

RECURSIVE SAMPLING FOR THE NYSTRÖM METHOD

Cameron Musco, Christopher Musco ♦ Massachusetts Institute of Technology

The Kernel Method

Adapt any **linear** data analysis method (regression, principal component analysis, support vector machines, etc.) to work with **nonlinear** similarity function

How?

These methods only depend on inner product information in the **Gram matrix**:

$$K_{i,j} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

(n data vectors, $\mathbf{x}_i = \{x_i^1, x_i^2, \dots, x_i^d\}$)

Replace Gram matrix and use as usual*

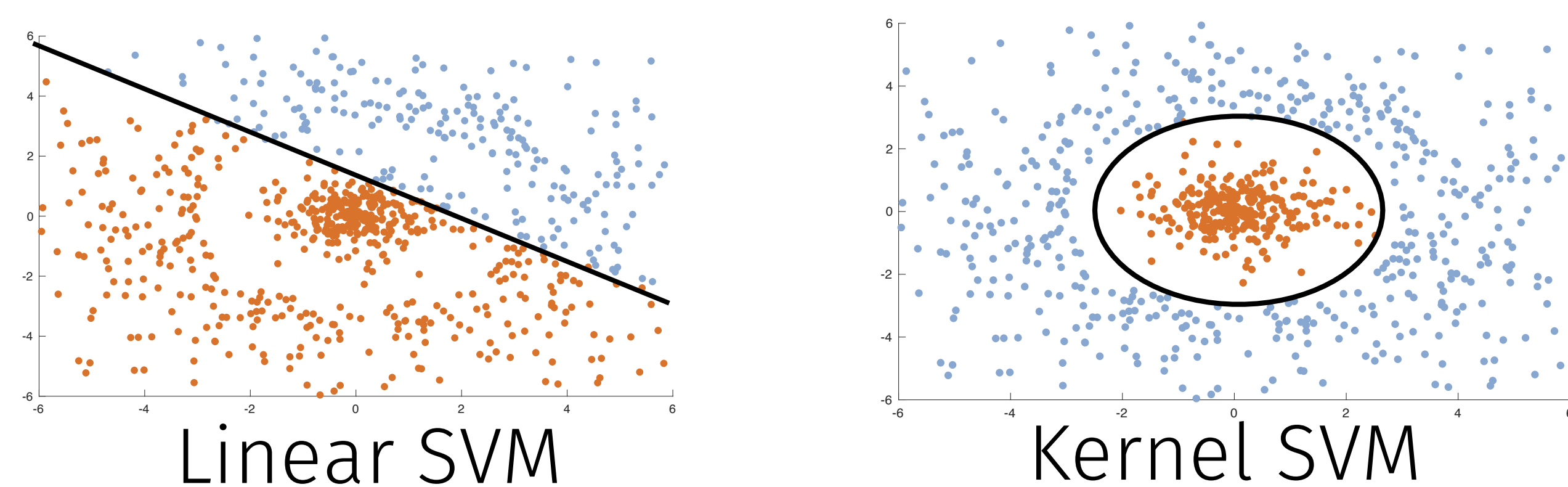
$$K_{i,j} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

linear inner product

*kernel function needs to be PSD

- $K_{i,j} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2)$
Gaussian kernel inner product
- $K_{i,j} = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^q$
Polynomial kernel inner product
- $K_{i,j} = \dots$
Custom nonlinear inner product

