

CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Feb 19, 2021

Outline

① Logistics

② Review of last lecture

③ Kernel methods

1 / 32

2 / 32

Logistics

Logistics

Outline

Logistics

① Logistics

② Review of last lecture

• HW 2 is due today, and HW 3 will be assigned!

③ Kernel methods

3 / 32

4 / 32

Outline

1 Logistics

2 Review of last lecture

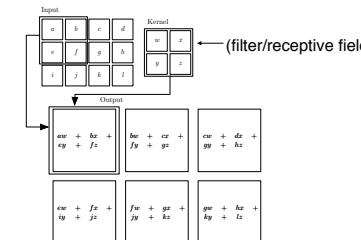
3 Kernel methods

Convolutional Neural Nets

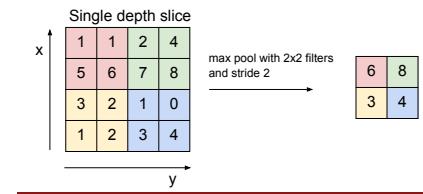
Typical architecture for CNNs:

Input \rightarrow [[Conv \rightarrow ReLU] $^N \rightarrow$ Pool?] $^M \rightarrow$ [FC \rightarrow ReLU] $^Q \rightarrow$ FC

2D Convolution



MAX POOLING



Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 73 April 18, 2017

Outline

Motivation

1 Logistics

2 Review of last lecture

3 Kernel methods

- Motivation
- Kernel Trick
- Dual formulation of linear regression

Recall the question: *how to choose nonlinear basis $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$?*

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

- neural network is one approach: learn ϕ from data
- **kernel method** is another one: sidestep the issue of choosing ϕ by using *kernel functions*

What are Kernels?

Consider the following example, where the data is not linearly separable in the **ambient space** but is separable **feature space**¹

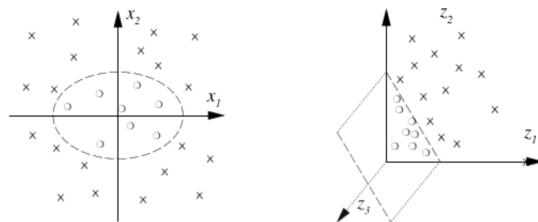


Figure 2.1 Toy example of a binary classification problem mapped into feature space. We assume that the true decision boundary is an ellipse in input space (left panel). The task of the learning process is to estimate this boundary based on empirical data consisting of training points in both classes (crosses and circles, respectively). When mapped into feature space via the nonlinear map $\Phi_2(x) = (z_1, z_2, z_3) = ([x]_1^2, [x]_2^2, \sqrt{2}[x]_1[x]_2)$ (right panel), the ellipse becomes a hyperplane (in the present simple case, it is parallel to the z_3 axis, hence all points are plotted in the (z_1, z_2) plane). This is due to the fact that ellipses can be written as linear equations in the entries of (z_1, z_2, z_3) . Therefore, in feature space, the problem reduces to that of estimating a hyperplane from the mapped data points. Note that via the polynomial kernel (see (2.12) and (2.13)), the dot product in the three-dimensional space can be computed without computing Φ_2 . Later in the book, we shall describe algorithms for constructing hyperplanes which are based on dot products (Chapter 7).

Schölkopf, Bernhard, and Alexander J. Smola. Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT press, 2002.

What are Kernels?

We observe that mapping data points to **higher dimension** feature spaces can help with separability, we can then use our favorite linear methods.

However, there are following issues:

- ① Computations in higher dimensions are cumbersome, and the
- ② Statistical issue of **curse of dimensionality** kicks-in, which means that as dimension increases we may require exponentially more data samples!

Kernel Trick

Wishlist: It would be great to have an *implicit* way to work in higher dimensions without having to do the computations there.

Kernels are special functions which allow us to accomplish exactly this!

How: Kernel functions allow us to compute inner-products in the **feature space** while operating on the data samples in the **ambient space**.

