

RANDOM FOURIER FEATURES  
FOR KERNEL RIDGE REGRESSION:  
APPROXIMATION BOUNDS AND STATISTICAL GUARANTEES

---

Haim Avron, Michael Kapralov, **Cameron Musco**, Christopher Musco, Ameya Velingker, and Amir Zandieh

Tel Aviv University, EPFL, and MIT.  
ICML 2017.

## **Our Contributions:**

## Our Contributions:

- Analyze the **random Fourier features** method (Rahimi Recht '07) for kernel approximation using **leverage score-based** techniques.

## Our Contributions:

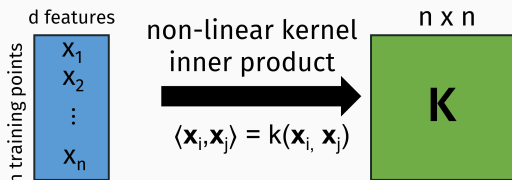
- Analyze the **random Fourier features** method (Rahimi Recht '07) for kernel approximation using **leverage score-based** techniques.
- **Concrete:** Introduce new sampling distribution that gives statistical guarantees for kernel ridge regression when used to approximate the Gaussian kernel.

## Our Contributions:

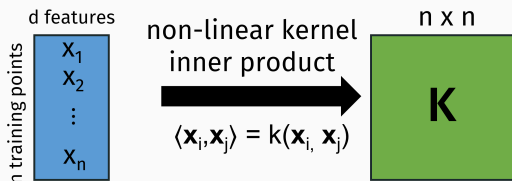
- Analyze the **random Fourier features** method (Rahimi Recht '07) for kernel approximation using **leverage score-based** techniques.
- **Concrete:** Introduce new sampling distribution that gives statistical guarantees for kernel ridge regression when used to approximate the Gaussian kernel.
- **High Level:** Hope that **Fourier leverage scores** will have further applications in kernel approximation, function approximation, and sparse Fourier transform methods.

# KERNEL APPROXIMATION

Kernel methods are expensive.



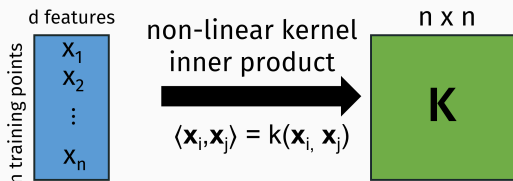
Kernel methods are expensive.



- Even writing down  $\mathbf{K}$  requires  $\Omega(n^2)$  time.

# KERNEL APPROXIMATION

Kernel methods are expensive.

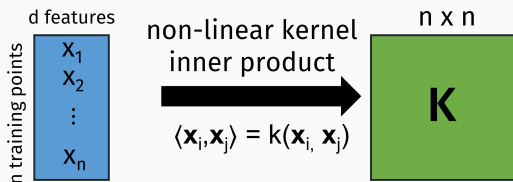


- Even writing down  $\mathbf{K}$  requires  $\Omega(n^2)$  time.
- Other operations require even more. A single iteration of a linear system solver takes  $\Omega(n^2)$  time.



# KERNEL APPROXIMATION

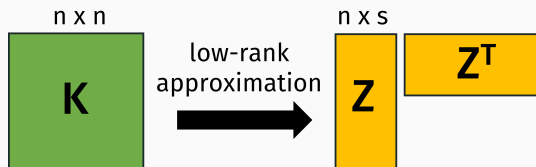
Kernel methods are expensive.



- Even writing down  $\mathbf{K}$  requires  $\Omega(n^2)$  time.
- Other operations require even more. A single iteration of a linear system solver takes  $\Omega(n^2)$  time.
- For  $n = 100,000$ ,  $\mathbf{K}$  has 10 billion entries. Takes 80 GB of storage if each is a double.

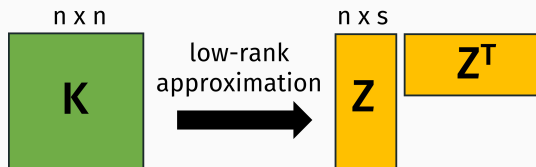
## SOLUTION: LOW-RANK APPROXIMATION.

Employ classic solution: low-rank approximation



## SOLUTION: LOW-RANK APPROXIMATION.

Employ classic solution: low-rank approximation



- Storing  $Z$  uses  $O(ns)$  space and computing  $ZZ^T x$  takes  $O(ns)$  time. Orthogonalization, eigendecomposition, and pseudo-inversion of  $ZZ^T$  all take just  $O(ns^2)$  time.