Pro: Learn nonlinear function classes

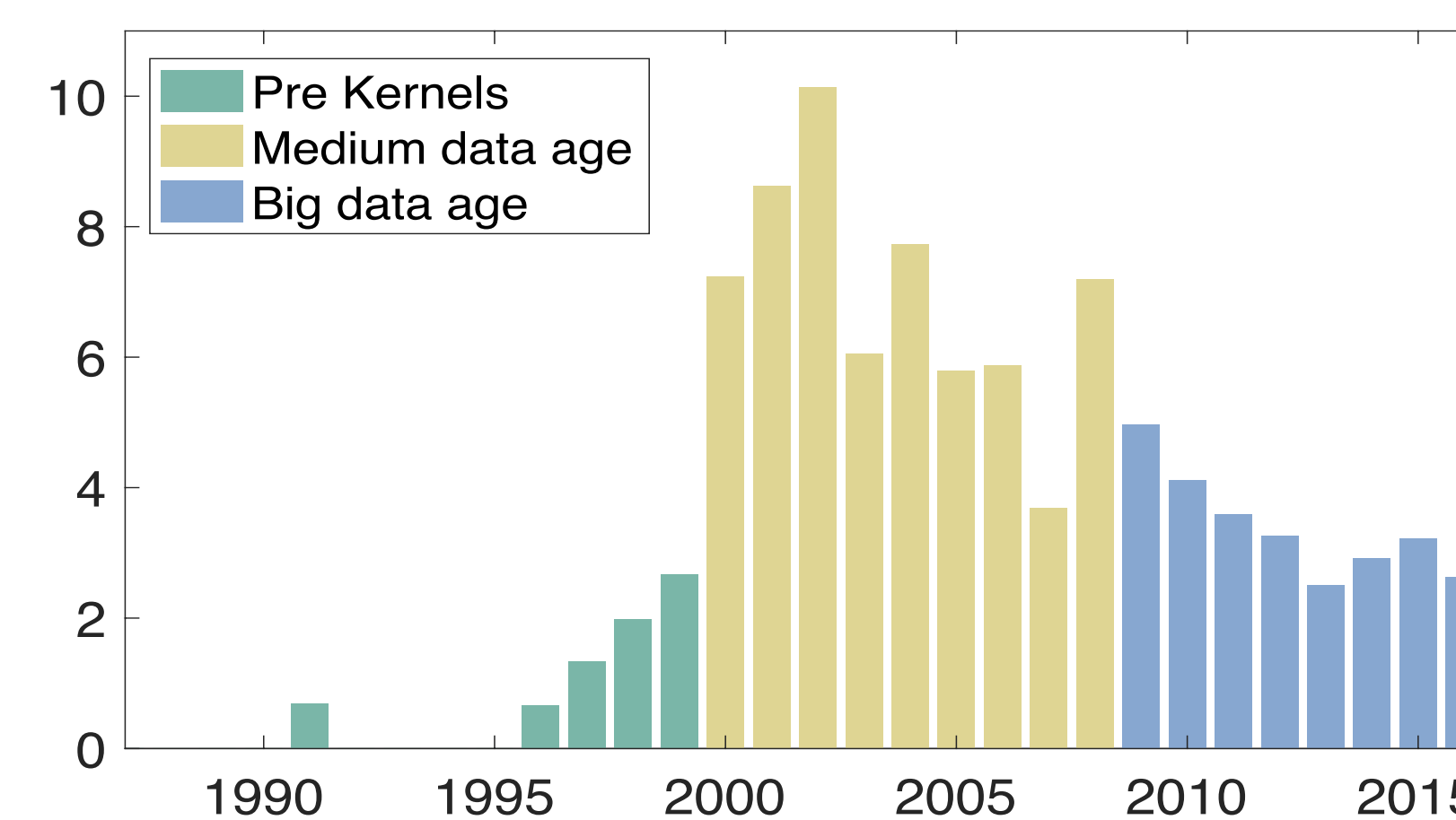


Theoretically sound, effective in practice.

