

Nonlinear Dimensionality Reduction for Faster Kernel Methods in Machine Learning

Christopher Musco, Massachusetts Institute of Technology
February 27, 2018

ICML 2017:

“Random Fourier Features for Kernel Ridge Regression:
Approximation Bounds and Statistical Guarantees”

Joint work with:

Haim Avron (TAU)
Michael Kapralov (EPFL)
Cameron Musco (MIT)
Ameya Velingker (EPFL)
Amir Zandieh (EPFL)

Main idea:

Study **Fourier kernel approximation methods** from a **matrix sampling point of view**.

Main idea:

Study **Fourier kernel approximation methods** from a **matrix sampling point of view**.

Specifics:

- Analyze Random Fourier Features method (Rahimi, Recht NIPS '07) using techniques based on leverage scores.

Main idea:

Study **Fourier kernel approximation methods** from a **matrix sampling point of view**.

Specifics:

- Analyze Random Fourier Features method (Rahimi, Recht NIPS '07) using techniques based on leverage scores.
- Develop an improved Random Fourier Features method based on this analysis (better in theory and experiments).

Main idea:

Study **Fourier kernel approximation methods** from a **matrix sampling point of view**.

Specifics:

- Analyze Random Fourier Features method (Rahimi, Recht NIPS '07) using techniques based on leverage scores.
- Develop an improved Random Fourier Features method based on this analysis (better in theory and experiments).

Lots of open questions and directions for future work.

Opportunities to combine techniques from randomized linear algebra and Fourier methods.

QUICK REFRESHER ON **KERNEL METHODS**

Adapt standard **linear learning methods** (least squares regression, support vector machines, PCA, k-means clustering) to **learn nonlinear relationships**.

Adapt standard **linear learning methods** (least squares regression, support vector machines, PCA, k-means clustering) to **learn nonlinear relationships**.

(theoretically well-understood, multipurpose, widely used)

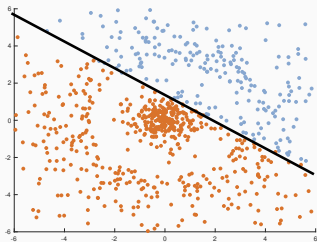
“Lift” data points to a higher dimensional feature space. E.g.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

“Lift” data points to a higher dimensional feature space. E.g.

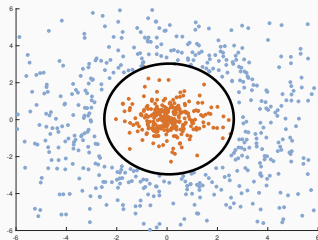
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \implies \phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ \vdots \\ x_d \\ x_1x_1 \\ x_1x_2 \\ \vdots \\ x_dx_d \end{bmatrix}$$

Linear Classifier



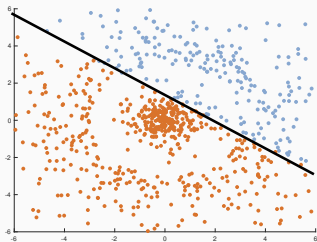
$$x_1 + 2x_2 \geq 6$$

Kernel Classifier



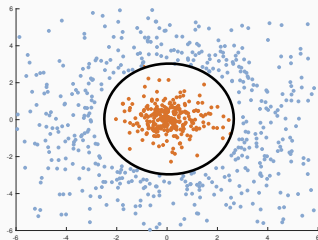
$$x_1^2 + x_2^2 \geq 10$$

Linear Classifier



$$x_1 + 2x_2 \geq 6$$

Kernel Classifier



$$x_1^2 + x_2^2 \geq 10$$

$$\phi(\mathbf{x})_3 + \phi(\mathbf{x})_4 \geq 10$$

Main computational issue: Forming $\phi(\mathbf{x})$ is intractable even for moderately complex kernels.

E.g. degree q polynomials $\implies O(d^q)$ dimensional vectors.

THE COST OF KERNEL METHODS

Main computational issue: Forming $\phi(\mathbf{x})$ is intractable even for moderately complex kernels.

E.g. degree q polynomials $\implies O(d^q)$ dimensional vectors.

Fix: For common linear learning methods, we only need the **kernel dot product** for each pair of data points \mathbf{x}, \mathbf{y} :

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle.$$

THE COST OF KERNEL METHODS

Main computational issue: Forming $\phi(\mathbf{x})$ is intractable even for moderately complex kernels.

E.g. degree q polynomials $\implies O(d^q)$ dimensional vectors.

Fix: For common linear learning methods, we only need the **kernel dot product** for each pair of data points \mathbf{x}, \mathbf{y} :

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle.$$

Can often be computed much more quickly than $\phi(\mathbf{x}_i), \phi(\mathbf{x}_j)$.

EXAMPLE: LEAST SQUARES REGRESSION

Input: Data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$, responses $\mathbf{b} = [b_1, \dots, b_n]^T$.

EXAMPLE: LEAST SQUARES REGRESSION

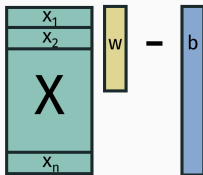
Input: Data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$, responses $\mathbf{b} = [b_1, \dots, b_n]^T$.

Solve:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

Predict:

$$b_{new} = \mathbf{w}^{*T} \mathbf{x}_{new}$$



EXAMPLE: LEAST SQUARES REGRESSION

Input: Data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$, responses $\mathbf{b} = [b_1, \dots, b_n]^T$.

Solve:

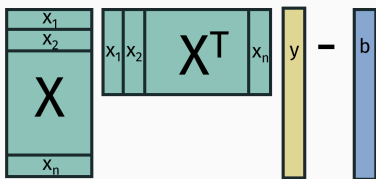
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$$\mathbf{y}^* = \arg \min_{\mathbf{y}} \|\mathbf{X}\mathbf{X}^T \mathbf{y} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{X}^T \mathbf{y}\|_2^2$$

Predict:

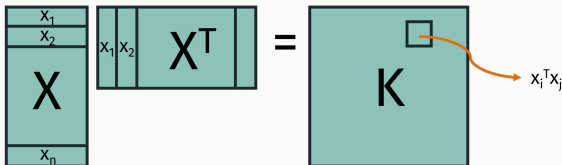
$$b_{new} = \mathbf{w}^{*T} \mathbf{x}_{new}$$

$$b_{new} = \mathbf{y}^{*T} \mathbf{X} \mathbf{x}_{new}$$



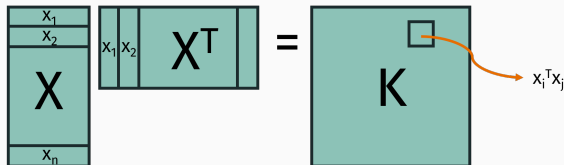
EXAMPLE: LEAST SQUARES REGRESSION

For training, we compute $(K + \lambda I)^{-1}$ for the **kernel matrix** K :



EXAMPLE: LEAST SQUARES REGRESSION

For training, we compute $(K + \lambda I)^{-1}$ for the **kernel matrix** K :



If we replace each x_i with $\phi(x_i)$ for nonlinear learning, we just need to alternatively compute:

$$K_{i,j} = \langle \phi(x_i), \phi(x_j) \rangle = k(x_i, x_j).$$

THE “KERNEL TRICK”

Kernel dot product can often be computed **implicitly** without forming $\phi(\mathbf{x})$ and $\phi(\mathbf{y})$. For example:

$$\begin{aligned}(1 + \langle \mathbf{x}, \mathbf{y} \rangle)^2 &= (1 + x_1 y_1 + \dots + x_d y_d)^2 \\&= (1 + x_1 y_1 + x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + \dots) \\&= \left\langle [1, x_1, x_1^2, \sqrt{2}x_1 x_2, \dots], [1, y_1, y_1^2, \sqrt{2}y_1 y_2, \dots] \right\rangle \\&= \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle\end{aligned}$$

THE “KERNEL TRICK”

Kernel dot product can often be computed **implicitly** without forming $\phi(\mathbf{x})$ and $\phi(\mathbf{y})$. For example:

$$\begin{aligned}(1 + \langle \mathbf{x}, \mathbf{y} \rangle)^2 &= (1 + x_1 y_1 + \dots + x_d y_d)^2 \\&= (1 + x_1 y_1 + x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + \dots) \\&= \left\langle [1, x_1, x_1^2, \sqrt{2}x_1 x_2, \dots], [1, y_1, y_1^2, \sqrt{2}y_1 y_2, \dots] \right\rangle \\&= \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle\end{aligned}$$

The **kernel function** $k(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^2$ provides an alternative **similarity metric** to the standard dot product.