Kernel Trick: To *kernelize* any given algorithm, our aim will be to write the computations as inner-products, and then utilize Kernel functions to do the computations.

Don't need to know $\phi(\cdot)$: Since we use **Kernel function** we actually don't need to know the mapping $\phi(\cdot)$. This means that $\phi(\cdot)$ may be infinite dimensional but we can still evaluate the inner-products in an infinite dimensional feature space!!

Example

Let's take a closer look at the example. Here, we consider the following polynomial basis $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$:

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

What is the inner product between $\phi(\mathbf{x})$ and $\phi(\mathbf{x}')$?

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1 x_2 x_1' x_2' + x_2^2 x_2'^2 \\ &= (x_1 x_1' + x_2 x_2')^2 = (\mathbf{x}^T \mathbf{x}')^2 \end{aligned}$$

Therefore, *the inner product in the new space is simply a function of the inner product in the original space.*

Another example

$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{2D}$ is parameterized by θ :

$$\phi_\theta(\mathbf{x}) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \vdots \\ \cos(\theta x_D) \\ \sin(\theta x_D) \end{pmatrix}$$

What is the inner product between $\phi_\theta(\mathbf{x})$ and $\phi_\theta(\mathbf{x}')$?

$$\begin{aligned} \phi_\theta(\mathbf{x})^\top \phi_\theta(\mathbf{x}') &= \sum_{d=1}^D \cos(\theta x_d) \cos(\theta x'_d) + \sin(\theta x_d) \sin(\theta x'_d) \\ &= \sum_{d=1}^D \cos(\theta(x_d - x'_d)) \end{aligned}$$

Once again, *the inner product in the new space is a simple function of the features in the original space.*

13 / 32

Infinite dimensional mapping

When $L \rightarrow \infty$, even if we cannot compute $\phi(\mathbf{x})$, a vector of *infinite dimension*, we can still compute inner product:

$$\begin{aligned} \phi_\infty(\mathbf{x})^\top \phi_\infty(\mathbf{x}') &= \int_0^{2\pi} \sum_{d=1}^D \cos(\theta(x_d - x'_d)) d\theta \\ &= \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d} \end{aligned}$$

Again, a simple function of the original features.

Note that using this mapping in linear regression, we are *learning a weight w^* with infinite dimension!*

More complicated example

Based on ϕ_θ , define $\phi_L : \mathbb{R}^D \rightarrow \mathbb{R}^{2D(L+1)}$ for some integer L :

$$\phi_L(\mathbf{x}) = \begin{pmatrix} \phi_0(\mathbf{x}) \\ \phi_{\frac{2\pi}{L}}(\mathbf{x}) \\ \phi_{2\frac{2\pi}{L}}(\mathbf{x}) \\ \vdots \\ \phi_{L\frac{2\pi}{L}}(\mathbf{x}) \end{pmatrix}$$

What is the inner product between $\phi_L(\mathbf{x})$ and $\phi_L(\mathbf{x}')$?

$$\begin{aligned} \phi_L(\mathbf{x})^\top \phi_L(\mathbf{x}') &= \sum_{\ell=0}^L \phi_{\frac{2\pi\ell}{L}}(\mathbf{x})^\top \phi_{\frac{2\pi\ell}{L}}(\mathbf{x}') \\ &= \sum_{\ell=0}^L \sum_{d=1}^D \cos\left(\frac{2\pi\ell}{L}(x_d - x'_d)\right) \end{aligned}$$

14 / 32

Kernel functions

Definition: a function $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ is called a *(positive semidefinite) kernel function* if there exists a function $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$ so that for any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

Examples we have seen

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^2$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d}$$

16 / 32

Using kernel functions

Choosing a nonlinear basis ϕ becomes choosing a kernel function.

As long as computing the kernel function is more efficient, we should apply the kernel trick.

Gram/kernel matrix needs to be *positive semi-definite* and *symmetric*

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

In fact, k is a kernel if and only if \mathbf{K} is positive semidefinite for *any N and any $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$* (formalized by the **Mercer theorem**).

