

Kernel Methods

Announcements

- HW3 due today, HW4 assigned

Outline

- Nonlinear classification and regression via nonlinear feature maps
- Kernels
- Kernel Methods

Nonlinear Feature Maps

- One way to create a nonlinear method for regression or classification is to first transform the feature vector via a *nonlinear feature map*

$$x \mapsto \Phi(x), \quad \Phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

and apply a linear method to the transformed data

- In regression, the nonlinear function model is

$$f(x) = w^T \Phi(x) + b$$

where $w \in \mathbb{R}^m, b \in \mathbb{R}$

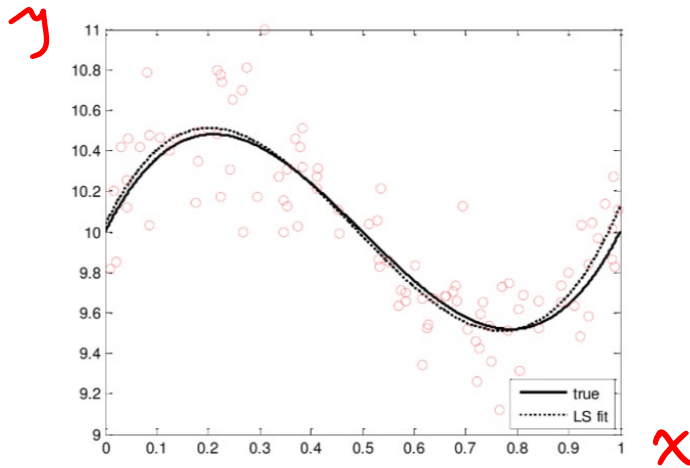
- In binary classification, the nonlinear classifier model is

$$h(x) = \text{sign}(w^T \Phi(x) + b)$$

where $w \in \mathbb{R}^m, b \in \mathbb{R}$

Nonlinear Feature Maps

- **Example:** cubic polynomial regression



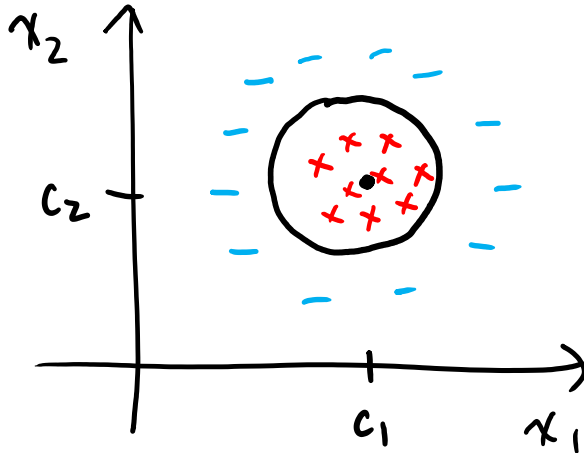
$$\Phi(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}, \quad w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$\min_{w, b} \frac{1}{n} \sum_{i=1}^n \left(y_i - \underbrace{[w_3 x_i^3 + w_2 x_i^2 + w_1 x_i + b]}_{w^T \Phi(x_i)} \right)^2$$

Nonlinear Feature Maps

- Example: circular classifier

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



$$x \mapsto \text{sign} \left\{ (x_1 - c_1)^2 + (x_2 - c_2)^2 - r^2 \right\}$$

$$= \text{sign} \left\{ x_1^2 - 2c_1x_1 + c_1^2 + x_2^2 - 2c_2x_2 + c_2^2 - r^2 \right\}$$

$$\Phi(x) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_1 \\ x_2 \\ x_1 x_2 \end{bmatrix}, w = \begin{bmatrix} 1 \\ 1 \\ -2c_1 \\ -2c_2 \\ 0 \end{bmatrix}, b = c_1^2 + c_2^2 - r^2$$

$$= \text{sign} \left\{ w^T \Phi(x) + b \right\}$$

Inner Product Kernels

- Problem with the above approach: m can explode as d increases
- This makes it prohibitive to compute/store/manipulate Φ directly
- Fortunately, the following facts allow us to use nonlinear feature maps for large m :
 - Many ML algorithms depend on $\Phi(\mathbf{x})$ only via *inner products*

$$\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$$

- For certain Φ , the function

$$k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$$

can be computed efficiently even if m is huge or possibly infinite!

