

Explicit Feature Methods for Accelerated Kernel Learning

Purushottam Kar

Quick Motivation

- Kernel Algorithms (SVM, SVR, KPCA) have output

$$h(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

- Number of “support vectors” is typically large
 - Provably a **constant fraction** of training set size*
 - Prediction time $\Omega(nd)$ where $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^d$
 - Slow for real time applications

*[Steinwart NIPS 03, Steinwart-Christmann. NIPS 08]

The General Idea

- Approximate kernel using explicit feature maps

$$\mathbf{Z}: \mathcal{X} \rightarrow \mathbb{R}^D \text{ s.t. } K(\mathbf{x}, \mathbf{x}_i) \approx \mathbf{Z}(\mathbf{x})^\top \mathbf{Z}(\mathbf{x}_i)$$

- Speeds up prediction time to $\mathcal{O}(Dd) \ll \mathcal{O}(nd)$

$$h(\mathbf{x}) \approx \sum_{i=1}^n \alpha_i \langle \mathbf{Z}(\mathbf{x}), \mathbf{Z}(\mathbf{x}_i) \rangle = \mathbf{Z}(\mathbf{x})^\top \mathbf{w}$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{Z}(\mathbf{x}_i)$$

- Speeds up training time as well

Why Should Such Maps Exist?

- **Mercer's theorem***

Every PSD kernel K has the following expansion

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{\infty} \lambda_i \Phi_i(\mathbf{x}) \Phi_i(\mathbf{y})$$

- The series **converges uniformly** to kernel

- For every $\epsilon > 0$, $\exists D_\epsilon$ such that if we construct the map

$$\mathbf{Z}_{D_\epsilon} = (\Phi_1, \Phi_2, \dots, \Phi_{D_\epsilon}) \in \mathbb{R}^{D_\epsilon},$$

- then for all $\mathbf{x}, \mathbf{y} \in X$

$$|K(\mathbf{x}, \mathbf{y}) - \mathbf{Z}_{D_\epsilon}(\mathbf{x})^\top \mathbf{Z}_{D_\epsilon}(\mathbf{y})| \leq \epsilon$$

- Call such maps **uniformly ϵ -approximate**

*[Mercer 09]

Today's Agenda

- **Some explicit feature map constructions**
 - **Randomized feature maps**
e.g. Translation invariant, rotation invariant
 - **Deterministic feature maps**
e.g. Intersection, scale invariant
- **Some “fast” random feature constructions**
 - **Translation invariant, dot product**
- **The BIG picture?**

Random Feature Maps

Approximate recovery of kernel values with high confidence

Translation Invariant Kernels*

- Kernels of the form $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{x} - \mathbf{y})$
 - Gaussian kernel, Laplacian kernel

- **Bochner's Theorem****

For every K there exists a positive function p

$$\begin{aligned} K(\mathbf{x} - \mathbf{y}) &= \int_{\omega \in \hat{X}} \cos(\omega^\top (\mathbf{x} - \mathbf{y})) p(\omega) d\omega \\ &= \mathbb{E}_{\omega \sim p} \left[\cos(\omega^\top (\mathbf{x} - \mathbf{y})) \right] \end{aligned}$$

- Finding p : take inverse Fourier transform of K
- Select $\omega_i \sim p$ for $i = 1, \dots, D$

$$\mathbf{Z}_i: \mathbf{x} \mapsto \left(\cos(\omega_i^\top \mathbf{x}), \sin(\omega_i^\top \mathbf{x}) \right)$$

*[Rahimi-Recht NIPS 07], ** Special case for $\mathcal{X} \subset \mathbb{R}^d$, [Bochner 33]

Translation Invariant Kernels

- *Empirical averages approximate expectations*

- Let $Z: x \mapsto (Z_1(x), Z_2(x), \dots, Z_D(x))$

$$\begin{aligned} Z(x)^\top Z(y) &= \frac{1}{D} \sum_{i=1}^D Z_i(x)^\top Z_i(y) \\ &= \frac{1}{D} \sum_{i=1}^D \cos(\omega_i^\top (x - y)) \\ &\approx \mathbb{E}_{\omega \sim p} [\cos(\omega^\top (x - y))] \\ &= K(x - y) \end{aligned}$$

