mathbf{w} - \mathbf{y}||_2^2$$

Equivalent $\min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^{(i)} - y^{(i)})^2$ by definition of Euclidean norm

Given n training points with k features, we define:

- $\mathbf{X} \in \mathbb{R}^{n \times k}$: matrix storing points
- $\mathbf{y} \in \mathbb{R}^n$: real-valued labels
- $\hat{\mathbf{y}} \in \mathbb{R}^n$: predicted labels, where $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- $\mathbf{w} \in \mathbb{R}^k$: regression parameters / model to learn



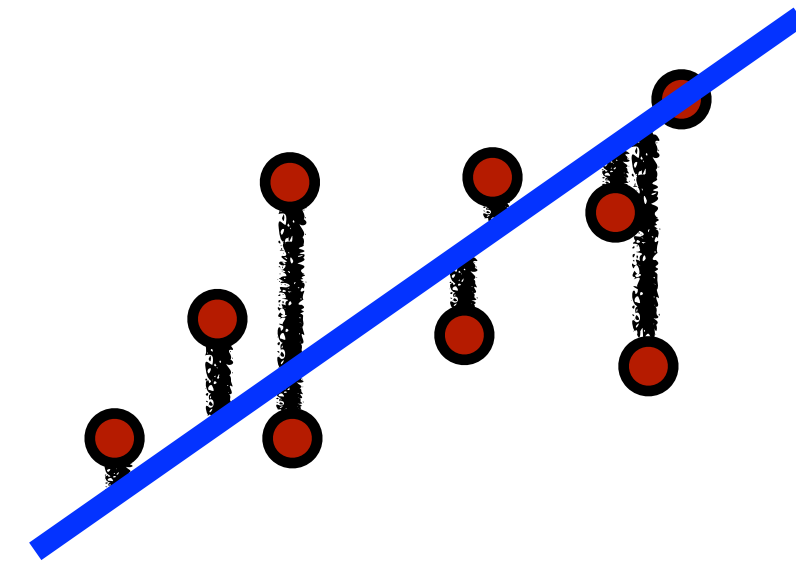
Ridge Regression: Learn mapping (\mathbf{w}) that minimizes residual sum of squares along with a regularization term:

$$\min_{\mathbf{w}} \underbrace{||\mathbf{X}\mathbf{w} - \mathbf{y}||_2^2}_{\text{Training Error}} + \underbrace{\lambda ||\mathbf{w}||_2^2}_{\text{Model Complexity}}$$

Closed-form solution: $\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_k)^{-1} \mathbf{X}^\top \mathbf{y}$

free parameter trades off
between training error and
model complexity

Why a linear mapping?



Simple

Often works well in practice

Can introduce complexity via feature extraction

KERNEL METHODS

Motivation

How to choose nonlinear basis function for regression?

$$\boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x})$$

- $\boldsymbol{\phi}(\cdot)$ maps the original feature vector \boldsymbol{x} to a *new* M -dimensional feature vector

KERNEL METHODS

Motivation

How to choose nonlinear basis function for regression?

$$w^T \phi(x)$$

- $\phi(\cdot)$ maps the original feature vector x to a *new* M -dimensional feature vector
- We can sidestep the issue of choosing which $\phi(\cdot)$ to use by *equivalently* choosing a *kernel function*
 - ➔ kernel function operates on inner products of the observations/data points

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $c \geq 0$ and d is a positive integer.

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $c \geq 0$ and d is a positive integer.

Gaussian kernel, RBF kernel, or Gaussian RBF kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

- Shift-invariant kernel (only depends on difference between two inputs)
- Corresponds to a feature space with *infinite* dimensions (but we can work directly with the original features)

‘Kernel trick’

- Many learning methods rely on training and test data only in the form of *inner products*, e.g., regularized least squares, nearest neighbors, SVMs
- We can use a kernel function to introduce nonlinearity, i.e., “*kernelizing*” the methods
 - ➔ allows us to consider complex, high dimensional feature spaces
 - ➔ but, requires expressing training and inference in terms of inner products, e.g., in terms of the kernel matrix, K

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

TODAY

Kernel methods at scale

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

- Kernel methods allow us to easily encode data in high-dimensional feature spaces
- But, they come at a cost: storing and operating on the $O(n^2)$ kernel matrix
- This may make kernel methods infeasible for large values of n

Question: Can we approximate the kernel matrix with something simpler and more scalable?

Outline

1. Background: Kernel Methods
2. Kernel Ridge Regression
3. SVD
4. Large-Scale Kernel Methods (via Approximate SVD)

EXAMPLE: KERNEL RIDGE REGRESSION

Motivation

How to choose nonlinear basis function for regression?

$$w^T \phi(x)$$

- $\phi(\cdot)$ maps the original feature vector x to a *new* M -dimensional feature vector
- We can sidestep the issue of choosing which $\phi(\cdot)$ to use by *equivalently* choosing a *kernel function*

EXAMPLE: KERNEL RIDGE REGRESSION

Motivation

How to choose nonlinear basis function for regression?

$$\boldsymbol{w}^T \phi(\boldsymbol{x})$$

- $\phi(\cdot)$ maps the original feature vector \boldsymbol{x} to a *new* M -dimensional feature vector
- We can sidestep the issue of choosing which $\phi(\cdot)$ to use by *equivalently* choosing a *kernel function*

Regularized least squares

$$J(\boldsymbol{w}) = \frac{1}{2} \sum_i (y_i - \boldsymbol{w}^T \phi(\boldsymbol{x}_i))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

Its solution is given by

$$\frac{\partial J(\boldsymbol{w})}{\partial \boldsymbol{w}} = \sum_i (y_i - \boldsymbol{w}^T \phi(\boldsymbol{x}_i))(-\phi(\boldsymbol{x}_i)) + \lambda \boldsymbol{w} = 0$$

EXAMPLE: KERNEL RIDGE REGRESSION

Solution

The optimal parameter vector is a linear combination of features

$$\boldsymbol{w}^* = \sum_i \frac{1}{\lambda} (y_i - \boldsymbol{w}^{*\top} \boldsymbol{\phi}(\boldsymbol{x}_i)) \boldsymbol{\phi}(\boldsymbol{x}_i)$$

EXAMPLE: KERNEL RIDGE REGRESSION

Solution

The optimal parameter vector is a linear combination of features

$$\boldsymbol{w}^* = \sum_i \frac{1}{\lambda} (y_i - \boldsymbol{w}^{*\top} \boldsymbol{\phi}(\boldsymbol{x}_i)) \boldsymbol{\phi}(\boldsymbol{x}_i) = \sum_i \alpha_i \boldsymbol{\phi}(\boldsymbol{x}_i) = \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$$

EXAMPLE: KERNEL RIDGE REGRESSION

Solution

The optimal parameter vector is a linear combination of features

$$\mathbf{w}^* = \sum_i \frac{1}{\lambda} (y_i - \mathbf{w}^{*\top} \phi(\mathbf{x}_i)) \phi(\mathbf{x}_i) = \sum_i \alpha_i \phi(\mathbf{x}_i) = \mathbf{\Phi}^\top \boldsymbol{\alpha}$$

- $\alpha_i = \frac{1}{\lambda} (y_i - \mathbf{w}^{*\top} \phi(\mathbf{x}_i))$
- $\mathbf{\Phi}$ is the *design matrix* of *transformed* features
- Its transpose is made of column vectors and is given by

$$\mathbf{\Phi}^\top = (\phi(\mathbf{x}_1) \ \phi(\mathbf{x}_2) \ \cdots \ \phi(\mathbf{x}_n)) \in \mathbb{R}^{M \times n}$$

EXAMPLE: KERNEL RIDGE REGRESSION

Solution

The optimal parameter vector is a linear combination of features

$$\boldsymbol{w}^* = \sum_i \frac{1}{\lambda} (y_i - \boldsymbol{w}^{*\top} \phi(\boldsymbol{x}_i)) \phi(\boldsymbol{x}_i) = \sum_i \alpha_i \phi(\boldsymbol{x}_i) = \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$$

- $\alpha_i = \frac{1}{\lambda} (y_i - \boldsymbol{w}^{*\top} \phi(\boldsymbol{x}_i))$
- $\boldsymbol{\Phi}$ is the *design matrix* of *transformed* features
- Its transpose is made of column vectors and is given by

$$\boldsymbol{\Phi}^\top = (\phi(\boldsymbol{x}_1) \ \phi(\boldsymbol{x}_2) \ \cdots \ \phi(\boldsymbol{x}_n)) \in \mathbb{R}^{M \times n}$$

Of course, we don't know what $\boldsymbol{\alpha}$ (the vector of α_i) corresponds to, given its dependence on \boldsymbol{w}^* !