# EFFICIENT LOW-RANK APPROXIMATION?

Low-rank approximation is itself an expensive task.

# EFFICIENT LOW-RANK APPROXIMATION?

Low-rank approximation is itself an expensive task.

- Optimal low-rank approximation via a direct eigendecomposition, or even approximation via Krylov subspace methods are out of the question since they at least require fully forming  $\mathbf{K}$ .

Low-rank approximation is itself an expensive task.

- Optimal low-rank approximation via a direct eigendecomposition, or even approximation via Krylov subspace methods are out of the question since they at least require fully forming  $\mathbf{K}$ .
- Many faster methods have been studied: incomplete Cholesky factorization (Fine & Scheinberg '02, Bach & Jordan '02), entrywise sampling (Achlioptas, McSherry, & Schölkopf '01), Nyström approximation (Williams & Seeger '01), random Fourier features (Rahimi & Recht '07)

Low-rank approximation is itself an expensive task.

- Optimal low-rank approximation via a direct eigendecomposition, or even approximation via Krylov subspace methods are out of the question since they at least require fully forming  $\mathbf{K}$ .
- Many faster methods have been studied: incomplete Cholesky factorization (Fine & Scheinberg '02, Bach & Jordan '02), entrywise sampling (Achlioptas, McSherry, & Schölkopf '01), Nyström approximation (Williams & Seeger '01), **random Fourier features (Rahimi & Recht '07)**

Rahimi & Recht NIPS '07:



Rahimi & Recht NIPS '07:


- For any shift-invariant  $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i - \mathbf{x}_j)$  let  $p(\cdot)$  be the Fourier transform of  $k(\cdot)$ . By Bochner's theorem,  $p(\boldsymbol{\eta}) \geq 0$  for all  $\boldsymbol{\eta}$ .

Rahimi & Recht NIPS '07:

- For any shift-invariant  $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i - \mathbf{x}_j)$  let  $\rho(\cdot)$  be the Fourier transform of  $k(\cdot)$ . By Bochner's theorem,  $\rho(\boldsymbol{\eta}) \geq 0$  for all  $\boldsymbol{\eta}$ .
- Sample  $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_s \in \mathbb{R}^d$  with probabilities proportional to  $\rho(\boldsymbol{\eta})$ .

Rahimi & Recht NIPS '07:

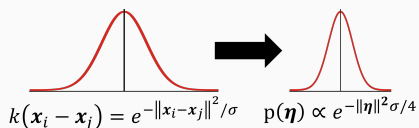
- For any shift-invariant  $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i - \mathbf{x}_j)$  let  $p(\cdot)$  be the Fourier transform of  $k(\cdot)$ . By Bochner's theorem,  $p(\boldsymbol{\eta}) \geq 0$  for all  $\boldsymbol{\eta}$ .
- Sample  $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_s \in \mathbb{R}^d$  with probabilities proportional to  $p(\boldsymbol{\eta})$ .



$k(\mathbf{x}_i - \mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / \sigma}$ 
 $\quad \rightarrow \quad$ 
 $p(\boldsymbol{\eta}) \propto e^{-\|\boldsymbol{\eta}\|^2 \sigma / 4}$

Rahimi & Recht NIPS '07:

- For any shift-invariant  $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i - \mathbf{x}_j)$  let  $p(\cdot)$  be the Fourier transform of  $k(\cdot)$ . By Bochner's theorem,  $p(\boldsymbol{\eta}) \geq 0$  for all  $\boldsymbol{\eta}$ .
- Sample  $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_s \in \mathbb{R}^d$  with probabilities proportional to  $p(\boldsymbol{\eta})$ .



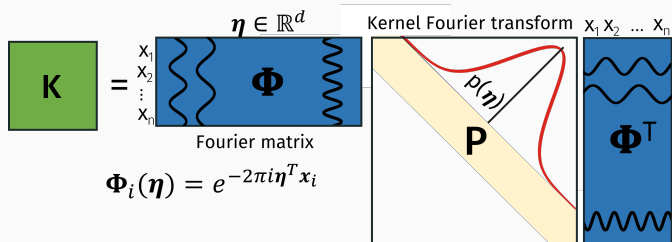
- Set  $\mathbf{z}_i = \frac{1}{\sqrt{s}} \left[ e^{-2\pi i \boldsymbol{\eta}_1^T \mathbf{x}_i}, \dots, e^{-2\pi i \boldsymbol{\eta}_s^T \mathbf{x}_i} \right]$ .
 