- k is called an *inner product kernel*
- Let's look at an example involving the *dot product*

Example

- The *inhomogeneous polynomial kernel of degree 2* is

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + 1)^2$$

where $\mathbf{u}, \mathbf{v} \in \mathbb{R}^2$. Determine a feature map Φ such that

$$k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle.$$

$$k(u, v) = (u_1 v_1 + u_2 v_2 + 1)^2$$

$$= u_1^2 v_1^2 + u_2^2 v_2^2 + 2u_1 u_2 v_1 v_2 + 2u_1 v_1 + 2u_2 v_2 + 1$$

$$= \langle \Phi(u), \Phi(v) \rangle, \quad \Phi(u) =$$

$$\begin{bmatrix} u_1^2 \\ u_2^2 \\ \sqrt{2} u_1 u_2 \\ \sqrt{2} u_1 \\ \sqrt{2} u_2 \\ 1 \end{bmatrix}$$

Important Kernels

- Homogeneous polynomial kernel

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^p,$$

- Inhomogeneous polynomial

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + c)^p, \quad c > 0$$

- Gaussian kernel

$$k(\mathbf{u}, \mathbf{v}) = \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{u} - \mathbf{v}\|^2 \right), \quad \sigma > 0$$

- For the Gaussian kernel, $m = \infty$ and the dot product is replaced by the ℓ_2 inner product: $\langle (a_i), (b_i) \rangle = \sum_{i=1}^{\infty} a_i b_i$.

SPD Kernels

- One way to determine an IP Kernel is to construct Φ explicitly as we did in the examples above.
- We can also verify that k is an IP kernel if it satisfies the following properties.
- Let $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$. We say k is *symmetric* if $k(\mathbf{u}, \mathbf{v}) = k(\mathbf{v}, \mathbf{u}) \forall \mathbf{u}, \mathbf{v}$.
- We say k is *positive definite* if

$$\begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

is a positive *semi*-definite matrix for all $n \in \mathbb{N}$ and $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$

- If k is both symmetric and positive definite, it is referred to as a *symmetric, positive definite (SPD) kernel*.
- **Theorem:** k is an SPD kernel $\iff k$ is an inner product kernel

$$\iff \exists \Phi \text{ s.t. } k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle$$

The Kernel Trick

- A machine learning algorithm is said to be *kernelizable* if it is possible to formulate the algorithm such that all training instances \mathbf{x}_i and any test instance \mathbf{x} occur exclusively in inner products of the form $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$, $\langle \mathbf{x}_i, \mathbf{x} \rangle$ or $\langle \mathbf{x}, \mathbf{x} \rangle$.
- Suppose Φ is a feature map associated to an inner product kernel k .
- If we apply a kernelizable algorithm to the training data

$$(\Phi(\mathbf{x}_1), y_1), \dots, (\Phi(\mathbf{x}_n), y_n)$$

then we can formulate the algorithm such that transformed feature vectors only appear via inner products $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$ with other transformed feature vectors.

- Can implement by evaluating $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$

which eliminates the need to ever compute $\Phi(\mathbf{x})$ explicitly.

Example: Nearest Neighbors

- Given training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, the 1-nearest neighbor classifier assigns a test point \mathbf{x} to the label

$$y_{\arg\min_{1 \leq i \leq n} \|\mathbf{x} - \mathbf{x}_i\|_2}$$

- Suppose Φ is a feature map associated to an inner product/SPD kernel k , i.e.,

$$k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle$$

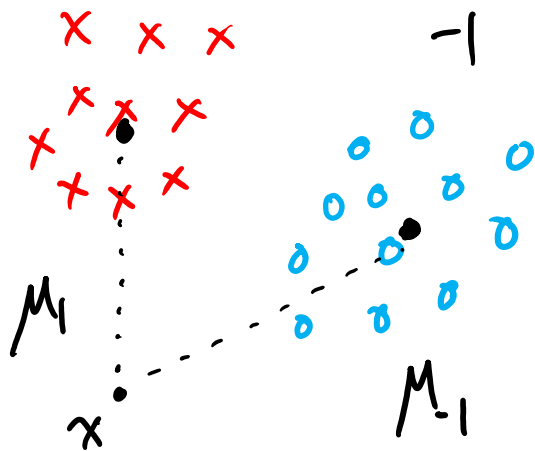
for all \mathbf{u}, \mathbf{v} .

- Now consider the 1-NN classifier after first applying Φ :

$$y_{\arg\min_{1 \leq i \leq n} \|\Phi(\mathbf{x}) - \Phi(\mathbf{x}_i)\|_2}$$

$$\begin{aligned} \|\Phi(\mathbf{x}) - \Phi(\mathbf{x}_i)\|_2^2 &= (\Phi(\mathbf{x}) - \Phi(\mathbf{x}_i))^T (\Phi(\mathbf{x}) - \Phi(\mathbf{x}_i)) \\ &= \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}) \rangle - 2\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}_i) \rangle + \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_i) \rangle \\ &= k(\mathbf{x}, \mathbf{x}) - 2k(\mathbf{x}, \mathbf{x}_i) + k(\mathbf{x}_i, \mathbf{x}_i) \end{aligned}$$