EXAMPLE: KERNEL RIDGE REGRESSION

Dual formulation

Regularized least squares: $J(\boldsymbol{w}) = \frac{1}{2} \sum_n (y_n - \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$

Substitute $\boldsymbol{w}^* = \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$ into $J(\boldsymbol{w})$ to obtain a function of $\boldsymbol{\alpha}$:

$$J(\boldsymbol{\alpha}) = \frac{1}{2} \boldsymbol{\alpha}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha} - (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{y})^\top \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$$

EXAMPLE: KERNEL RIDGE REGRESSION

Dual formulation

Regularized least squares: $J(\boldsymbol{w}) = \frac{1}{2} \sum_n (y_n - \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$

Substitute $\boldsymbol{w}^* = \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$ into $J(\boldsymbol{w})$ to obtain a function of $\boldsymbol{\alpha}$:

$$J(\boldsymbol{\alpha}) = \frac{1}{2} \boldsymbol{\alpha}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha} - (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{y})^\top \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$$

The *Gram matrix* or *kernel matrix* — $\boldsymbol{\Phi} \boldsymbol{\Phi}^\top$ — appears multiple times

$$\begin{aligned} \boldsymbol{K} &= \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \\ &= \begin{pmatrix} \boldsymbol{\phi}(\boldsymbol{x}_1)^\top \boldsymbol{\phi}(\boldsymbol{x}_1) & \boldsymbol{\phi}(\boldsymbol{x}_1)^\top \boldsymbol{\phi}(\boldsymbol{x}_2) & \cdots & \boldsymbol{\phi}(\boldsymbol{x}_1)^\top \boldsymbol{\phi}(\boldsymbol{x}_N) \\ \boldsymbol{\phi}(\boldsymbol{x}_2)^\top \boldsymbol{\phi}(\boldsymbol{x}_1) & \boldsymbol{\phi}(\boldsymbol{x}_2)^\top \boldsymbol{\phi}(\boldsymbol{x}_2) & \cdots & \boldsymbol{\phi}(\boldsymbol{x}_2)^\top \boldsymbol{\phi}(\boldsymbol{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ \boldsymbol{\phi}(\boldsymbol{x}_N)^\top \boldsymbol{\phi}(\boldsymbol{x}_1) & \boldsymbol{\phi}(\boldsymbol{x}_N)^\top \boldsymbol{\phi}(\boldsymbol{x}_2) & \cdots & \boldsymbol{\phi}(\boldsymbol{x}_N)^\top \boldsymbol{\phi}(\boldsymbol{x}_N) \end{pmatrix} \in \mathbb{R}^{n \times n} \end{aligned}$$

EXAMPLE: KERNEL RIDGE REGRESSION

The derivation of $J(\boldsymbol{\alpha})$

$$J(\boldsymbol{w}) = \frac{1}{2} \sum_n (y - \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

EXAMPLE: KERNEL RIDGE REGRESSION

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\boldsymbol{w}) &= \frac{1}{2} \sum_n (y - \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2 \\ &= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{w}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2 \end{aligned}$$

$$(\boldsymbol{y}, \boldsymbol{\Phi} \boldsymbol{w} \in \mathbb{R}^N)$$

EXAMPLE: KERNEL RIDGE REGRESSION

The derivation of $J(\boldsymbol{\alpha})$

$$J(\boldsymbol{w}) = \frac{1}{2} \sum_n (y - \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{w}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2$$

$$(\boldsymbol{y}, \boldsymbol{\Phi} \boldsymbol{w} \in \mathbb{R}^N)$$

$$(\boldsymbol{w}^* = \boldsymbol{\Phi}^T \boldsymbol{\alpha})$$

EXAMPLE: KERNEL RIDGE REGRESSION

The derivation of $J(\boldsymbol{\alpha})$

$$J(\boldsymbol{w}) = \frac{1}{2} \sum_n (y - \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{w}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

$$(\boldsymbol{y}, \boldsymbol{\Phi} \boldsymbol{w} \in \mathbb{R}^N)$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^\top \boldsymbol{\alpha}\|_2^2$$

$$(\boldsymbol{w}^* = \boldsymbol{\Phi}^\top \boldsymbol{\alpha})$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{K} \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$$

$$(\boldsymbol{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^\top, \|\boldsymbol{v}\|_2^2 = \boldsymbol{v}^\top \boldsymbol{v})$$

EXAMPLE: KERNEL RIDGE REGRESSION

The derivation of $J(\boldsymbol{\alpha})$

$$J(\boldsymbol{w}) = \frac{1}{2} \sum_n (y - \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{w}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

$$(\boldsymbol{y}, \boldsymbol{\Phi} \boldsymbol{w} \in \mathbb{R}^N)$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^\top \boldsymbol{\alpha}\|_2^2$$

$$(\boldsymbol{w}^* = \boldsymbol{\Phi}^\top \boldsymbol{\alpha})$$

$$= \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{K} \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\alpha}$$

$$(\boldsymbol{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^\top, \|\boldsymbol{v}\|_2^2 = \boldsymbol{v}^\top \boldsymbol{v})$$

$$\propto \frac{1}{2} \boldsymbol{\alpha}^\top \boldsymbol{K}^\top \boldsymbol{K} \boldsymbol{\alpha} - \boldsymbol{y}^\top \boldsymbol{K} \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{K} \boldsymbol{\alpha}$$

EXAMPLE: KERNEL RIDGE REGRESSION

Optimal α

$$\frac{\partial J(\alpha)}{\partial \alpha} = \mathbf{K}^2 \alpha - \mathbf{K} y + \lambda \mathbf{K} \alpha = 0$$

EXAMPLE: KERNEL RIDGE REGRESSION

Optimal α

$$\frac{\partial J(\alpha)}{\partial \alpha} = \mathbf{K}^2 \alpha - \mathbf{K} \mathbf{y} + \lambda \mathbf{K} \alpha = 0$$

which leads to (assuming that \mathbf{K} is invertible)

$$\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

EXAMPLE: KERNEL RIDGE REGRESSION

Optimal α

$$\frac{\partial J(\alpha)}{\partial \alpha} = \mathbf{K}^2 \alpha - \mathbf{K} \mathbf{y} + \lambda \mathbf{K} \alpha = 0$$

which leads to (assuming that \mathbf{K} is invertible)

$$\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

- We only need to know \mathbf{K} in order to compute α !
- Solution doesn't involve $\phi(\cdot)$, but instead inner products $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$
- This observation will give rise to the use of *kernel functions*

EXAMPLE: KERNEL RIDGE REGRESSION

Optimal α

$$\frac{\partial J(\alpha)}{\partial \alpha} = \mathbf{K}^2 \alpha - \mathbf{K} \mathbf{y} + \lambda \mathbf{K} \alpha = 0$$

which leads to (assuming that \mathbf{K} is invertible)

$$\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

- We only need to know \mathbf{K} in order to compute α !
- Solution doesn't involve $\phi(\cdot)$, but instead inner products $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$
- This observation will give rise to the use of *kernel functions*

Note that the initial parameter vector does require knowledge of Φ

$$\mathbf{w}^* = \Phi^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

EXAMPLE: KERNEL RIDGE REGRESSION

Computing prediction needs only inner products too!