THE “KERNEL TRICK”

Kernel learning pipeline for data points $\mathbf{x}_1, \dots, \mathbf{x}_n$:

1. Choose kernel function $k(\mathbf{x}_i, \mathbf{x}_j)$:

$$(1 + \mathbf{x}_i^T \mathbf{x}_j)^q, e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}, e^{-\|\mathbf{x}_i - \mathbf{x}_j\|_1}, \text{ etc.}$$

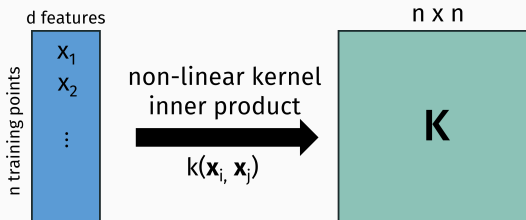
2. Form $n \times n$ **kernel matrix** \mathbf{K} with:

$$\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$$

3. Compute model using \mathbf{K} : compute $(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{b}$ for kernel regression, eigendecomposition of \mathbf{K} for kernel PCA, etc.

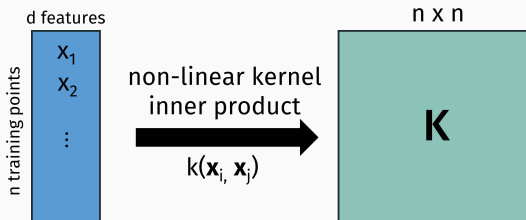
ALGORITHMIC CHALLENGE

Even if we avoid explicit feature expansion, kernel methods are slow. Quadratic dependence on the number of data points.



ALGORITHMIC CHALLENGE

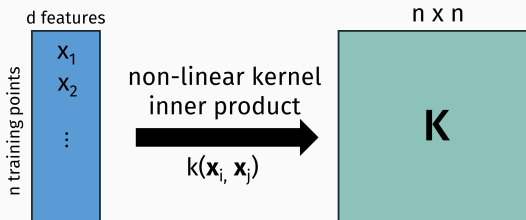
Even if we avoid explicit feature expansion, kernel methods are slow. Quadratic dependence on the number of data points.



- $n = 100,000 \Rightarrow 10$ billion entries $\Rightarrow 80$ GB to store K .

ALGORITHMIC CHALLENGE

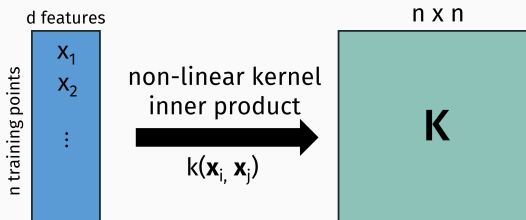
Even if we avoid explicit feature expansion, kernel methods are slow. Quadratic dependence on the number of data points.



- $n = 100,000 \Rightarrow 10$ billion entries $\Rightarrow 80$ GB to store K .
- Just writing down K requires $\Omega(n^2)$ time.

ALGORITHMIC CHALLENGE

Even if we avoid explicit feature expansion, kernel methods are slow. Quadratic dependence on the number of data points.

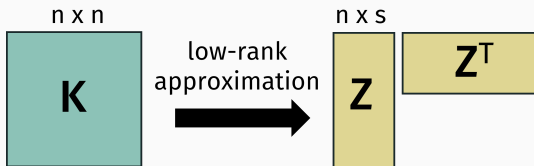


- $n = 100,000 \Rightarrow 10$ billion entries $\Rightarrow 80$ GB to store K .
- Just writing down K requires $\Omega(n^2)$ time.
- Other operations require even more. A single iteration for a linear system solver takes $\Omega(n^2)$ time.

New algorithmic ideas are needed to scale kernel methods.
Even for moderately large datasets.

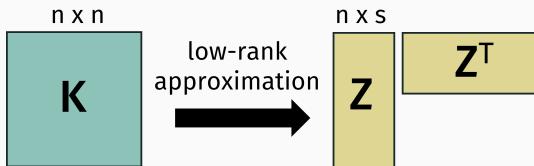
STANDARD APPROACH: LOW-RANK APPROXIMATION

Find approximation $\tilde{K} = ZZ^T$ for K (which is symmetric and positive semidefinite):



STANDARD APPROACH: LOW-RANK APPROXIMATION

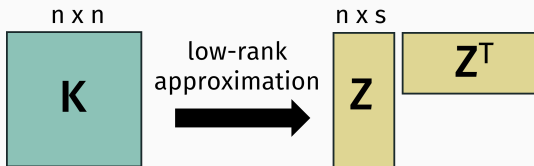
Find approximation $\tilde{\mathbf{K}} = \mathbf{Z}\mathbf{Z}^T$ for \mathbf{K} (which is symmetric and positive semidefinite):



We can typically set $s \ll n$.

STANDARD APPROACH: LOW-RANK APPROXIMATION

Find approximation $\tilde{K} = ZZ^T$ for K (which is symmetric and positive semidefinite):

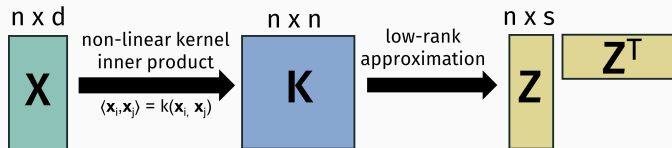


We can typically set $s \ll n$.

- Z takes $O(ns)$ space to store.
- Orthogonalization, eigendecomposition, and inversion of ZZ^T all take just $O(ns^2)$ time.
- $ZZ^T x$ can be computed in $O(ns)$ time.

IMPLICIT LOW-RANK APPROXIMATION

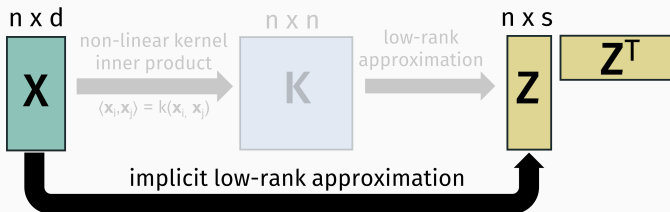
All standard low-rank approximation algorithms take $\Omega(n^2)$ time for kernel matrices: we at least need to compute K .



Goal: Find methods that avoid explicit access to the K .

IMPLICIT LOW-RANK APPROXIMATION

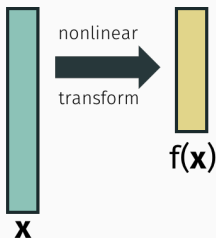
All standard low-rank approximation algorithms take $\Omega(n^2)$ time for kernel matrices: we at least need to compute K .



Goal: Find methods that avoid explicit access to the K .

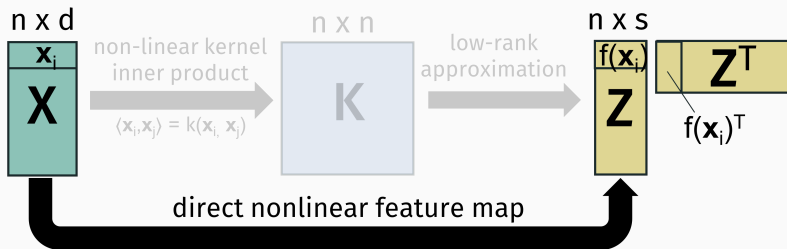
Goal: Find a compact nonlinear function $f(\cdot)$ such that

$$\langle f(\mathbf{x}), f(\mathbf{y}) \rangle \approx k(\mathbf{x}, \mathbf{y})$$

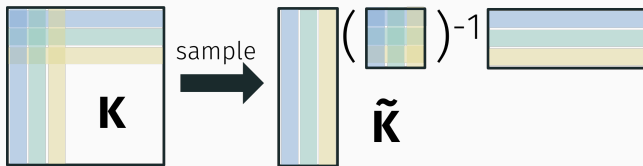


Goal: Find a compact nonlinear function $f(\cdot)$ such that

$$\langle f(\mathbf{x}), f(\mathbf{y}) \rangle \approx k(\mathbf{x}, \mathbf{y})$$

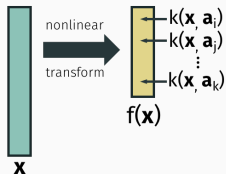


Example: Nyström approximation.



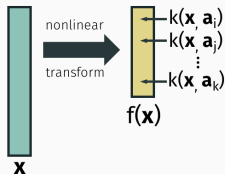
Construct $\tilde{\mathbf{K}}$ from subsample of \mathbf{K} 's columns and rows.

Example: Nyström approximation.



Construct $\tilde{\mathbf{K}}$ from subsample of \mathbf{K} 's columns and rows.

Example: Nyström approximation.



Construct $\tilde{\mathbf{K}}$ from subsample of \mathbf{K} 's columns and rows.

Can match optimal k -rank approximation to $(1 + \epsilon)$ factor with $O(nk/\epsilon)$ total samples [Musco, Musco, NIPS 2017].

Suppose $k(\mathbf{x}, \mathbf{y})$ is shift invariant.

