

Kernel Methods

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Outline

- 1 A general view of supervised learning
- 2 Kernel methods
- 3 Example
- 4 Kernels
- 5 Another example

Supervised learning

We aim to build a function $h(\boldsymbol{x})$ to predict the true value y associated with \boldsymbol{x} . If we make a mistake, we incur a *loss*

$$\ell(y, h(\boldsymbol{x}))$$

Supervised learning

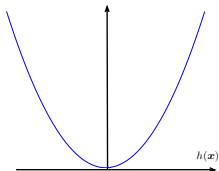
We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

$$\ell(y, h(\mathbf{x}))$$

Example: squared loss function for regression when y is continuous

$$\ell(y, h(\mathbf{x})) = [h(\mathbf{x}) - y]^2$$

Ex: when $y = 0$



Other types of loss functions

For classification: 0/1 loss

$$\ell(y, h(\mathbf{x})) = \mathbf{1}\{y \neq h(\mathbf{x})\}$$

Assumes h outputs values in $\{-1, +1\}$.

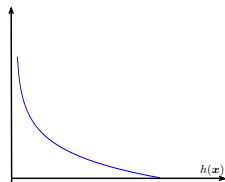
Other types of loss functions

For classification: *logistic* loss

$$\ell(y, h(\mathbf{x})) = -y \log h(\mathbf{x}) - (1-y) \log[1-h(\mathbf{x})]$$

Assumes h outputs values in $[0, 1]$.

Ex: when $y = 1$



Measure how good our hypothesis h is

Risk/ expected test loss: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$\mathcal{R}[h(\mathbf{x})] = \sum_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y)$$

Measure how good our hypothesis h is

Risk/ expected test loss: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$\mathcal{R}[h(\mathbf{x})] = \sum_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y)$$

However, we cannot compute $\mathcal{R}[h(\mathbf{x})]$, so we use *empirical risk/training error*, given a training dataset \mathcal{D}

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Measure how good our hypothesis h is

Risk/ expected test loss: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$\mathcal{R}[h(\mathbf{x})] = \sum_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y)$$

However, we cannot compute $\mathcal{R}[h(\mathbf{x})]$, so we use *empirical risk/training error*, given a training dataset \mathcal{D}

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Intuitively, as $N \rightarrow +\infty$,

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] \rightarrow \mathcal{R}[h(\mathbf{x})]$$

Pick a good hypothesis h

A good hypothesis h is one that minimizes risk/expected test loss $\mathcal{R}[h(\mathbf{x})]$ but we cannot even compute it!

Instead pick h that minimizes empirical risk/training error $\mathcal{R}^{\text{EMP}}[h(\mathbf{x})]$

This strategy is known as **empirical risk minimization**.

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, and we use squared loss
- For logistic regression, $h_{\mathbf{w},b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$, and we use logistic loss
- Finding the best h is achieved by searching for (\mathbf{w}, b) that minimizes the training error/empirical risk.

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, and we use squared loss
- For logistic regression, $h_{\mathbf{w},b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$, and we use logistic loss
- Finding the best h is achieved by searching for (\mathbf{w}, b) that minimizes the training error/empirical risk.

ERM might be problematic

- If $h(\mathbf{x})$ is complicated enough,

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] \rightarrow 0$$

- But then $h(\mathbf{x})$ is unlikely to do well in predicting things out of the training dataset \mathcal{D}
- This is called *poor generalization* or *overfitting*. We have just discussed approaches to address this issue.

Regularizer

Instead of \mathcal{R}^{emp} , use

$$\begin{aligned} & \arg \min_{\mathbf{w}, b} \mathcal{R}^{\text{EMP}}[h_{\mathbf{w}, b}(\mathbf{x})] + \lambda R(\mathbf{w}, b) \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \ell(y_n, h_{\mathbf{w}, b}(\mathbf{x}_n)) + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$

Loss functions ℓ , $\hat{y} = \mathbf{w}^T \mathbf{x} + b$

- Zero/One: $\ell(y, h(\mathbf{x})) = \mathbf{1}\{y \neq h(\mathbf{x})\}$
- Squared loss: $\ell(y, h(\mathbf{x})) = [h(\mathbf{x}) - y]^2$
- Logistic loss: $\ell(y, h(\mathbf{x})) = -y \log h(\mathbf{x}) - (1 - y) \log[1 - h(\mathbf{x})]$

