

Random Projections

Lopez-Paz & Duvenaud

November 7, 2013

Random Outline

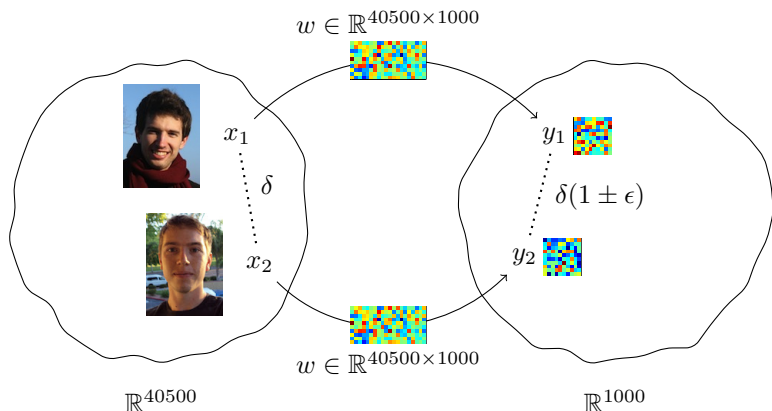
The Johnson-Lindenstrauss Lemma (1984)

Random Kitchen Sinks (Rahimi and Recht, NIPS 2008)

Fastfood (Le et al., ICML 2013)

Why random projections?

Fast, efficient and \mathcal{O} distance-preserving **dimensionality reduction!**

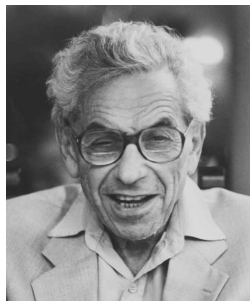


$$(1 - \epsilon)\|x_1 - x_2\|^2 \leq \|y_1 - y_2\|^2 \leq (1 + \epsilon)\|x_1 - x_2\|^2$$

This result is formalized in the *Johnson-Lindenstrauss Lemma*

The Johnson-Lindenstrauss Lemma

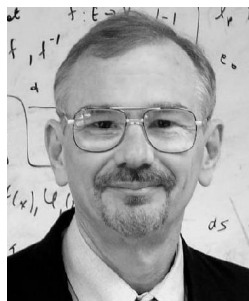
The proof is a great example of Erdős' *probabilistic method* (1947).



Paul Erdős
1913-1996



Joram Lindenstrauss
1936-2012



William B. Johnson
1944-

§12.5 of *Foundations of Machine Learning* (Mohri et al., 2012)

Auxiliary Lemma 1

Let Q be a random variable following a χ^2 distribution with k degrees of freedom. Then, for any $0 < \epsilon < 1/2$:

$$\Pr[(1 - \epsilon)k \leq Q \leq (1 + \epsilon)k] \geq 1 - 2e^{-(\epsilon^2 - \epsilon^3)k/4}.$$

Proof: we start by using *Markov's inequality* ($\Pr[X > a] \leq \frac{E[X]}{a}$):

$$\Pr[Q \geq (1 + \epsilon)k] = \Pr[e^{\lambda Q} \geq e^{\lambda(1+\epsilon)k}] \leq \frac{E[e^{\lambda Q}]}{e^{\lambda(1+\epsilon)k}} = \frac{(1 - 2\lambda)^{-k/2}}{e^{\lambda(1+\epsilon)k}},$$

where $E[e^{\lambda Q}] = (1 - 2\lambda)^{-k/2}$ is the m.g.f. of a χ^2 distribution, $\lambda < \frac{1}{2}$.

To tight the bound we minimize the right-hand side with $\lambda = \frac{\epsilon}{2(1+\epsilon)}$:

$$\Pr[Q \geq (1 + \epsilon)k] \leq \frac{(1 - \frac{\epsilon}{1+\epsilon})^{-k/2}}{e^{\epsilon k/2}} = \frac{(1 + \epsilon)^{k/2}}{(e^\epsilon)^{k/2}} = \left(\frac{1 + \epsilon}{e^\epsilon}\right)^{k/2}.$$

Auxiliary Lemma 1

Using $1 + \epsilon \leq e^{\epsilon - (\epsilon^2 - \epsilon^3)/2}$ yields

$$\Pr[Q \geq (1 + \epsilon)k] \leq \left(\frac{1 + \epsilon}{e^\epsilon}\right)^{k/2} \leq \left(\frac{e^{\epsilon - \frac{\epsilon^2 - \epsilon^3}{2}}}{e^\epsilon}\right) = e^{-\frac{k}{4}(\epsilon^2 - \epsilon^3)}.$$

$\Pr[Q \leq (1 - \epsilon)k]$ is bounded similarly, and the lemma follows by applying the union bound:

$$\begin{aligned} \Pr[\overline{(1 - \epsilon)k \leq Q \leq (1 + \epsilon)k}] &\leq \\ \Pr[Q \geq (1 + \epsilon)k \cup Q \leq (1 - \epsilon)k] &\leq \\ \Pr[Q \geq (1 + \epsilon)k] + \Pr[Q \leq (1 - \epsilon)k] &= \\ &2e^{-\frac{k}{4}(\epsilon^2 - \epsilon^3)} \end{aligned}$$

Then,

$$\Pr[(1 - \epsilon)k \leq Q \leq (1 + \epsilon)k] \geq 1 - 2e^{-\frac{k}{4}(\epsilon^2 - \epsilon^3)}$$



Auxiliary Lemma 2

Let $\mathbf{x} \in \mathbb{R}^N$, $k < N$ and $\mathbf{A} \in \mathbb{R}^{k \times N}$ with $A_{ij} \sim \mathcal{N}(0, 1)$. Then, for any $0 \leq \epsilon \leq 1/2$:

$$\Pr[(1 - \epsilon)\|\mathbf{x}\|^2 \leq \|\frac{1}{\sqrt{k}}\mathbf{A}\mathbf{x}\|^2 \leq (1 + \epsilon)\|\mathbf{x}\|^2] \geq 1 - 2e^{-(\epsilon^2 - \epsilon^3)k/4}.$$

Proof: let $\hat{\mathbf{x}} = \mathbf{A}\mathbf{x}$. Then,

$$E[\hat{x}_j^2] = E\left[\left(\sum_i^N A_{ji}x_i\right)^2\right] = E\left[\sum_i^N A_{ji}^2x_i^2\right] = \sum_i^N x_i^2 = \|\mathbf{x}\|^2.$$

Note that $T_j = \hat{x}_j/\|\mathbf{x}\| \sim \mathcal{N}(0, 1)$. Then, $Q = \sum_i^k T_j^2 \sim \chi_k^2$.

Remember the previous lemma?