- Let us assume points $x, y \in \mathcal{B}(\mathbf{0}, R) \subset \mathbb{R}^d$

Then we require $D \geq \frac{30d}{\epsilon^2} \log\left(\frac{RC_K}{\delta\epsilon}\right)$

C_K depends on spectrum of kernel K

Translation Invariant Kernels

- For the RBF Kernel

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(\frac{-\|\mathbf{x} - \mathbf{y}\|_2^2}{2}\right)$$
$$p(\boldsymbol{\omega}) = \frac{1}{(2\pi)^{d/2}} \exp\left(\frac{-\|\boldsymbol{\omega}\|_2^2}{2}\right)$$

- If kernel K offers a γ margin, then we should require

$$D \gtrsim \frac{30d}{\gamma^2} \log\left(\frac{Rd}{\delta\gamma}\right)$$

Here $C_K \approx d$ where $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$

Rotation Invariant Kernels*

- Kernels of the form $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{x}^\top \mathbf{y})$
 - Polynomial kernels, exponential kernel
- **Schoenberg's theorem****

$$K(\mathbf{x}^\top \mathbf{y}) = \sum_{p \geq 0} a_p (\mathbf{x}^\top \mathbf{y})^p, \quad a_p \geq 0$$

- Select $p_i \sim \mu \in \mathbb{N}$ for $i = 1, \dots, D$
 - Approx. $(\mathbf{x}^\top \mathbf{y})^{p_i}$: select $\omega_1, \dots, \omega_{p_i} \sim \{-1, 1\}^d$

$$Z_i: \mathbf{x} \mapsto \sqrt{a_{p_i}} \prod_{j=1}^{p_i} \omega_j^\top \mathbf{x}$$

- Similar approximation guarantees as earlier

*[K.-Karnick AISTATS 12], **[Schoenberg 42]

Deterministic Feature Maps

Exact/approximate recovery of kernel values
with certainty

Intersection Kernel*

- Kernel of the form $K(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^d \min\{x^j, y^j\}$
- Exploit **additive separability** of the kernel

$$h(\mathbf{x}) = \sum_{i=1}^n \alpha_i \sum_{j=1}^d \min\{x^j, x_i^j\} = \sum_j h_j(\mathbf{x})$$

$$h_j(\mathbf{x}) = \sum_i \alpha_i \min\{x^j, x_i^j\}$$

- Each $h_j(\mathbf{x})$ can be calculated in $O(\log n)$ time !
 - Requires $O(n \log n)$ preprocessing time per dimension
- Prediction time *almost* independent of n
 - However, **deterministic** and **exact** method – no ϵ or δ

*[Maji-Berg-Malik CVPR 08]

Scale Invariant Kernels*

- Kernels of the form $K(x, y) = \sum_{j=1}^d K_j(x^j, y^j)$ where

$$K_j(x^j, y^j) = (x^j)^\gamma K_j\left(\frac{x^j}{y^j}\right) (y^j)^\gamma, \quad \gamma \geq 0$$

- Bochner's theorem still applies**

- Involves working with $\tilde{K}_j(x^j, y^j) = \tilde{K}(\log |x^j| - \log |y^j|)$
- Restrict domain so that we have a Fourier series

$$\tilde{K}_j(\lambda) = \sum_{k=-\infty}^{\infty} \tilde{\mu}_k e^{ij\Delta\lambda}$$

- Use only lower frequencies $k \in \{-A, \dots, A\}$
- **Deterministic** ϵ -approximate maps

*[Vedaldi-Zisserman CVPR 10], **[K. 12]