Regularizer

Instead of \mathcal{R}^{emp} , use

$$\begin{aligned} & \arg \min_{\mathbf{w}, b} \mathcal{R}^{\text{EMP}}[h_{\mathbf{w}, b}(x)] + \lambda R(\mathbf{w}, b) \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \ell(y_n, h_{\mathbf{w}, b}(\mathbf{x}_n)) + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$

Regularizer

- Squared 2-norm: $R(\mathbf{w}, b) = \|\mathbf{w}\|_2^2 = \sum_{d=1}^D w_d^2$
- 1-norm: $R(\mathbf{w}, b) = \|\mathbf{w}\|_1 = \sum_{d=1}^D |w_d|$.
- 0-norm: $R(\mathbf{w}, b) = \sum_{d=1}^D \mathbf{1}\{w_d \neq 0\}$.
- p -norm: $R(\mathbf{w}, b) = \|\mathbf{w}\|_p = (\sum_{d=1}^D |w_d|^p)^{\frac{1}{p}}$

Framework for supervised learning

- Pick:

- ▶ Model/hypotheses
- ▶ Loss function
- ▶ Regularizer
- ▶ Algorithm to solve optimization problem

These choices lead to different learning algorithms

Framework for machine learning

- Application
 - ▶ Labeled vs unlabeled data
 - ▶ Labeled: supervised learning. Type of label: categorical (classification), quantitative (regression)
 - ▶ Unlabeled: unsupervised learning.
- Model/hypotheses
- Optimization problem
- Algorithm to solve optimization problem

Outline

- 1 A general view of supervised learning
- 2 Kernel methods
 - Motivation
- 3 Example
- 4 Kernels
- 5 Another example

Motivation

- Linear models are convenient.
 - ▶ Computationally efficient for learning (training) and prediction.
- We would like our models to be “expressive”.
 - ▶ If it is not expressive enough, it will underfit.
 - ▶ Too expressive models can overfit.

Motivation

How to increase the expressive power of linear models?

- Map the feature vector \mathbf{x} to an expanded version $\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$.
- Use a non-linear basis function $\phi(\mathbf{x})$ as input to linear model.

Motivation

How to increase the expressive power of linear models?

- Map the feature vector \mathbf{x} to an expanded version $\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$.
- Use a non-linear basis function $\phi(\mathbf{x})$ as input to linear model.
- **Difficulty:** When M is large, computational difficulty.

Motivation

How to choose nonlinear basis function for regression?

$$\mathbf{w}^T \phi(\mathbf{x})$$

where $\phi(\cdot)$ maps the original feature vector \mathbf{x} to a M -dimensional *new* feature vector. In the following, we will show that we can sidestep the issue of choosing which $\phi(\cdot)$ to use — instead, we will choose *equivalently* a *kernel function* that are often computationally efficient to work with.

Example

$$\phi(\mathbf{x}) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \vdots \\ \sqrt{2}x_D \\ x_1^2 \\ x_1x_2 \\ \vdots \\ x_1x_D \\ x_2x_1 \\ x_2^2 \\ \vdots \\ x_2x_D \\ \vdots \\ x_Dx_1 \\ \vdots \\ x_D^2 \end{pmatrix}$$

Example

Computing $\mathbf{w}^T \phi(\mathbf{x})$ is $O(D^2)$.

Example

Many learning algorithms can be rewritten to depend on the instances $\mathbf{x}_i, \mathbf{x}_j$ only through inner products $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$.

Why is this helpful ?

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$$

Inner product can be computed in $O(D)$.

Outline

- 1 A general view of supervised learning
- 2 Kernel methods
- 3 **Example**
 - Kernelized ridge regression
- 4 Kernels
- 5 Another example

Kernelized ridge regression

Ridge regression

$$J(\mathbf{w}) = \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2$$

Kernelized ridge regression

Ridge regression

$$J(\mathbf{w}) = \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2$$

Its solution is given by

$$\nabla J(\mathbf{w}) = 2 \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))(-\phi(\mathbf{x}_n)) + 2\lambda \mathbf{w} = 0$$

Solution

The optimal parameter vector is a linear combination of features

$$\hat{\mathbf{w}} = \sum_n \frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) = \sum_n \alpha_n \phi(\mathbf{x}_n) = \Phi^T \boldsymbol{\alpha}$$

where we have designated $\frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))$ as α_n . And the matrix Φ is the *design matrix* made of *transformed* features.

