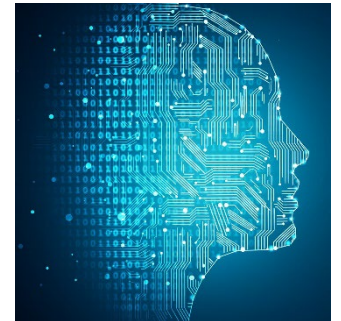


Principles of Machine Learning

Kernels



Kevin Moon (kevin.moon@usu.edu)
STAT/CS 5810/6655



One weird kernel trick



- <http://oneweirdkerneltrick.com/>

Overview



- So far we have mainly focused on linear methods
- There are several nonlinear methods that we will consider in this course
- The first class of nonlinear methods that we will study are “kernel methods”
- We start here because kernel methods build directly on the linear methods we have just studied

Outline



1. Nonlinear classification via nonlinear feature maps
2. Inner product kernels
3. Symmetric positive definite kernels
4. The kernel trick

Nonlinear feature maps



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

$$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

and apply a linear method to the transformed data

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

- Nonlinear regression:

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

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

- Nonlinear classification:

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

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

Polynomial Regression



- Determine the least squares cubic polynomial fit to training data $(x_1, y_1), \dots, (x_n, y_n)$ where $x_i, y_i \in \mathbb{R}$

- Can write

$$f(x) = a + bx + cx^2 + dx^3 = \mathbf{w}^T \mathbf{\Phi}(x) + a$$

where

$$\mathbf{w} = \begin{bmatrix} b \\ c \\ d \end{bmatrix}, \quad \mathbf{\Phi}(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}$$

- Note that $m = 3$ and $\mathbf{\Phi}(x)$ is a nonlinear function

Polynomial Regression



- The empirical risk (using squared error loss):

$$\sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n (y_i - \mathbf{w}^T \boldsymbol{\Phi}(x) - a)^2$$

- The minimizer is:

$$\begin{bmatrix} a \\ \mathbf{w} \end{bmatrix} = (X^T X)^{-1} X^T \mathbf{y}$$

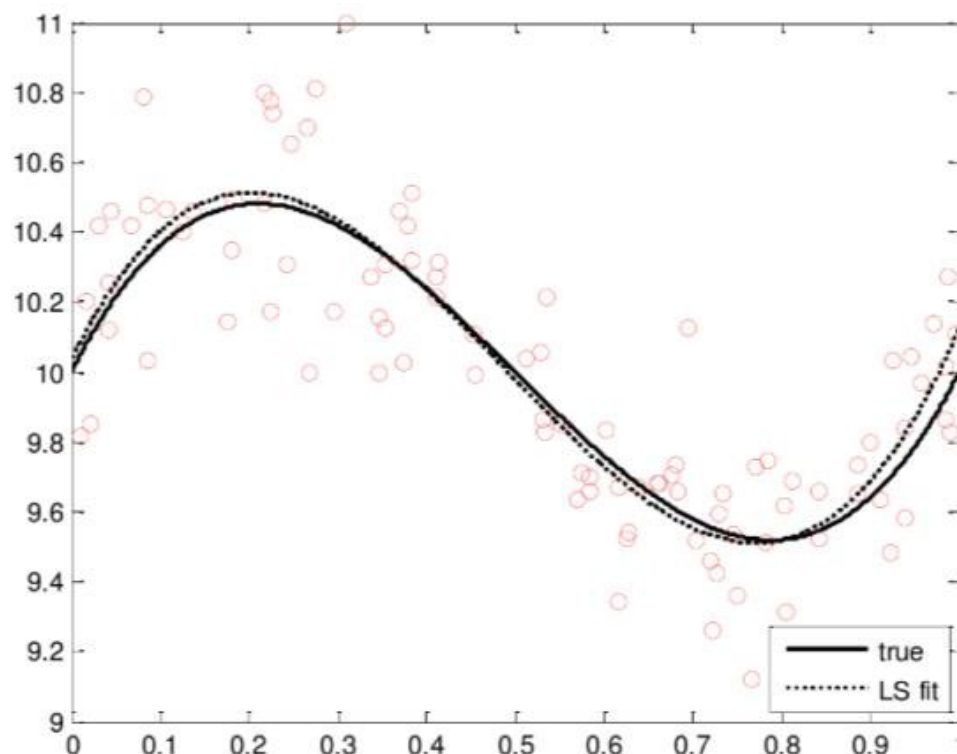
where

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{bmatrix}$$

Polynomial Regression



- For large m (i.e. polynomial degree) relative to n , regularization becomes necessary
- Otherwise, the matrix $X^T X$ becomes ill-conditioned and regularization is needed to avoid overfitting



