

ECS289: Scalable Machine Learning

Cho-Jui Hsieh
UC Davis

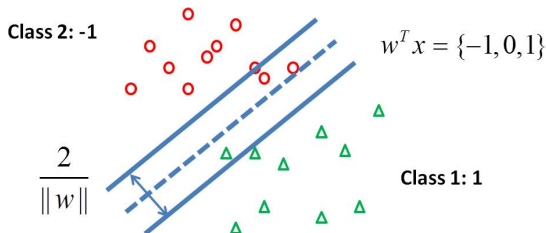
Oct 18, 2016

Outline

- Kernel Methods
- Solving the dual problem
- Kernel approximation

Support Vector Machines (SVM)

- SVM is a widely used classifier.
- Given:
 - Training data points $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Each $\mathbf{x}_i \in \mathbb{R}^d$ is a feature vector:
 - Consider a simple case with two classes: $y_i \in \{+1, -1\}$.
- Goal: Find a hyperplane to separate these two classes of data:
if $y_i = 1$, $\mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i$; $y_i = -1$, $\mathbf{w}^T \mathbf{x}_i \leq -1 + \xi_i$.



Support Vector Machines (SVM)

- Given training data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ with labels $y_i \in \{+1, -1\}$.
- SVM primal problem (find optimal $\mathbf{w} \in \mathbb{R}^d$):

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i$$
$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, n,$$

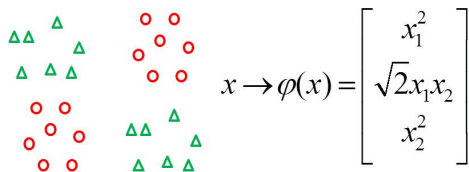
- SVM dual problem (find optimal $\alpha \in \mathbb{R}^n$):

$$\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha$$
$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, n$$

- Q : n -by- n kernel matrix, $Q_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$
- Each α_i corresponds to one training data point.
- Primal-dual relationship: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$

Non-linearly separable problems

- What if the data is not linearly separable?



Solution: map data x_i to higher dimensional (maybe infinite) feature space $\varphi(x_i)$, where they are linearly separable.

Support Vector Machines (SVM)

- SVM primal problem:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \varphi(\mathbf{x}_i)) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, n, \end{aligned}$$

- The dual problem for SVM:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha, \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad \text{for } i = 1, \dots, n, \end{aligned}$$

where $Q_{ij} = y_i y_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$ and $\mathbf{e} = [1, \dots, 1]^T$.

- Kernel trick: define $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$.
- At optimum: $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \varphi(\mathbf{x}_i)$,

Various types of kernels

- Gaussian kernel: $K(\mathbf{x}_i, \mathbf{y}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|_2^2}$;
- Polynomial kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i^T \mathbf{x}_j + c)^d$.
- Other kernels for specific problems:
 - Graph kernels
(Vishwanathan et al., "Graph Kernels", JMLR, 2010)
 - Pyramid kernel for image matching
(Grauman and Darrell, "The Pyramid Match Kernel: Discriminative Classification with Sets of Image Features". In ICCV, 2005)
 - String kernel
(Lodhi et al., "Text classification using string kernels". JMLR, 2002).

General Kernelized ERM

- L2-Regularized Empirical Risk Minimization:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \ell_i(\mathbf{w}^T \varphi(\mathbf{x}_i))$$

$\mathbf{x}_1, \dots, \mathbf{x}_n$: training samples

$\varphi(\mathbf{x}_i)$: nonlinear mapping to a higher dimensional space

- Dual problem:

$$\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha + \sum_{i=1}^n \ell_i^*(-\alpha_i)$$

where $Q \in \mathbb{R}^{n \times n}$ and $Q_{ij} = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$.

Kernel Ridge Regression

- Given training samples $(\mathbf{x}_i, y_i), i = 1, \dots, n$.