Con: Limited Scalability

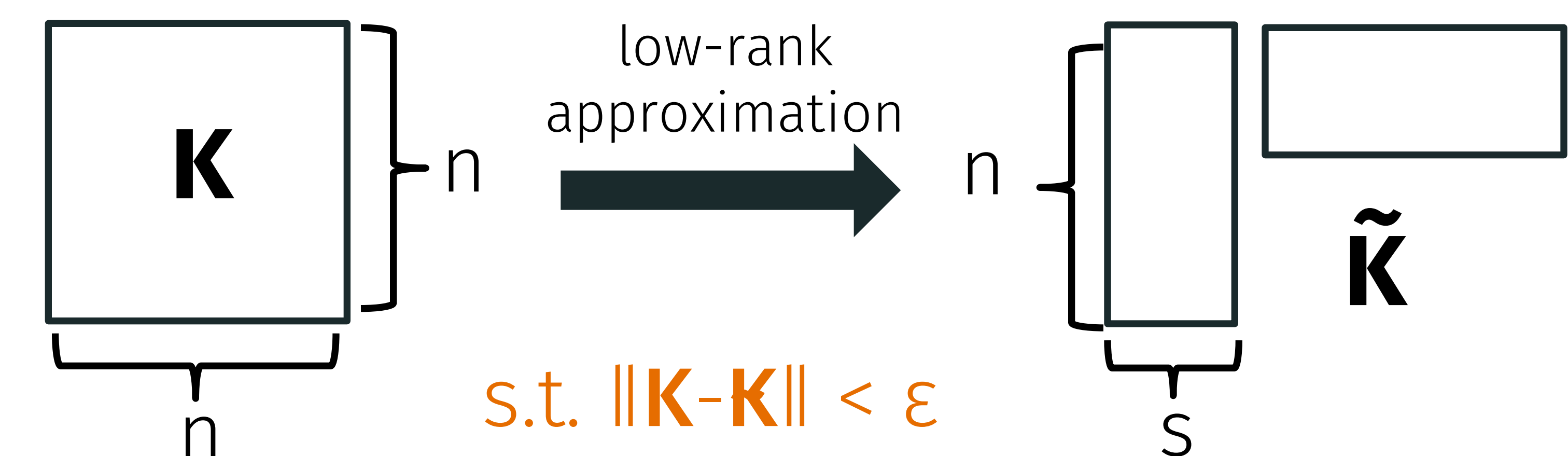
Just writing down \mathbf{K} takes $O(n^2)$ time!

% of NIPS titles containing "kernel"



Kernel methods can handle high dimensional data, but not large training sets.

Can we use an approximation to \mathbf{K} ?

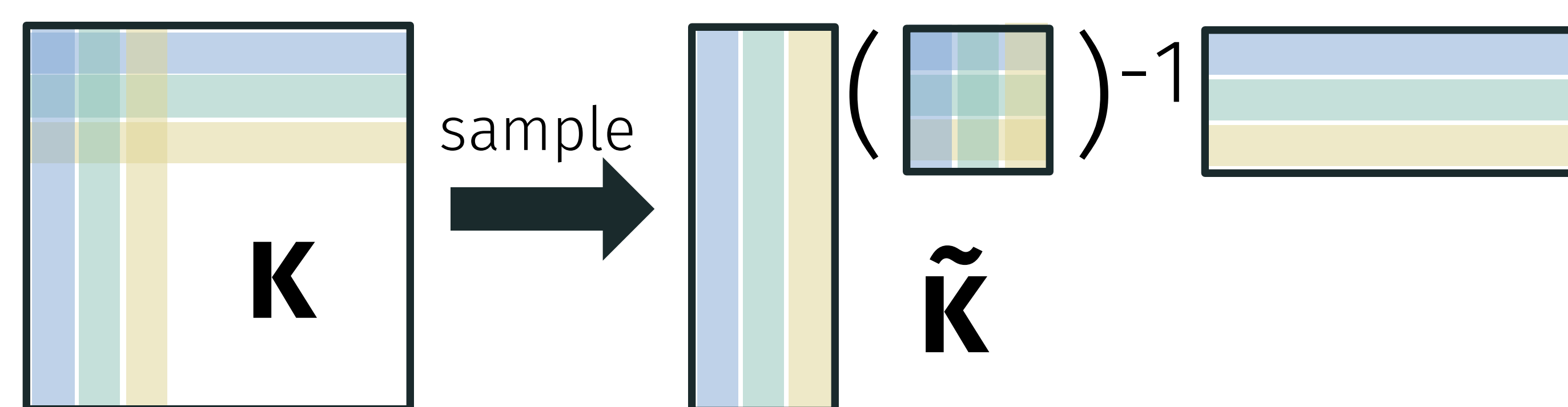


Rank s approximation $\tilde{\mathbf{K}}$ can be stored in $O(ns)$ space, $\tilde{\mathbf{K}}^{-1}$ can be computed in $O(ns^2)$ time, eigendecomposition in $O(ns^2)$ time, etc.