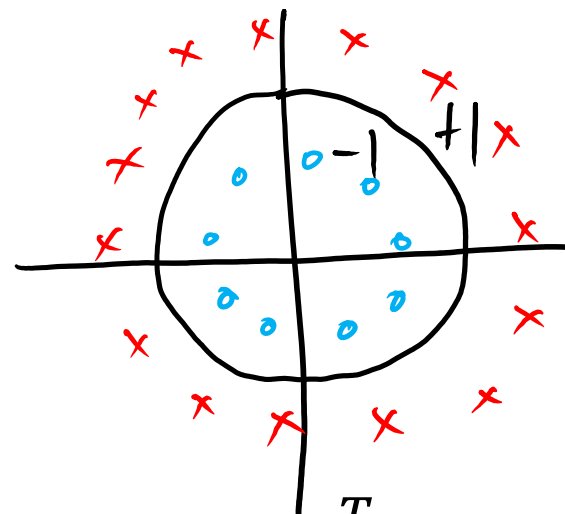
Binary Classification



- Write $\mathbf{x} = [x^{(1)} \quad x^{(2)}]^T \in \mathbb{R}^2$

- Consider

$$\Phi(\mathbf{x}) = \left[(x^{(1)})^2 \quad (x^{(2)})^2 \quad x^{(1)} \quad x^{(2)} \right]^T$$



Binary Classification



- Training data are separated by a circular classifier
$$\mathbf{x} \mapsto \text{sign} \left\{ \left(x^{(1)} - c^{(1)} \right)^2 + \left(x^{(2)} - c^{(2)} \right)^2 - r^2 \right\}$$
for a certain radius r and center

$$\mathbf{c} = \begin{bmatrix} c^{(1)} \\ c^{(2)} \end{bmatrix}$$

- This is a linear classifier in the transformed space where

$$\mathbf{w} = \begin{bmatrix} 1 \\ 1 \\ -2c^{(1)} \\ -2c^{(2)} \end{bmatrix} \quad b = \left(c^{(1)} \right)^2 + \left(c^{(2)} \right)^2 - r^2$$

- In this example, you can also absorb the offset b into Φ

Inner Product Kernels



- Problem with previous approach: m can explode as d increases
 - Thus it can be difficult to store/compute/manipulate Φ directly
- Fortunately, the following facts allow us to use nonlinear feature maps for large d :
 - 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 infinite!
- k is called an *inner product kernel*
- Let's look at some examples using the *dot product*



Homogeneous Polynomial Kernel



- Degree $p = 2, d = 2$

$$\begin{aligned}k(\mathbf{u}, \mathbf{v}) &= (\mathbf{u}^T \mathbf{v})^2 \\&= (u^{(1)}v^{(1)} + u^{(2)}v^{(2)})^2 \\&= (u^{(1)})^2(v^{(1)})^2 + 2u^{(1)}u^{(2)}v^{(1)}v^{(2)} + (u^{(2)})^2(v^{(2)})^2 \\&= \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle\end{aligned}$$

where

$$\Phi(\mathbf{u}) = \left[(u^{(1)})^2, \sqrt{2}u^{(1)}u^{(2)}, (u^{(2)})^2 \right]^T.$$



Homogeneous Polynomial Kernel



- Degree $p = 2$, d arbitrary

$$\begin{aligned}k(\mathbf{u}, \mathbf{v}) &= (\mathbf{u}^T \mathbf{v})^2 \\&= \left(\sum_{i=1}^d u^{(i)} v^{(i)} \right)^2 \\&= \sum_{i=1}^d (u^{(i)})^2 (v^{(i)})^2 + \sum_{i < j} 2u^{(i)} u^{(j)} v^{(i)} v^{(j)} \\&= \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle\end{aligned}$$

where

$$\Phi(\mathbf{u}) = \left[(u^{(1)})^2, \dots, (u^{(d)})^2, \sqrt{2}u^{(1)}u^{(2)}, \dots, \sqrt{2}u^{(d-1)}u^{(d)} \right]^T.$$