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^T \varphi(\mathbf{x}_i) - y_i)^2$$

- Dual problem:

$$\min_{\boldsymbol{\alpha}} \boldsymbol{\alpha}^T Q \boldsymbol{\alpha} + \|\boldsymbol{\alpha}\|^2 - 2\boldsymbol{\alpha}^T \mathbf{y}$$

Scalability

- Challenge for solving kernel SVMs (for dataset with n samples):
 - **Space:** $O(n^2)$ for storing the n -by- n kernel matrix (can be reduced in some cases);
 - **Time:** $O(n^3)$ for computing the exact solution.

Greedy Coordinate Descent for Kernel SVM

Nonlinear SVM

- SVM primal problem:

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i$$

s.t. $y_i(\mathbf{w}^T \varphi(\mathbf{x}_i)) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, n,$

w: an infinite dimensional vector, very hard to solve

- The dual problem for SVM:

$$\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha,$$

s.t. $0 \leq \alpha_i \leq C, \quad \text{for } i = 1, \dots, n,$

where $Q_{ij} = y_i y_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$ and $\mathbf{e} = [1, \dots, 1]^T$.

- Kernel trick: define $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$.
- Example: Gaussian kernel $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}$

Nonlinear SVM

- Can we solve the problem by dual coordinate descent?
- The vector $\mathbf{w} = \sum_j y_j \alpha_j \varphi(\mathbf{x}_j)$ may have infinite dimensionality:

Cannot maintain \mathbf{w}

- Closed form solution needs $O(n)$ computational time:

$$\delta^* = \max(-\alpha_j \min(C - \alpha_j, \frac{1 - (Q\alpha)_j}{Q_{jj}}))$$

(Assume Q is stored in memory)

- Can we improve coordinate descent using the same $O(n)$ time complexity?

Greedy Coordinate Descent

The Greedy Coordinate Descent (GCD) algorithm:

For $t = 1, 2, \dots$

1. Compute $\delta_i^* := \operatorname{argmin}_{\delta} f(\boldsymbol{\alpha} + \delta \mathbf{e}_i)$ for all $i = 1, \dots, n$
2. Find the best i^* according to the following criterion:

$$i^* = \operatorname{argmax}_i |\delta_i^*|$$

3. $\alpha_{i^*} \leftarrow \alpha_{i^*} + \delta_{i^*}^*$

Greedy Coordinate Descent

Other variable selection criterion:

- The coordinate with the maximum step size:

$$i^* = \operatorname{argmax}_i |\delta_i^*|$$

- The coordinate with maximum objective function reduction:

$$i^* = \operatorname{argmax}_i (f(\boldsymbol{\alpha}) - f(\boldsymbol{\alpha} + \delta_i^* \mathbf{e}_i))$$

- The coordinate with the maximum projected gradient.
- ...

Greedy Coordinate Descent

- How to compute the optimal coordinate?

- Closed form solution of best δ :

$$\delta_i^* = \max \left(-\alpha_i, \min \left(C - \alpha_i, \frac{1 - (Q\alpha)_i}{Q_{ii}} \right) \right)$$

- Observations:

- 1 Computing all δ_i^* needs $O(n^2)$ time
 - 2 If $Q\alpha$ is stored in memory, computing all δ_i^* only needs $O(n)$ time
- Maintaining $Q\alpha$ also needs $O(n)$ time after each update

Greedy Coordinate Descent

Initial: α , $\mathbf{z} = Q\alpha$

For $t = 1, 2, \dots$

For all $i = 1, \dots, n$, compute

$$\delta_i^* = \max \left(-\alpha_i, \min \left(C - \alpha_i, \frac{1 - z_i}{Q_{ii}} \right) \right)$$

Let $i^* = \operatorname{argmax}_i |\delta_i^*|$

$\alpha \leftarrow \alpha + \delta_{i^*}^*$

$\mathbf{z} \leftarrow \mathbf{z} + \mathbf{q}_{i^*} \delta_{i^*}^*$ (\mathbf{q}_{i^*} is the i^* -th column of Q)

(This is a simplified version of the Sequential Minimal Optimization (SMO) algorithm proposed in Platt et al., 1998)

(A similar version is implemented in LIBSVM)

How to solve problems with millions of samples?

