

Random features for large-scale kernel machines

Rahimi, Recht

(NIPS 2007)

April 25, 2013

Introduction

In kernel methods, learned functions take the form

$$f(x) = \sum_i \alpha_i k(x, x_i) = \sum_i \alpha_i \langle \phi(x), \phi(x_i) \rangle_{\mathcal{H}}$$

for training points x_i .

- 1 Advantage: can work with infinite feature spaces.
- 2 Disadvantage: need to store all the training points.

Introduction

In kernel methods, learned functions take the form

$$f(x) = \sum_i \alpha_i k(x, x_i) = \sum_i \alpha_i \langle \phi(x), \phi(x_i) \rangle_{\mathcal{H}}$$

for training points x_i .

- 1 Advantage: can work with infinite feature spaces.
- 2 Disadvantage: need to store all the training points.

Ways to get around this:

- 1 Throw points away (incomplete Cholesky, sparse methods,...)
- 2 This paper: finite random feature spaces

Method 1: Fourier space

Bochner's theorem: a **continuous** kernel $k(\mathbf{x} - \mathbf{y})$ on \mathbb{R}^d is positive definite iff

$$k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{R}^d} p(\omega) e^{i\omega^\top (\mathbf{x} - \mathbf{y})} d\omega$$

for a **probability measure** $p(\omega)$ (actually a finite non-negative Borel measure: prob. measure with appropriate normalization)

Method 1: Fourier space

Bochner's theorem: a **continuous** kernel $k(\mathbf{x} - \mathbf{y})$ on \mathbb{R}^d is positive definite iff

$$k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{R}^d} p(\omega) e^{i\omega^\top (\mathbf{x} - \mathbf{y})} d\omega$$

for a **probability measure** $p(\omega)$ (actually a finite non-negative Borel measure: prob. measure with appropriate normalization)

Define $\zeta_\omega := e^{i\omega^\top \mathbf{x}}$. Then

$$\begin{aligned} k(\mathbf{x} - \mathbf{y}) &= \mathbb{E}_\omega \left[\left(e^{i\omega^\top \mathbf{x}} \right) \left(e^{i\omega^\top \mathbf{y}} \right)^* \right] \\ &= \mathbb{E}_\omega \left(\cos(\omega^\top (\mathbf{x} - \mathbf{y})) \right) + \underbrace{i \mathbb{E}_\omega \left(\sin(\omega^\top (\mathbf{x} - \mathbf{y})) \right)}_{=0}. \end{aligned}$$

Method 1: Fourier space

Because $k(\mathbf{x} - \mathbf{y})$ is real and $p(\omega)$ is real, can replace this with cosine features:

$$z_{\omega,b}(\mathbf{x}) = \sqrt{2} \cos(\omega^\top \mathbf{x} + b)$$

where b uniform on $[0, 2\pi)$

Then

$$k(\mathbf{x} - \mathbf{y}) = \mathbb{E}_{\omega,b} [z_{\omega,b}(\mathbf{x})z_{\omega,b}(\mathbf{y})]$$

Method 1: Fourier space

Because $k(\mathbf{x} - \mathbf{y})$ is real and $p(\omega)$ is real, can replace this with cosine features:

$$z_{\omega,b}(\mathbf{x}) = \sqrt{2} \cos(\omega^\top \mathbf{x} + b)$$

where b uniform on $[0, 2\pi)$

Then

$$k(\mathbf{x} - \mathbf{y}) = \mathbb{E}_{\omega,b} [z_{\omega,b}(\mathbf{x})z_{\omega,b}(\mathbf{y})]$$

Proof:

$$2 \cos(\omega^\top \mathbf{x} + b) \cos(\omega^\top \mathbf{y} + b) = \underbrace{\cos(\omega^\top (\mathbf{x} + \mathbf{y}) + 2b)}_{\text{expectation zero}} + \cos(\omega^\top (\mathbf{x} - \mathbf{y}))$$

Method 1: Fourier space

Generate D random features to decrease variance. Then

$$k(\mathbf{x} - \mathbf{y}) \approx \frac{1}{D} \sum_{j=1}^D z_{\omega,b}^{(j)}(\mathbf{x}) z_{\omega,b}^{(j)}(\mathbf{y}).$$