Suppose $k(\mathbf{x}, \mathbf{y})$ is **shift invariant**. I.e.

$$k(\mathbf{x}, \mathbf{y}) = k(\mathbf{\Delta}) \quad \text{where} \quad \mathbf{\Delta} = \mathbf{x} - \mathbf{y}.$$

Suppose $k(\mathbf{x}, \mathbf{y})$ is **shift invariant**. I.e.

$$k(\mathbf{x}, \mathbf{y}) = k(\mathbf{\Delta}) \quad \text{where} \quad \mathbf{\Delta} = \mathbf{x} - \mathbf{y}.$$

Suppose $k(\mathbf{x}, \mathbf{y})$ is **shift invariant**. I.e.

$$k(\mathbf{x}, \mathbf{y}) = k(\mathbf{\Delta}) \quad \text{where} \quad \mathbf{\Delta} = \mathbf{x} - \mathbf{y}.$$

- Gaussian (RBF) kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2}$

Suppose $k(\mathbf{x}, \mathbf{y})$ is **shift invariant**. I.e.

$$k(\mathbf{x}, \mathbf{y}) = k(\mathbf{\Delta}) \quad \text{where} \quad \mathbf{\Delta} = \mathbf{x} - \mathbf{y}.$$

- Gaussian (RBF) kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2}$
- Laplace kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_1}$

Suppose $k(\mathbf{x}, \mathbf{y})$ is **shift invariant**. I.e.

$$k(\mathbf{x}, \mathbf{y}) = k(\mathbf{\Delta}) \quad \text{where} \quad \mathbf{\Delta} = \mathbf{x} - \mathbf{y}.$$

- Gaussian (RBF) kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2}$
- Laplace kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_1}$
- Matern kernel, Cauchy kernel, rational quadratic, etc.

Suppose $k(\mathbf{x}, \mathbf{y})$ is **shift invariant**. I.e.

$$k(\mathbf{x}, \mathbf{y}) = k(\mathbf{\Delta}) \quad \text{where} \quad \mathbf{\Delta} = \mathbf{x} - \mathbf{y}.$$

- Gaussian (RBF) kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2}$
- Laplace kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_1}$
- Matern kernel, Cauchy kernel, rational quadratic, etc.

RANDOM FOURIER FEATURES (RFF) ALGORITHM

Write $k(\mathbf{\Delta})$ using its inverse Fourier transform,

$$p(\boldsymbol{\eta}) = \mathcal{F}^{-1}k(\mathbf{\Delta}) \qquad k(\mathbf{\Delta}) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-i\pi \boldsymbol{\eta}^\top \mathbf{\Delta}} d\boldsymbol{\eta}$$

RANDOM FOURIER FEATURES (RFF) ALGORITHM

Write $k(\Delta)$ using its inverse Fourier transform,

$$p(\eta) = \mathcal{F}^{-1}k(\Delta) \quad k(\Delta) = \int_{\eta \in \mathbb{R}^d} p(\eta) e^{-i\pi \eta^\top \Delta} d\eta$$

Approximate with finite sum for any Δ :

$$k(\Delta) \approx \frac{1}{S} \sum_{j=1}^S c_j e^{-i\pi \eta_j^\top \Delta}$$

RANDOM FOURIER FEATURES (RFF) ALGORITHM

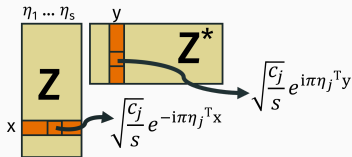
Write $k(\Delta)$ using its inverse Fourier transform,

$$p(\eta) = \mathcal{F}^{-1}k(\Delta) \quad k(\Delta) = \int_{\eta \in \mathbb{R}^d} p(\eta) e^{-i\pi \eta^T \Delta} d\eta$$

Approximate with finite sum for any Δ :

$$k(\Delta) \approx \frac{1}{S} \sum_{j=1}^S c_j e^{-i\pi \eta_j^T \Delta}$$

Immediately gives low-rank approximation:



RANDOM FOURIER FEATURES (RFF) ALGORITHM

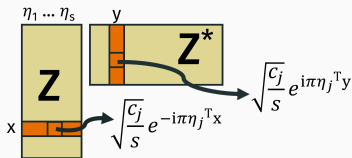
Write $k(\Delta)$ using its inverse Fourier transform,

$$p(\eta) = \mathcal{F}^{-1}k(\Delta) \quad k(\Delta) = \int_{\eta \in \mathbb{R}^d} p(\eta) e^{-i\pi \eta^T \Delta} d\eta$$

Approximate with finite sum for any Δ :

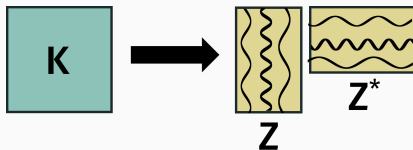
$$k(\Delta) \approx \frac{1}{S} \sum_{j=1}^S c_j e^{-i\pi \eta_j^T \Delta}$$

Immediately gives low-rank approximation:



$$\begin{aligned} Z(x)Z(y)^* &= \sum_j \frac{c_j}{s} e^{-i\pi \eta_j^T (x-y)} \\ &\approx k(x-y) \end{aligned}$$

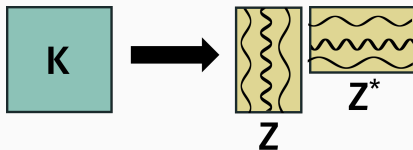
RANDOM FOURIER FEATURES (RFF) ALGORITHM



How do we choose the frequencies in the finite sum?

$$\int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-i\pi \boldsymbol{\eta}^T \boldsymbol{\Delta}} d\boldsymbol{\eta} \approx \frac{1}{s} \sum_{j=1}^s c_j e^{-i\pi \boldsymbol{\eta}_j^T \boldsymbol{\Delta}}$$

RANDOM FOURIER FEATURES (RFF) ALGORITHM



How do we choose the frequencies in the finite sum?

$$\int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-i\pi \boldsymbol{\eta}^T \boldsymbol{\Delta}} d\boldsymbol{\eta} \approx \frac{1}{s} \sum_{j=1}^s c_j e^{-i\pi \boldsymbol{\eta}_j^T \boldsymbol{\Delta}}$$

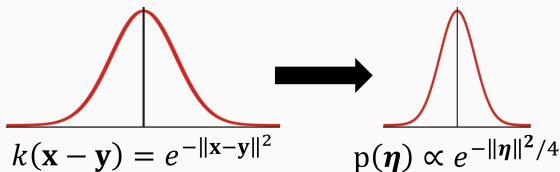
Bochner's Theorem: For shift-invariant, positive semidefinite kernel functions, $p(\boldsymbol{\eta}) \geq 0$ for all $\boldsymbol{\eta}$.

Bochner's Theorem: For shift-invariant, positive semidefinite kernel functions, $p(\boldsymbol{\eta}) \geq 0$ for all $\boldsymbol{\eta}$.

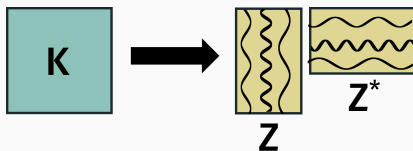
RANDOM FOURIER FEATURES (RFF) ALGORITHM

Bochner's Theorem: For shift-invariant, positive semidefinite kernel functions, $p(\boldsymbol{\eta}) \geq 0$ for all $\boldsymbol{\eta}$.

Example (Gaussian kernel):

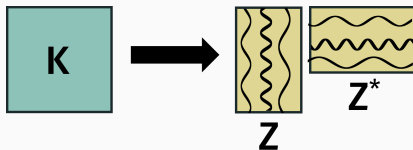


RANDOM FOURIER FEATURES (RFF) ALGORITHM



Sample frequencies from distribution $p(\eta)$.

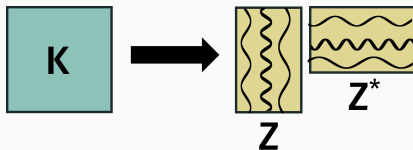
RANDOM FOURIER FEATURES (RFF) ALGORITHM



Sample frequencies from distribution $p(\eta)$.

$$\mathbb{E} \left[\frac{1}{s} \sum_{j=1}^s e^{-i\pi \eta_j^T \Delta} \right] = \mathbb{E} \left[e^{-i\pi \eta^T \Delta} \right]$$

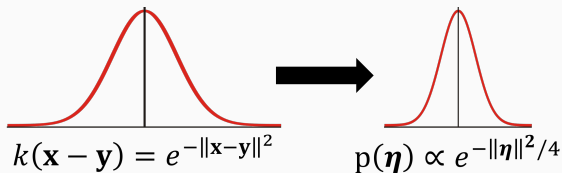
RANDOM FOURIER FEATURES (RFF) ALGORITHM



Sample frequencies from distribution $p(\eta)$.