Auxiliary Lemma 2

Remember: $\hat{\mathbf{x}} = \mathbf{A}\mathbf{x}$, $T_j = \hat{x}_j/\|\mathbf{x}\| \sim \mathcal{N}(0, 1)$, $Q = \sum_i^k T_j^2 \sim \chi_k^2$:

$$\Pr[(1 - \epsilon)\|\mathbf{x}\|^2 \leq \|\frac{1}{\sqrt{k}}\mathbf{A}\mathbf{x}\|^2 \leq (1 + \epsilon)\|\mathbf{x}\|^2] =$$

$$\Pr[(1 - \epsilon)\|\mathbf{x}\|^2 \leq \frac{\|\hat{\mathbf{x}}\|^2}{k} \leq (1 + \epsilon)\|\mathbf{x}\|^2] =$$

$$\Pr[(1 - \epsilon)k \leq \frac{\|\hat{\mathbf{x}}\|^2}{\|\mathbf{x}\|^2} \leq (1 + \epsilon)k] =$$

$$\Pr \left[(1 - \epsilon)k \leq \sum_i^k T_j^2 \leq (1 + \epsilon)k \right] =$$

$$\Pr [(1 - \epsilon)k \leq Q \leq (1 + \epsilon)k] \geq$$

$$1 - 2e^{-(\epsilon^2 - \epsilon^3)k/4}$$



The Johnson-Lindenstrauss Lemma

For any $0 < \epsilon < 1/2$ and any integer $m > 4$, let $k = \frac{20 \log m}{\epsilon^2}$. Then, for any set V of m points in $\mathbb{R}^N \exists f : \mathbb{R}^N \rightarrow \mathbb{R}^k$ s.t. $\forall \mathbf{u}, \mathbf{v} \in V$:

$$(1 - \epsilon)\|\mathbf{u} - \mathbf{v}\|^2 \leq \|f(\mathbf{u}) - f(\mathbf{v})\|^2 \leq (1 + \epsilon)\|\mathbf{u} - \mathbf{v}\|^2.$$

Proof: Let $f = \frac{1}{\sqrt{k}}\mathbf{A}$, $\mathbf{A} \in \mathbb{R}^{k \times N}$, $k < N$ and $A_{ij} \sim \mathcal{N}(0, 1)$.

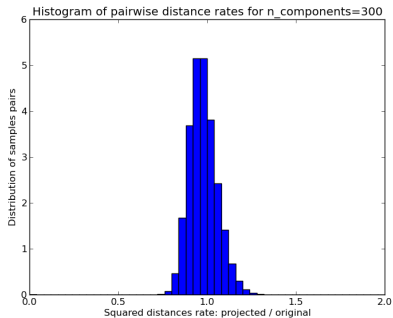
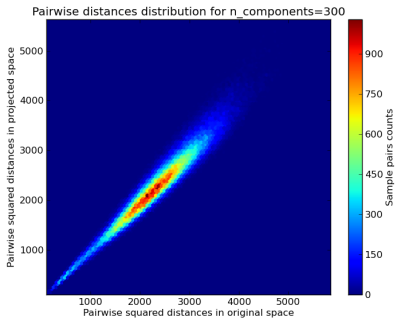
- For fixed $\mathbf{u}, \mathbf{v} \in V$, we apply the previous lemma with $\mathbf{x} = \mathbf{u} - \mathbf{v}$ to lower bound the *success* probability by $1 - 2e^{-(\epsilon^2 - \epsilon^3)k/4}$.
- Union bound again! This time over the m^2 pairs in V with $k = \frac{20 \log m}{\epsilon^2}$ and $\epsilon < 1/2$ to obtain:

$$\Pr[\text{success}] \geq 1 - 2m^2 e^{-(\epsilon^2 - \epsilon^3)k/4} = 1 - 2m^{5\epsilon - 3} > 1 - 2m^{-1/2} > 0.$$



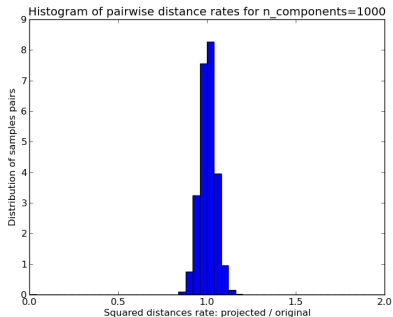
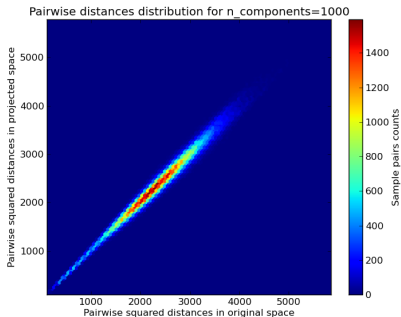
JL Experiments

Data: 20-newsgroups, from 100.000 features to 300 (0.3%)



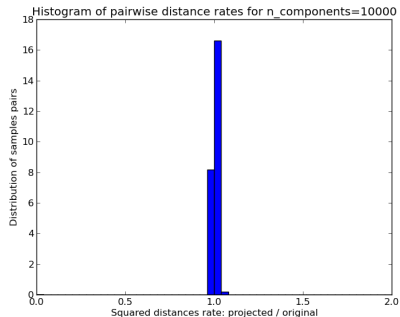
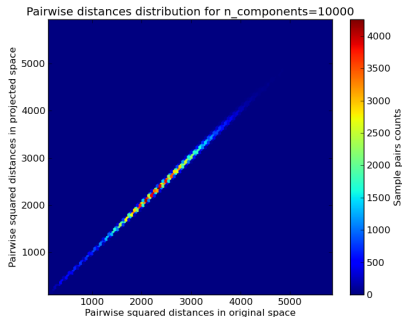
JL Experiments

Data: 20-newsgroups, from 100.000 features to 1.000 (1%)



JL Experiments

Data: 20-newsgroups, from 100.000 features to 10.000 (10%)



MATLAB implementation: `1/sqrt(k) .* randn(k, N) % %X.`

JL Conclusions

- Do you have a huge feature space?
- Are you wasting too much time with PCA?
- Random Projections are fast, compact & efficient!
- Monograph (Vampala, 2004)
- Sparse Random Projections (Achlioptas, 2003)
- Random Projections can Improve MoG! (Dasgupta, 2000)
- Code for previous experiments: <http://bit.ly/17FTfbH>

But... What about **non-linear** random projections?