- $Q \in \mathbb{R}^{n \times n}$ cannot be fully stored
- Have a fixed size of memory to “cache” the computed columns of Q
- For each coordinate update:

If \mathbf{q}_i is in memory, directly use it to update

If \mathbf{q}_i is not in memory

- 1 Kick out the “Least Recent Used” column
- 2 Recompute \mathbf{q}_i and store it in memory
- 3 Update α_i

LIBSVM

- Implemented in LIBSVM:

<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>

- Other functionalities:
 - Multi-class classification
 - Support vector regression
 - Cross-validation

Kernel Approximation

Kernel Approximation

- Kernel methods are hard to scale up because
 - Kernel matrix G is an n -by- n matrix
- Can we approximate the kernel using a low-rank representation?
- Two main algorithms:
 - Nystrom approximation: Approximate the [kernel matrix](#)
 - Random Features: Approximate the [kernel function](#)

Kernel Matrix Approximation

- We want to form a low rank approximation of $G \in \mathbb{R}^{n \times n}$,
where $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- Can we do SVD?

Kernel Matrix Approximation

- We want to form a low rank approximation of $G \in \mathbb{R}^{n \times n}$,
where $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- Can we do SVD?
No, SVD needs to form the n -by- n matrix
 $O(n^2)$ space and $O(n^3)$ time

Kernel Matrix Approximation

- We want to form a low rank approximation of $G \in \mathbb{R}^{n \times n}$,
where $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- Can we do SVD?
No, SVD needs to form the n -by- n matrix
 $O(n^2)$ space and $O(n^3)$ time
- Can we do matrix completion or matrix sketching?

Kernel Matrix Approximation

- We want to form a low rank approximation of $G \in \mathbb{R}^{n \times n}$,
where $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- Can we do SVD?
No, SVD needs to form the n -by- n matrix
 $O(n^2)$ space and $O(n^3)$ time
- Can we do matrix completion or matrix sketching?
Main problem: need to generalize to new points

Nystrom Approximation

- Nystrom approximation:
 - Randomly sample m columns of G :

$$G = \begin{bmatrix} W & G_{12} \\ G_{21} & G_{22} \end{bmatrix}$$

W : m -by- m square matrix (observed)

G_{21} : $(n - m)$ -by- m matrix (observed)

$G_{12} = G_{21}^T$ (observed)

G_{22} : $(n - m)$ -by- $(n - m)$ matrix (not observed)

- Form a kernel approximation based on these mn elements

$$G \approx \tilde{G} = \begin{bmatrix} W \\ G_{21} \end{bmatrix} W^\dagger [W \quad G_{12}]$$

W^\dagger : pseudo-inverse of W

Nystrom Approximation

- Why W^\dagger ?

Exact recover the top left m -by- m matrix

- The kernel approximation:

$$\begin{bmatrix} W & G_{12} \\ G_{21} & G_{22} \end{bmatrix} \approx \begin{bmatrix} W \\ G_{21} \end{bmatrix} W^\dagger \begin{bmatrix} W & G_{12} \end{bmatrix} = \begin{bmatrix} W & G_{12} \\ G_{21} W^\dagger G_{12} \end{bmatrix}$$

Nystrom Approximation

Algorithm:

- 1 Sample m "landmark points": $\mathbf{v}_1, \dots, \mathbf{v}_m$
- 2 Compute the kernel values between all training data to landmark points:

$$C = \begin{bmatrix} W \\ G_{21} \end{bmatrix} = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{v}_1) & \cdots & K(\mathbf{x}_1, \mathbf{v}_m) \\ K(\mathbf{x}_2, \mathbf{v}_1) & \cdots & K(\mathbf{x}_2, \mathbf{v}_m) \\ \vdots & \vdots & \vdots \\ K(\mathbf{x}_n, \mathbf{v}_1) & \cdots & K(\mathbf{x}_n, \mathbf{v}_m) \end{bmatrix}$$

- 3 Form the kernel approximation $\tilde{G} = CW^\dagger C^T$
(No need to explicitly form the n -by- n matrix)
- Time complexity: mn kernel evaluations ($O(mnd)$ time if use classical kernels)
 - How to choose m ?
Trade-off between accuracy v.s computational time and memory space