Method 1: Fourier space

Generate D random features to decrease variance. Then

$$k(\mathbf{x} - \mathbf{y}) \approx \frac{1}{D} \sum_{j=1}^D z_{\omega,b}^{(j)}(\mathbf{x}) z_{\omega,b}^{(j)}(\mathbf{y}).$$

Convergence result:

Claim 1 (Uniform convergence of Fourier features). *Let \mathcal{M} be a compact subset of \mathcal{R}^d with diameter $\text{diam}(\mathcal{M})$. Then, for the mapping \mathbf{z} defined in Algorithm 1, we have*

$$\Pr \left[\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{M}} |\mathbf{z}(\mathbf{x})' \mathbf{z}(\mathbf{y}) - k(\mathbf{x}, \mathbf{y})| \geq \epsilon \right] \leq 2^8 \left(\frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \right)^2 \exp \left(-\frac{D\epsilon^2}{4(d+2)} \right),$$

where $\sigma_p^2 \equiv E_p[\omega' \omega]$ is the second moment of the Fourier transform of k . Further, $\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{M}} |\mathbf{z}(\mathbf{x})' \mathbf{z}(\mathbf{y}) - k(\mathbf{y}, \mathbf{x})| \leq \epsilon$ with any constant probability when $D = \Omega \left(\frac{d}{\epsilon^2} \log \frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \right)$.

Method 2: randomly shifted grid

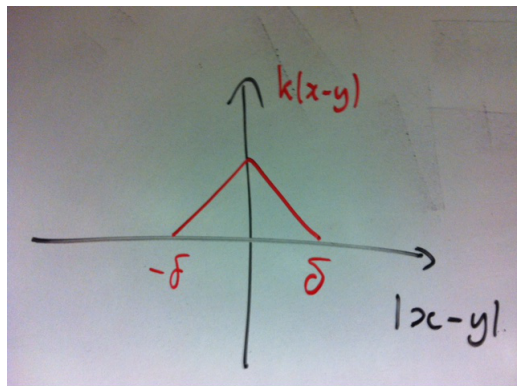


Figure: Kernel $k_{\text{hat}}(x - y) = \max\left(0, 1 - \frac{|x-y|}{\delta}\right)$.

Method 2: randomly shifted grid

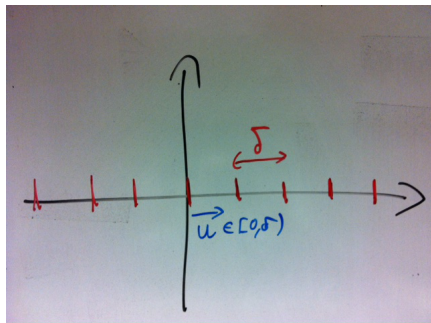


Figure: Randomly shifted grid. $u \sim \mathcal{U}(0, \delta)$.

Probability of x, y falling in the same bin:

$$\Pr_u(\hat{x} = \hat{y} \mid \delta) = k_{\text{hat}}(x - y) \quad \hat{x} = \left\lfloor \frac{x - u}{\delta} \right\rfloor.$$

Method 2: randomly shifted grid

As before, take distributions over features to get more advanced kernels:

$$k(x, y) = \int_0^\infty k_{\text{hat}}(x, y; \delta) p(\delta) d\delta.$$

Given a kernel, how to compute $p(\delta)$?

Method 2: randomly shifted grid

As before, take distributions over features to get more advanced kernels:

$$k(x, y) = \int_0^\infty k_{\text{hat}}(x, y; \delta) p(\delta) d\delta.$$

Given a kernel, how to compute $p(\delta)$?

$$\begin{aligned} k(|x - y|) &=: k(\Delta) \\ &= \int_0^\infty \max\left(0, 1 - \frac{\Delta}{\delta}\right) p(\delta) d\delta \\ &= \int_\Delta^\infty p(\delta) d\delta - \Delta \int_\Delta^\infty \frac{p(\delta)}{\delta} d\delta. \end{aligned}$$

Take 2nd derivative wrt Δ :

$$\frac{d^2 k}{d\Delta^2} = \frac{p(\Delta)}{\Delta} \quad \implies \quad p(\Delta) = \Delta \frac{d^2 k}{d\Delta^2}$$

Method 2: randomly shifted grid

Example:

$$k_{\text{lap}} = \exp(-|x - y|) = \exp(-\Delta)$$

then $p(\delta) = \delta \exp(-\delta)$ (Gamma distribution).