Since $\boldsymbol{w}^* = \boldsymbol{\Phi}^T (\boldsymbol{K} + \lambda \boldsymbol{I})^{-1} \boldsymbol{y}$, at test time we compute:

$$\boldsymbol{w}^{*\top} \boldsymbol{\phi}(x) = \boldsymbol{y}^T (\boldsymbol{K} + \lambda \boldsymbol{I})^{-1} \boldsymbol{\Phi} \boldsymbol{\phi}(x)$$

EXAMPLE: KERNEL RIDGE REGRESSION

Computing prediction needs only inner products too!

Since $\boldsymbol{w}^* = \boldsymbol{\Phi}^T(\boldsymbol{K} + \lambda\boldsymbol{I})^{-1}\boldsymbol{y}$, at test time we compute:

$$\begin{aligned} \boldsymbol{w}^{*\top}\boldsymbol{\phi}(\boldsymbol{x}) &= \boldsymbol{y}^T(\boldsymbol{K} + \lambda\boldsymbol{I})^{-1}\boldsymbol{\Phi}\boldsymbol{\phi}(\boldsymbol{x}) \\ &= \boldsymbol{y}^T(\boldsymbol{K} + \lambda\boldsymbol{I})^{-1} \begin{pmatrix} \boldsymbol{\phi}(\boldsymbol{x}_1)^\top \boldsymbol{\phi}(\boldsymbol{x}) \\ \boldsymbol{\phi}(\boldsymbol{x}_2)^\top \boldsymbol{\phi}(\boldsymbol{x}) \\ \vdots \\ \boldsymbol{\phi}(\boldsymbol{x}_N)^\top \boldsymbol{\phi}(\boldsymbol{x}) \end{pmatrix} \end{aligned}$$

EXAMPLE: KERNEL RIDGE REGRESSION

Computing prediction needs only inner products too!

Since $\mathbf{w}^* = \Phi^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$, at test time we compute:

$$\begin{aligned}\mathbf{w}^{*\top}\phi(\mathbf{x}) &= \mathbf{y}^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\Phi\phi(\mathbf{x}) \\ &= \mathbf{y}^T(\mathbf{K} + \lambda\mathbf{I})^{-1} \begin{pmatrix} \phi(\mathbf{x}_1)^\top\phi(\mathbf{x}) \\ \phi(\mathbf{x}_2)^\top\phi(\mathbf{x}) \\ \vdots \\ \phi(\mathbf{x}_N)^\top\phi(\mathbf{x}) \end{pmatrix} = \mathbf{y}^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{k}_x\end{aligned}$$

- We used the property that $(\mathbf{K} + \lambda\mathbf{I})^{-1}$ is symmetric (as \mathbf{K} is)
- \mathbf{k}_x is shorthand notation for the column vector of inner products between training set and test point

EXAMPLE: KERNEL RIDGE REGRESSION

Computing prediction needs only inner products too!

Since $\mathbf{w}^* = \Phi^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$, at test time we compute:

$$\begin{aligned}\mathbf{w}^{*\top}\phi(\mathbf{x}) &= \mathbf{y}^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\Phi\phi(\mathbf{x}) \\ &= \mathbf{y}^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\begin{pmatrix} \phi(\mathbf{x}_1)^\top\phi(\mathbf{x}) \\ \phi(\mathbf{x}_2)^\top\phi(\mathbf{x}) \\ \vdots \\ \phi(\mathbf{x}_N)^\top\phi(\mathbf{x}) \end{pmatrix} = \mathbf{y}^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{k}_x\end{aligned}$$

- We used the property that $(\mathbf{K} + \lambda\mathbf{I})^{-1}$ is symmetric (as \mathbf{K} is)
- \mathbf{k}_x is shorthand notation for the column vector of inner products between training set and test point
- To make a prediction we *only need* $\phi(\mathbf{x}_n)^T\phi(\mathbf{x})$!

Kernel PCA

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Given: $n \times k$ matrix of uncentered raw data

Goal: Compute $r \ll k$ dimensional representation

Step 1: Compute kernel matrix, \mathbf{K}

Step 2: Center \mathbf{K} via $\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{A}\mathbf{K} - \mathbf{K}\mathbf{A} + \mathbf{A}\mathbf{K}\mathbf{A}$, where \mathbf{A} is matrix containing $\frac{1}{n}$

Step 3: Find top r eigenvectors of $\tilde{\mathbf{K}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ and normalize to find \mathbf{P}

Step 4: Compute PCA Scores $\mathbf{Z} = \tilde{\mathbf{K}}\mathbf{P}$

‘Kernel trick’

- Many learning methods rely on training and test data only in the form of *inner products*, e.g., regularized least squares, nearest neighbors, SVMs
- We can use a kernel function to introduce nonlinearity, i.e., “*kernelizing*” the methods

Why use kernel functions?

Can define kernel matrix without specifying $\phi(\cdot)$

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Question: Can we approximate the kernel matrix with something simpler and more scalable?

Why do we care?

Kernel-based algorithms

- $n \times n$ similarity matrix (\mathbf{K})
- flexibility in choice of kernel (similarity measure)
- e.g., SVMs, Kernel Ridge Regression, Kernel PCA
- Recall for KRR: $\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$

Why do we care?

Kernel-based algorithms

- $n \times n$ similarity matrix (K)
- flexibility in choice of kernel (similarity measure)
- e.g., SVMs, Kernel Ridge Regression, Kernel PCA
- Recall for KRR: $\alpha = (K + \lambda I)^{-1}y$

Expensive for large-scale problems

- $O(n^2)$ space
- at least $O(n^2)$ time to operate on K
- $O(n^3)$ to invert

Example: we want to invert a large matrix, e.g., $n = 18M$

- K requires roughly 1300 TB
- Requires 40K 32GB RAM machines
- Even if we can store it, $O(n^3)$ time is not possible

Possible solution: Encourage implicit sparsity via **low-rank approximation**

- Some data are naturally approximately low rank (e.g., images)
- Low-rankedness can be viewed as a form of regularization
- Ultimately it's a modeling assumption; common in practice, but not guaranteed to be a good one

Key technical tool: SVD!

Outline

1. Background: Kernel Methods
2. Kernel Ridge Regression
3. SVD
4. Large-Scale Kernel Methods (via Approximate SVD)

Singular value decomposition (SVD)

Every matrix has the following decomposition:

SVD

Let $X \in \mathbb{R}^{n \times m}$ then

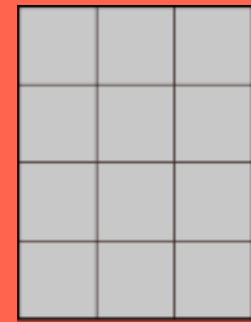
$$X = U \Sigma V^{\top},$$

where $U \in \mathbb{R}^{n \times n}$, $V \in \mathbb{R}^{m \times m}$ are orthogonal matrices (i.e. $U^{\top} = U^{-1}$) and $\Sigma \in \mathbb{R}^{n \times m}$ is a diagonal matrix with *singular values* of X denoted by σ_i appearing by non-increasing order: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(m,n)} \geq 0$.