Nystrom Approximation: Training

- Solve the dual problem using low-rank representation.
- Example: Kernel ridge regression

$$\min_{\alpha} \alpha^T \tilde{G} \alpha + \lambda \|\alpha\|^2 - \mathbf{y}^T \alpha$$

Solve a linear system $(\tilde{G} + \lambda I)\alpha = \mathbf{y}$

- Use iterative methods (such as CG), with fast matrix-vector multiplication

$$(\tilde{G} + \lambda I)\mathbf{p} = CW^\dagger C^T \mathbf{p} + \lambda \mathbf{p}$$

$O(nm)$ time complexity per iteration

Nystrom Approximation: Training

- Another approach: reduce to linear ERM problems!
- Rewrite as

$$\tilde{G} = CW^\dagger C^T = C(W^\dagger)^{1/2}(W^\dagger)^{1/2}C^T$$

- So the approximate kernel can be represent by linear kernel with feature matrix $C(W^\dagger)^{1/2}$:

$$\tilde{K}(\mathbf{x}_i, \mathbf{x}_j) = \tilde{G}_{ij} = \mathbf{u}_i^T \mathbf{u}_j,$$

where $\mathbf{u}_j = (W^\dagger)^{1/2} \begin{bmatrix} K(\mathbf{x}_j, \mathbf{v}_1) \\ \vdots \\ K(\mathbf{x}_j, \mathbf{v}_m) \end{bmatrix}$

- The problem is equivalent to a linear models with features $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbb{R}^m$
 - Kernel SVM \Rightarrow Linear SVM with m features
 - Kernel Ridge Regression \Rightarrow Linear Ridge regression with m features

Nystrom Approximation: Prediction

- Given the dual solution α .
- Prediction for testing sample \mathbf{x} :

$$\sum_{i=1}^n \alpha_i \tilde{K}(\mathbf{x}_i, \mathbf{x}) \quad \text{or} \quad \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}) ?$$

Nystrom Approximation: Prediction

- Given the dual solution α .
- Prediction for testing sample \mathbf{x} :

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

- **Need to use approximate kernel instead of original kernel!**
 - Approximate kernel gives much better accuracy!
 - Because training & testing should use the same kernel.
- Generalization bound:
 - (Alaoui and Mahoney, “Fast Randomized Kernel Methods With Statistical Guarantees”. 2014.)
 - (Rudi et al., “Less is More: Nystrom Computational Regularization”. NIPS 2015.)
 - (using small number of landmark points can be viewed as a regularization)

Nystrom Approximation: Prediction

- How to define $\tilde{K}(\mathbf{x}, \mathbf{x}_i)$?

$$\tilde{G} = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{v}_1) & \cdots & K(\mathbf{x}_1, \mathbf{v}_m) \\ \vdots & \vdots & \vdots \\ K(\mathbf{x}_n, \mathbf{v}_1) & \cdots & K(\mathbf{x}_n, \mathbf{v}_m) \\ K(\mathbf{x}, \mathbf{v}_1) & \cdots & K(\mathbf{x}, \mathbf{v}_m) \end{bmatrix} W^\dagger \begin{bmatrix} K(\mathbf{v}_1, \mathbf{x}_1) & \cdots & K(\mathbf{v}_1, \mathbf{x}) \\ \vdots & \vdots & \vdots \\ K(\mathbf{v}_m, \mathbf{x}_1) & \cdots & K(\mathbf{v}_m, \mathbf{x}) \end{bmatrix}$$

- Compute $\mathbf{u} = (W^\dagger)^{1/2} \begin{bmatrix} K(\mathbf{v}_1, \mathbf{x}) \\ \vdots \\ K(\mathbf{v}_m, \mathbf{x}) \end{bmatrix}$ (need m kernel evaluations)
- Approximate kernel can be defined by

$$\tilde{K}(\mathbf{x}, \mathbf{x}_i) = \mathbf{u}^T \mathbf{u}_i$$

Nystrom Approximation: Prediction

- Prediction:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i \mathbf{u}^T \mathbf{u}_i = \mathbf{u}^T \left(\sum_{i=1}^n \alpha_i \mathbf{u}_i \right)$$