Note: for a Gaussian, $p(\delta)$ not a valid prob. density.

Method 2: randomly shifted grid

Example:

$$k_{\text{lap}} = \exp(-|x - y|) = \exp(-\Delta)$$

then $p(\delta) = \delta \exp(-\delta)$ (Gamma distribution).

Note: for a Gaussian, $p(\delta)$ not a valid prob. density.

Reduce variance by averaging over P independent grids (u, δ) .

Method 2: randomly shifted grid

Example:

$$k_{\text{lap}} = \exp(-|x - y|) = \exp(-\Delta)$$

then $p(\delta) = \delta \exp(-\delta)$ (Gamma distribution).

Note: for a Gaussian, $p(\delta)$ not a valid prob. density.

Reduce variance by averaging over P independent grids (u, δ) .

Multiple dimensions: use independent grids in each dimension, and

$$k(\mathbf{x} - \mathbf{y}) = \prod_{k=1}^m k_m(x^m - y^m).$$

The feature is an m -dimensional binary tensor with a single one at coordinate $\left[\left[\frac{x_1 - u_1}{\delta_1} \right] \quad \dots \quad \left[\frac{x_m - u_m}{\delta_m} \right] \right]$

Method 2: randomly shifted grid

In practice: use a hash of the binary vector as a feature map.

Convergence result:

Claim 2. Let \mathcal{M} be a compact subset of \mathcal{R}^d with diameter $\text{diam}(\mathcal{M})$. Let $\alpha = E[1/\delta]$ and let L_k denote the Lipschitz constant of k with respect to the L_1 norm. With \mathbf{z} as above, we have

$$\Pr \left[\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{M}} |\mathbf{z}(\mathbf{x})' \mathbf{z}(\mathbf{y}) - k(\mathbf{x}, \mathbf{y})| \leq \epsilon \right] \geq 1 - 36dP\alpha \text{diam}(\mathcal{M}) \exp \left(\frac{- \left(\frac{P\epsilon^2}{8} + \ln \frac{\epsilon}{L_k} \right)}{d+1} \right),$$

Results

| Dataset | Fourier+LS | Binning+LS | CVM | Exact SVM |
|---|--------------------------------|-------------------------------|----------------------------------|---|
| CPU regression 6500 instances 21 dims | 3.6% 20 secs $D = 300$ | 5.3% 3 mins $P = 350$ | 5.5% 51 secs | 11% 31 secs ASVM |
| Census regression 18,000 instances 119 dims | 5% 36 secs $D = 500$ | 7.5% 19 mins $P = 30$ | 8.8% 7.5 mins | 9% 13 mins SVMTorch |
| Adult classification 32,000 instances 123 dims | 14.9% 9 secs $D = 500$ | 15.3% 1.5 mins $P = 30$ | 14.8% 73 mins | 15.1% 7 mins SVM ^{light} |
| Forest Cover classification 522,000 instances 54 dims | 11.6% 71 mins $D = 5000$ | 2.2% 25 mins $P = 50$ | 2.3% 7.5 hrs | 2.2% 44 hrs libSVM |
| KDDCUP99 (see footnote) classification 4,900,000 instances 127 dims | 7.3% 1.5 min $D = 50$ | 7.3% 35 mins $P = 10$ | 6.2% (18%) 1.4 secs (20 secs) | 8.3% < 1 s SVM+sampling |

Interpretation: for data where interpolation is needed, use Fourier kernels.
For data where “memorization” is needed, use binning features.

Caveat: the Gaussian kernel was used for Fourier+LS, the Laplace for Binning+LS