- The square singular values of X are the eigenvalues of the matrix XX^{\top} or $X^{\top}X$, i.e., $\sigma_i(X) = \sqrt{\lambda_i(XX^{\top})} = \sqrt{\lambda_i(X^{\top}X)}$

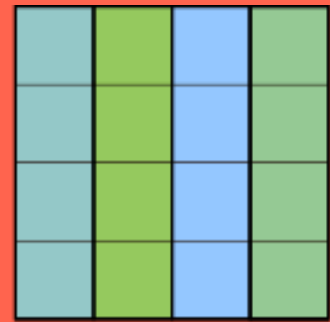
RECALL

SVD vs Compact SVD vs Truncated SVD



M

$m \times n$



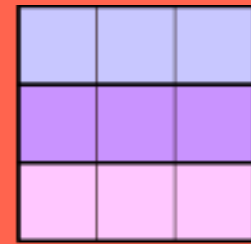
U

$m \times m$



Σ

$m \times n$



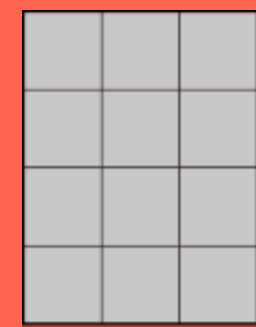
V*

$n \times n$

SVD: U, V orthogonal,
 Σ rectangular diagonal matrix

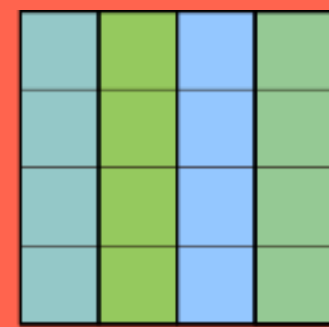
RECALL

SVD vs Compact SVD vs Truncated SVD



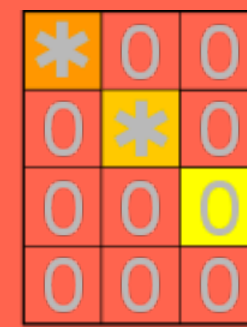
M

$m \times n$



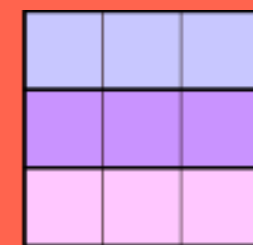
U

$m \times m$



Σ

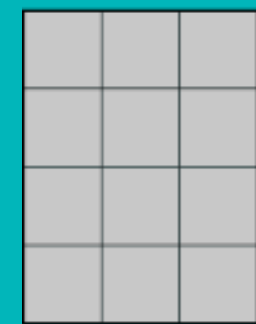
$m \times n$



V*

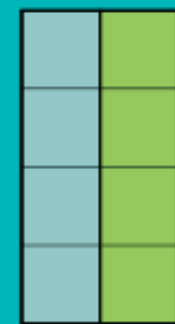
$n \times n$

SVD: U, V orthogonal,
 Σ rectangular diagonal matrix



M

$m \times n$



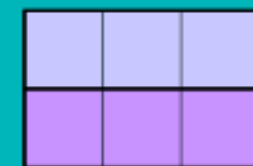
U_r

$m \times r$



Σ_r

$r \times r$



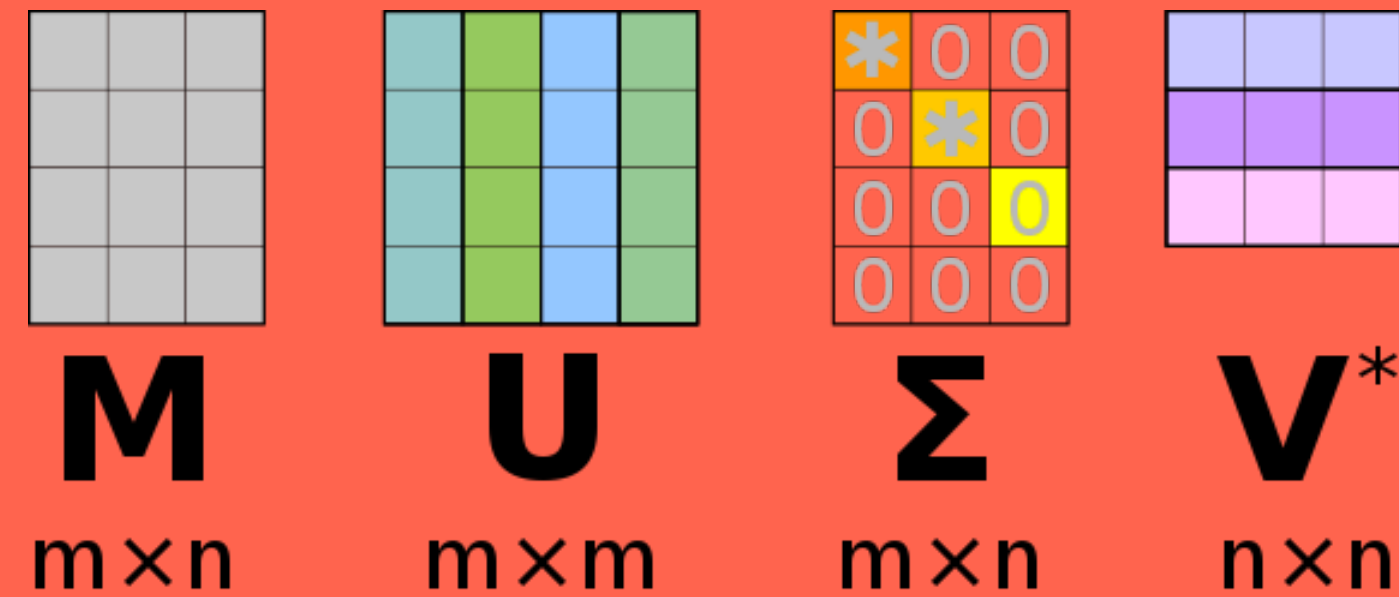
V_r*

$r \times n$

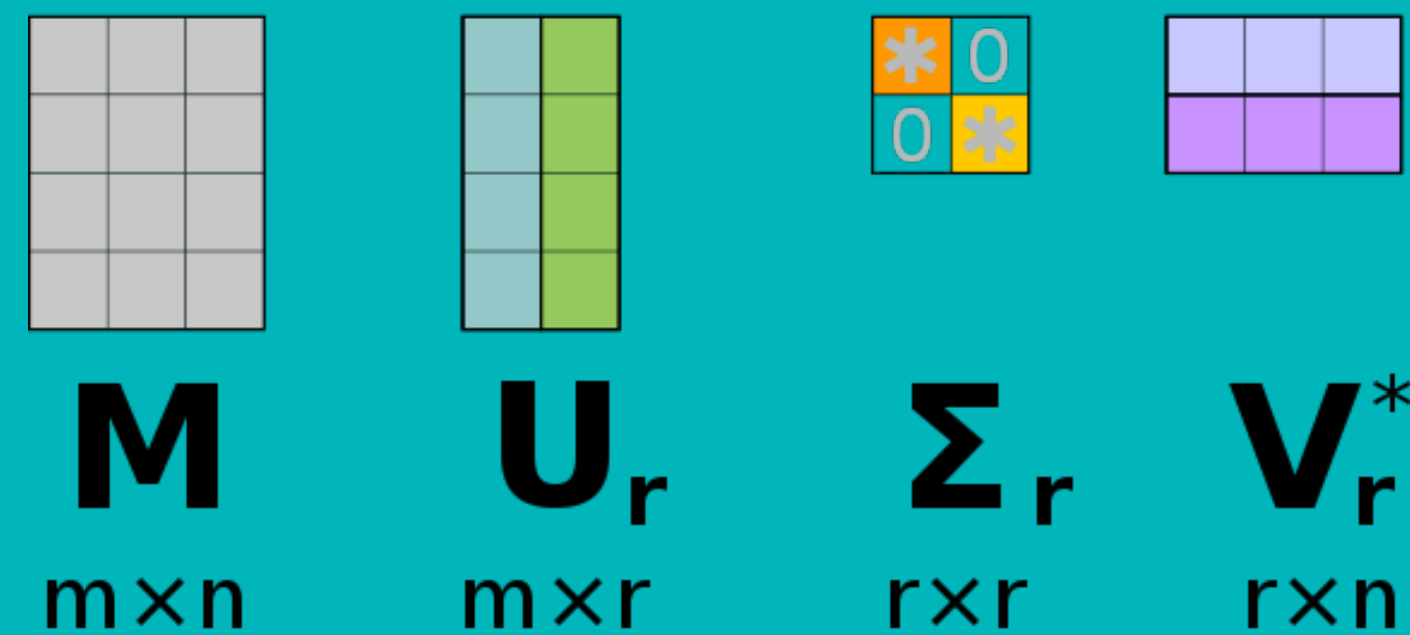
Compact SVD: U, V semi-orthogonal,
 Σ square diagonal matrix with non-zero SVs