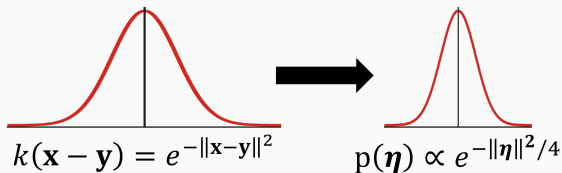
$$\mathbb{E} \left[\frac{1}{s} \sum_{j=1}^s e^{-i\pi \eta_j^T \Delta} \right] = \mathbb{E} \left[e^{-i\pi \eta^T \Delta} \right] = \int_{\eta \in \mathbb{R}^d} p(\eta) e^{-i\pi \eta^T \Delta} d\eta$$

RANDOM FOURIER FEATURES (RFF) ALGORITHM



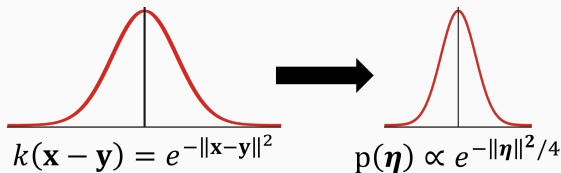
- Sampling frequencies by $p(\boldsymbol{\eta})$ is correct in expectation.

RANDOM FOURIER FEATURES (RFF) ALGORITHM



- Sampling frequencies by $p(\boldsymbol{\eta})$ is correct in expectation.
- Final sum $\frac{1}{S} \sum_{j=1}^S e^{-i\pi \boldsymbol{\eta}_j^T \boldsymbol{\Delta}} \approx k(\boldsymbol{\Delta})$ only has terms with magnitude $|e^{-i\pi \boldsymbol{\eta}_j^T \boldsymbol{\Delta}}| = 1$. Real part ≤ 1 . Imaginary part ≤ 1 .

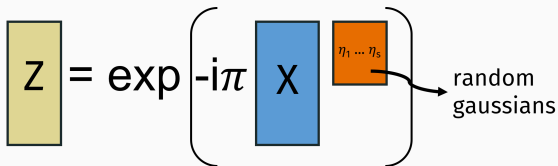
RANDOM FOURIER FEATURES (RFF) ALGORITHM



- Sampling frequencies by $p(\boldsymbol{\eta})$ is correct in expectation.
- Final sum $\frac{1}{S} \sum_{j=1}^S e^{-i\pi \boldsymbol{\eta}_j^T \boldsymbol{\Delta}} \approx k(\boldsymbol{\Delta})$ only has terms with magnitude $|e^{-i\pi \boldsymbol{\eta}_j^T \boldsymbol{\Delta}}| = 1$. Real part ≤ 1 . Imaginary part ≤ 1 .
- By Chernoff bound, if we take $\mathcal{O}\left(\frac{\log n}{\epsilon^2}\right)$ samples, we approximate every entry in \mathbf{K} to error $\pm \epsilon$ w.h.p.

RANDOM FOURIER FEATURES (RFF) ALGORITHM

Super simple algorithm. For Gaussian kernel $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_2^2}$:

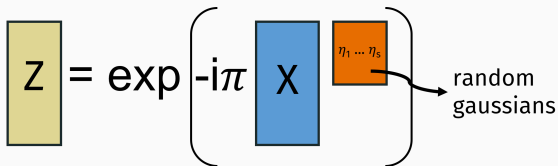


The diagram illustrates the RFF algorithm equation: $\mathbf{Z} = \exp \left(-i\pi \mathbf{X} \boldsymbol{\eta} \right)$. The variable \mathbf{Z} is represented by a yellow rectangle. The variable \mathbf{X} is represented by a blue rectangle. The vector $\boldsymbol{\eta}$ is represented by an orange rectangle containing the text $\eta_1 \dots \eta_s$. An arrow points from the orange rectangle to the text "random gaussians". The entire expression is enclosed in large parentheses.

$$z_{x,j} = e^{-i\pi \mathbf{x}^T \boldsymbol{\eta}_j} \text{ for } \boldsymbol{\eta}_j \sim \mathcal{N}$$

RANDOM FOURIER FEATURES (RFF) ALGORITHM

Super simple algorithm. For Gaussian kernel $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_2^2}$:



The diagram illustrates the RFF algorithm formula. It shows a yellow box labeled 'Z' followed by an equals sign and an exponential function. The argument of the exponential is $-i\pi$ multiplied by a blue box labeled 'X' and an orange box labeled $\eta_1 \dots \eta_s$. These three components are enclosed in large parentheses. An arrow points from the orange box to the text 'random gaussians'.

$$\mathbf{Z} = \exp \left(-i\pi \mathbf{X} \begin{matrix} \eta_1 \dots \eta_s \end{matrix} \right)$$

random gaussians

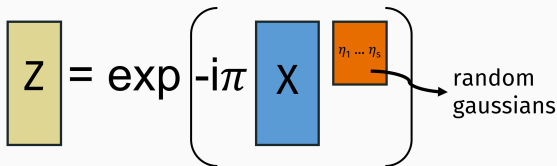
$$Z_{\mathbf{x},j} = e^{-i\pi \mathbf{x}^T \boldsymbol{\eta}_j} \text{ for } \boldsymbol{\eta}_j \sim \mathcal{N}$$

```
G = randn(d, s);
```

```
Z = exp(-sqrt(-1)*pi*X*G)/sqrt(s);
```

RANDOM FOURIER FEATURES (RFF) ALGORITHM

Super simple algorithm. For Gaussian kernel $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_2^2}$:



The diagram illustrates the RFF algorithm equation: $\mathbf{Z} = \exp \left(-i\pi \mathbf{X} \begin{bmatrix} \eta_1 \\ \vdots \\ \eta_s \end{bmatrix} \right)$. A yellow box labeled \mathbf{Z} is on the left. To its right is the expression $\exp \left(-i\pi \mathbf{X} \begin{bmatrix} \eta_1 \\ \vdots \\ \eta_s \end{bmatrix} \right)$. Inside the exponent, there is a blue box labeled \mathbf{X} followed by an orange box containing $\eta_1 \dots \eta_s$. An arrow points from the orange box to the text "random gaussians".

$$Z_{x,j} = e^{-i\pi \mathbf{x}^T \boldsymbol{\eta}_j} \text{ for } \boldsymbol{\eta}_j \sim \mathcal{N}$$

```
G = randn(d, s);
```

```
Z = exp(-sqrt(-1)*pi*X*G)/sqrt(s);
```

This is a so-called oblivious sketch.

STARTING POINT

Want to improve the approximation guarantee:

$$\mathbf{Z} \mathbf{Z}^* = \mathbf{K} + \mathbf{E}$$

STARTING POINT

Want to improve the approximation guarantee:

$$\mathbf{Z} \mathbf{Z}^* = \mathbf{K} + \mathbf{E}$$

- Does not give bounds on $\|\mathbf{Z}\mathbf{Z}^* - \mathbf{K}\|$ unless we take $\Omega(n^2)$ samples.

STARTING POINT

Want to improve the approximation guarantee:

$$\mathbf{Z} \mathbf{Z}^* = \mathbf{K} + \mathbf{E}$$

- Does not give bounds on $\|\mathbf{Z}\mathbf{Z}^* - \mathbf{K}\|$ unless we take $\Omega(n^2)$ samples.
- No clear implications for downstream learning tasks. E.g., does $(\mathbf{Z}\mathbf{Z}^* + \lambda\mathbf{I})^{-1}$ approximate $(\mathbf{K} + \lambda\mathbf{I})^{-1}$?

STARTING POINT

Want to improve the approximation guarantee:

$$\mathbf{Z} \mathbf{Z}^* = \mathbf{K} + \mathbf{E}$$

- Does not give bounds on $\|\mathbf{Z}\mathbf{Z}^* - \mathbf{K}\|$ unless we take $\Omega(n^2)$ samples.
- No clear implications for downstream learning tasks. E.g., does $(\mathbf{Z}\mathbf{Z}^* + \lambda\mathbf{I})^{-1}$ approximate $(\mathbf{K} + \lambda\mathbf{I})^{-1}$?
- Faster but less accurate than a Nyström approximation with the same number of samples in practice.

We want matrix-like error bounds, not entrywise bounds.

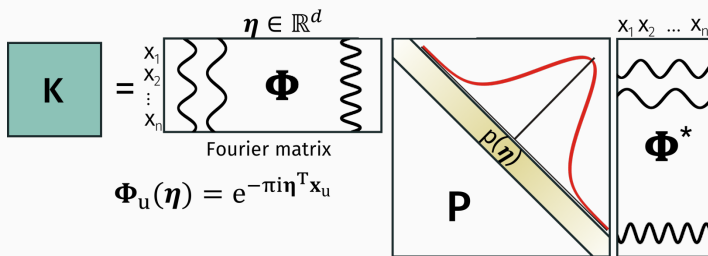
We want matrix-like error bounds, not entrywise bounds.