Ali Rahimi



Ben Recht

Random Outline

The Johnson-Lindenstrauss Lemma (1984)

Random Kitchen Sinks (Rahimi and Recht, NIPS 2008)

Fastfood (Le et al., ICML 2013)

A Familiar Creature

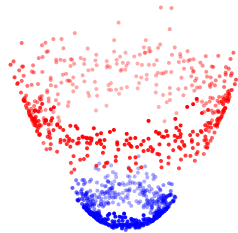
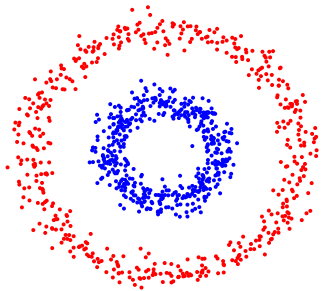
$$f(\mathbf{x}) = \sum_{i=1}^T \alpha_i \phi(\mathbf{x}; \mathbf{w}_i)$$

- Gaussian Process
- Kernel Regression
- SVM
- AdaBoost
- Multilayer Perceptron
- ...

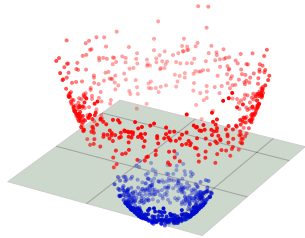
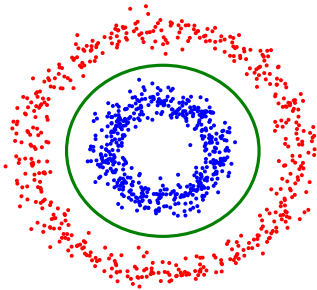
Same model, different training approaches!

Things get interesting when ϕ is non-linear...

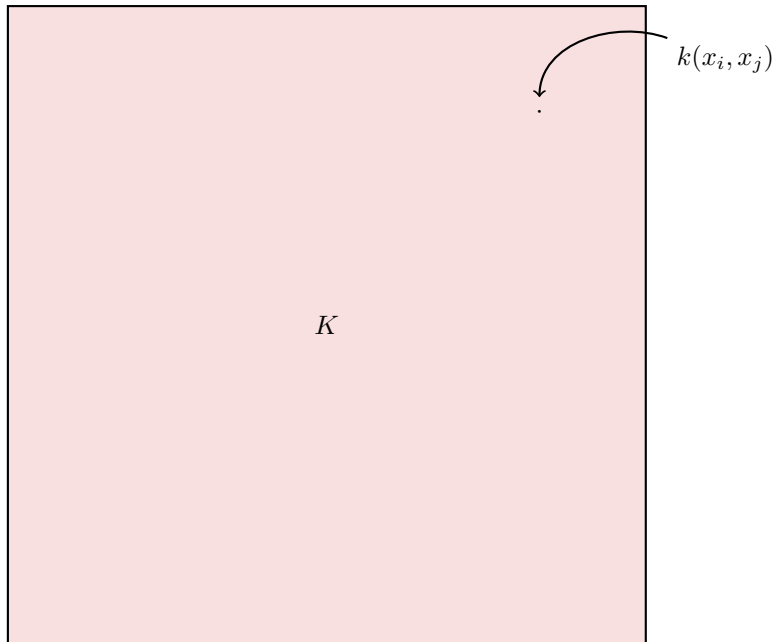
A Familiar Solution



A Familiar Solution



A Familiar Monster: The Kernel Trap



Greedy Approximation of Functions

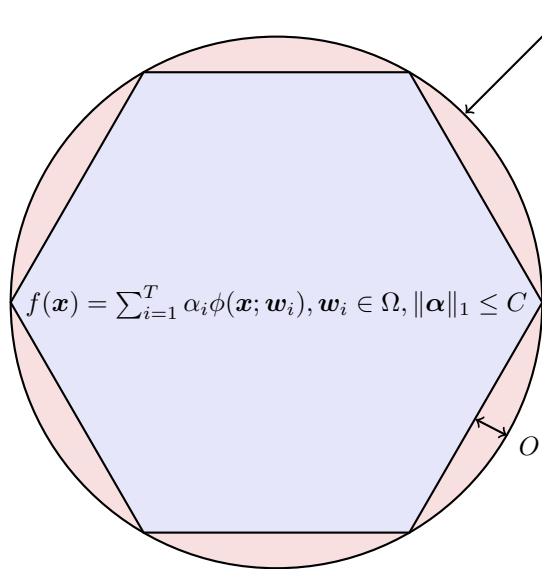
Approx. $f(\mathbf{x}) = \sum_{i=1}^{\infty} \alpha_i \phi(\mathbf{x}; \mathbf{w}_i)$ with $f_T(\mathbf{x}) = \sum_{i=1}^T \alpha_i \phi(\mathbf{x}; \mathbf{w}_i)$.

Greedy Fitting

$$(\boldsymbol{\alpha}^*, \mathbf{W}^*) = \min_{\boldsymbol{\alpha}, \mathbf{W}} \left\| \sum_{i=1}^T \alpha_i \phi(\cdot; \mathbf{w}_i) - f \right\|_{\mu}$$

Greedy Approximation of Functions

$$\mathcal{F} \equiv \{f(\mathbf{x}) = \sum_{i=1}^{\infty} \alpha_i \phi(\mathbf{x}; \mathbf{w}_i), \mathbf{w}_i \in \Omega, \|\alpha\|_1 \leq C\}$$



$$f(\mathbf{x}) = \sum_{i=1}^T \alpha_i \phi(\mathbf{x}; \mathbf{w}_i), \mathbf{w}_i \in \Omega, \|\alpha\|_1 \leq C$$

$$\begin{aligned} \|f_T - f\|_{\mu} &= \\ \sqrt{\int_{\mathcal{X}} (f_T(x) - f(x))^2 \mu(dx)} &= \\ O\left(\frac{C}{\sqrt{T}}\right) & \text{ (Jones, 1992)} \end{aligned}$$

RKS Approximation of Functions

Approx. $f(\mathbf{x}) = \sum_{i=1}^{\infty} \alpha_i \phi(\mathbf{x}; \mathbf{w}_i)$ with $f_T(\mathbf{x}) = \sum_{i=1}^T \alpha_i \phi(\mathbf{x}; \mathbf{w}_i)$.

Greedy Fitting

$$(\boldsymbol{\alpha}^*, \mathbf{W}^*) = \min_{\boldsymbol{\alpha}, \mathbf{W}} \left\| \sum_{i=1}^T \alpha_i \phi(\cdot; \mathbf{w}_i) - f \right\|_{\mu}$$

Random Kitchen Sinks Fitting

$$\mathbf{w}_i^*, \dots, \mathbf{w}_T^* \sim p(\mathbf{w}), \quad \boldsymbol{\alpha}^* = \min_{\boldsymbol{\alpha}} \left\| \sum_{i=1}^T \alpha_i \phi(\cdot; \mathbf{w}_i^*) - f \right\|_{\mu}$$

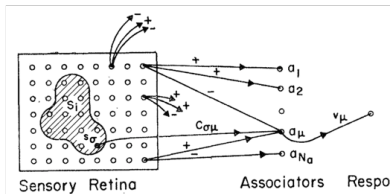
The Perceptron: A Model for Brain Functioning. I*

H. D. BLOCK

Cornell University, Ithaca, New York

THE Perceptron is a self-organizing or adaptive system proposed by Rosenblatt.¹ Its primary purpose is to shed some light on the problem of explaining brain function in terms of brain structure. It also has technological applications as a pattern-recognizing device, but here our emphasis is on the brain function-structure problem. The technological aspects are not completely irrelevant however, since a model, no matter how appealing it may appear from the point of view of structural similarity, must also be judged on the basis of its performance.