Kernel Approximations	Speed	Accuracy
Incomplete Cholesky / explicit low-rank approximation	$O(n^3)$	High
Random Sketching	$O(n^2)$	High
Random Fourier features	$O(n)$	Variable
Standard Nyström method	$O(n)$	Variable
Our Recursive Nyström	$O(n)$	High

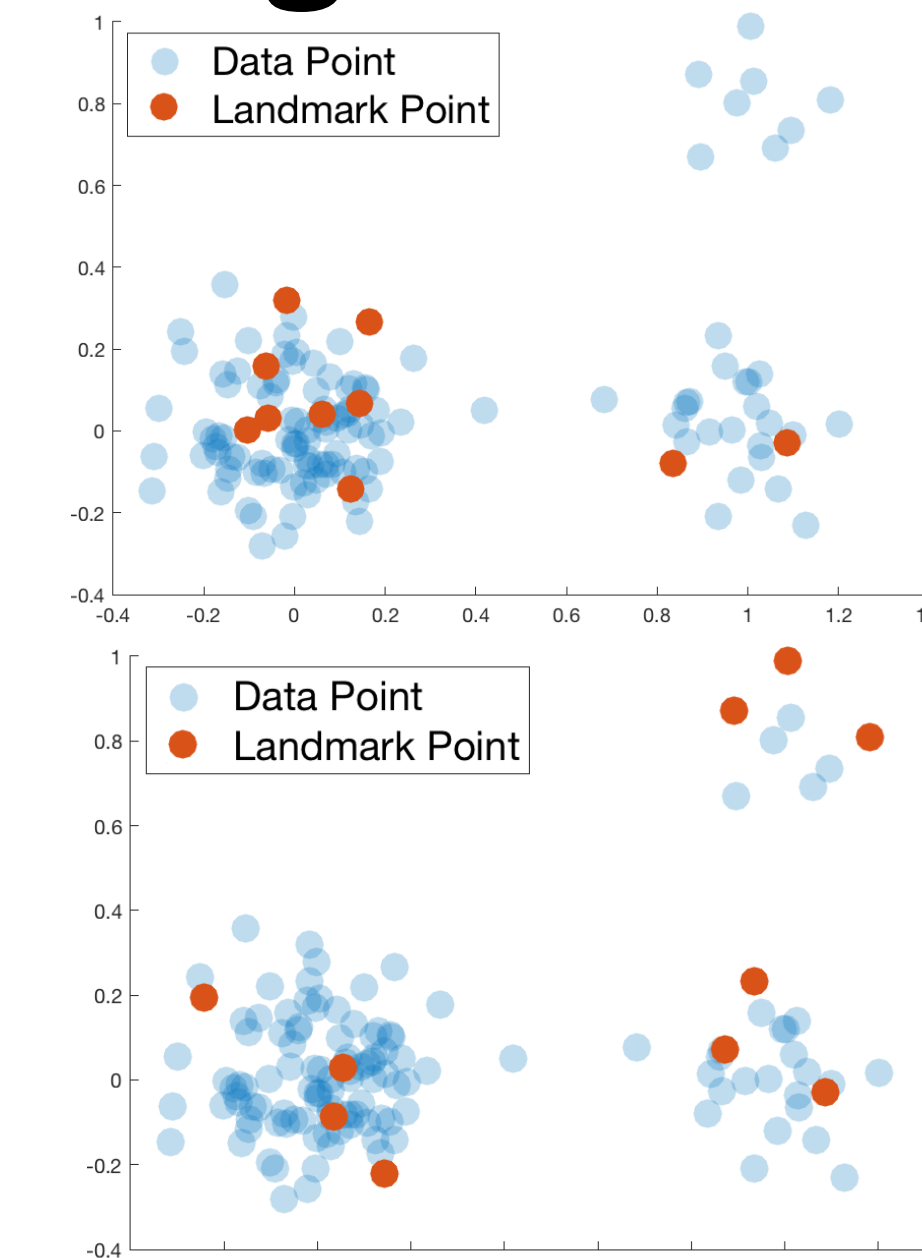
Nyström Method

Low-rank approximation from random sample of the "landmark" data points



Time linear in n. Does not require all of \mathbf{K} !