RECALL

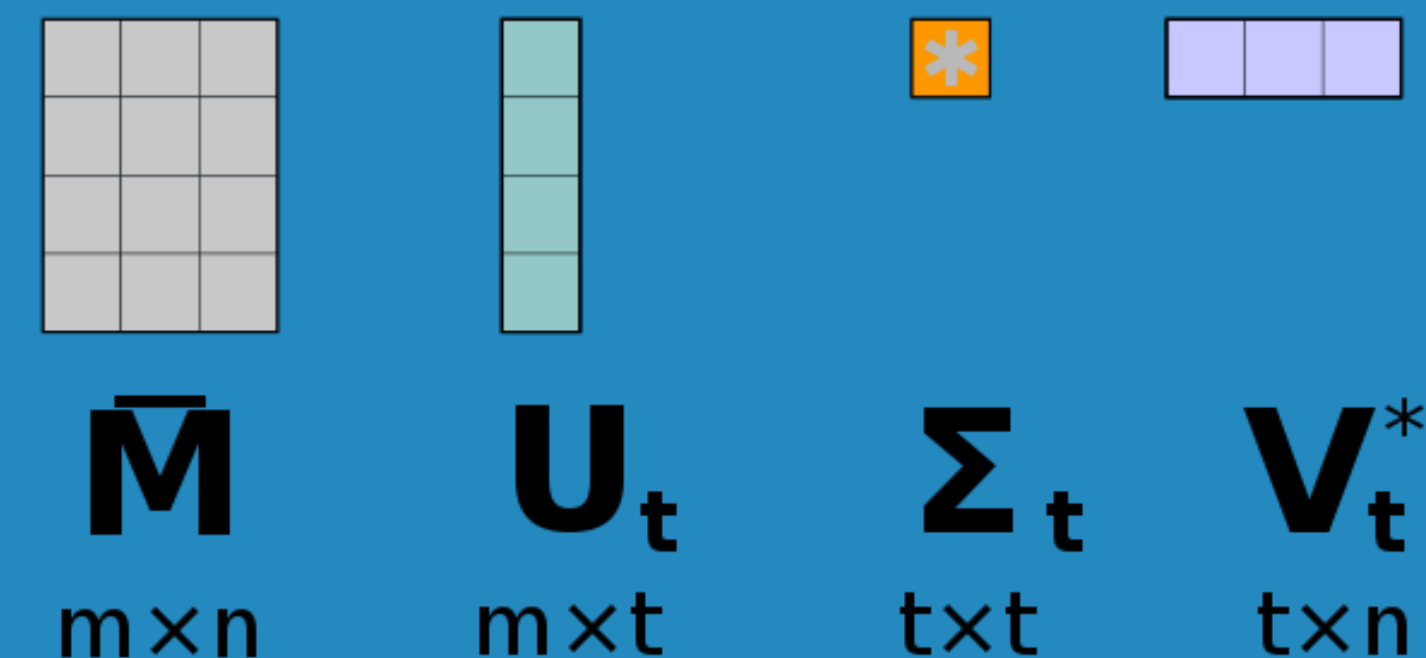
SVD vs Compact SVD vs Truncated SVD



SVD: U, V orthogonal,
 Σ rectangular diagonal matrix



Compact SVD: U, V semi-orthogonal,
 Σ square diagonal matrix with non-zero SVs



Truncated SVD: U, V semi-orthogonal,
 Σ square diagonal matrix with t largest SVs

Properties of the matrix \mathbf{K}

- Symmetric $K_{ij} = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_i) = K_{ji}$

- Positive semidefinite: for any vector \mathbf{a}

$$\mathbf{a}^\top \mathbf{K} \mathbf{a} = (\Phi^\top \mathbf{a})^\top (\Phi^\top \mathbf{a}) \geq 0$$

Connections between Φ , C , and K

Assuming $\Phi \in \mathbb{R}^{n \times M}$ is mean centered:

- Covariance matrix: $C = \Phi^\top \Phi \in \mathbb{R}^{M \times M}$
- Kernel matrix: $K = \Phi \Phi^\top \in \mathbb{R}^{n \times n}$
- Rank (number of non-zero singular values) of all three matrices are the same

Why? Let's consider SVD

- $\Phi = U \Sigma^{1/2} V^\top$
- $C = \Phi^\top \Phi = V \Sigma^{1/2} U^\top U \Sigma^{1/2} V^\top = V \Sigma V^\top$
- $K = \Phi \Phi^\top = U \Sigma^{1/2} V^\top V \Sigma^{1/2} U^\top = U \Sigma U^\top$
- C and K have the same non-zero singular values

Note: variants of the SVD do bookkeeping differently for zero-valued singular values (see HW2 for details on compact SVD)

Sanity Check: How does this relate to PCA?

SVD implications

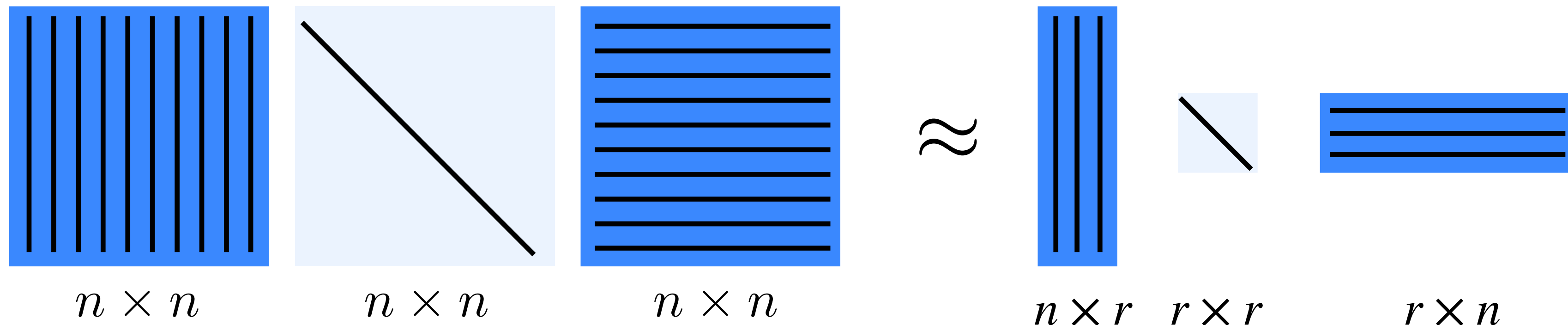
- $\Phi = U\Sigma^{1/2}V^T$
- $C = \Phi^T\Phi = V\Sigma^{1/2}U^T U\Sigma^{1/2}V^T = V\Sigma V^T$
- $K = \Phi\Phi^T = U\Sigma^{1/2}V^T V\Sigma^{1/2}U^T = U\Sigma U^T$