In brief a perceptron consists of a *retina* of *sensory units* (for example photocells); these are connected (for example by wires) to *associator units*. The connections are **many to many and random**. The associator units may be connected to each other or to *response units*. When a *stimulus* is presented to the retina



Just an old idea?

Liquid state machine

From Wikipedia, the free encyclopedia

A **liquid state machine (LSM)** is a computational **construct like a neural network**. An LSM consists of a large collection of units (called *nodes*, or *neurons*). Each node receives time varying input from external sources (the **inputs**) as well as from other nodes. **Nodes are randomly connected to each other**. The **recurrent** nature of the connections turns the time varying input into a spatio-temporal pattern of activations in the network nodes. The spatio-temporal patterns of activation are read out by **linear discriminant** units.

The soup of recurrently connected nodes will end up computing a large variety of **nonlinear functions** on the input. Given a **large enough variety** of such nonlinear functions, it is theoretically possible to obtain linear combinations (using the read out units) to perform whatever mathematical operation is needed to perform a certain task, such as **speech recognition** or **computer vision**.

W. Maass, T. Natschlaeger, H. Markram, Real-time computing without stable states: A new framework for neural computation based on perturbations, *Neural Computation*. **14**(11), 2531-2560, (2002)

How does RKS work?

For functions $f(\mathbf{x}) = \int_{\Omega} \alpha(\mathbf{w})\phi(\mathbf{x}; \mathbf{w})d\mathbf{w}$, define the p -norm as:

$$\|f\|_p = \sup_{\mathbf{w} \in \Omega} \frac{|\alpha(\mathbf{w})|}{p(\mathbf{w})}$$

and let \mathcal{F}_p be all f with $\|f\|_p \leq C$. Then, for $\mathbf{w}_1, \dots, \mathbf{w}_T \sim p(\mathbf{w})$, w.p. at least $1 - \delta$, $\delta > 0$, there exist some α s.t. f_T satisfies:

$$\|f_T - f\|_{\mu} = O\left(\frac{C}{\sqrt{T}} \sqrt{1 + 2 \log \frac{1}{\delta}}\right)$$

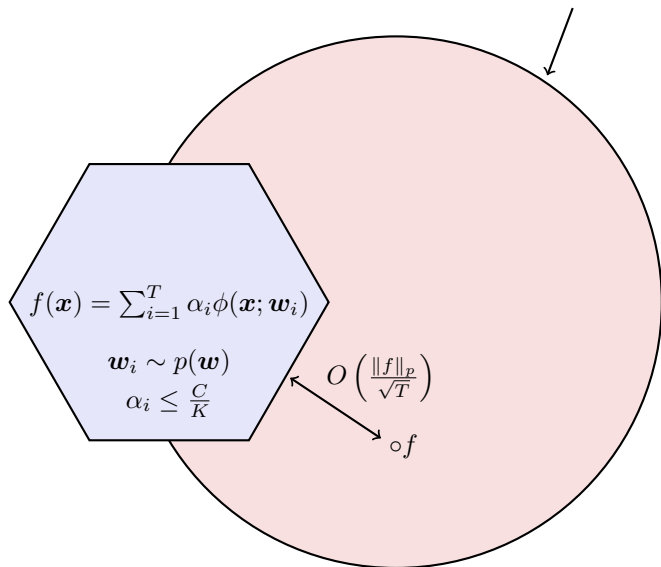
Why? Set $\alpha_i = \frac{1}{T}\alpha(\mathbf{w}_i)$. Then, the discrepancy is given by $\|f\|_p$:

$$f_T(\mathbf{x}) = \frac{1}{T} \sum_i^T a_i(\mathbf{w}_i)\phi(\mathbf{x}; \mathbf{w}_i)$$

With a dataset of size N , an error $O\left(\frac{1}{\sqrt{N}}\right)$ is added to all bounds.

RKS Approximation of Functions

$$\mathcal{F}_p \equiv \{f(\mathbf{x}) = \int_{\Omega} \alpha(\mathbf{w})\phi(\mathbf{x}; \mathbf{w})d\mathbf{w}, |\alpha(\mathbf{w})| \leq Cp(\mathbf{w})\}$$



Random Kitchen Sinks: Auxiliary Lemma

Let $\mathbf{X} = \{x_1, \dots, x_K\}$ be i.i.d. r.v. drawn from a centered C -radius ball of a Hilbert Space \mathcal{H} . Let $\bar{\mathbf{X}} = \frac{1}{K} \sum_k x_k$. Then, for any $\delta > 0$, with probability at least $1 - \delta$:

$$\|\bar{\mathbf{X}} - E[\bar{\mathbf{X}}]\| \leq \frac{C}{\sqrt{K}} \left(1 + \sqrt{2 \log \frac{1}{\delta}} \right)$$

Proof: Show that $f(\mathbf{X}) = \|\bar{\mathbf{X}} - E[\bar{\mathbf{X}}]\|$ is stable w.r.t. perturbations:

$$|f(\mathbf{X}) - f(\tilde{\mathbf{X}})| \leq \|\bar{\mathbf{X}} - \tilde{\bar{\mathbf{X}}}\| \leq \frac{\|x_i - x'_i\|}{K} \leq \frac{2C}{K}.$$

Second, the variance of the average of i.i.d. random variables is:

$$E[\|\bar{\mathbf{X}} - E[\bar{\mathbf{X}}]\|^2] = \frac{1}{K} (E[\|x\|^2] - \|E[x]\|^2).$$

Third, using Jensen's inequality and given that $\|x\| \leq C$:

$$E[f(\mathbf{X})] \leq \sqrt{E[f^2(\mathbf{X})]} = \sqrt{E[\|\bar{\mathbf{X}} - E[\bar{\mathbf{X}}]\|^2]} \leq \frac{C}{\sqrt{K}}$$

Fourth, use McDiarmid's inequality and rearrange.