High accuracy requires better landmarks



Uniformly Random (standard Nyström)
Better Landmarks (importance sampled)

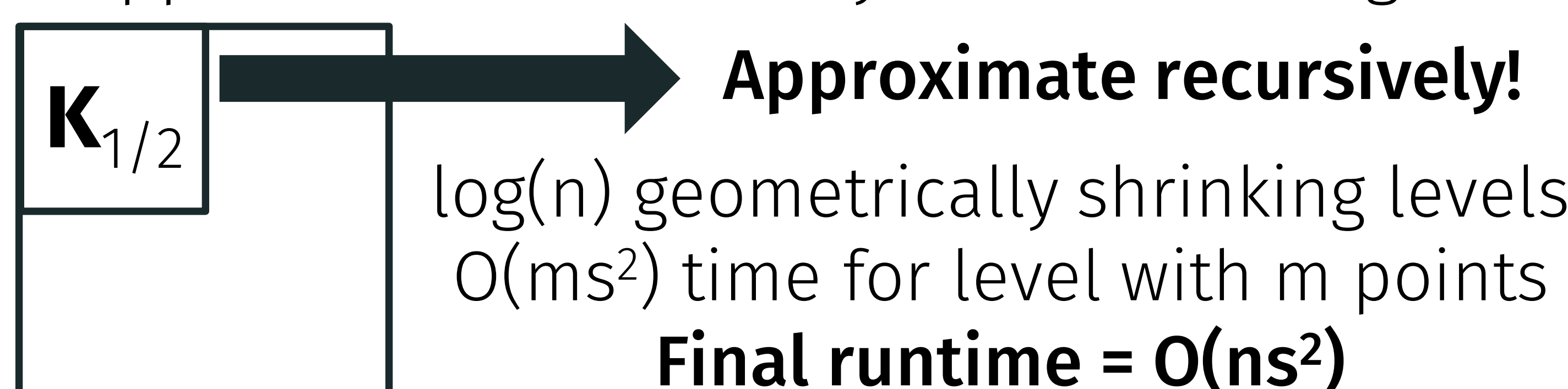
The Nyström method is like "triangulation with noise". For data in s dimensions, we need s landmark points to determine all distances in \mathbf{K} . If data nearly lies in s dimensions we need $O(s)$ "well conditioned" points spread throughout data.

Fast "leverage score" sampling

All good importance sampling probabilities require a good approximation to \mathbf{K} to compute!

$$\text{leverage score}(\mathbf{x}_i) = \mathbf{k}_i^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{k}_i$$

Main technique: Uniform sampling gives good approximation with many landmarks – e.g. $n/2$.