Fast Feature Maps

Accelerated Random Feature Constructions

Fast Fourier Features

- Special case of $K(\mathbf{x}, \mathbf{y}) = \exp\left(\frac{-\|\mathbf{x}-\mathbf{y}\|^2}{2\sigma^2}\right)$
 - Old method: $\mathbf{W} \in \mathbb{R}^{D \times d}$, $W_{ij} \sim \mathcal{N}(0, \sigma^{-2})$, $\mathcal{O}(dD)$ time
 - Instead use $\tilde{\mathbf{Z}}: \mathbf{x} \mapsto \cos(\mathbf{V}\mathbf{x})$ where $\mathbf{V} = \mathbf{S}\mathbf{H}\mathbf{G}\mathbf{\Pi}\mathbf{H}\mathbf{B}$
 - $\mathbf{\Pi}$ is the **Hadamard transform**, $\mathbf{\Pi}$ is a random permutation
 $\mathbf{S}, \mathbf{G}, \mathbf{B}$ random diagonal scaling, Gaussian and sign matrices
- Prediction time $\mathcal{O}(D \log d)$, $\mathbb{E}[\tilde{\mathbf{Z}}(\mathbf{x})^\top \tilde{\mathbf{Z}}(\mathbf{y})] = K(\mathbf{x}, \mathbf{y})$
 - Rows of \mathbf{V} are (non independent) Gaussian vectors
 - Correlations are sufficiently low $\text{Var}[\tilde{\mathbf{Z}}(\mathbf{x})^\top \tilde{\mathbf{Z}}(\mathbf{y})] \leq \mathcal{O}\left(\frac{1}{D}\right)$
 - However, exponential convergence (for now) only for $D = d$

Fast Taylor Features

- **Special case of $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + c)^p$**
 - Earlier method $Z: \mathbf{x} \mapsto \prod_{j=1}^p \omega_j^\top \mathbf{x}$, takes $O(pdD)$ time
 - New method* takes $O(p(d + D \log D))$ time
 - Earlier method works (a bit) better for $c > 0$
 - Should be possible to improve new method as well
- **Crucial idea $(\mathbf{x}^\top \mathbf{y})^p = \langle \mathbf{x}^{\otimes p}, \mathbf{y}^{\otimes p} \rangle$**
 - **Count Sketch**** $C: \mathbf{x} \mapsto C(\mathbf{x})$ such that $C(\mathbf{x})^\top C(\mathbf{y}) \approx \mathbf{x}^\top \mathbf{y}$
 - Create sketch $C(X) \in \mathbb{R}^D$ of tensor $X = \mathbf{x}^{\otimes p}$
 - Create p independent count sketches $C_1(\mathbf{x}), \dots, C_p(\mathbf{x})$
 - Can show that $C(\mathbf{x}^{\otimes p}) \sim \prod_{j=1}^p P(C_j(\mathbf{x}))$
 - Can be done in time $O(p(d + D \log D))$ time using FFT

*[Pham-Pagh KDD 13], **[Charikar et al 02]

The BIG Picture

An Overview of Explicit Feature Methods

Other Feature Construction Methods

- Efficiently evaluable maps for efficient prediction
 - **Fidelity** to a particular kernel not an objective
 - Hard(er?) to give **generalization guarantees**
- **Local Deep Kernel Learning (LDKL)***
 - Sparse features speed up evaluation time to $O(d \log D)$
 - Training phase more involved
- **Pairwise Piecewise Linear Embedding (PL2)****
 - Encodes (discretization of) individual and pairs of features
 - Construct a $D = O\left(K(d + d^2)\right)$ dimensional feature map
 - Features are $O(d + d^2)$ -sparse

*[Jose et al ICML 13], **[Pele et al ICML 13]

A Taxonomy of Feature Methods

Data Dependence

Kernel Dependence

	Yes	No
Yes	Nystrom Methods <ul style="list-style-type: none">- Slow training- Data aware- Problem oblivious	Explicit Maps <ul style="list-style-type: none">- Fast training- Data oblivious- Problem oblivious
No	LDKL, PL2 <ul style="list-style-type: none">- Slow(er) training- Data aware- Problem aware	

Discussion

The next big thing in accelerated kernel learning ?