4.1 Nearest Centroid Classifier



$$f(x) = \text{sign} \left\{ \|x - \mu_{-1}\|^2 - \|x - \mu_1\|^2 \right\}$$

$$\mu_1 = \frac{1}{n_+} \sum_{i: y_i = 1} x_i$$

$$\mu_{-1} = \frac{1}{n_-} \sum_{i: y_i = -1} x_i$$

$$f(x) = \text{sign} \left\{ \frac{1}{n_+} \sum_{i: y_i = 1} \langle x_i, x \rangle - \frac{1}{n_-} \sum_{i: y_i = -1} \langle x_i, x \rangle + \frac{1}{2} \left(\frac{1}{n_+^2} \sum_{\substack{i: y_i = 1 \\ j: y_j = 1}} \langle x_i, x_j \rangle - \frac{1}{n_-^2} \sum_{\substack{i: y_i = -1 \\ j: y_j = -1}} \langle x_i, x_j \rangle \right) \right\}$$

\uparrow $k(x_i, x)$ \uparrow $k(x_i, x_j)$

Poll

True or false: To apply the kernel trick using a kernel k , we need to know the output dimension of the feature map Φ associated to k .

(A) True

(B) False

Kernel Ridge Regression

- We will argue that ridge regression is kernelizable, and use the kernel trick to extend it to a nonlinear regression method called *kernel ridge regression*.
- To simplify the presentation, we'll consider ridge regression *without offset*

$$\min_w \frac{1}{n} \sum (y_i - w^T x_i)^2 + \lambda \|w\|^2$$

-
- Since some Φ contain a constant term, as with the inhomogeneous polynomial kernel, the offset is not always needed.
 - KRR with offset is described in my lecture notes.

Poll

- Introduce the notation

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}$$

- The solution to

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

is

(A) $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

(B) $\hat{\mathbf{w}} = (\mathbf{X} \mathbf{X}^T + n\lambda \mathbf{I})^{-1} \mathbf{X} \mathbf{y}$

(C) $\hat{\mathbf{w}} = (\mathbf{X} \mathbf{X}^T + n\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

(D) $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X} \mathbf{y}$

Kernel Ridge Regression

- Note $\mathbf{X}^T \mathbf{X}$ is *not* the Gram matrix of the training data.
- Idea: apply matrix inversion lemma:

$$(\mathbf{P} + \mathbf{QRS})^{-1} = \mathbf{P}^{-1} - \mathbf{P}^{-1} \mathbf{Q} (\mathbf{R}^{-1} + \mathbf{S} \mathbf{P}^{-1} \mathbf{Q})^{-1} \mathbf{S} \mathbf{P}^{-1}$$

- After simplification, the solution of RR without offset is

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \hat{\mathbf{w}}^T \mathbf{x} \\ &= \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{x} \\ &= \mathbf{y}^T (\mathbf{X} \mathbf{X}^T + n\lambda \mathbf{I})^{-1} \underline{\mathbf{X} \mathbf{x}} \\ &= \mathbf{y}^T (\mathbf{K} + n\lambda \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}) \end{aligned} \quad \text{MIL + algebra}$$

where

$$\mathbf{K} = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \cdots & \langle \mathbf{x}_1, \mathbf{x}_n \rangle \\ \vdots & \ddots & \vdots \\ \langle \mathbf{x}_n, \mathbf{x}_1 \rangle & \cdots & \langle \mathbf{x}_n, \mathbf{x}_n \rangle \end{bmatrix} \quad \mathbf{k}(\mathbf{x}) = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{x}_n, \mathbf{x} \rangle \end{bmatrix}$$

Kernel Ridge Regression

- If we redefine

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} \quad \mathbf{k}(\mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}) \\ \vdots \\ k(\mathbf{x}_n, \mathbf{x}) \end{bmatrix}$$

then the function estimated by KRR is

$$\hat{f}(\mathbf{x}) = \underbrace{\mathbf{y}^T (\mathbf{K} + n\lambda \mathbf{I})^{-1}}_{\boldsymbol{\alpha}^T} \mathbf{k}(\mathbf{x})$$

- The form of the estimated function is

$$\hat{f}(\mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{k}(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

Summary

- Kernels: Convenient way to incorporate nonlinear feature maps into an algorithm
- Kernel ridge regression: powerful nonlinear regression method, useful element of the ML toolbox
- See lecture notes for KRR with offset
- Neural tangent kernel: interesting connections to neural networks
- Next two lectures: We will kernelize the optimal soft-margin hyperplane classifier, leading to the support vector machine

$$k(x, x') = \exp(-\gamma \|x - x'\|^2)$$

$d=1$

$$\sum \alpha_i k(x, x_i)$$