$\mathbf{Z}$	$\mathbf{Z}^T$	$\mathbf{z}_i$	$\approx$	$\mathbf{K}$
$\mathbf{z}_i$				

Fourier transform  $k(\mathbf{z}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-2\pi i \boldsymbol{\eta}^T \mathbf{z}} d\boldsymbol{\eta}$  gives:

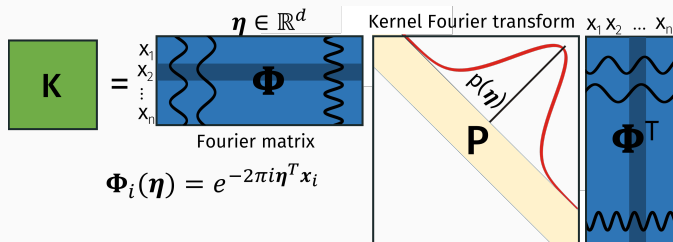
# VIEW AS MATRIX SAMPLING METHOD

Fourier transform  $k(\mathbf{z}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-2\pi i \boldsymbol{\eta}^T \mathbf{z}} d\boldsymbol{\eta}$  gives:



# VIEW AS MATRIX SAMPLING METHOD

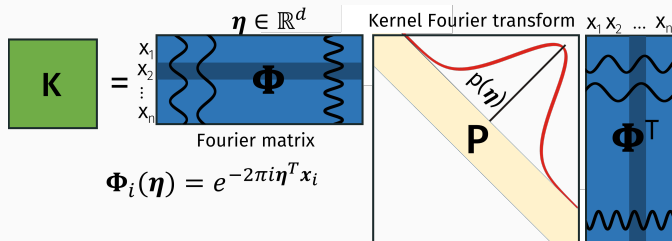
Fourier transform  $k(\mathbf{z}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-2\pi i \boldsymbol{\eta}^T \mathbf{z}} d\boldsymbol{\eta}$  gives:



- $\int_{\boldsymbol{\eta}} \Phi_i(\boldsymbol{\eta}) p(\boldsymbol{\eta}) \Phi_j(\boldsymbol{\eta})^* d\boldsymbol{\eta}$

# VIEW AS MATRIX SAMPLING METHOD

Fourier transform  $k(\mathbf{z}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-2\pi i \boldsymbol{\eta}^T \mathbf{z}} d\boldsymbol{\eta}$  gives:

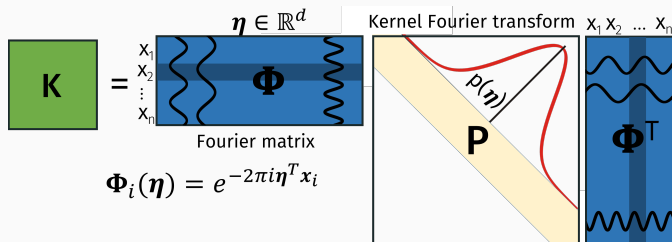


- $\int_{\boldsymbol{\eta}} \Phi_i(\boldsymbol{\eta}) p(\boldsymbol{\eta}) \Phi_j(\boldsymbol{\eta})^* d\boldsymbol{\eta} = \int_{\boldsymbol{\eta}} e^{-2\pi i \boldsymbol{\eta}^T (x_i - x_j)} p(\boldsymbol{\eta}) d\boldsymbol{\eta}$



# VIEW AS MATRIX SAMPLING METHOD

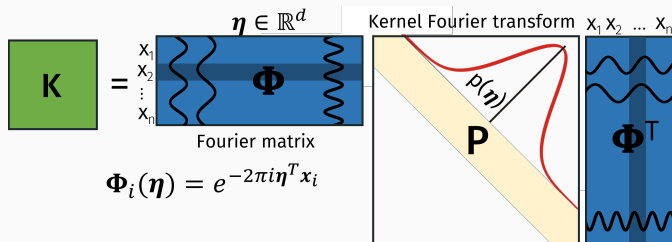
Fourier transform  $k(\mathbf{z}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-2\pi i \boldsymbol{\eta}^T \mathbf{z}} d\boldsymbol{\eta}$  gives:



