

CSCI567 Machine Learning (Spring 2023)

Week 6: Kernel Methods

Prof. Yan Liu

University of Southern California

Outline

1 Kernel methods

Motivation

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

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

Motivation

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

Motivation

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*

Case study: regularized linear regression

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

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:

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

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:

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

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

A closer look at the least square solution

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

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

A closer look at the least square solution

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

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

we know

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

A closer look at the least square solution

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

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

we know

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

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

A closer look at the least square solution

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

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

we know

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

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

Note this is true for perceptron and many other problems.

A closer look at the least square solution

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

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

we know

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

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 $\boldsymbol{\alpha}$ is.

Why is this helpful?

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

$$w^{*\text{T}}\phi(x) = \sum_{n=1}^N \alpha_n \phi(x_n)^{\text{T}} \phi(x)$$

Why is this helpful?

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

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

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

Why is this helpful?

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

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

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

Kernel methods are exactly about computing inner products *without knowing ϕ* .

Why is this helpful?

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

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

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

Kernel methods are exactly about computing inner products *without knowing ϕ* .

But we need to figure out what α is first!

How to find α ?

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

$$G(\alpha) = F(\Phi^T \alpha)$$

How to find α ?

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

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

How to find α ?

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

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

How to find α ?

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

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

How to find α ?

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

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

How to find α ?

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

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

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

How to find α ?

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

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

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

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

$$\phi(x_i)^T \phi(x_j)$$

Examples of kernel matrix

3 data points in \mathbb{R}

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

ϕ is polynomial basis with degree 4:

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

Examples of kernel matrix

3 data points in \mathbb{R}

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

ϕ is polynomial basis with degree 4:

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

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

Calculation of the Gram matrix

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

Gram/Kernel matrix

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

Gram matrix vs covariance matrix

	dimensions	entry (i, j)	property
$\Phi\Phi^T$			
$\Phi^T\Phi$			

Gram matrix vs covariance matrix

	dimensions	entry (i, j)	property
$\Phi\Phi^T$	$N \times N$		
$\Phi^T\Phi$			

Gram matrix vs covariance matrix

	dimensions	entry (i, j)	property
$\Phi\Phi^T$	$N \times N$		
$\Phi^T\Phi$	$M \times M$		

Gram matrix vs covariance matrix

	dimensions	entry (i, j)	property
$\Phi\Phi^T$	$N \times N$	$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$	
$\Phi^T\Phi$	$M \times M$		

Gram matrix vs covariance matrix

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

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
$\Phi^T\Phi$	$M \times M$	$\sum_{n=1}^N \phi(\mathbf{x}_n)_i \phi(\mathbf{x}_n)_j$	

How to find α ?

Minimize the dual formulation

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

How to find α ?

Minimize the dual formulation

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

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

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

How to find α ?

Minimize the dual formulation

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

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

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

How to find α ?

Minimize the dual formulation

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

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

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

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

How to find α ?

Minimize the dual formulation

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

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

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

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

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

How to find α ?

Minimize the dual formulation

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

Setting the derivative to 0 we have

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

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

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

Exercise: *are there other minimizers?*

Comparing two solutions

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

Comparing two solutions

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

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

Comparing two solutions

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

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

Note I has different dimensions in two formulas.

Comparing two solutions

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

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

Note I has different dimensions in two formulas.

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

Comparing two solutions

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

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

Note I has different dimensions in two formulas.

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

They are the same

Comparing two solutions

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

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

Note I has different dimensions in two formulas.

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

They are the same

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

Comparing two solutions

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

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

Note I has different dimensions in two formulas.

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

They are the same

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

Comparing two solutions

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

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

Note I has different dimensions in two formulas.

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

They are the same

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

Comparing two solutions

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

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

Note I has different dimensions in two formulas.

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

They are the same

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

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 $N \leq M$.