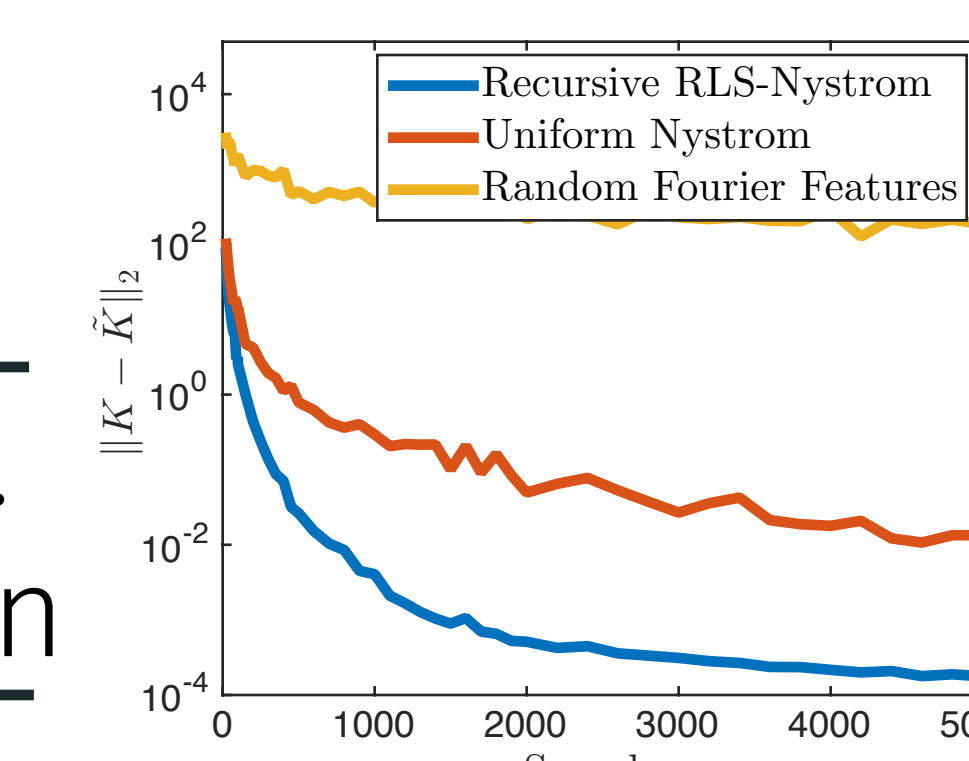


Strongest theoretical guarantees for approximate kernel learning

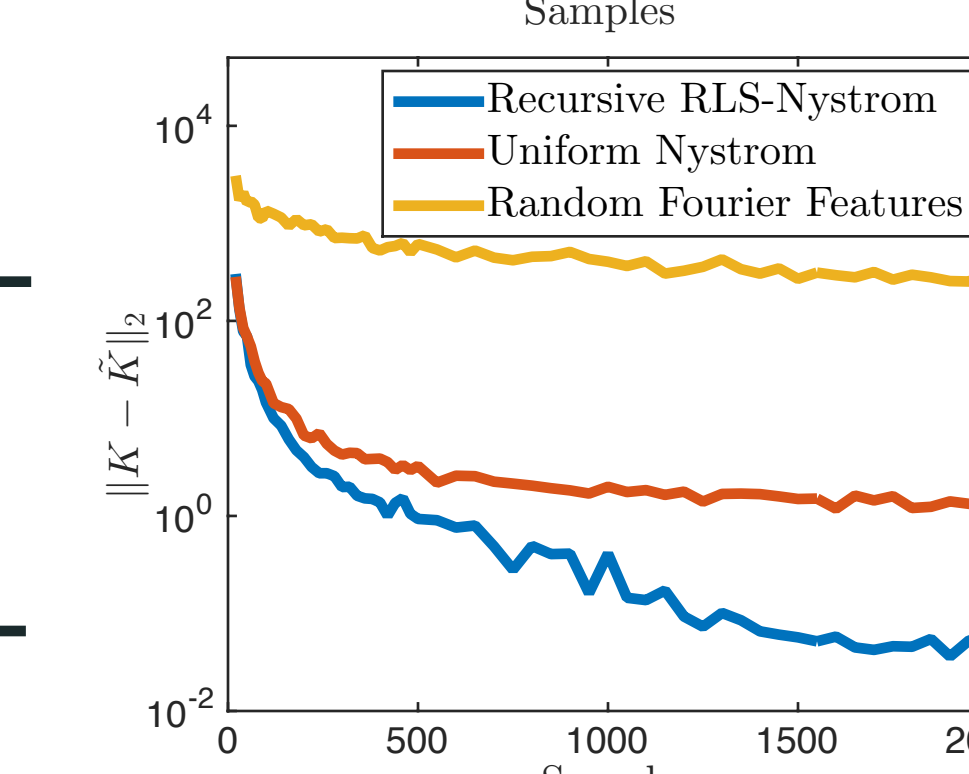
# of samples	guarantee
$O(\text{statistical dimension})$	relative error approx. kernel ridge regression
$O(k/\epsilon)$	$(1+\epsilon)$ error rank k kernel PCA
$O(k/\epsilon)$	$(1+\epsilon)$ error kernel k-means clustering