- $\int_{\boldsymbol{\eta}} \Phi_i(\boldsymbol{\eta}) p(\boldsymbol{\eta}) \Phi_j(\boldsymbol{\eta})^* d\boldsymbol{\eta} = \int_{\boldsymbol{\eta}} e^{-2\pi i \boldsymbol{\eta}^T (\mathbf{x}_i - \mathbf{x}_j)} p(\boldsymbol{\eta}) d\boldsymbol{\eta} = k(\mathbf{x}_i - \mathbf{x}_j)$

# VIEW AS MATRIX SAMPLING METHOD

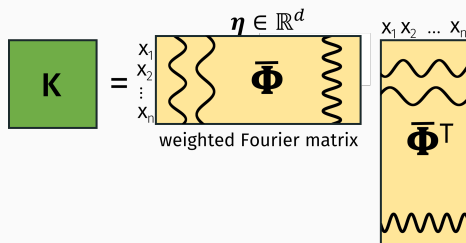
Fourier transform  $k(\mathbf{z}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-2\pi i \boldsymbol{\eta}^T \mathbf{z}} d\boldsymbol{\eta}$  gives:



- $$\int_{\boldsymbol{\eta}} \Phi_i(\boldsymbol{\eta}) p(\boldsymbol{\eta}) \Phi_j(\boldsymbol{\eta})^* d\boldsymbol{\eta} = \int_{\boldsymbol{\eta}} e^{-2\pi i \boldsymbol{\eta}^T (\mathbf{x}_i - \mathbf{x}_j)} p(\boldsymbol{\eta}) d\boldsymbol{\eta} = k(\mathbf{x}_i - \mathbf{x}_j) = \mathbf{K}_{i,j}.$$

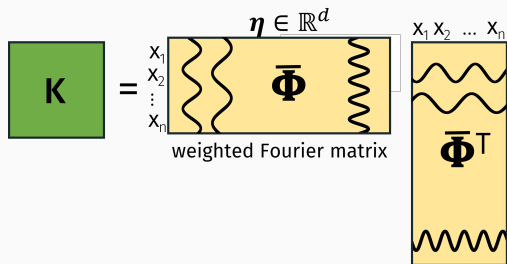
# VIEW AS MATRIX SAMPLING METHOD

Fourier transform  $k(\mathbf{z}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-2\pi i \boldsymbol{\eta}^T \mathbf{z}} d\boldsymbol{\eta}$  gives:

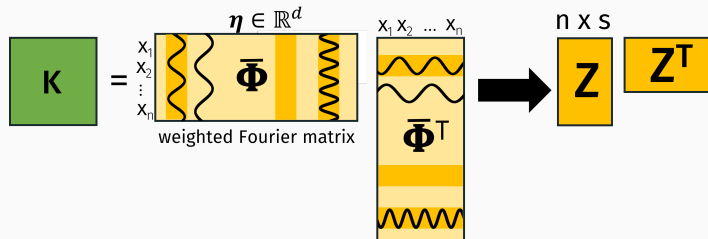


- $\int_{\boldsymbol{\eta}} \boldsymbol{\Phi}_i(\boldsymbol{\eta}) p(\boldsymbol{\eta}) \boldsymbol{\Phi}_j(\boldsymbol{\eta})^* d\boldsymbol{\eta} = \int_{\boldsymbol{\eta}} e^{-2\pi i \boldsymbol{\eta}^T (\mathbf{x}_i - \mathbf{x}_j)} p(\boldsymbol{\eta}) d\boldsymbol{\eta} = k(\mathbf{x}_i - \mathbf{x}_j) = \mathbf{K}_{i,j}$ .
- Set  $\bar{\boldsymbol{\Phi}} = \boldsymbol{\Phi} \mathbf{P}^{1/2}$ . So  $\mathbf{K} = \bar{\boldsymbol{\Phi}} \bar{\boldsymbol{\Phi}}^T$ .