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

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 $K = (\Phi\Phi^T + \lambda I)^{-1}$ 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 $\phi(x)^T\phi(x')$.*

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 $K = (\Phi\Phi^T + \lambda I)^{-1}$ 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 $\phi(x)^T\phi(x')$.*

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

Example

Consider the following polynomial basis $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$:

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

Example

Consider the following polynomial basis $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$:

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

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

$$\phi(\mathbf{x})^T \phi(\mathbf{x}') = x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2$$

Example

Consider the following polynomial basis $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$:

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

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

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 \end{aligned}$$

Example

Consider the following polynomial basis $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$:

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

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

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (\mathbf{x}^T \mathbf{x}')^2 \end{aligned}$$

Example

Consider the following polynomial basis $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$:

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

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

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (\mathbf{x}^T \mathbf{x}')^2 \end{aligned}$$

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 :

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

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 $\phi_{\theta}(\mathbf{x}')$?

$$\phi_{\theta}(\mathbf{x})^T \phi_{\theta}(\mathbf{x}') = \sum_{d=1}^D \cos(\theta x_d) \cos(\theta x'_d) + \sin(\theta x_d) \sin(\theta x'_d)$$

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 $\phi_{\theta}(\mathbf{x}')$?

$$\begin{aligned} \phi_{\theta}(\mathbf{x})^T \phi_{\theta}(\mathbf{x}') &= \sum_{d=1}^D \cos(\theta x_d) \cos(\theta x'_d) + \sin(\theta x_d) \sin(\theta x'_d) \\ &= \sum_{d=1}^D \cos(\theta(x_d - x'_d)) \end{aligned}$$

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 $\phi_{\theta}(\mathbf{x}')$?

$$\begin{aligned} \phi_{\theta}(\mathbf{x})^T \phi_{\theta}(\mathbf{x}') &= \sum_{d=1}^D \cos(\theta x_d) \cos(\theta x'_d) + \sin(\theta x_d) \sin(\theta x'_d) \\ &= \sum_{d=1}^D \cos(\theta(x_d - x'_d)) \end{aligned}$$

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

More complicated example

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

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

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 $\phi_L(\mathbf{x}')$?

$$\begin{aligned} \phi_L(\mathbf{x})^T \phi_L(\mathbf{x}') &= \sum_{\ell=0}^L \phi_{\frac{2\pi\ell}{L}}(\mathbf{x})^T \phi_{\frac{2\pi\ell}{L}}(\mathbf{x}') \\ &= \sum_{\ell=0}^L \sum_{d=1}^D \cos\left(\frac{2\pi\ell}{L}(x_d - x'_d)\right) \end{aligned}$$

Infinite dimensional mapping

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

Infinite dimensional mapping

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

$$\begin{aligned}\phi_{\infty}(\mathbf{x})^T \phi_{\infty}(\mathbf{x}') &= \int_0^{2\pi} \sum_{d=1}^D \cos(\theta(x_d - x'_d)) d\theta \\ &= \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d}\end{aligned}$$

Again, a simple function of the original features.

Infinite dimensional mapping

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

$$\begin{aligned}\phi_{\infty}(\mathbf{x})^T \phi_{\infty}(\mathbf{x}') &= \int_0^{2\pi} \sum_{d=1}^D \cos(\theta(x_d - x'_d)) d\theta \\ &= \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d}\end{aligned}$$

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!*

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 $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

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 $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

Examples we have seen

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$$

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

Using kernel functions

Choosing a nonlinear basis ϕ becomes choosing a kernel function.

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.

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 becomes:

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

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 becomes:

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

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

Using kernel functions

The prediction on a new example \mathbf{x} is

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

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.

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 $\sigma > 0$.

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 $\sigma > 0$.

Think about *what the corresponding ϕ is* for each kernel.

Composing kernels

Creating more kernel functions using the following rules:

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

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

- **linear combination**: $\alpha k_1(\cdot, \cdot) + \beta k_2(\cdot, \cdot)$ if $\alpha, \beta \geq 0$

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

- **linear combination**: $\alpha k_1(\cdot, \cdot) + \beta k_2(\cdot, \cdot)$ if $\alpha, \beta \geq 0$
- **product**: $k_1(\cdot, \cdot)k_2(\cdot, \cdot)$

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

- **linear 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:** $e^{k(\cdot, \cdot)}$

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

- **linear 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!

Examples that are not kernels

Function

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

is *not a kernel*, why?

Examples that are not kernels

Function

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

is *not a kernel*, why?

Example: 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,

Examples that are not kernels

Function

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

is *not a kernel*, why?

Example: 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?*

Kernelizing other ML algorithms

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

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

Example: Kernelized NNC

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

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

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

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

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**

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