$(\sum_{i=1}^n \alpha_i \mathbf{u}_i)$ can be pre-computed

⇒ prediction time is m kernel evaluations

- Original kernel method: need n kernel evaluations for prediction.
- Summary:

	Quality	Training time	Prediction time
Original kernel	exact	$n^{1 \cdot x}$ kernel + kernel training	n kernel
Nystrom (m)	rank m	nm kernel + linear training	m kernel

Nystrom Approximation

- How to select landmark points ($\mathbf{v}_1, \dots, \mathbf{v}_m$)
 - Traditional approach: uniform random sampling from training data
 - Importance sampling (leverage score):
 - Drineas and Mahoney “On the Nystrom method for approximating a Gram matrix for improved kernel-based learning”. JMLR 2015.
 - Gittens and Mahoney “Revisiting the Nystrom method for improved large-scale machine learning”. ICML 2013
 - Kmeans sampling (performs extremely well) :
Run kmeans clustering and set $\mathbf{v}_1, \dots, \mathbf{v}_m$ to be cluster centers
 - Zhang et al., “Improved Nystrom low rank approximation and error analysis”. ICML 2008.
 - Subspace distance:
 - Lim et al., “Double Nystrom method: An efficient and accurate Nystrom scheme for large-scale data sets”. ICML 2015.

Nystrom Approximation (other related papers)

- Distributed setting: (Kumar et al., “Ensemble Nystrom method”. NIPS 2009).
- Block diagonal + low-rank approximation: (Si et al., “Memory efficient kernel approximation”. ICML 2014.)
- Nystrom method for fast prediction: (Hsieh et al., “Fast prediction for large-scale kernel machines. ” NIPS, 2014.)
- Structured landmark points: (Si et al., “Computational Efficient Nystrom Approximation using Fast Transforms”. ICML 2016.)

Kernel Approximation (random features)

Random Features

- Directly approximation the kernel function
(Data-independent)
- Generate nonlinear “features”
⇒ reduce the problem to linear model training
- Several examples:
 - Random Fourier Features:
(Rahimi and Recht, “Random features for large-scale kernel machines”. NIPS 2007)
 - Random features for polynomial kernel:
(Kar and Karnick, “Random feature maps for dot product kernels”. AISTATS 2012.)

Random Fourier Features

- Random features for “Shift-invariant kernel”:

A continuous kernel is shift-invariant if $K(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y})$.

- Gaussian kernel: $K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|^2}$
- Laplacian kernel: $K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_1}$
- Shift invariant kernel (if positive definite) can be written as:

$$k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{R}^d} P(\mathbf{w}) e^{j\mathbf{w}^T(\mathbf{x}-\mathbf{y})} d\mathbf{w}$$

for some probability distribution $P(\mathbf{w})$.

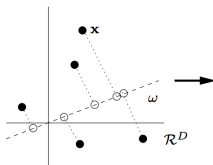
Random Fourier Features

- Taking the real part we have

$$\begin{aligned}k(\mathbf{x} - \mathbf{y}) &= \int_{\mathbb{R}^d} P(\mathbf{w}) \cos(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{y}) \\ &= \int_{\mathbb{R}^d} P(\mathbf{w}) (\cos(\mathbf{w}^T \mathbf{x}) \cos(\mathbf{w}^T \mathbf{y}) + \sin(\mathbf{w}^T \mathbf{x}) \sin(\mathbf{w}^T \mathbf{y})) \\ &= E_{\mathbf{w} \sim P(\mathbf{w})} [z_{\mathbf{w}}(\mathbf{x})^T z_{\mathbf{w}}(\mathbf{y})]\end{aligned}$$

where $z_{\mathbf{w}}(\mathbf{x})^T = [\cos(\mathbf{w}^T \mathbf{x}), \sin(\mathbf{w}^T \mathbf{x})]$.