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}$$

where M is the dimensionality of $\phi(\mathbf{x})$.

Solution

The optimal parameter vector is a linear combination of features

$$\hat{\mathbf{w}} = \sum_n \frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) = \sum_n \alpha_n \phi(\mathbf{x}_n) = \Phi^T \boldsymbol{\alpha}$$

where we have designated $\frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))$ as α_n . And the matrix Φ is the *design matrix* made of *transformed* features.

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}$$

where M is the dimensionality of $\phi(\mathbf{x})$.

Of course, we do not know what $\boldsymbol{\alpha}$ (the vector of all α_n) corresponds to $\hat{\mathbf{w}}$!

Dual formulation

We substitute $\hat{\mathbf{w}} = \Phi^T \alpha$ into $J(\mathbf{w})$, and obtain the following function of α

$$J(\alpha) = \alpha^T \Phi \Phi^T \Phi \Phi^T \alpha - 2(\Phi \Phi^T \mathbf{y})^T \alpha + \lambda \alpha^T \Phi \Phi^T \alpha + \text{const}$$

Dual formulation

We substitute $\hat{\mathbf{w}} = \Phi^T \alpha$ into $J(\mathbf{w})$, and obtain the following function of α

$$J(\alpha) = \alpha^T \Phi \Phi^T \Phi \Phi^T \alpha - 2(\Phi \Phi^T \mathbf{y})^T \alpha + \lambda \alpha^T \Phi \Phi^T \alpha + \text{const}$$

Before we show how $J(\alpha)$ is derived, we make an important observation. We see repeated structures $\Phi \Phi^T$, to which we refer as *Gram matrix* or *kernel matrix*

$$\begin{aligned} K &= \Phi \Phi^T \\ &= \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times N} \end{aligned}$$

Properties of the matrix \mathbf{K}

- Symmetric

$$K_{mn} = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = K_{nm}$$

- Positive semidefinite: for any vector \mathbf{a}

$$\mathbf{a}^T \mathbf{K} \mathbf{a} = (\Phi^T \mathbf{a})^T (\Phi^T \mathbf{a}) \geq 0$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2 \\ &= \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2 \\ &= \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \\ J(\boldsymbol{\alpha}) &= \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 \end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2 \\ &= \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \\ J(\boldsymbol{\alpha}) &= \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 \\ &= \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} \end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned}J(\mathbf{w}) &= \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2 \\&= \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \\J(\boldsymbol{\alpha}) &= \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 \\&= \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} \\&= (\mathbf{y} - \mathbf{K} \boldsymbol{\alpha})^T (\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}) + \lambda \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned}J(\mathbf{w}) &= \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2 \\&= \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \\J(\boldsymbol{\alpha}) &= \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 \\&= \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} \\&= (\mathbf{y} - \mathbf{K} \boldsymbol{\alpha})^T (\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}) + \lambda \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} \\&= \boldsymbol{\alpha}^T \mathbf{K}^T \mathbf{K} \boldsymbol{\alpha} - 2\mathbf{y}^T \mathbf{K} \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} + \text{const}\end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned}J(\mathbf{w}) &= \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|_2^2 \\&= \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \\J(\boldsymbol{\alpha}) &= \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 \\&= \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} \\&= (\mathbf{y} - \mathbf{K} \boldsymbol{\alpha})^T (\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}) + \lambda \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} \\&= \boldsymbol{\alpha}^T \mathbf{K}^T \mathbf{K} \boldsymbol{\alpha} - 2\mathbf{y}^T \mathbf{K} \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} + \text{const} \\&= \boldsymbol{\alpha}^T \mathbf{K}^2 \boldsymbol{\alpha} - 2(\mathbf{K} \mathbf{y})^T \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} = J(\boldsymbol{\alpha})\end{aligned}$$

where we have used the property that \mathbf{K} is symmetric.

Optimal α

$$\nabla J(\alpha) = 2\mathbf{K}^2\alpha - 2\mathbf{K}\mathbf{y} + 2\lambda\mathbf{K}\alpha = 0$$

which leads to (assuming \mathbf{K} is invertible).

$$\hat{\alpha} = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$$

Note that we only need to know \mathbf{K} in order to compute $\hat{\alpha}$ — the exact form of $\phi(\cdot)$ is not essential — as long as we know how to get inner products $\phi(\mathbf{x}_m)^T\phi(\mathbf{x}_n)$. That observation will give rise to the use of *kernel function*.