Key Idea: Analyze the random Fourier features algorithm as a **matrix sampling process**.

$$k(\mathbf{x}_u - \mathbf{x}_v) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-i\pi \boldsymbol{\eta}^T (\mathbf{x}_u - \mathbf{x}_v)} d\boldsymbol{\eta}$$

HIGH LEVEL APPROACH

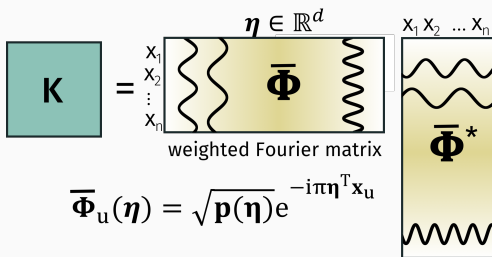
$$k(\mathbf{x}_u - \mathbf{x}_v) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-i\pi \boldsymbol{\eta}^T (\mathbf{x}_u - \mathbf{x}_v)} d\boldsymbol{\eta}$$



Kernel Fourier Transform

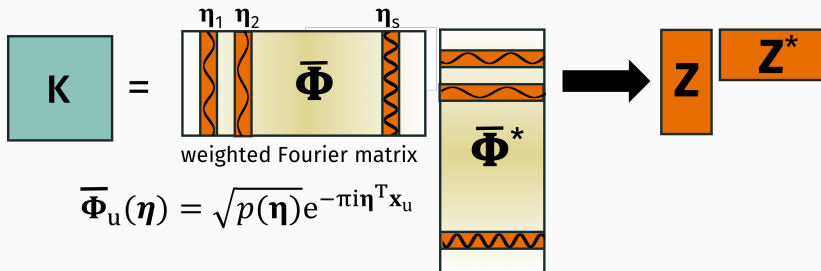
HIGH LEVEL APPROACH

$$k(\mathbf{x}_u - \mathbf{x}_v) = \int_{\boldsymbol{\eta} \in \mathbb{R}^d} p(\boldsymbol{\eta}) e^{-i\pi \boldsymbol{\eta}^T (\mathbf{x}_u - \mathbf{x}_v)} d\boldsymbol{\eta}$$

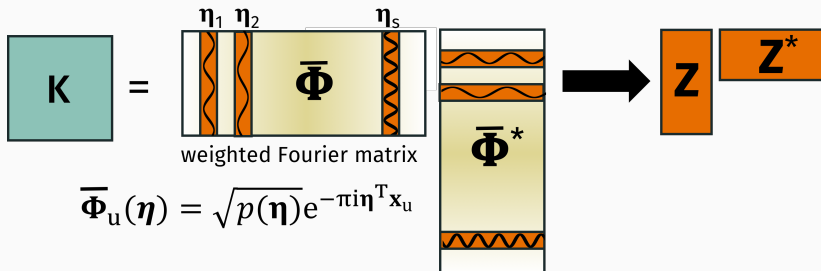


Kernel Fourier Transform

STANDARD RFF = COLUMN NORM SAMPLING

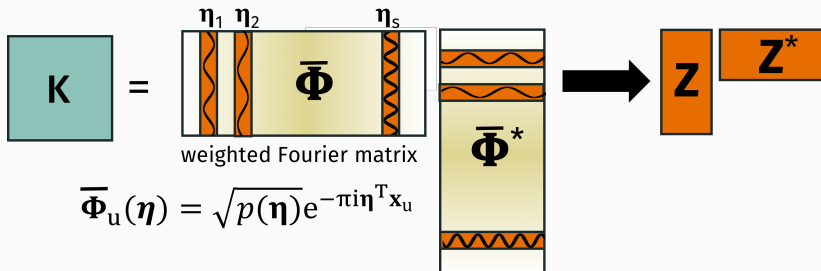


STANDARD RFF = COLUMN NORM SAMPLING



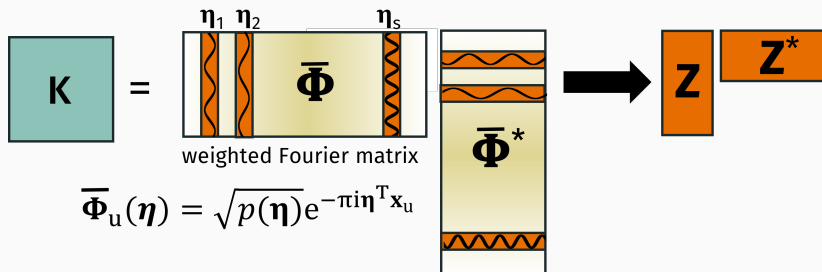
Standard RFF selects column $\bar{\Phi}(\eta)$ with probability $\propto p(\eta)$.

STANDARD RFF = COLUMN NORM SAMPLING



Standard RFF selects column $\bar{\Phi}(\eta)$ with probability $\propto p(\eta)$.
 (with probability proportional to the column's **squared norm**).

STANDARD RFF = COLUMN NORM SAMPLING



Standard RFF selects column $\bar{\Phi}(\eta)$ with probability $\propto p(\eta)$.
(with probability proportional to the column's **squared norm**).

Simple matrix Chernoff already gives better bounds:

$$\|K - ZZ^*\|_2 \leq \epsilon \text{ with } \tilde{O}(n/\epsilon^2) \text{ samples.}$$

For matrix approximation, norm based sampling probabilities are known to be suboptimal!

LEVERAGE SCORE SAMPLING

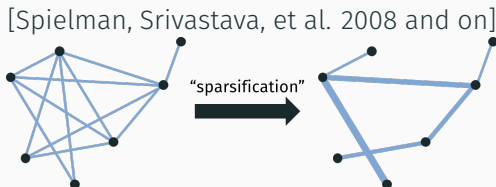
For matrix approximation, norm based sampling probabilities are known to be suboptimal!

We can obtain better theoretical bounds and empirical performance if we sample by column **leverage scores**.

LEVERAGE SCORE SAMPLING

For matrix approximation, norm based sampling probabilities are known to be suboptimal!

We can obtain better theoretical bounds and empirical performance if we sample by column **leverage scores**.



(see work on effective resistances for spectral graph sparsification, randomized linear algebra, etc.)

λ -Ridge Leverage Score Sampling:

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

λ -Ridge Leverage Score Sampling:

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

- Since $A \preceq B \Rightarrow B^{-1} \preceq A^{-1}$, $(\mathbf{Z}\mathbf{Z}^* + \lambda\mathbf{I})^{-1} \approx (\mathbf{K} + \lambda\mathbf{I})^{-1}$

λ -Ridge Leverage Score Sampling:

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

- Since $A \preceq B \Rightarrow B^{-1} \preceq A^{-1}$, $(\mathbf{Z}\mathbf{Z}^* + \lambda\mathbf{I})^{-1} \approx (\mathbf{K} + \lambda\mathbf{I})^{-1}$
- Spectral approximation also gives bounds for kernel PCA, k-means, CCA, etc. (Cohen, Musco, Musco '16,'17)

λ -Ridge Leverage Score Sampling:

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

- Since $A \preceq B \Rightarrow B^{-1} \preceq A^{-1}$, $(\mathbf{Z}\mathbf{Z}^* + \lambda\mathbf{I})^{-1} \approx (\mathbf{K} + \lambda\mathbf{I})^{-1}$
- Spectral approximation also gives bounds for kernel PCA, k-means, CCA, etc. (Cohen, Musco, Musco '16,'17)

Need $\tilde{O}(s_\lambda/\epsilon^2)$ samples where s_λ is the **statistical dimension**:

$$s_\lambda = \text{tr}(\mathbf{K}(\mathbf{K} + \lambda\mathbf{I})^{-1})$$

λ -Ridge Leverage Score Sampling:

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

- Since $A \preceq B \Rightarrow B^{-1} \preceq A^{-1}$, $(\mathbf{Z}\mathbf{Z}^* + \lambda\mathbf{I})^{-1} \approx (\mathbf{K} + \lambda\mathbf{I})^{-1}$
- Spectral approximation also gives bounds for kernel PCA, k-means, CCA, etc. (Cohen, Musco, Musco '16,'17)

Need $\tilde{O}(s_\lambda/\epsilon^2)$ samples where s_λ is the **statistical dimension**:

$$s_\lambda = \text{tr}(\mathbf{K}(\mathbf{K} + \lambda\mathbf{I})^{-1}) = \sum_{i=1}^n \frac{\sigma_i(\mathbf{K})}{\sigma_i(\mathbf{K}) + \lambda}$$

λ -Ridge Leverage Score Sampling:

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

- Since $A \preceq B \Rightarrow B^{-1} \preceq A^{-1}$, $(\mathbf{Z}\mathbf{Z}^* + \lambda\mathbf{I})^{-1} \approx (\mathbf{K} + \lambda\mathbf{I})^{-1}$
- Spectral approximation also gives bounds for kernel PCA, k-means, CCA, etc. (Cohen, Musco, Musco '16,'17)

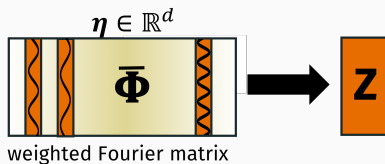
Need $\tilde{O}(s_\lambda/\epsilon^2)$ samples where s_λ is the **statistical dimension**:

$$s_\lambda = \text{tr}(\mathbf{K}(\mathbf{K} + \lambda\mathbf{I})^{-1}) = \sum_{i=1}^n \frac{\sigma_i(\mathbf{K})}{\sigma_i(\mathbf{K}) + \lambda}$$

Roughly the number of singular values $\geq \lambda$ plus a term depending on the tail singular values.