PCA implications

- Top r PCA scores (via C): $Z = \Phi V_r \in \mathbb{R}^{n \times r}$
 - V_r are top r eigenvectors of C
- Top r PCA scores (via kernel matrix):
$$Z = \Phi V_r = U\Sigma^{1/2}V^T V_r = U_r \Sigma_r^{1/2}$$

Low-rank Approximation

- $\mathbf{K} = \Phi\Phi^\top \in \mathbb{R}^{n \times n}$ is an SPSP matrix
- $\mathbf{K} = \mathbf{U}\Sigma\mathbf{U}^\top$: Singular value decomposition (SVD)
- $\mathbf{K}_r = \mathbf{U}_r\Sigma_r\mathbf{U}_r^\top$: ‘optimal’ low-rank matrix ($r \ll n, M$)
- $\mathbf{K}_r = \min_{rank(\tilde{\mathbf{K}})=r} \|\mathbf{K} - \tilde{\mathbf{K}}\|_F$



Low-rank Approximation

- $\mathbf{K} = \Phi\Phi^\top \in \mathbb{R}^{n \times n}$ is an SPSSD matrix
- $\mathbf{K} = \mathbf{U}\Sigma\mathbf{U}^\top$: Singular value decomposition (SVD)
- $\mathbf{K}_r = \mathbf{U}_r\Sigma_r\mathbf{U}_r^\top$: ‘optimal’ low-rank matrix ($r \ll n, M$)

How does factorized \mathbf{K}_r help speed up kernel methods?

- KRR: substituting it into $\alpha = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$ allows for linear time computation (instead of cubic) via “Woodbury Inversion Lemma”
- PCA: it directly gives us the PCA solution (see previous slide)

Low-rank Approximation

- $\mathbf{K} = \Phi\Phi^\top \in \mathbb{R}^{n \times n}$ is an SPSD matrix
- $\mathbf{K} = \mathbf{U}\Sigma\mathbf{U}^\top$: Singular value decomposition (SVD)
- $\mathbf{K}_r = \mathbf{U}_r\Sigma_r\mathbf{U}_r^\top$: ‘optimal’ low-rank matrix ($r \ll n, M$)

Goal: Fast (linear in n), Approximate SVD for large-scale kernel methods

- Minimize reconstruction error $\|\mathbf{K} - \tilde{\mathbf{K}}\|_F$
- Achieve $\mathcal{L}(\mathbf{K}) \approx \mathcal{L}(\tilde{\mathbf{K}})$ for learning algorithm $\mathcal{L}(\cdot)$

Low-rank Approximation

- $\mathbf{K} = \Phi\Phi^\top \in \mathbb{R}^{n \times n}$ is an SPSP matrix
- $\mathbf{K} = \mathbf{U}\Sigma\mathbf{U}^\top$: Singular value decomposition (SVD)
- $\mathbf{K}_r = \mathbf{U}_r\Sigma_r\mathbf{U}_r^\top$: ‘optimal’ low-rank matrix ($r \ll n, M$)

Can we just compute \mathbf{K}_r and call it a day?

- No, SVD is expensive! — $O(n^2r)$ for top r singular vectors/values
- Note that the distributed techniques from last lecture don’t apply b/c we have access to \mathbf{K} but not Φ (and even storing \mathbf{K} for very large n may be infeasible)

Outline

1. Background: Kernel Methods
2. Kernel Ridge Regression
3. SVD
4. Large-Scale Kernel Methods (via Approximate SVD)

Key Question: Can we approximate the SVD w/o looking at all the entries in the matrix?

What approximation should we use?

What entries should we look at?

Key Question: Can we approximate the SVD w/o looking at all the entries in the matrix?

What approximation should we use?

What entries should we look at?

Column-sampling Method

[Frieze et al., '98]

$$\mathbf{K} = \left[\right]_{n \times n}$$

Column-sampling Method

[Frieze et al., '98]

- Algorithm
 - Sample r columns of \mathbf{K}
 - SVD of \mathbf{C} : $O(nr^2)$

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{C} \end{array} \right]_{n \times n}$$

The diagram illustrates the column-sampling method. It shows a large matrix \mathbf{K} of size $n \times n$. A submatrix \mathbf{C} is highlighted within \mathbf{K} , representing the sampled columns. The width of \mathbf{C} is labeled r (number of sampled columns), and its height is labeled n (number of rows). The overall dimensions of \mathbf{K} are indicated as $n \times n$ at the bottom right.

Column-sampling Method

[Frieze et al., '98]

- Algorithm

- Sample r columns of \mathbf{K}
- SVD of \mathbf{C} : $O(nr^2)$

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{C} \end{array} \right]_{n \times n}$$

The diagram shows a large square matrix \mathbf{K} of size $n \times n$. Inside \mathbf{K} , a smaller rectangle represents the matrix \mathbf{C} . The width of \mathbf{C} is labeled r (in blue) above the rectangle, and its height is labeled n (in blue) to the right of the rectangle.

- Approximation: $\tilde{\mathbf{K}}_{col} = \tilde{\mathbf{U}}_{col} \tilde{\Sigma}_{col} \tilde{\mathbf{U}}_{col}^T$

- $\tilde{\Sigma}_{col} = \sqrt{\frac{n}{r}} \Sigma^C$
- $\tilde{\mathbf{U}}_{col} = \mathbf{U}^C$

Intuition: SVD on \mathbf{C} provides good approximation of full SVD

Nystrom Method

[Williams & Seeger, '00]

[Drineas & Mahoney, '05]

$$\mathbf{K} = \begin{bmatrix} & \end{bmatrix}_{n \times n}$$

Nystrom Method

[Williams & Seeger, '00]

[Drineas & Mahoney, '05]

- Algorithm
 - Sample r columns of \mathbf{K}
 - Sample the same rows of \mathbf{C}
 - SVD / inverse of $\mathbf{W} : O(r^3)$

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{C} \end{array} \right]_{n \times n}$$

Diagram illustrating the Nystrom Method. A large matrix \mathbf{K} of size $n \times n$ is shown. A blue box highlights a submatrix \mathbf{C} of size $r \times n$, representing the first r columns of \mathbf{K} .

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{W} \end{array} \right]_{n \times n}$$

Diagram illustrating the Nystrom Method. A large matrix \mathbf{K} of size $n \times n$ is shown. A blue box highlights a submatrix \mathbf{W} of size $r \times r$, representing the first r rows of the first r columns of \mathbf{K} .

Nystrom Method

[Williams & Seeger, '00]

[Drineas & Mahoney, '05]

- Algorithm
 - Sample r columns of \mathbf{K}
 - Sample the same rows of \mathbf{C}
 - SVD / inverse of $\mathbf{W} : O(r^3)$
 - Extrapolate: $O(nr^2)$
- Approximation: $\tilde{\mathbf{K}}_{nys} = \mathbf{C}\mathbf{W}^{-1}\mathbf{C}^\top$
 - What's going on here?

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{C} \\ \hline \end{array} \right]_{n \times n}$$

The diagram shows a large square matrix \mathbf{K} of size $n \times n$. A blue rectangular box highlights a submatrix \mathbf{C} in the top-left corner. The width of this box is labeled with a blue r , and its height is labeled with a blue n .

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{W} \\ \hline \end{array} \right]_{n \times n}$$

The diagram shows a large square matrix \mathbf{K} of size $n \times n$. A blue rectangular box highlights a submatrix \mathbf{W} in the top-left corner. The width of this box is labeled with a blue r , and its height is labeled with a blue r .