- useful for proving that a function is not a kernel

More examples of kernel functions

Two most commonly used kernel functions in practice:

Polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d$$

for $c \geq 0$ and d is a positive integer.

Gaussian kernel or Radial basis function (RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2\sigma^2}}$$

for some $\sigma > 0$.

Examples that are not kernels

Function

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

is *not a kernel*, why?

If it is a kernel, the kernel matrix for two data points \mathbf{x}_1 and \mathbf{x}_2 :

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

must be positive semidefinite, *but is it?*

Composing kernels

Creating more kernel functions using the following rules:

If $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$ are kernels, the followings are kernels too

- conical combination:** $\alpha k_1(\cdot, \cdot) + \beta k_2(\cdot, \cdot)$ if $\alpha, \beta \geq 0$
- product:** $k_1(\cdot, \cdot)k_2(\cdot, \cdot)$
- exponential:** $e^{k(\cdot, \cdot)}$
- ...

Verify using the definition of kernel!

Case study: regularized linear regression

Kernel methods work for *many problems* and we take **regularized linear regression** as an example.

Recall the regularized least square solution, where $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$:

$$\begin{aligned} \mathbf{w}^* &= \underset{\mathbf{w}}{\operatorname{argmin}} F(\mathbf{w}) \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} (\|\Phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_2^2) \\ &= (\Phi^T\Phi + \lambda I)^{-1}\Phi^T\mathbf{y} \end{aligned} \quad \left| \begin{array}{l} \Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \end{array} \right.$$

Issue: *operate in space \mathbb{R}^M and M could be huge or even infinity!*

Our aim: pose these computation as inner-products between $\phi(\cdot)$.

Why is this helpful?

Assuming we know α , the prediction of \mathbf{w}^* on a new example \mathbf{x} is

$$\mathbf{w}^{*\top}\phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

Therefore we do not really need to know \mathbf{w}^* . *Only inner products in the new feature space matter!*

As we've seen, we can now use Kernels to compute inner products *without knowing ϕ* .

Also, *this is a non-parametric method!*

But we need to figure out what α is first!

A closer look at the least square solution, where

By setting the gradient of $F(\mathbf{w}) = \|\Phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_2^2$ to be **0**:

$$\Phi^T(\Phi\mathbf{w}^* - \mathbf{y}) + \lambda\mathbf{w}^* = \mathbf{0}$$

we know

$$\mathbf{w}^* = \frac{1}{\lambda}\Phi^T(\mathbf{y} - \Phi\mathbf{w}^*) = \Phi^T\alpha = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$$

Thus the least square solution is **a linear combination of features!**

Note this is true for perceptron and many other problems.

Of course, the above calculation does not show what α is.

How to find α ?

Plugging in $\mathbf{w} = \Phi^T\alpha$ into $F(\mathbf{w})$ gives

$$\begin{aligned} G(\alpha) &= F(\Phi^T\alpha) \\ &= \|\Phi\Phi^T\alpha - \mathbf{y}\|_2^2 + \lambda\|\Phi^T\alpha\|_2^2 \\ &= \|\mathbf{K}\alpha - \mathbf{y}\|_2^2 + \lambda\alpha^T\mathbf{K}\alpha \quad (\mathbf{K} = \Phi\Phi^T) \\ &= \alpha^T\mathbf{K}^T\mathbf{K}\alpha - 2\mathbf{y}^T\mathbf{K}\alpha + \lambda\alpha^T\mathbf{K}\alpha + \text{cnt.} \\ &= \alpha^T(\mathbf{K}^2 + \lambda\mathbf{K})\alpha - 2\mathbf{y}^T\mathbf{K}\alpha + \text{cnt.} \quad (\mathbf{K}^T = \mathbf{K}) \end{aligned}$$

This is sometime called the *dual formulation* of linear regression.