COMPUTING RIDGE LEVERAGE SCORES

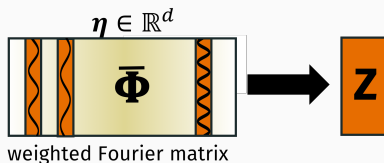
What are the ridge leverage scores?



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

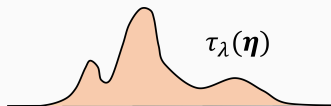
COMPUTING RIDGE LEVERAGE SCORES

What are the ridge leverage scores?



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

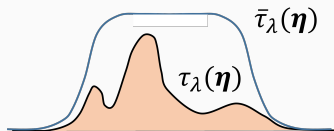
Expensive to invert $K + \lambda I$. Even if you could, not at all clear how to efficiently sample from the leverage score distribution.



Goal: Upper bound Fourier ridge leverage scores for common kernels with simple distributions.

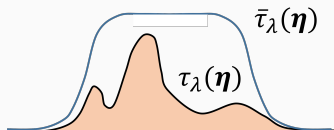
BOUNDING THE FOURIER LEVERAGE SCORES

Goal: Upper bound Fourier ridge leverage scores for common kernels with simple distributions.



BOUNDING THE FOURIER LEVERAGE SCORES

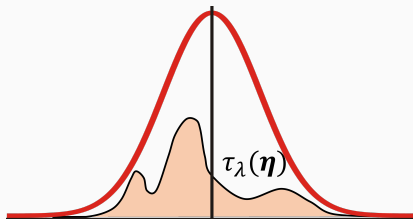
Goal: Upper bound Fourier ridge leverage scores for common kernels with simple distributions.



More closely match leverage score sampling to improve random Fourier features.

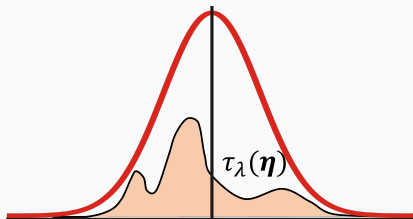
BOUNDING THE FOURIER LEVERAGE SCORES

First observation: Scaled by $\frac{n}{\lambda}$, the standard Rahimi Recht distribution upper bounds the λ -ridge leverage scores.



BOUNDED THE FOURIER LEVERAGE SCORES

First observation: Scaled by $\frac{n}{\lambda}$, the standard Rahimi Recht distribution upper bounds the λ -ridge leverage scores.



Basic Result: Sampling $O(\frac{n}{\lambda} \cdot \frac{1}{\epsilon^2})$ frequencies with RFF gives spectral guarantees.

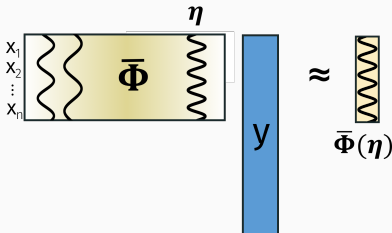
Ridge leverage score $\bar{\Phi}(\boldsymbol{\eta})^*(\mathbf{K} + \lambda\mathbf{I})^{-1}\bar{\Phi}(\boldsymbol{\eta})$ solves:

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

ALTERNATIVE CHARACTERIZATION OF LEVERAGE SCORES

Ridge leverage score $\bar{\Phi}(\eta)^*(K + \lambda I)^{-1}\bar{\Phi}(\eta)$ solves:

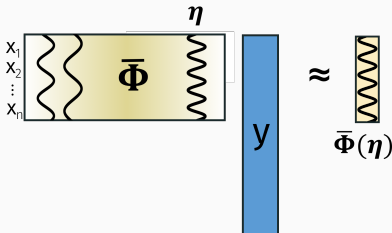
$$\tau_{\lambda}(\eta) = \min_y \left[\frac{1}{\lambda} \|\bar{\Phi}y - \bar{\Phi}(\eta)\|_2^2 + \|y\|_2^2 \right].$$



ALTERNATIVE CHARACTERIZATION OF LEVERAGE SCORES

Ridge leverage score $\bar{\Phi}(\boldsymbol{\eta})^*(\mathbf{K} + \lambda\mathbf{I})^{-1}\bar{\Phi}(\boldsymbol{\eta})$ solves:

$$\tau_{\lambda}(\boldsymbol{\eta}) = \min_y \left[\frac{1}{\lambda} \|\bar{\Phi}\mathbf{y} - \bar{\Phi}(\boldsymbol{\eta})\|_2^2 + \|\mathbf{y}\|_2^2 \right].$$

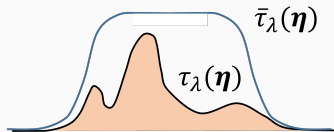


Intuition: \mathbf{y} reconstructs frequency $\boldsymbol{\eta}$ from other frequencies.

If $\boldsymbol{\eta}$ is “easy” to reconstruct, it is less important to sample.

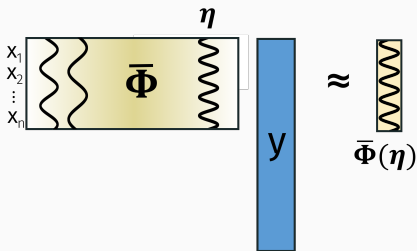
Approach: Obtain closed form upper bound on leverage scores by exhibiting simple candidate vector $\tilde{\mathbf{y}}$ and noting that:

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



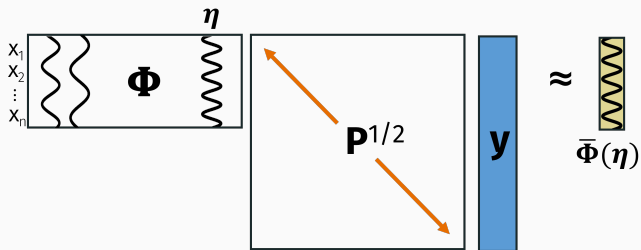
ALTERNATIVE CHARACTERIZATION OF LEVERAGE SCORES

$\bar{\Phi}y$ is just the Fourier transform of $\sqrt{P}y$ evaluated at x_1, \dots, x_n !



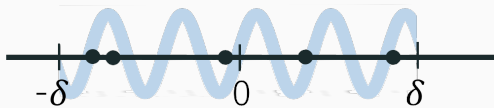
ALTERNATIVE CHARACTERIZATION OF LEVERAGE SCORES

$\bar{\Phi}y$ is just the Fourier transform of $\sqrt{P}y$ evaluated at x_1, \dots, x_n !



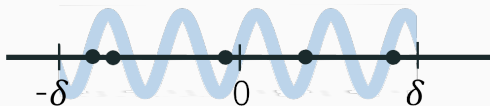
SIMPLIFYING ASSUMPTION

For remainder of talk: Assume Gaussian kernel, n points in 1-dimension, bounded between $[-\delta, \delta]$.



SIMPLIFYING ASSUMPTION

For remainder of talk: Assume Gaussian kernel, n points in 1-dimension, bounded between $[-\delta, \delta]$.



We want a function whose Fourier transform matches frequency η on these data points.

Need to bound: $\frac{1}{\lambda} \|\bar{\Phi} \tilde{\mathbf{y}} - \bar{\Phi}(\boldsymbol{\eta})\|_2^2 + \|\tilde{\mathbf{y}}\|_2^2$.

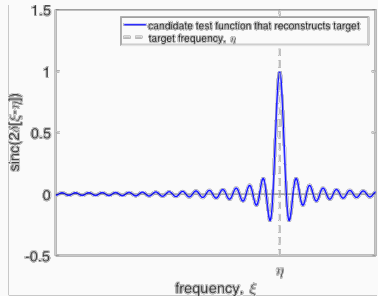
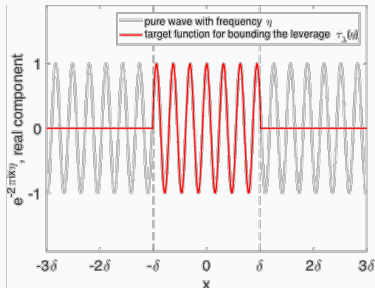
Need to bound: $\frac{1}{\lambda} \|\bar{\Phi} \tilde{\mathbf{y}} - \bar{\Phi}(\boldsymbol{\eta})\|_2^2 + \|\tilde{\mathbf{y}}\|_2^2$.