Random Kitchen Sinks Proof

Let μ be a measure on \mathcal{X} , and $f^* \in \mathcal{F}_p$. Let $\mathbf{w}_1, \dots, \mathbf{w}_T \sim p(\mathbf{w})$. Then, w.p. at least $1 - \delta$, with $\delta > 0$, $\exists f_T(x) = \sum_i^T \beta_i \phi(\mathbf{x}; \mathbf{w}_i)$ s.t.:

$$\|f_T - f^*\|_\mu \leq \frac{C}{\sqrt{T}} \left(1 + \sqrt{2 \log \frac{1}{\delta}} \right)$$

Proof:

Let $f_i = \beta_i \phi(\mathbf{x}; \mathbf{w}_i)$, $1 \leq k \leq T$ and $\beta_i = \frac{\alpha(\mathbf{w}_i)}{p(\mathbf{w}_i)}$. Then, $E[f_i] = f^*$:

$$E[f_i] = E_{\mathbf{w}} \left[\frac{\alpha(\mathbf{w}_i)}{p(\mathbf{w}_i)} \phi(\cdot; \mathbf{w}_i) \right] = \int_{\Omega} p(\mathbf{w}) \frac{\alpha(\mathbf{w})}{p(\mathbf{w})} \phi(\cdot; \mathbf{w}) d\mathbf{w} = f^*$$

The claim is mainly completed by describing the concentration of the average $f_T = \frac{1}{T} \sum f_i$ around f^* with the previous lemma.

Approximating Kernels with RKS

Bochner's Theorem: A kernel $k(\mathbf{x} - \mathbf{y})$ on \mathbb{R}^d is PSD if and only if $k(\mathbf{x} - \mathbf{y})$ is the Fourier transform of a non-negative measure $p(\mathbf{w})$.

$$\begin{aligned}k(\mathbf{x} - \mathbf{y}) &= \int_{\mathbb{R}^d} p(\mathbf{w}) e^{j\mathbf{w}'(\mathbf{x} - \mathbf{y})} d\mathbf{w} \\ &\approx \frac{1}{T} \sum_{i=1}^T e^{j\mathbf{w}'_i(\mathbf{x} - \mathbf{y})} \quad (\text{Monte-Carlo, } O(T^{-1/2})) \\ &= \frac{1}{T} \sum_{i=1}^T \underbrace{e^{j\mathbf{w}'_i \mathbf{x}}}_{\phi(\mathbf{x}; \mathbf{W})} \underbrace{e^{-j\mathbf{w}'_i \mathbf{y}}}_{\phi(\mathbf{y}; \mathbf{W})} \\ &= \frac{1}{\sqrt{T}} \phi(\mathbf{x}; \mathbf{W})^* \frac{1}{\sqrt{T}} \phi(\mathbf{y}; \mathbf{W})\end{aligned}$$

Now solve least squares in the primal in $O(n)$ time!

Random Kitchen Sinks: Implementation

```
% Training
function ytest = kitchen_sinks( X, y, Xtest, T, noise)

Z = randn(T, size(X,1)); % Sample feature frequencies
phi = exp(i*Z*X); % Compute feature matrix

% Linear regression with observation noise.
w = (phi*phi' + eye(T)*noise)\(phi*y);

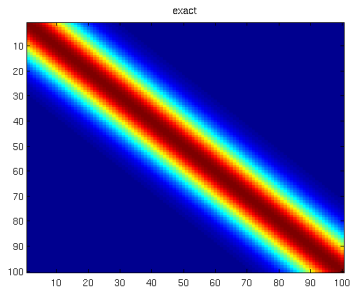
% testing
ytest = w'*exp(i*Z*xtest);
```

from <http://www.keysduplicated.com/~ali/random-features/>

- That's fast, approximate GP regression! (with a sq-exp kernel)
- Or linear regression with $[\sin(zx), \cos(zx)]$ feature pairs
- Show demo

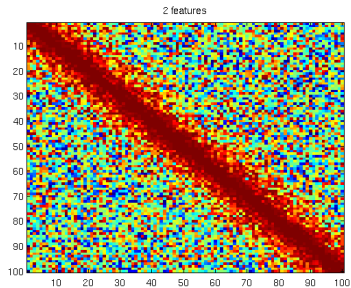
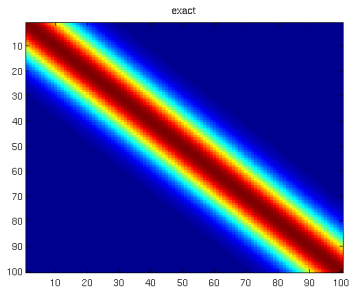
Random Kitchen Sinks and Gram Matrices

- How fast do we approach the exact Gram matrix?
- $k(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^\top \Phi(\mathbf{X})$



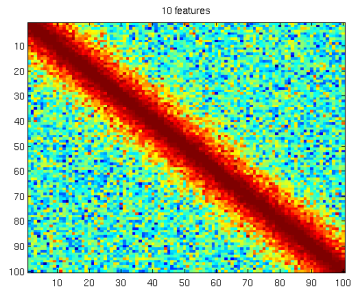
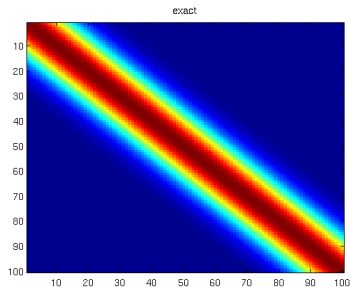
Random Kitchen Sinks and Gram Matrices

- How fast do we approach the exact Gram matrix?
- $k(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^\top \Phi(\mathbf{X})$



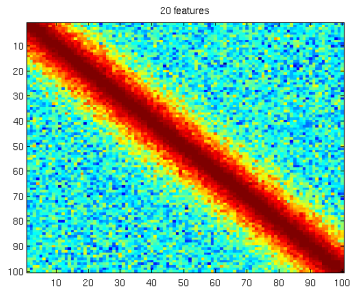
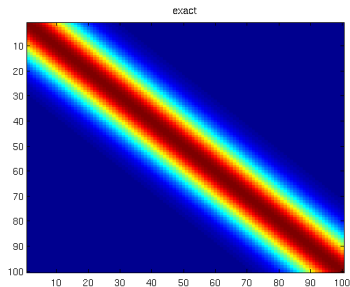
Random Kitchen Sinks and Gram Matrices

- How fast do we approach the exact Gram matrix?
- $k(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^\top \Phi(\mathbf{X})$



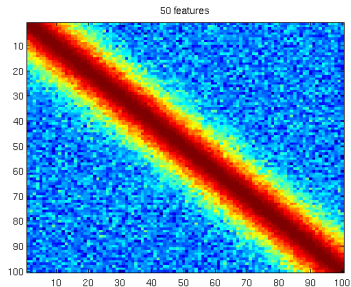
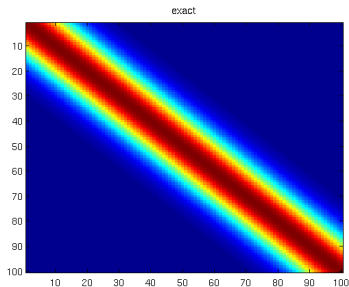
Random Kitchen Sinks and Gram Matrices