# VIEW AS MATRIX SAMPLING METHOD

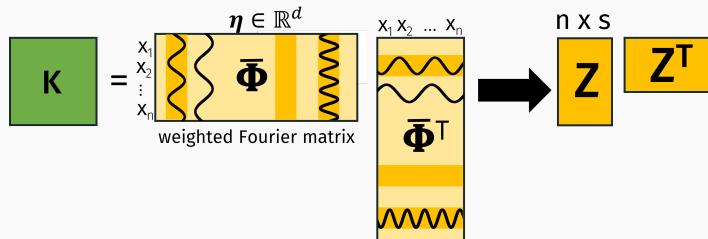


# VIEW AS MATRIX SAMPLING METHOD



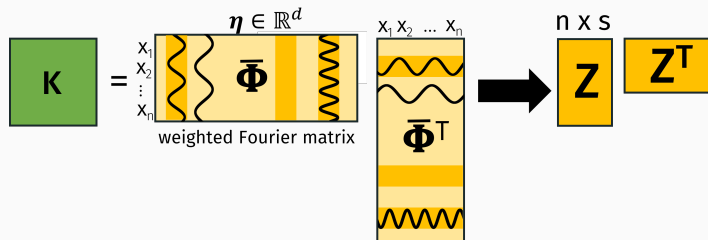
- $\mathbf{Z}(j) = \frac{1}{\sqrt{sp(\eta)}} \bar{\Phi}(\eta)$  with probability  $p(\eta)$ . So  $\mathbb{E}[\mathbf{Z}\mathbf{Z}^T] = \mathbf{K}$ .

# VIEW AS MATRIX SAMPLING METHOD

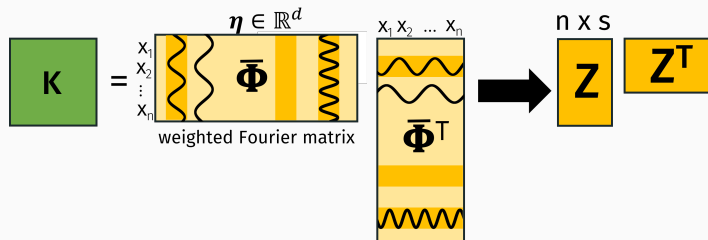


- $\mathbf{Z}(j) = \frac{1}{\sqrt{sp(\eta)}} \bar{\Phi}(\eta)$  with probability  $p(\eta)$ . So  $\mathbb{E}[\mathbf{Z}\mathbf{Z}^T] = \mathbf{K}$ .
- $\mathbf{z}_i = \frac{1}{\sqrt{s}} \left[ e^{-2\pi i \eta_1^T \mathbf{x}_i}, \dots, e^{-2\pi i \eta_s^T \mathbf{x}_i} \right]$  for  $\eta_1, \dots, \eta_s$  sampled according to  $p(\eta)$ .

# VIEW AS MATRIX SAMPLING METHOD



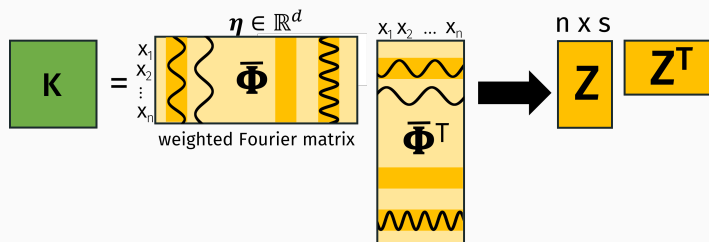
# VIEW AS MATRIX SAMPLING METHOD



- $\mathbf{Z}$  is a sample of  $\bar{\Phi} = \Phi \mathbf{P}^{1/2}$ . Columns are sampled with probability  $\propto p(\eta)$ , i.e., their **squared column norms**.

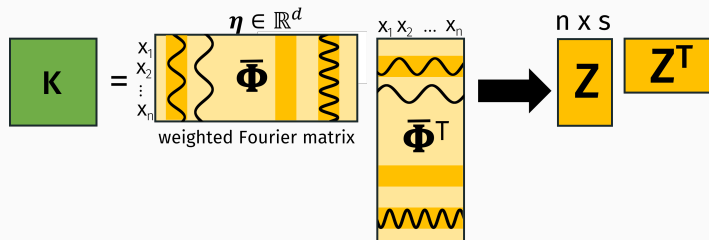


# VIEW AS MATRIX SAMPLING METHOD



- $Z$  is a sample of  $\bar{\Phi} = \Phi P^{1/2}$ . Columns are sampled with probability  $\propto p(\eta)$ , i.e., their **squared column norms**.
- It is well known from work on randomized methods in linear algebra that there are better sampling probabilities (in both theory and practice): the column **leverage scores**.