Set $\sqrt{\mathbf{P}}\tilde{\mathbf{y}}$ to be a shifted sinc function.

FIRST ATTEMPT

Need to bound: $\frac{1}{\lambda} \|\bar{\Phi} \tilde{y} - \bar{\Phi}(\eta)\|_2^2 + \|\tilde{y}\|_2^2$.

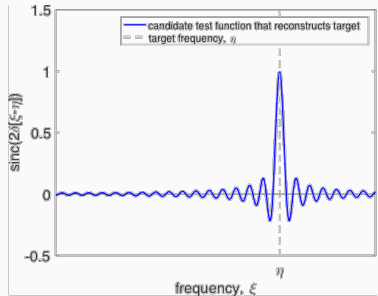
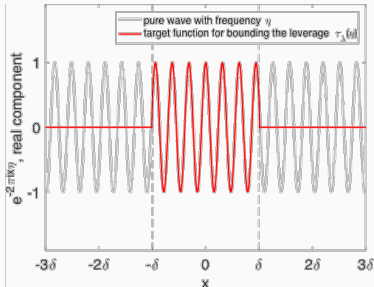
Set $\sqrt{P} \tilde{y}$ to be a **shifted sinc function**.



FIRST ATTEMPT

Need to bound: $\frac{1}{\lambda} \|\bar{\Phi}\tilde{y} - \bar{\Phi}(\eta)\|_2^2 + \|\tilde{y}\|_2^2$.

Set $\sqrt{P}\tilde{y}$ to be a **shifted sinc function**.

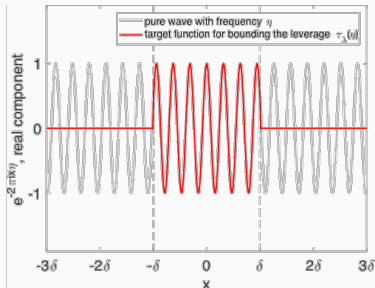


$$\frac{1}{\lambda} \|\bar{\Phi}\tilde{y} - \bar{\Phi}(\eta)\|_2^2 = 0!$$

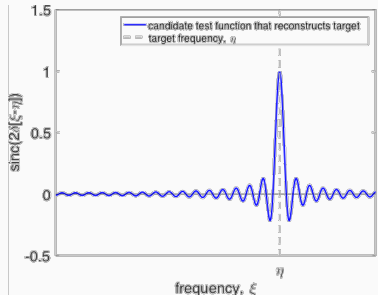
FIRST ATTEMPT

Need to bound: $\frac{1}{\lambda} \|\bar{\Phi} \tilde{\mathbf{y}} - \bar{\Phi}(\eta)\|_2^2 + \|\tilde{\mathbf{y}}\|_2^2$.

Set $\sqrt{\mathbf{P}}\tilde{\mathbf{y}}$ to be a **shifted sinc function**.



$$\frac{1}{\lambda} \|\bar{\Phi} \tilde{\mathbf{y}} - \bar{\Phi}(\eta)\|_2^2 = 0!$$

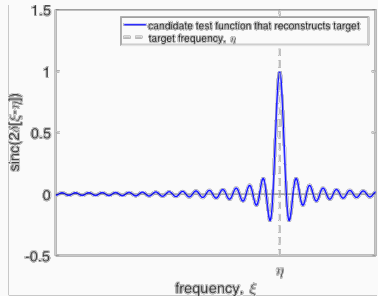
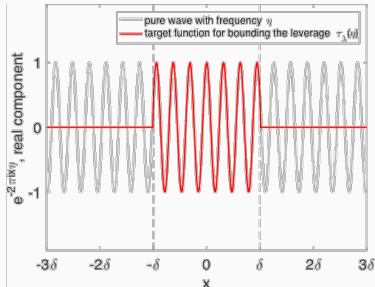


$$\|\sqrt{\mathbf{P}}\tilde{\mathbf{y}}\|_2^2 = O(\delta).$$

FIRST ATTEMPT

Need to bound: $\frac{1}{\lambda} \|\bar{\Phi}\tilde{\mathbf{y}} - \bar{\Phi}(\eta)\|_2^2 + \|\tilde{\mathbf{y}}\|_2^2$.

Set $\sqrt{\mathbf{P}}\tilde{\mathbf{y}}$ to be a **shifted sinc function**.

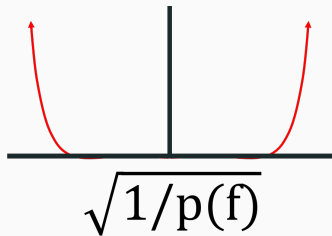


$$\frac{1}{\lambda} \|\bar{\Phi}\tilde{\mathbf{y}} - \bar{\Phi}(\eta)\|_2^2 = 0!$$

$$\|\sqrt{\mathbf{P}}\tilde{\mathbf{y}}\|_2^2 = O(\delta).$$

$\|\tilde{\mathbf{y}}\|_2^2$ is much larger!

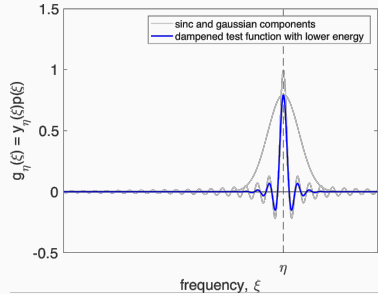
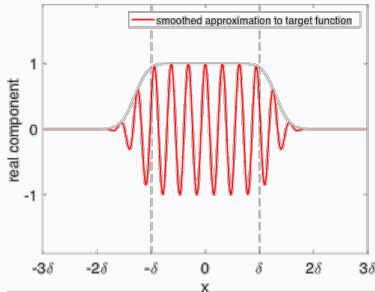
$$\tilde{\mathbf{y}} = \mathbf{P}^{-1/2} \mathbf{P}^{1/2} \tilde{\mathbf{y}} = \mathbf{P}^{-1/2} \cdot \text{sinc function}$$



Sinc falls off as $O(1/f)$, but $\frac{1}{\sqrt{p(f)}}$ grows as $e^{O(f^2)}$,
so $\|\mathbf{y}\|_2^2$ explodes.

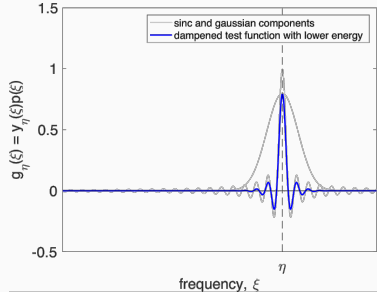
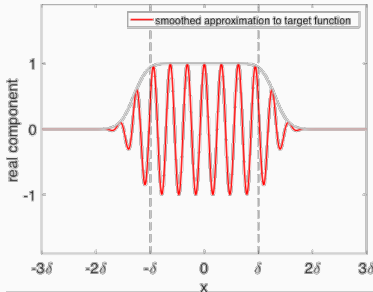
FINAL TEST FUNCTION

Solution: Dampen sinc with a narrow Gaussian.



FINAL TEST FUNCTION

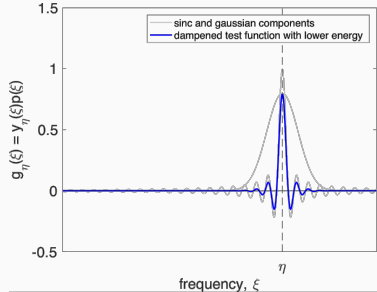
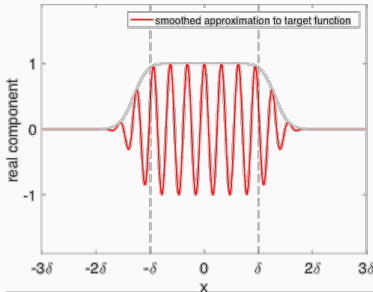
Solution: Dampen sinc with a narrow Gaussian.



$\frac{1}{\lambda} \|\bar{\Phi}\tilde{y} - \bar{\Phi}(\eta)\|_2^2$ remains very small, $\|y\|_2^2 \approx O(\delta)$.

FINAL TEST FUNCTION

Solution: Dampen sinc with a narrow Gaussian.



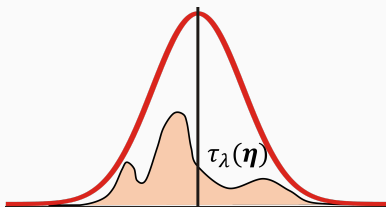
$\frac{1}{\lambda} \|\bar{\Phi}\tilde{y} - \bar{\Phi}(\eta)\|_2^2$ remains very small, $\|y\|_2^2 \approx O(\delta)$.