- How fast do we approach the exact Gram matrix?
- $k(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^\top \Phi(\mathbf{X})$



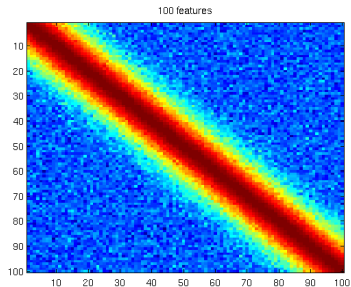
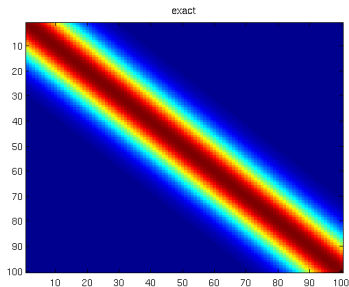
Random Kitchen Sinks and Gram Matrices

- How fast do we approach the exact Gram matrix?
- $k(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^\top \Phi(\mathbf{X})$



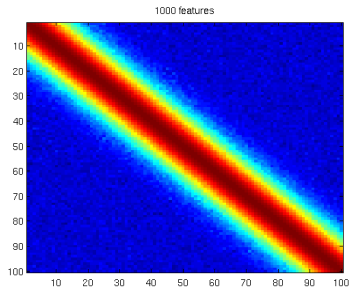
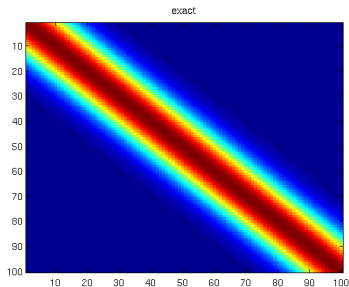
Random Kitchen Sinks and Gram Matrices

- How fast do we approach the exact Gram matrix?
- $k(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^T \Phi(\mathbf{X})$



Random Kitchen Sinks and Gram Matrices

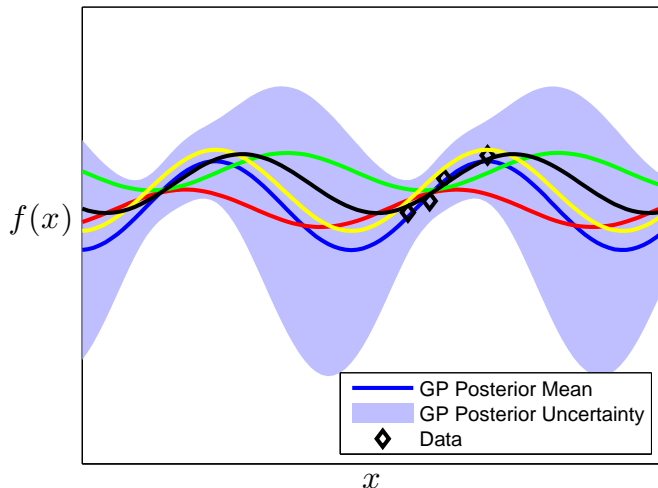
- How fast do we approach the exact Gram matrix?
- $k(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^\top \Phi(\mathbf{X})$



Random Kitchen Sinks and Posterior

- How fast do we approach the exact Posterior?

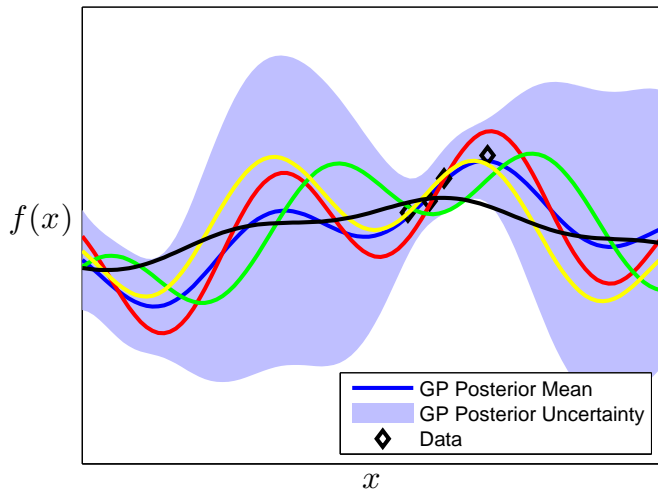
3 features



Random Kitchen Sinks and Posterior

- How fast do we approach the exact Posterior?

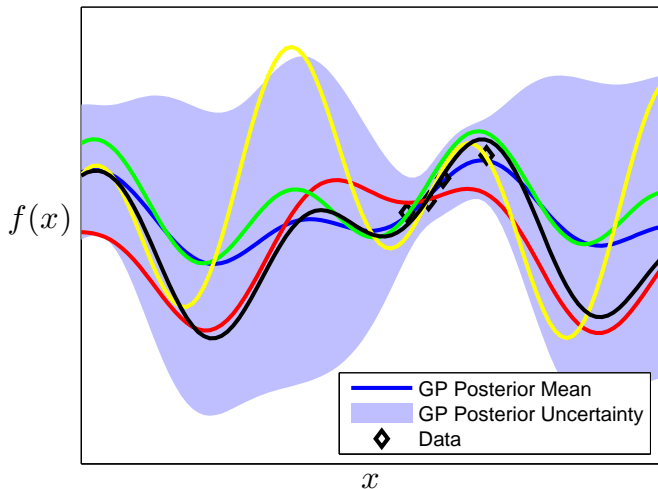
4 features



Random Kitchen Sinks and Posterior

- How fast do we approach the exact Posterior?

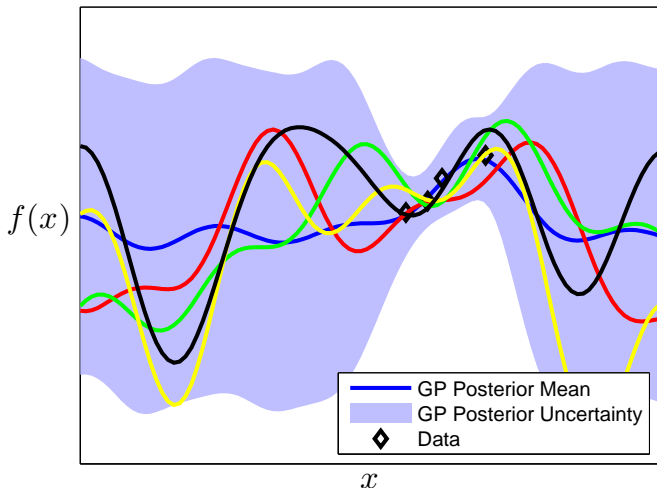
5 features



Random Kitchen Sinks and Posterior

- How fast do we approach the exact Posterior?