# VIEW AS MATRIX SAMPLING METHOD



- $Z$  is a sample of  $\bar{\Phi} = \Phi P^{1/2}$ . Columns are sampled with probability  $\propto p(\eta)$ , i.e., their **squared column norms**.
- It is well known from work on randomized methods in linear algebra that there are better sampling probabilities (in both theory and practice): the column **leverage scores**.
- Also noted by Bach '17, implicit in Rudi et al. '16.

**Column Norm Sampling:**  $s = \tilde{O}(d/\epsilon^2)$  samples ensure that  $(\mathbf{Z}\mathbf{Z}^T)_{i,j} = \mathbf{K}_{i,j} \pm \epsilon$  for all  $i, j$  with high probability [RR07].

**Column Norm Sampling:**  $s = \tilde{O}(d/\epsilon^2)$  samples ensure that  $(\mathbf{ZZ}^T)_{i,j} = \mathbf{K}_{i,j} \pm \epsilon$  for all  $i, j$  with high probability [RR07].

**Ridge Leverage Score Sampling:**  $s = \tilde{O}(s_\lambda/\epsilon^2)$  samples gives spectral approximation:

$$(1 - \epsilon)(\mathbf{ZZ}^T + \lambda\mathbf{I}) \preceq \mathbf{K} + \lambda\mathbf{I} \preceq (1 + \epsilon)(\mathbf{ZZ}^T + \lambda\mathbf{I}).$$

where  $s_\lambda = \text{tr}(\mathbf{K}(\mathbf{K} + \lambda\mathbf{I})^{-1})$  is the **statistical dimension**.

**Column Norm Sampling:**  $s = \tilde{O}(d/\epsilon^2)$  samples ensure that  $(\mathbf{ZZ}^T)_{i,j} = \mathbf{K}_{i,j} \pm \epsilon$  for all  $i, j$  with high probability [RR07].

**Ridge Leverage Score Sampling:**  $s = \tilde{O}(s_\lambda/\epsilon^2)$  samples gives spectral approximation:

$$(1 - \epsilon)(\mathbf{ZZ}^T + \lambda\mathbf{I}) \preceq \mathbf{K} + \lambda\mathbf{I} \preceq (1 + \epsilon)(\mathbf{ZZ}^T + \lambda\mathbf{I}).$$

where  $s_\lambda = \text{tr}(\mathbf{K}(\mathbf{K} + \lambda\mathbf{I})^{-1})$  is the **statistical dimension**.

- Spectral approximation gives statistical guarantees for kernel ridge regression (this work), and approximation bounds for kernel PCA and k-means clustering (Cohen, Musco, Musco '16, '17)

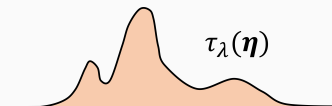
The ridge leverage score for frequency  $\eta$  is given by:

$$\tau_\lambda(\eta) = \bar{\Phi}(\eta)^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \bar{\Phi}(\eta).$$

The ridge leverage score for frequency  $\eta$  is given by:

$$\tau_\lambda(\eta) = \bar{\Phi}(\eta)^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \bar{\Phi}(\eta).$$

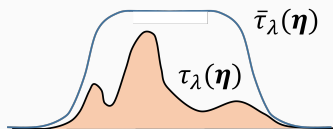
- Expensive to invert  $\mathbf{K} + \lambda \mathbf{I}$ . But even if you could do that efficiently, it is not at all clear you could efficiently sample from the leverage score distribution.



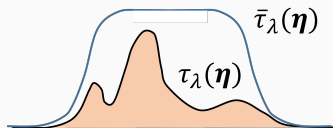
**Main goal:** Get a handle on the Fourier ridge leverage scores for common kernels by upper bounding them with simple distributions.



**Main goal:** Get a handle on the Fourier ridge leverage scores for common kernels by upper bounding them with simple distributions.



**Main goal:** Get a handle on the Fourier ridge leverage scores for common kernels by upper bounding them with simple distributions.



1. Improve random Fourier features.
2. Bound statistical dimension by the sum of leverage scores.
3. Connections with sparse Fourier transforms, Fourier interpolation, and other problems.