$\mathbf{K} = \Phi\Phi^T \in \mathbb{R}^{N \times N}$ is called **Gram matrix** or **kernel matrix** where the (i, j) entry is

$$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

Examples of kernel matrix \mathbf{K}

3 data points in \mathbb{R}

$$x_1 = -1, x_2 = 0, x_3 = 1$$

ϕ is polynomial basis with degree 4:

$$\phi(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ x^3 \end{pmatrix}$$

$$\phi(x_1) = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \phi(x_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

25 / 32

Calculation of the Gram matrix \mathbf{K}

$$\phi(x_1) = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \phi(x_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

Gram/Kernel matrix

$$\begin{aligned} \mathbf{K} &= \begin{pmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \phi(x_1)^T \phi(x_3) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \phi(x_2)^T \phi(x_3) \\ \phi(x_3)^T \phi(x_1) & \phi(x_3)^T \phi(x_2) & \phi(x_3)^T \phi(x_3) \end{pmatrix} \\ &= \begin{pmatrix} 4 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 4 \end{pmatrix} \end{aligned}$$

26 / 32

Gram matrix vs covariance matrix

	dimensions	entry (i, j)	property
$\Phi\Phi^T$	$N \times N$	$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$	both are symmetric and positive semidefinite
$\Phi^T\Phi$	$M \times M$	$\sum_{n=1}^N \phi(\mathbf{x}_n)_i \phi(\mathbf{x}_n)_j$	

27 / 32

How to find α ?

Minimize the dual formulation

$$G(\alpha) = \alpha^T (\mathbf{K}^2 + \lambda \mathbf{K}) \alpha - 2y^T \mathbf{K} \alpha + \text{cnt.}$$

Setting the derivative to $\mathbf{0}$ we have

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

Thus $\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} y$ is a **minimizer** and we obtain

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

Exercise: *are there other minimizers? and are there other \mathbf{w}^* 's?*

28 / 32

Comparing two solutions

Minimizing $F(\mathbf{w})$ gives $\mathbf{w}^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y}$

Minimizing $G(\alpha)$ gives $\mathbf{w}^* = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y}$

Note I has different dimensions in these two formulas.

Natural question: *are they the same or different?*

They have to be the same because $F(\mathbf{w})$ has a unique minimizer!

And they are:

$$\begin{aligned} & (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y} \\ &= (\Phi^T \Phi + \lambda I)^{-1} \Phi^T (\Phi \Phi^T + \lambda I) (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y} \\ &= (\Phi^T \Phi + \lambda I)^{-1} (\Phi^T \Phi \Phi^T + \lambda \Phi^T) (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y} \\ &= (\Phi^T \Phi + \lambda I)^{-1} (\Phi^T \Phi + \lambda I) \Phi^T (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y} \\ &= \Phi^T (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y} \end{aligned}$$

29 / 32

Then what is the difference?

First, computing $(\Phi \Phi^T + \lambda I)^{-1}$ can be more efficient than computing $(\Phi^T \Phi + \lambda I)^{-1}$ when $N \leq M$.

More importantly, computing $\alpha = (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y}$ also *only requires computing inner products in the new feature space!*

Now we can conclude that the exact form of $\phi(\cdot)$ is not essential; *all we need is computing inner products $\phi(\mathbf{x})^T \phi(\mathbf{x}')$.*

For some ϕ it is indeed possible to compute $\phi(\mathbf{x})^T \phi(\mathbf{x}')$ without computing/knowing ϕ . This is the *kernel trick*.

30 / 32

Kernelizing other ML algorithms

Kernel trick is applicable to **many ML algorithms**:

- nearest neighbor classifier
- perceptron
- logistic regression
- SVM
- ...

31 / 32

Example: Kernelized NNC

For NNC with **L2 distance**, the key is to compute for any two points \mathbf{x}, \mathbf{x}'

$$d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2^2 = \mathbf{x}^T \mathbf{x} + \mathbf{x}'^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}'$$

With a kernel function k , we simply compute

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

which by definition is the **L2 distance in a new feature space**

$$d^{\text{KERNEL}}(\mathbf{x}, \mathbf{x}') = \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2^2$$

32 / 32