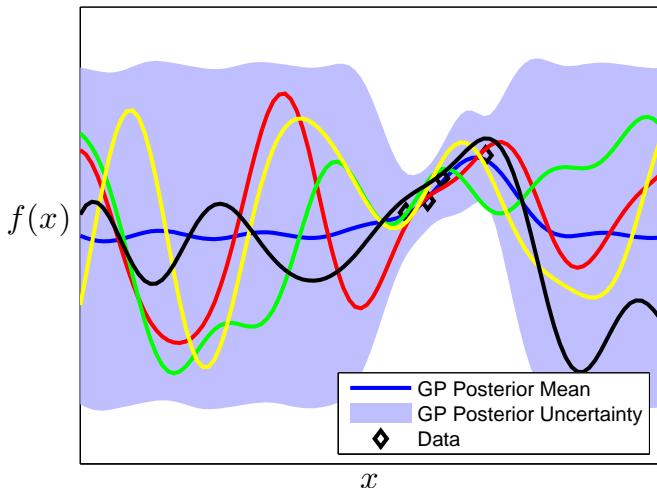
10 features



Random Kitchen Sinks and Posterior

- How fast do we approach the exact Posterior?

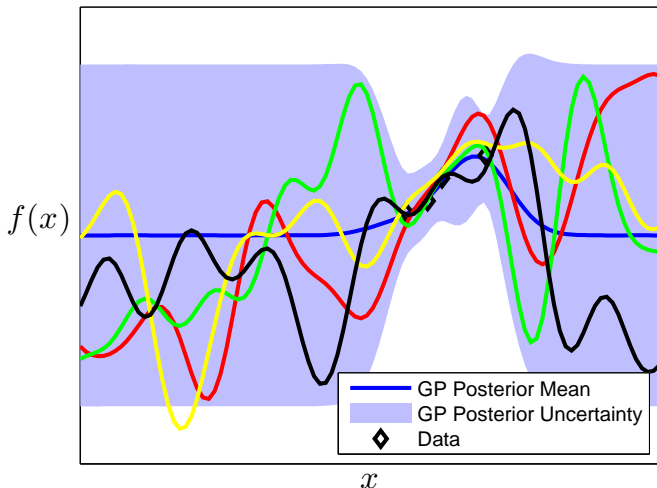
20 features



Random Kitchen Sinks and Posterior

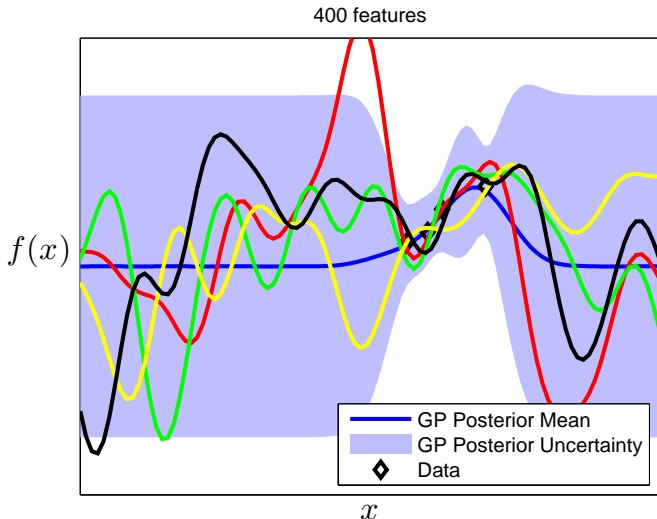
- How fast do we approach the exact Posterior?

200 features



Random Kitchen Sinks and Posterior

- How fast do we approach the exact Posterior?

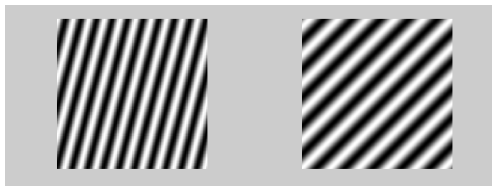


Kitchen Sinks in multiple dimensions

D dimensions, T random features, N datapoints

- First: Sample $T \times D$ random numbers: $Z \sim \mathcal{N}(0, \sigma^{-2})$
- each $\phi_j(\mathbf{x}) = \exp(i[Zx]_j)$
- or: $\Phi(\mathbf{x}) = \exp(iZx)$

Each feature $\phi_j(\cdot)$ is a sine, cos wave varying along the direction given by one row of Z , with varying periods.



Can be slow for many features in high dimensions. For example, computing $\Phi(\mathbf{x})$ is $\mathcal{O}(NTD)$.

But isn't linear regression already fast?

D dimensions, T features, N Datapoints

- Computing features: $\Phi(\mathbf{x}) = \frac{1}{\sqrt{T}} \exp(iZx)$
- Regression: $\mathbf{w} = [\Phi(\mathbf{x})^T \Phi(\mathbf{x})]^{-1} \Phi(\mathbf{x})^T \mathbf{y}$
- Train time complexity:

$$\mathcal{O}\left(\underbrace{DTN}_{\text{Computing Features}} + \underbrace{T^3}_{\text{Inverting Covariance Matrix}} + \underbrace{T^2N}_{\text{Multiplication}} \right)$$

- Prediction: $\mathbf{y}^* = \mathbf{w}^T \Phi(\mathbf{x}^*)$
- Test time complexity: $\mathcal{O}\left(\underbrace{DTN^*}_{\text{Computing Features}} + \underbrace{T^2N^*}_{\text{Multiplication}} \right)$
- For images, D is often $> 10,000$.

Random Outline

The Johnson-Lindenstrauss Lemma (1984)

Random Kitchen Sinks (Rahimi and Recht, NIPS 2008)

Fastfood (Le et al., ICML 2013)

Fastfood (Q. Le, T. Sarlós, A. Smola, 2013)

Main idea: Approximately compute Zx quickly, by replacing Z with some easy-to-compute matrices.

- Uses the Subsampled Randomized Hadamard Transform (SRHT) (Sarlós, 2006)

- Train time complexity:

$$\mathcal{O}\left(\underbrace{\log(D)TN}_{\text{Computing Features}} + \underbrace{T^3}_{\text{Inverting Feature Matrix}} + \underbrace{T^2N}_{\text{Multiplication}} \right)$$

- Test time complexity: $\mathcal{O}\left(\underbrace{\log(D)TN^*}_{\text{Computing Features}} + \underbrace{T^2N^*}_{\text{Multiplication}} \right)$

- For images, if $D = 10,000$, $\log_{10}(D) = 4$.