Ridge leverage score  $\tau_\lambda(\boldsymbol{\eta}) = \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})$  also equals:

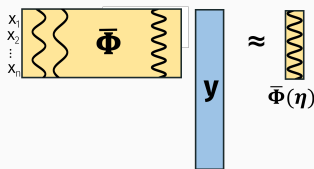
$$\tau_\lambda(\boldsymbol{\eta}) = \min_{\mathbf{y}} \lambda^{-1} \|\bar{\boldsymbol{\Phi}}\mathbf{y} - \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})\|_2^2 + \|\mathbf{y}\|_2^2.$$

# ALTERNATIVE LEVERAGE SCORE CHARACTERIZATION

Ridge leverage score  $\tau_\lambda(\boldsymbol{\eta}) = \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})$  also equals:

$$\tau_\lambda(\boldsymbol{\eta}) = \min_{\mathbf{y}} \lambda^{-1} \|\bar{\boldsymbol{\Phi}}\mathbf{y} - \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})\|_2^2 + \|\mathbf{y}\|_2^2.$$

**Intuition:**  $\tau_\lambda(\boldsymbol{\eta})$  is small iff there exists a function  $\mathbf{y} : \mathbb{R}^d \rightarrow \mathbb{C}$  with low energy ( $\|\mathbf{y}\|_2^2$  small) whose ( $\sqrt{\rho(\boldsymbol{\eta})}$  weighted) Fourier transform is close to the frequency  $e^{-2\pi i \mathbf{x}_j^T \boldsymbol{\eta}}$  at each data point.

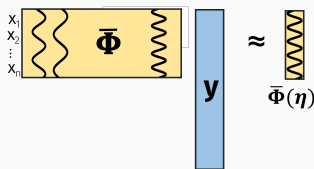


# ALTERNATIVE LEVERAGE SCORE CHARACTERIZATION

Ridge leverage score  $\tau_\lambda(\boldsymbol{\eta}) = \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})$  also equals:

$$\tau_\lambda(\boldsymbol{\eta}) = \min_{\mathbf{y}} \lambda^{-1} \|\bar{\boldsymbol{\Phi}} \mathbf{y} - \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta})\|_2^2 + \|\mathbf{y}\|_2^2.$$

**Intuition:**  $\tau_\lambda(\boldsymbol{\eta})$  is small iff there exists a function  $\mathbf{y} : \mathbb{R}^d \rightarrow \mathbb{C}$  with low energy ( $\|\mathbf{y}\|_2^2$  small) whose ( $\sqrt{p(\boldsymbol{\eta})}$  weighted) Fourier transform is close to the frequency  $e^{-2\pi i \mathbf{x}_j^T \boldsymbol{\eta}}$  at each data point.



- $\mathbf{y}$  reconstructs frequency  $\boldsymbol{\eta}$  from other frequencies. The easier it is to reconstruct, the less important it is to sample.

Assume data points are 1-dimensional and bounded:

$x_1, \dots, x_n \in [-\delta, \delta]$ . One possibility is to choose  $\mathbf{y}$  with ( $\sqrt{p(\eta)}$  weighted) Fourier transform equal to  $e^{-2\pi i x \eta}$  for all  $x \in [-\delta, \delta]$ .

Assume data points are 1-dimensional and bounded:

$x_1, \dots, x_n \in [-\delta, \delta]$ . One possibility is to choose  $\mathbf{y}$  with ( $\sqrt{p(\eta)}$  weighted) Fourier transform equal to  $e^{-2\pi i x \eta}$  for all  $x \in [-\delta, \delta]$ .

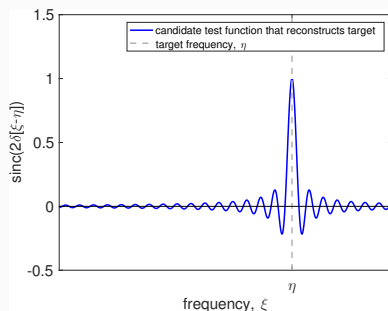
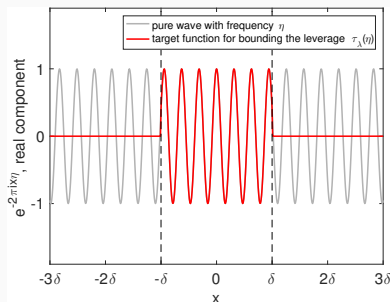
- Achieved by the shifted sinc function weighted by  $1/\sqrt{p(\eta)}$ .

# FREQUENCY RECONSTRUCTION FOR BOUNDED DATA

Assume data points are 1-dimensional and bounded:

$x_1, \dots, x_n \in [-\delta, \delta]$ . One possibility is to choose  $\mathbf{y}$  with ( $\sqrt{p(\eta)}$  weighted) Fourier transform equal to  $e^{-2\pi i x \eta}$  for all  $x \in [-\delta, \delta]$ .