- $m = d + \frac{d(d-1)}{2}$

Group Exercise



1. Let $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^3$ where $\mathbf{u}, \mathbf{v} \in \mathbb{R}^2$. Find Φ such that

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

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



- General homogeneous polynomial kernel:

$$\begin{aligned} k(\mathbf{u}, \mathbf{v}) &= (\mathbf{u}^T \mathbf{v})^p \\ &= \sum_{(j_1, \dots, j_d): \sum j_i = p} \binom{p}{j_1 \dots j_d} (u^{(1)})^{j_1} \dots (u^{(d)})^{j_d} (v^{(1)})^{j_1} \dots (v^{(d)})^{j_d} \\ &\Rightarrow \Phi(\mathbf{u}) = [\dots, \sqrt{\binom{p}{j_1 \dots j_d}} (u^{(1)})^{j_1} \dots (u^{(d)})^{j_d}, \dots]^T \end{aligned}$$

- For the general inhomogeneous polynomial kernel $\Phi =$ all monomials in d variables up to degree p
- All the preceding examples involve the dot product, but some important kernels involve other inner products

Inner Products



- A (*real*) *inner product space* is a vector space V on which we can define a function $\langle \mathbf{u}, \mathbf{v} \rangle$ (called an *inner product*) such that

1. $\forall \alpha_1, \alpha_2 \in \mathbb{R}, \mathbf{u}_1, \mathbf{u}_2, \mathbf{v} \in V$

$$\langle \alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2, \mathbf{v} \rangle = \alpha_1 \langle \mathbf{u}_1, \mathbf{v} \rangle + \alpha_2 \langle \mathbf{u}_2, \mathbf{v} \rangle$$

2. $\forall \mathbf{u}, \mathbf{v} \in V$

$$\langle \mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{u} \rangle$$

3. $\forall \mathbf{u} \in V,$

$$\langle \mathbf{u}, \mathbf{u} \rangle \geq 0,$$

with equality iff $\mathbf{u} = \mathbf{0}$.

- Recall: We say $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is an *inner product kernel* if \exists an inner product space V and a feature map $\Phi : \mathbb{R}^d \rightarrow V$ such that

$$k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$$

- *Note:* Φ and V are not unique for a given k

Symmetric Positive Definite 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 $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$

- 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

Important Kernels



- Homogeneous polynomial kernel

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

- Inhomogeneous polynomial kernel

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

- Gaussian kernel

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

- For the Gaussian kernel, V is infinite dimensional!
- This is ok, though. In kernel methods we don't have to work with Φ , just k

Big Picture: The Kernel Trick



- Using kernels, we can obtain nonlinear methods from linear methods as follows:
 1. Select an IP/SPD kernel k
 2. Formulate your linear method such that feature vectors (i.e., the training data and an arbitrary test instance) only appear via inner products $\langle \mathbf{x}, \mathbf{x}' \rangle$
 3. Replace $\langle \mathbf{x}, \mathbf{x}' \rangle$ with
$$k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$$
throughout the algorithm
- This idea is called the *kernel trick*
- The resulting method is equivalent to applying the original linear method to the non-linearly transformed data $(\Phi(\mathbf{x}_1), y_1), \dots, (\Phi(\mathbf{x}_n), y_n)$. With the kernel trick, we never have to compute Φ , just k .

Applications of the kernel trick



- Many standard methods can be kernelized
 - Ridge regression (kernel ridge regression)
 - Learns a nonlinear regression function
 - PCA (kernel PCA)
 - Learns a nonlinear dimensionality reduction
 - Optimal soft-margin hyperplane (the support vector machine)
 - Learns a nonlinear decision boundary
 - Others...
- We'll cover kernel ridge regression and the support vector machine in the next few lectures

Reproducing Kernel Hilbert Spaces (RKHS)



- Can define a space (set) of real-valued functions defined in terms of a positive definite kernel.
 - Referred to as an RKHS
- Optimizing over an RKHS can give many ML algorithms
 - Kernel ridge regression
 - SVM
 - Kernel logistic regression
 - Semi-supervised learning
- Provides a very rich theory that is often used to develop new methods
- If time, we'll go into more detail later in the course

Further reading



- ESL Section 5.8