Fastfood (Q. Le, T. Sarlós, A. Smola, 2013)

Main idea: Compute $Vx \approx Zx$ quickly, by building V out of easy-to-compute matrices.

V has similar properties to the Gaussian matrix Z .

$$V = \frac{1}{\sigma\sqrt{d}} SHG\Pi HB \quad (1)$$

- Π is a $D \times D$ permutation matrix
- G is diagonal random Gaussian.
- B is diagonal random $\{+1, -1\}$.
- S is diagonal random scaling.
- H is Walsh-Hadamard Matrix.

Fastfood (Q. Le, T. Sarlós, A. Smola, 2013)

Main idea: Compute $Vx \approx Zx$ quickly, by building V out of easy-to-compute matrices.

$$V = \frac{1}{\sigma\sqrt{d}} SHG\Pi HB \quad (2)$$

- H is Walsh-Hadamard Matrix. Multiplication in $\mathcal{O}(T \log(D))$. H is orthogonal, can be built recursively.

$$H_3 = \frac{1}{2^{\frac{3}{2}}} \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 \\ 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 \\ 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 \\ 1 & -1 & -1 & 1 & -1 & 1 & 1 & -1 \end{pmatrix} \quad (3)$$

Fastfood (Q. Le, T. Sarlós, A. Smola, 2013)

Main idea: Compute $Vx \approx Zx$ quickly, by building V out of easy-to-compute matrices.

$$Vx = \frac{1}{\sigma\sqrt{d}} SHG\Pi HBx \quad (4)$$

Intuition: Scramble a single Gaussian random vector many different ways, to create a matrix with similar properties to Z . [Draw on board]

- $HG\Pi HB$ produces pseudo-random Gaussian vectors (of identical length)
- S fixes the lengths to have the correct distribution.

Fastfood results

Computing Vx :

d	T	Fastfood	RKS	Speedup	RAM
1024	16384	0.00058s	0.0139s	24x	256x
4096	32768	0.00136s	0.1224s	90x	1024x
8192	65536	0.00268s	0.5360s	200x	2048x

We never store V !

Fastfood results

Regression MSE:

Dataset	m	d	Exact RBF	Nystrom RBF
Insurance Company	5,822	85	0.231	0.232
Wine Quality	4,080	11	0.819	0.797
Parkinson Telemonitor	4,700	21	0.059	0.058
CPU	6,554	21	7.271	6.758
Location of CT slices (axial)	42,800	384	n.a.	60.683
KEGG Metabolic Network	51,686	27	n.a.	17.872
Year Prediction MSD	463,715	90	n.a.	0.113
Forest	522,910	54	n.a.	0.837

Random Kitchen Sinks (RBF)	Fastfood FFT	Fastfood RBF	Exact Matern	Fastfood Matern
0.266	0.266	0.264	0.234	0.235
0.740	0.721	0.740	0.753	0.720
0.054	0.052	0.054	0.053	0.052
7.103	4.544	7.366	4.345	4.211
49.491	58.425	43.858	n.a.	14.868
17.837	17.826	17.818	n.a.	17.846
0.123	0.106	0.115	n.a.	0.116
0.840	0.838	0.840	n.a.	0.976

“The Trouble with Kernels” (Smola)

- Kernel Expansion: $f(x) = \sum_{i=1}^m \alpha_i k(x_i, x)$
- Feature Expansion: $f(x) = \sum_{i=1}^T \mathbf{w}_i \phi_i(\mathbf{x})$

D dimensions, T features, N samples

Method	Train Time	Test Time	Train Mem	Test Mem
Naive	$\mathcal{O}(N^2 D)$	$\mathcal{O}(ND)$	$\mathcal{O}(ND)$	$\mathcal{O}(ND)$
Low Rank	$\mathcal{O}(NTD)$	$\mathcal{O}(TD)$	$\mathcal{O}(TD)$	$\mathcal{O}(TD)$
Kitchen Sinks	$\mathcal{O}(NTD)$	$\mathcal{O}(TD)$	$\mathcal{O}(TD)$	$\mathcal{O}(TD)$
Fastfood	$\mathcal{O}(NT \log(D))$	$\mathcal{O}(T \log(D))$	$\mathcal{O}(T \log(D))$	$\mathcal{O}(T)$

“Can run on a phone”

Feature transforms of more interesting kernels

In high dimensions, all Euclidian distances become the same, unless data lie on manifold. Usually need structured kernel.

- can do any stationary kernel (Matérn, rational-quadratic) with Fastfood

- sum of kernels is concatenation of features:

$$k^{(+)}(\mathbf{x}, \mathbf{x}') = k^{(1)}(\mathbf{x}, \mathbf{x}') + k^{(2)}(\mathbf{x}, \mathbf{x}') \implies \Phi^{(+)} = \begin{bmatrix} \Phi^{(1)}(\mathbf{x}) \\ \Phi^{(2)}(\mathbf{x}) \end{bmatrix}$$

- product of kernels is outer product of features:

$$k^{(\times)}(\mathbf{x}, \mathbf{x}') = k^{(1)}(\mathbf{x}, \mathbf{x}')k^{(2)}(\mathbf{x}, \mathbf{x}') \implies \Phi^{(\times)}(\mathbf{x}) = \Phi^{(1)}(\mathbf{x})\Phi^{(2)}(\mathbf{x})^{\top}$$

For example, can build translation-invariant kernel:

$$f\left(\begin{array}{|c|c|c|c|c|} \hline \square & \square & \square & \square & \square \\ \hline \square & \blacksquare & \square & \square & \square \\ \hline \square & \blacksquare & \blacksquare & \square & \square \\ \hline \square & \square & \blacksquare & \blacksquare & \square \\ \hline \square & \square & \square & \square & \square \\ \hline \end{array}\right) = f\left(\begin{array}{|c|c|c|c|c|} \hline \square & \square & \square & \square & \square \\ \hline \square & \square & \square & \blacksquare & \square \\ \hline \square & \square & \square & \blacksquare & \square \\ \hline \square & \square & \square & \square & \blacksquare \\ \hline \square & \square & \square & \square & \square \\ \hline \end{array}\right) \quad (5)$$

$$k((x_1, x_2, \dots, x_D), (x'_1, x'_2, \dots, x'_D)) = \sum_{i=1}^D \prod_{j=1}^D k(x_j, x'_{i+j \bmod D}) \quad (6)$$

Takeaways

Random Projections

- Preserve Euclidian distances
- while reducing dimensionality
- Allow for nonlinear mappings

Random Features

- RKS can approximate GP posterior quickly
- Fastfood can compute T nonlinear basis functions in $\mathcal{O}(T \log D)$ time.
- Can operate on structured kernels.

Gourmet cuisine (exact inference) is nice,
but often fastfood is good enough.