- $z_{\mathbf{w}}(\mathbf{x})^T z_{\mathbf{w}}(\mathbf{y})$ is an unbiased estimator of $K(\mathbf{x}, \mathbf{y})$ if \mathbf{w} is sampled from $P(\cdot)$



Kernel Name	$k(\Delta)$	$p(\omega)$
Gaussian	$e^{-\frac{\ \Delta\ _2^2}{2}}$	$(2\pi)^{-\frac{D}{2}} e^{-\frac{\ \omega\ _2^2}{2}}$
\mathcal{R}^2 Laplacian	$e^{-\ \Delta\ _1}$	$\prod_d \frac{1}{\pi(1+\omega_d^2)}$
Cauchy	$\prod_d \frac{2}{1+\Delta^2}$	$e^{-\ \Delta\ _1}$

Random Fourier Features

- Sample m vectors $\mathbf{w}_1, \dots, \mathbf{w}_m$ from $P(\cdot)$ distribution
- Generate random features for each sample:

$$\mathbf{u}_i = \begin{bmatrix} \sin(\mathbf{w}_1^T \mathbf{x}_i) \\ \cos(\mathbf{w}_1^T \mathbf{x}_i) \\ \sin(\mathbf{w}_2^T \mathbf{x}_i) \\ \cos(\mathbf{w}_2^T \mathbf{x}_i) \\ \vdots \\ \sin(\mathbf{w}_m^T \mathbf{x}_i) \\ \cos(\mathbf{w}_m^T \mathbf{x}_i) \end{bmatrix}$$

- $K(\mathbf{x}_i, \mathbf{x}_j) \approx \mathbf{u}_i^T \mathbf{u}_j$
(larger m leads to better approximation)

Random Features

- $G \approx UU^T$, where $U = \begin{bmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_n^T \end{bmatrix}$
- Rank $2m$ approximation.
- The problem can be reduced to linear classification/regression for $\mathbf{u}_1, \dots, \mathbf{u}_n$.
- Time complexity: $O(nmd)$ time + linear training
(Nystrom: $O(nmd + m^3)$ time + linear training)
- Prediction time: $O(md)$ per sample
(Nystrom: $O(md)$ per sample)

Fastfood

- Random Fourier features require $O(md)$ time to generate m features:
Main bottleneck: Computing $\mathbf{w}_i^T \mathbf{x}$ for all $i = 1, \dots, m$
- Can we compute this in $O(m \log d)$ time?

Fastfood

- Random Fourier features require $O(md)$ time to generate m features:
Main bottleneck: Computing $\mathbf{w}_i^T \mathbf{x}$ for all $i = 1, \dots, m$
- Can we compute this in $O(m \log d)$ time?
- Fastfood: proposed in (Le et al., “Fastfood Approximating Kernel Expansions in Loglinear Time”. ICML 2013.)
- Reduce time by using “structured random features”

- Tool: fast matrix-vector multiplication for Hadamard matrix ($H_d \in \mathbb{R}^{d \times d}$)

$H_d \mathbf{x}$ requires $O(d \log d)$ time

- Hadamard matrix:

$$H_1 = [1]$$

$$H_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$H_{2^k} = \begin{bmatrix} H_{2^{k-1}} & H_{2^{k-1}} \\ H_{2^{k-1}} & -H_{2^{k-1}} \end{bmatrix}$$

- $H_d \mathbf{x}$ can be computed efficiently by Fast Hadamard Transform (dynamic programming)

- Sample $m = d$ random features by

$$W = [H_d] \begin{bmatrix} v_1 & 0 & \cdots & 0 \\ 0 & v_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & v_d \end{bmatrix}$$

So $W\mathbf{x}$ only needs $O(d \log d)$ time.

- Instead of sampling W (d^2 elements), we just sample v_1, \dots, v_d from $N(0, \sigma)$
- Each row of W is still $N(0, \sigma)$ (but rows are not independent)
- Faster computation, but performance is slightly worse than original random features.

Extensions

- Fastfood: structured random features to improve computational speed:
(Le et al., “Fastfood Approximating Kernel Expansions in Loglinear Time”. ICML 2013.)
- Structured landmark points for Nystrom approximation:
(Si et al., “Computational Efficient Nystrom Approximation using Fast Transforms”. ICML 2016.)
- Doubly stochastic gradient with random features:
(Dai et al., “Scalable Kernel Methods via Doubly Stochastic Gradients”. NIPS 2015.)
- Generalization bound (?)

Coming up

- Choose the paper for presentation (due this Sunday, Oct 23).

Questions?