Optimal α

$$\nabla J(\alpha) = 2\mathbf{K}^2\alpha - 2\mathbf{K}\mathbf{y} + 2\lambda\mathbf{K}\alpha = 0$$

which leads to (assuming \mathbf{K} is invertible).

$$\hat{\alpha} = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$$

Note that we only need to know \mathbf{K} in order to compute $\hat{\alpha}$ — the exact form of $\phi(\cdot)$ is not essential — as long as we know how to get inner products $\phi(\mathbf{x}_m)^T\phi(\mathbf{x}_n)$. That observation will give rise to the use of *kernel function*.

Note that computing the parameter vector does require knowledge of Φ

$$\hat{\mathbf{w}} = \Phi^T\hat{\alpha} = \Phi^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$$

Computing prediction needs only inner products too!

Given a test point $\phi(x)$, for prediction, we must compute:

$$\hat{\mathbf{w}}^T \phi(x) = \mathbf{y}^T (\mathbf{K} + \lambda I)^{-1} \Phi \phi(x)$$

Computing prediction needs only inner products too!

Given a test point $\phi(x)$, for prediction, we must compute:

$$\begin{aligned}\hat{\mathbf{w}}^T \phi(x) &= \mathbf{y}^T (\mathbf{K} + \lambda I)^{-1} \Phi \phi(x) \\ &= \mathbf{y}^T (\mathbf{K} + \lambda I)^{-1} \begin{pmatrix} \phi(x_1)^T \phi(x) \\ \phi(x_2)^T \phi(x) \\ \vdots \\ \phi(x_N)^T \phi(x) \end{pmatrix} = \mathbf{y}^T (\mathbf{K} + \lambda I)^{-1} \mathbf{k}_x\end{aligned}$$

where we have used the property that $(\mathbf{K} + \lambda I)^{-1}$ is symmetric (as \mathbf{K} is) and use \mathbf{k}_x as a shorthand notation for the column vector.

Note that, to make a prediction, once again, we *only need to know how to get $\phi(x_n)^T \phi(x)$* .

Outline

- 1 A general view of supervised learning
- 2 Kernel methods
- 3 Example
- 4 Kernels**
 - Kernel matrix and kernel functions
 - Kernelized machine learning methods
- 5 Another example

Inner products between features

Let us examine more closely the inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$ for a pair of data points \mathbf{x}_m and \mathbf{x}_n .

Polynomial-based nonlinear basis functions consider the following $\phi(\mathbf{x})$:

$$\phi : \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

Inner products between features

Let us examine more closely the inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$ for a pair of data points \mathbf{x}_m and \mathbf{x}_n .

Polynomial-based nonlinear basis functions consider the following $\phi(\mathbf{x})$:

$$\phi : \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

This gives rise to an inner product in a special form,

$$\begin{aligned} \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) &= x_{m1}^2 x_{n1}^2 + 2x_{m1}x_{m2}x_{n1}x_{n2} + x_{m2}^2 x_{n2}^2 \\ &= (x_{m1}x_{n1} + x_{m2}x_{n2})^2 = (\mathbf{x}_m^T \mathbf{x}_n)^2 \end{aligned}$$

Namely, the inner product can be computed by a function $(\mathbf{x}_m^T \mathbf{x}_n)^2$ defined in terms of the original features, *without computing $\phi(\cdot)$* .

Kernel functions

Some intuition

- To use kernelized ridge regression, we need to be able to compute an inner product between a test point and any training point.
- Since we need this for any possible pair of points, we need a function that takes a pair of points and computes an inner product.
- This is the kernel function $k(\cdot, \cdot)$.
- Given two inputs, $k(\cdot, \cdot)$ tells us how “similar” or “close” these inputs are in the space defined by the function ϕ .

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $\mathbf{x}_m, \mathbf{x}_n \in \mathbb{R}^D$, $c \geq 0$ and d is a positive integer.

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $\mathbf{x}_m, \mathbf{x}_n \in \mathbb{R}^D$, $c \geq 0$ and d is a positive integer.

Gaussian kernel, RBF kernel, or Gaussian RBF kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

- Only depends on difference between two inputs
- Corresponds to a feature space with *infinite* dimensions (but we can work directly with the original features)!