Nystrom Method

[Williams & Seeger, '00]

[Drineas & Mahoney, '05]

$$\tilde{\mathbf{K}}_{nys} = \mathbf{C}\mathbf{W}^{-1}\mathbf{C}^{\top}$$

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{C} \end{array} \right]_{n \times n}$$

r

$$\mathbf{K} = \left[\begin{array}{cc} \mathbf{W} & \mathbf{K}_{21}^{\top} \\ \mathbf{K}_{21} & \mathbf{K}_{22} \end{array} \right]_{n \times n}$$

r

Nystrom Method

[Williams & Seeger, '00]

[Drineas & Mahoney, '05]

$$\tilde{\mathbf{K}}_{nys} = \mathbf{C} \mathbf{W}^{-1} \mathbf{C}^{\top}$$

$$= \begin{matrix} & \overset{n \times r}{\mathbf{W}} & \overset{r \times r}{\mathbf{W}^{-1}} \left[\overset{r \times n}{\mathbf{W}} \quad \mathbf{K}_{21}^{\top} \right] \\ = & \left[\begin{array}{c} \mathbf{W} \\ \mathbf{K}_{21} \end{array} \right] \end{matrix}$$

$$\mathbf{K} = \left[\begin{array}{c} \mathbf{C} \end{array} \right]_{n \times n}$$

r (width of C), *n* (height of C)

$$\mathbf{K} = \left[\begin{array}{cc} \mathbf{W} & \mathbf{K}_{21}^{\top} \\ \mathbf{K}_{21} & \mathbf{K}_{22} \end{array} \right]_{n \times n}$$

r (width of W), *r* (height of W), *K*₂₂ is red

Nystrom Method

[Williams & Seeger, '00]

[Drineas & Mahoney, '05]

$$\begin{aligned}
 \tilde{\mathbf{K}}_{nys} &= \mathbf{C} \mathbf{W}^{-1} \mathbf{C}^T \\
 &= \begin{matrix} & \begin{matrix} n \times r & r \times r & r \times n \end{matrix} \\ \begin{bmatrix} \mathbf{W} \\ \mathbf{K}_{21} \end{bmatrix} & \mathbf{W}^{-1} \begin{bmatrix} \mathbf{W} & \mathbf{K}_{21}^T \end{bmatrix} \end{matrix} \\
 &= \begin{bmatrix} \mathbf{W} & \mathbf{K}_{21}^T \\ \mathbf{K}_{21} & \mathbf{K}_{21} \mathbf{W}^{-1} \mathbf{K}_{21}^T \end{bmatrix}_{n \times n}
 \end{aligned}$$

$$\mathbf{K} = \begin{bmatrix} \mathbf{C} \end{bmatrix}_{n \times n}$$

(A blue box highlights the top r rows of \mathbf{C} , and the label n is to the right of the box.)

$$\mathbf{K} = \begin{bmatrix} \mathbf{W} & \mathbf{K}_{21}^T \\ \mathbf{K}_{21} & \mathbf{K}_{22} \end{bmatrix}_{n \times n}$$

(A blue box highlights the top-left $r \times r$ block \mathbf{W} , and the label r is to the left of the box.)

note: $\mathbf{K}_{22} - \mathbf{K}_{21} \mathbf{W}^{-1} \mathbf{K}_{21}^T$ is a well studied quantity called the “Schur complement”

How Do They Compare?

Column Sampling

- Estimate SVD via subset of columns
- $\tilde{\mathbf{K}}_{col} = \mathbf{C}((r/n)\mathbf{C}^\top\mathbf{C})^{-1/2}\mathbf{C}^\top$

$$\mathbf{K} = \left[\begin{array}{c} \boxed{\mathbf{C}} \\ \end{array} \right]_{n \times n}$$

The diagram shows a large square matrix \mathbf{K} of size $n \times n$. Inside the matrix, a vertical rectangular sub-matrix \mathbf{C} is highlighted with a blue border. The width of this sub-matrix is labeled r in blue above the box, and the height is labeled n in blue to the right of the box.

Nystrom Method

- Estimate SVD via smaller sub-matrix
- $\tilde{\mathbf{K}}_{nys} = \mathbf{C}\mathbf{W}^{-1}\mathbf{C}^\top$

$$\mathbf{K} = \left[\begin{array}{c} \boxed{\mathbf{W}} \\ \end{array} \right]_{n \times n}$$

The diagram shows a large square matrix \mathbf{K} of size $n \times n$. Inside the matrix, a small square sub-matrix \mathbf{W} is highlighted with a blue border. The width and height of this sub-matrix are both labeled r in blue above and to the right of the box, respectively.

How Do They Compare?

Assume $\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{K})$

How Do They Compare?

Assume $\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{K})$

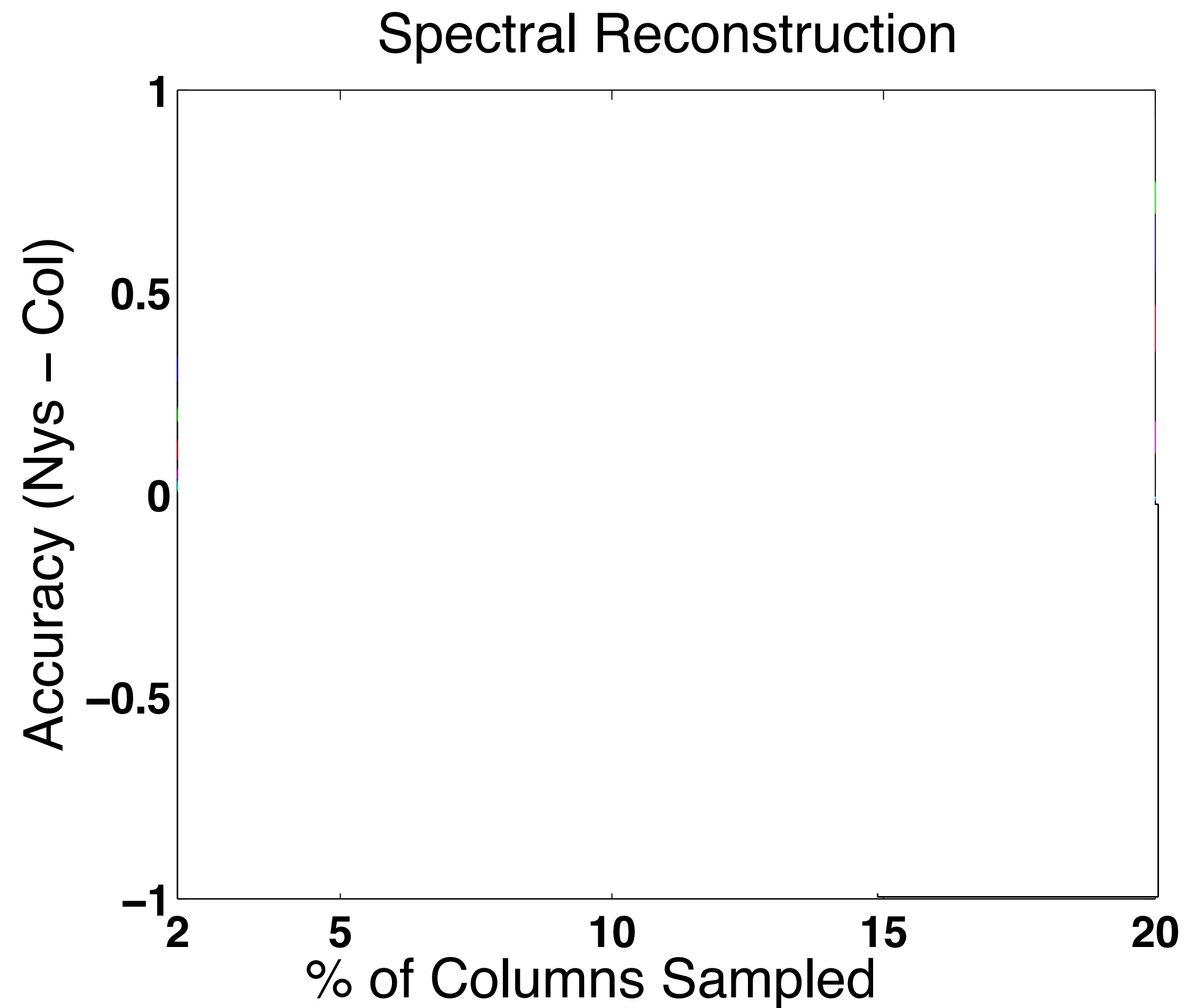
- *Nyström is exact*
- *Column-sampling exact iff it trivially reduces to Nyström, i.e., $\mathbf{W} = ((r/n)\mathbf{C}^\top \mathbf{C})^{1/2}$*

How Do They Compare?