+ general theorems to apply to other kernel problems

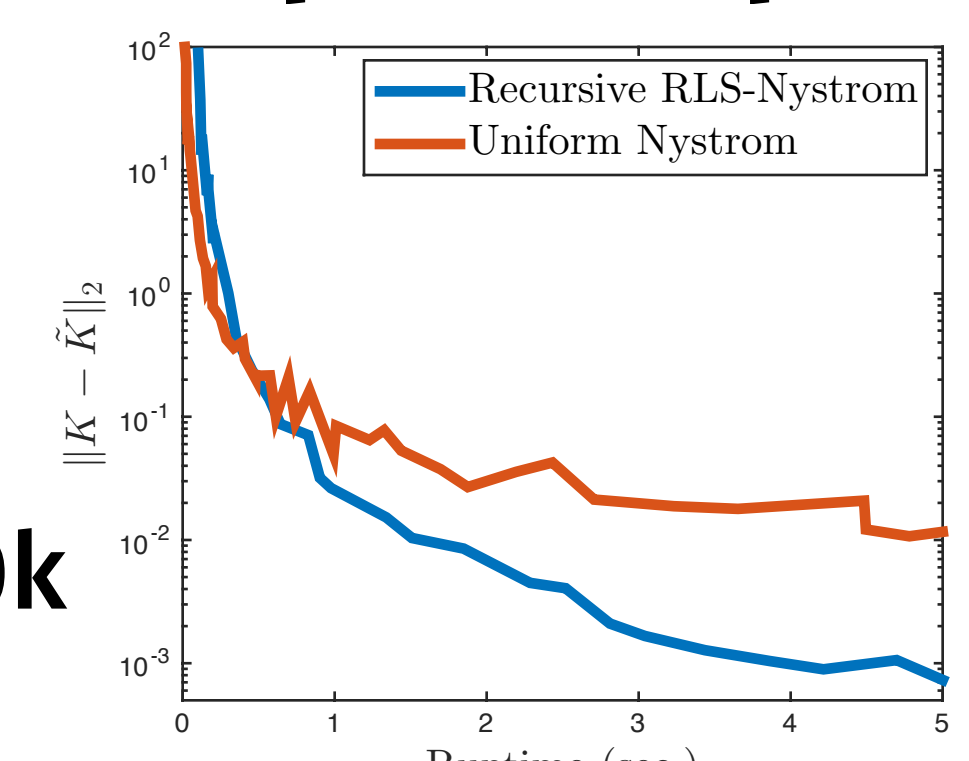
State-of-the-art empirical performance



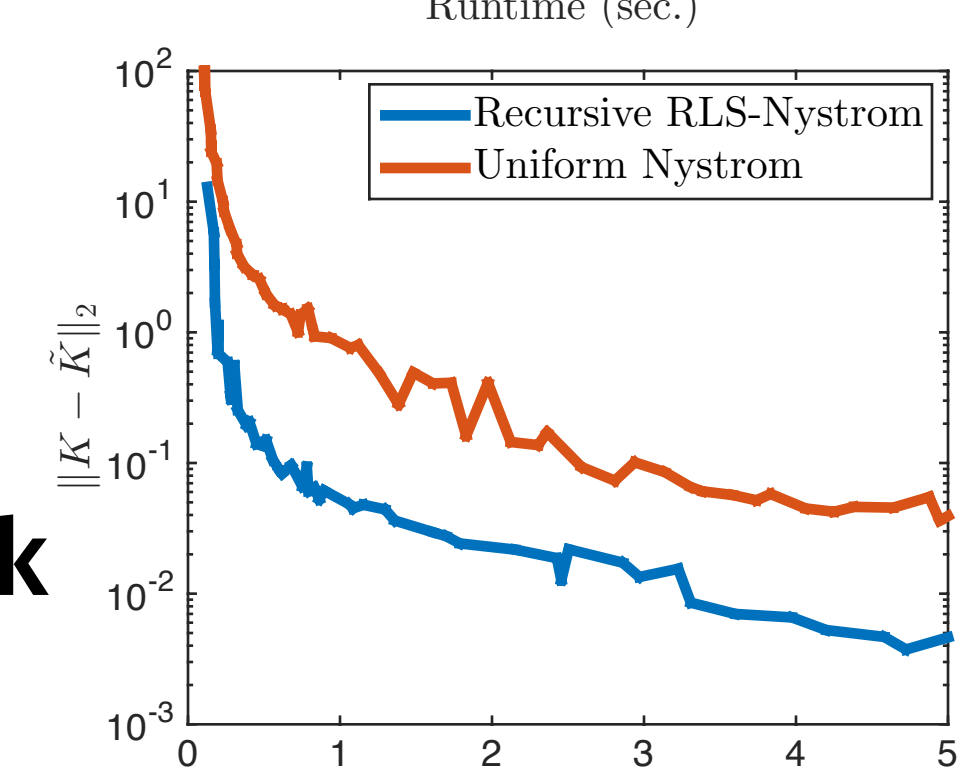
Forest cover data, n ~ 600k



RNA data, n ~ 350k



Many more experiments available in the paper.



Simple MATLAB code available at chrismusco.com

No tuning required! Works for any kernel.

BACKGROUND

THE PROBLEM

OUR APPROACH

RESULTS