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $\mathbf{x}_m, \mathbf{x}_n \in \mathbb{R}^D$, $c \geq 0$ and d is a positive integer.

Gaussian kernel, RBF kernel, or Gaussian RBF kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

- Only depends on difference between two inputs
- Corresponds to a feature space with *infinite* dimensions (but we can work directly with the original features)!

These kernels have hyperparameters to be tuned: d , c , σ^2

Kernel functions

Definition: a kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Not very useful though

Kernel functions

Definition: a kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Not very useful though

Examples we have seen

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2$$

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

Kernel functions

Definition: a kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Not very useful though

Examples we have seen

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2$$

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

Example that is not a kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2$$

(we'll see why later)

Conditions for being a positive semidefinite kernel function

Mercer theorem (loosely), a bivariate function $k(\cdot, \cdot)$ is a kernel function, if and only if, for *any N and any $\mathbf{x}_1, \mathbf{x}_2, \dots, \text{and } \mathbf{x}_N$* , the matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is positive semidefinite.

Why $\|\mathbf{x}_m - \mathbf{x}_n\|_2^2$ is not a positive semidefinite kernel?

Use the definition of positive semidefinite kernel function. We choose $N = 2$, and compute the matrix

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

This matrix cannot be positive semidefinite as it has both *negative* and positive eigenvalues.

Recap: why use kernel functions?

Without specifying $\phi(\cdot)$, the kernel matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is exactly the same as

$$\begin{aligned} \mathbf{K} &= \Phi \Phi^T \\ &= \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_N) \end{pmatrix} \end{aligned}$$

‘Kernel trick’

Many learning methods depend on computing *inner products* between features — we have seen the example of regularized least squares. For those methods, we can use a kernel function in the place of the inner products, i.e., “*kernelizing*” the methods, thus, introducing nonlinear features.

We will present one more to illustrate this “trick” by kernelizing the nearest neighbor classifier.

When we talk about support vector machines, we will see the trick one more time.

Outline

- 1 A general view of supervised learning
- 2 Kernel methods
- 3 Example
- 4 Kernels
- 5 Another example
 - Kernelized nearest neighbors

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

We replace all the inner products in the distance with a kernel function $k(\cdot, \cdot)$, arriving at the kernel distance

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_m) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}_m, \mathbf{x}_n)$$

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

We replace all the inner products in the distance with a kernel function $k(\cdot, \cdot)$, arriving at the kernel distance

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_m) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}_m, \mathbf{x}_n)$$

The distance is equivalent to compute the distance between $\phi(\mathbf{x}_m)$ and $\phi(\mathbf{x}_n)$

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = d(\phi(\mathbf{x}_m), \phi(\mathbf{x}_n))$$

where the $\phi(\cdot)$ is the nonlinear mapping function implied by the kernel function.

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

We replace all the inner products in the distance with a kernel function $k(\cdot, \cdot)$, arriving at the kernel distance

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_m) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}_m, \mathbf{x}_n)$$

The distance is equivalent to compute the distance between $\phi(\mathbf{x}_m)$ and $\phi(\mathbf{x}_n)$

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = d(\phi(\mathbf{x}_m), \phi(\mathbf{x}_n))$$

where the $\phi(\cdot)$ is the nonlinear mapping function implied by the kernel function. The nearest neighbor of a point \mathbf{x} is thus found with

$$\arg \min_n d^{\text{KERNEL}}(\mathbf{x}, \mathbf{x}_n)$$

There are infinite numbers of kernels to use!

Rules of composing kernels (this is just a partial list)

- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $ck(\mathbf{x}_m, \mathbf{x}_n)$ is also if $c > 0$.
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $\alpha k_1(\mathbf{x}_m, \mathbf{x}_n) + \beta k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also if $\alpha, \beta \geq 0$
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $k_1(\mathbf{x}_m, \mathbf{x}_n)k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also.
- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $e^{k(\mathbf{x}_m, \mathbf{x}_n)}$ is also.
- ...

In practice, choosing an appropriate kernel is an “art”

People typically start with polynomial and Gaussian RBF kernels or incorporate domain knowledge.

Summary

- Kernels allow us to design algorithms that use rich set of features while being computationally efficient.
- Many machine learning algorithms can be “kernelized”.
- Picking kernels is an art.
- We still need to tune hyperparameters.