- Achieved by the shifted sinc function weighted by  $1/\sqrt{p(\eta)}$ .





Unfortunately the sinc function falls off too slowly.

Unfortunately the sinc function falls off too slowly.

- For the Gaussian kernel, the  $\frac{1}{\sqrt{\rho(\eta)}} \approx e^{\eta^2/4}$  weighting, will grow faster than  $\text{sinc}(2\delta\eta) = \frac{\sin(2\delta\eta)}{\eta}$  falls off. So  $\|\mathbf{y}\|_2$  is unbounded.

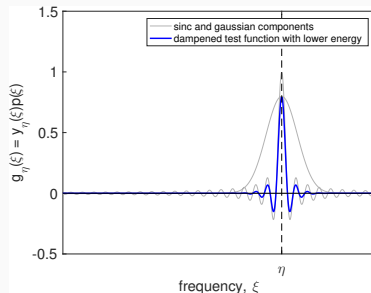
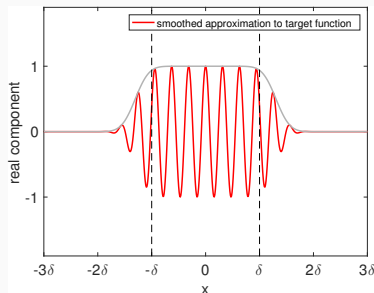
Unfortunately the sinc function falls off too slowly.

- For the Gaussian kernel, the  $\frac{1}{\sqrt{\rho(\eta)}} \approx e^{\eta^2/4}$  weighting, will grow faster than  $\text{sinc}(2\delta\eta) = \frac{\sin(2\delta\eta)}{\eta}$  falls off. So  $\|\mathbf{y}\|_2$  is unbounded.
- **Solution:** Dampen the sinc by multiplying with a Gaussian, keeping Fourier transform nearly identical.

# IMPROVED TEST FUNCTION

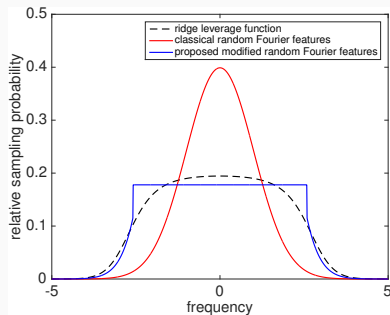
Unfortunately the sinc function falls off too slowly.

- For the Gaussian kernel, the  $\frac{1}{\sqrt{\rho(\eta)}} \approx e^{\eta^2/4}$  weighting, will grow faster than  $\text{sinc}(2\delta\eta) = \frac{\sin(2\delta\eta)}{\eta}$  falls off. So  $\|\mathbf{y}\|_2$  is unbounded.
- **Solution:** Dampen the sinc by multiplying with a Gaussian, keeping Fourier transform nearly identical.



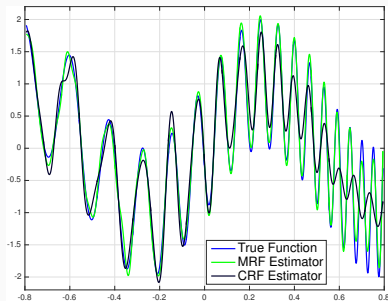
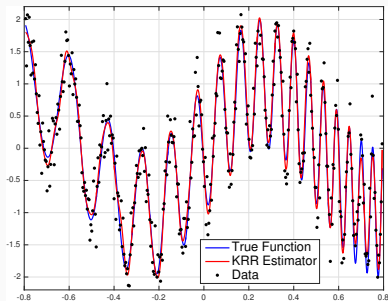
**Upshot:** easy to sample from approximate leverage distribution for the Gaussian kernel with  $x_1, \dots, x_n \in [-\delta, \delta]^d$ :

$$\bar{\tau}_\lambda(\boldsymbol{\eta}) \begin{cases} \tilde{O}(\delta^d) \text{ when } \|\boldsymbol{\eta}\|_\infty \leq \sqrt{\log n/\lambda} \\ \rho(\boldsymbol{\eta}) = e^{-\|\boldsymbol{\eta}\|_2^2/2} \text{ otherwise.} \end{cases}$$



# EXPERIMENTAL RESULTS

Example of approximating a synthetic 'wiggly function':



CRF = classic random Fourier features 'column norm' sampling,  
MRF = our modified sampling distribution.

Questions?