(as long as η is not too large)

FINAL TEST FUNCTION

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

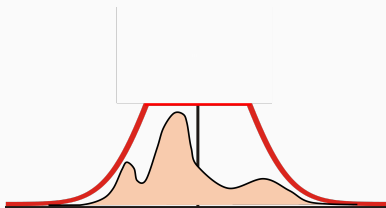
$$\bar{\tau}_\lambda(\boldsymbol{\eta}) = \begin{cases} \tilde{O}(\delta) & \text{when } \boldsymbol{\eta} \leq \sqrt{\log n / \lambda} \\ p(\boldsymbol{\eta}) = e^{-\|\boldsymbol{\eta}\|_2^2 / 2} & \text{otherwise.} \end{cases}$$



FINAL TEST FUNCTION

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

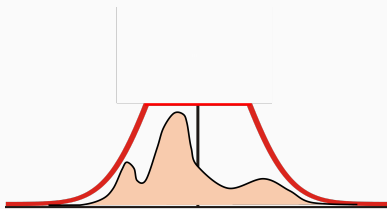
$$\bar{\tau}_\lambda(\boldsymbol{\eta}) = \begin{cases} \tilde{O}(\delta) & \text{when } \boldsymbol{\eta} \leq \sqrt{\log n / \lambda} \\ p(\boldsymbol{\eta}) = e^{-\|\boldsymbol{\eta}\|_2^2 / 2} & \text{otherwise.} \end{cases}$$



FINAL TEST FUNCTION

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

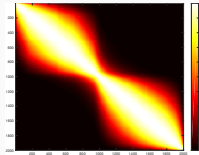
$$\bar{\tau}_\lambda(\boldsymbol{\eta}) = \begin{cases} \tilde{O}(\delta) & \text{when } \boldsymbol{\eta} \leq \sqrt{\log n / \lambda} \\ p(\boldsymbol{\eta}) = e^{-\|\boldsymbol{\eta}\|_2^2 / 2} & \text{otherwise.} \end{cases}$$



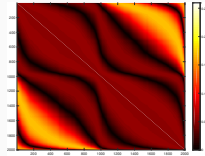
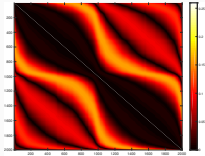
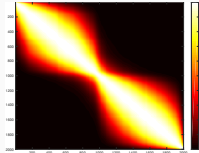
Requires $O(\delta \sqrt{\log(n/\lambda)}) \cdot \frac{1}{\epsilon^2}$ samples for spectral guarantee.

(vs. $O(n/\lambda) \cdot \frac{1}{\epsilon^2}$ for standard random Fourier features.)

Gaussian kernel for two clusters:



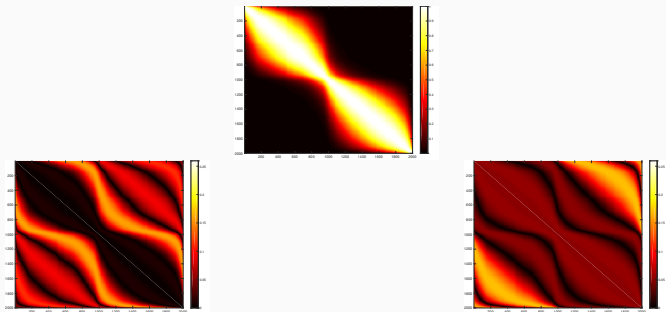
Gaussian kernel for two clusters:



Standard RFF element error.

Modified RFF element error.

Gaussian kernel for two clusters:



Standard RFF element error.

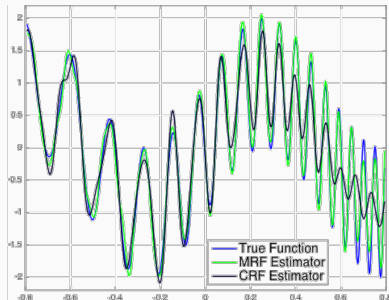
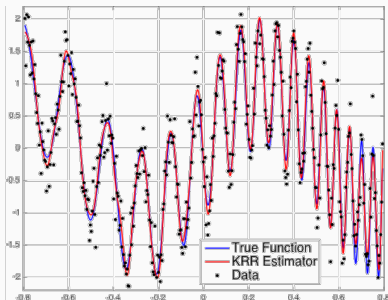
Modified RFF element error.

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

Sampling low frequencies relatively less biases error to align with large eigenvectors.

EXPERIMENTAL RESULTS

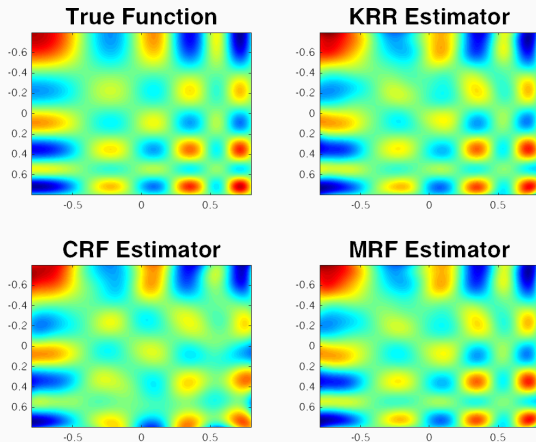
Example of approximate kernel ridge regression to interpolate a synthetic function:



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

EXPERIMENTAL RESULTS

In higher dimensions:



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

1. Viewed random Fourier features methods as a **matrix sampling problem**.

1. Viewed random Fourier features methods as a **matrix sampling problem**.
2. Used optimization perspective on leverage scores to certify upper bounds on these scores. Reduced score computation to **Fourier approximation problem**.

1. Viewed random Fourier features methods as a **matrix sampling problem**.
2. Used optimization perspective on leverage scores to certify upper bounds on these scores. Reduced score computation to **Fourier approximation problem**.
3. New sampling distribution under-samples lower frequencies to obtain better kernel approximations.

Major open question: Can we achieve our spectral guarantee with $O(s_\lambda)$ samples in high dimensions for any data set.

Major open question: Can we achieve our spectral guarantee with $O(s_\lambda)$ samples in high dimensions for any data set.

Conjecture: Yes, although maybe with polynomial loss (i.e. $\text{poly}(s_\lambda)$ samples).

Major open question: Can we achieve our spectral guarantee with $O(s_\lambda)$ samples in high dimensions for any data set.

Conjecture: Yes, although maybe with polynomial loss (i.e. $\text{poly}(s_\lambda)$ samples).

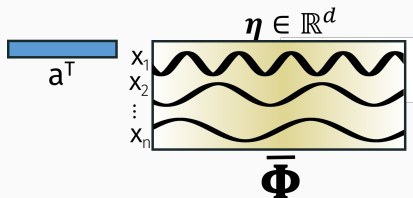
I.e. can we match data-adaptive methods like Nyström obviously?

Dual view of leverage scores:

$$\tau_{\lambda}(\boldsymbol{\eta}) = \max_{\mathbf{a}} \frac{(\mathbf{a}^T \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta}))^2}{\|\mathbf{a}^T \bar{\boldsymbol{\Phi}}\|_2^2 + \lambda \|\mathbf{a}\|_2^2}.$$

Dual view of leverage scores:

$$\tau_\lambda(\boldsymbol{\eta}) = \max_{\mathbf{a}} \frac{(\mathbf{a}^T \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta}))^2}{\|\mathbf{a}^T \bar{\boldsymbol{\Phi}}\|_2^2 + \lambda \|\mathbf{a}\|_2^2}.$$

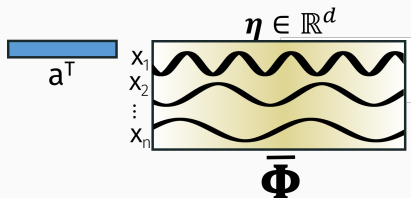


$\mathbf{a}^T \bar{\boldsymbol{\Phi}}$ is an n sparse Fourier function weighted by a Gaussian.

OPEN QUESTIONS

Dual view of leverage scores:

$$\tau_\lambda(\boldsymbol{\eta}) = \max_{\mathbf{a}} \frac{(\mathbf{a}^T \bar{\boldsymbol{\Phi}}(\boldsymbol{\eta}))^2}{\|\mathbf{a}^T \bar{\boldsymbol{\Phi}}\|_2^2 + \lambda \|\mathbf{a}\|_2^2}.$$



$\mathbf{a}^T \bar{\boldsymbol{\Phi}}$ is an n sparse Fourier function weighted by a Gaussian.

Immediately get leverage score bounds from bounds on smoothness of sparse Fourier functions, e.g. [Chen, Kane, Price, Song FOCS 2016].

Vague open question: Why does this all actually matter for function fitting?

Vague open question: Why does this all actually matter for function fitting?

We were surprised to beat random Fourier features on a kernel regression task.

THANK YOU!