Assume $\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{K})$

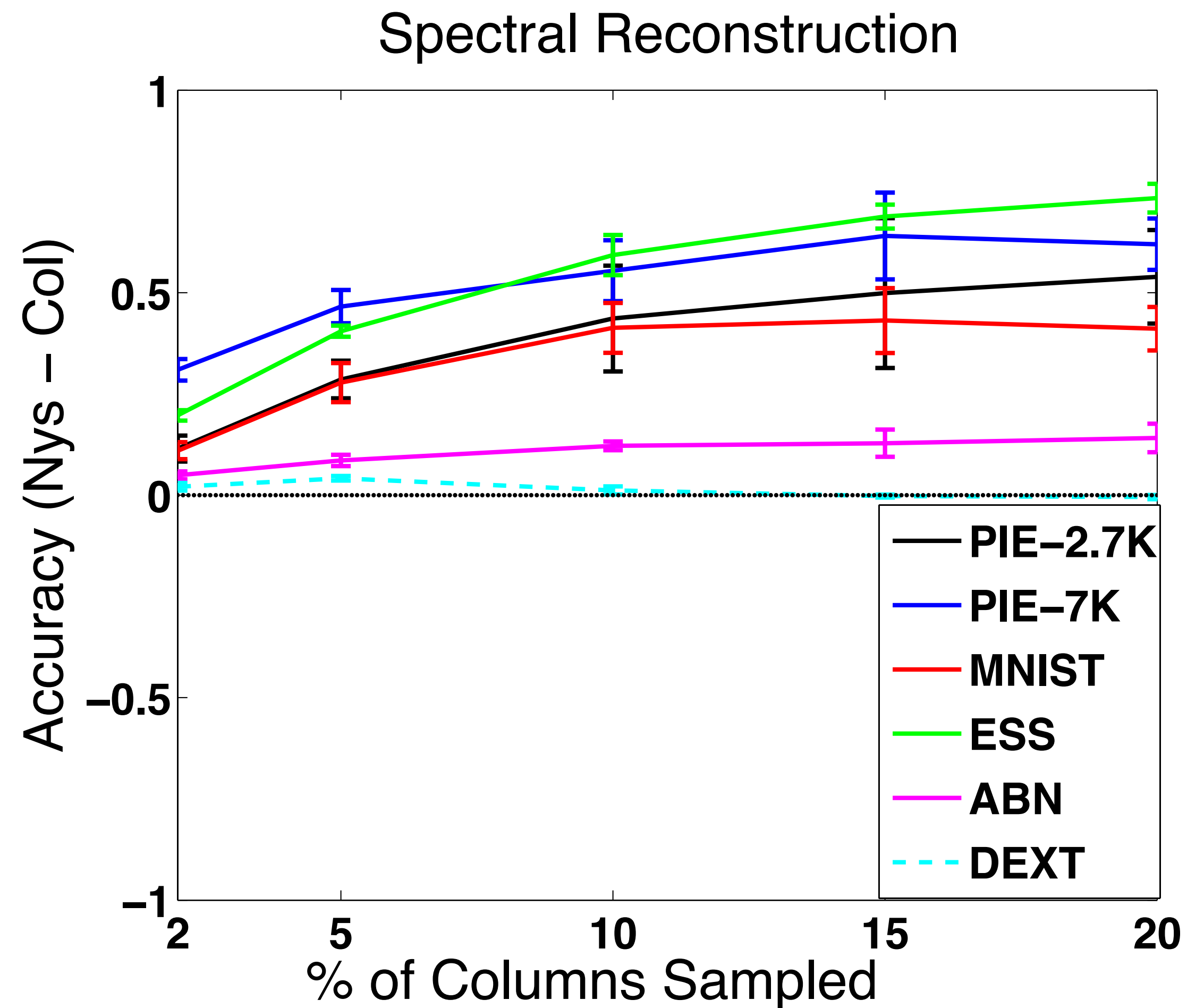
- *Nyström is exact*
 - *Column-sampling exact iff it trivially reduces to Nyström, i.e., $\mathbf{W} = ((r/n)\mathbf{C}^\top \mathbf{C})^{1/2}$*
- Under certain conditions, Nyström is optimal, Column-sampling is suboptimal
 - Suggests that if \mathbf{K} has low-rank structure ($\mathbf{K} \approx \mathbf{K}_r$), then Nyström will be better

How Do They Compare?



- Accuracy defined as $\frac{\| \mathbf{K} - \mathbf{K}_r \|_F}{\| \mathbf{K} - \tilde{\mathbf{K}} \|_F}$
- Above 0 signifies Nyström is better

How Do They Compare?



- Accuracy defined as $\frac{\| \mathbf{K} - \mathbf{K}_r \|_F}{\| \mathbf{K} - \tilde{\mathbf{K}} \|_F}$
- Above 0 signifies Nyström is better
- Nyström better on all 5 dense matrices with quickly decaying eigenvalues

- Suggests that if \mathbf{K} has low-rank structure ($\mathbf{K} \approx \mathbf{K}_r$), then Nyström will be better

Key Question: Can we approximate the SVD w/o looking at all the entries in the matrix?

What approximation should we use?

What entries should we look at?

Limitations of Random Sampling

$$\mathbf{K} = \begin{bmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & 1 & & & & \\ & & & 0 & & & \\ & & & & 0 & & \\ & & & & & 0 & \\ & & & & & & 0 \\ & 0 & & & & & & 0 \end{bmatrix}_{n \times n}$$

Optimal Choice

Leverage scores

- function of exact eigenvectors
- Very expensive: requires computing SVD!

Interesting theoretically, but not for practical, large-scale settings

How to sample in practice?

- Distribution over the columns remains **fixed**
 - uniform: $O(1)$
 - diagonal: $O(n)$
 - column-norm: $O(n^2)$
- **Adaptively** pick columns or perform pre-processing
 - eg. via K -means [Zhang et al., '08]
- Combine approximations (**ensemble algorithm**)

Take-aways

- **Kernel methods** are a popular approach for encoding complex features
- But, kernel methods don't scale well with n
- A natural approach is **low-rank approximations** of the kernel matrix, K
- **SVD** is a powerful tool
- But, calculating the optimal low-rank approximation (eg., via SVD) is expensive:
 - **Column sampling** or **Nystrom method** can make this feasible

Notes

- Random Fourier Features: another popular alternative for fast, approximate kernel methods, focusing on approximating $\phi(\cdot)$
- Theory connecting reconstruction error with generalization error
 - [On the Impact of Kernel Approximation on Learning Accuracy](#). Cortes, Mohri, Talwalkar. AISTATS 2010.
 - [Sharp analysis of low-rank kernel matrix approximations](#). Bach. COLT, 2013.
- Other References
 - [Randomized algorithms for matrices and data](#). Mahoney. Foundations & Trends in Machine Learning, 2011.
 - [Random Features for Large-Scale Kernel Machines](#). Rahimi and Recht. NeurIPS 2007.
 - [Matrix Approximation for Large-scale Learning](#). Talwalkar. PhD Thesis. 2010.