

Kernel

Announcements

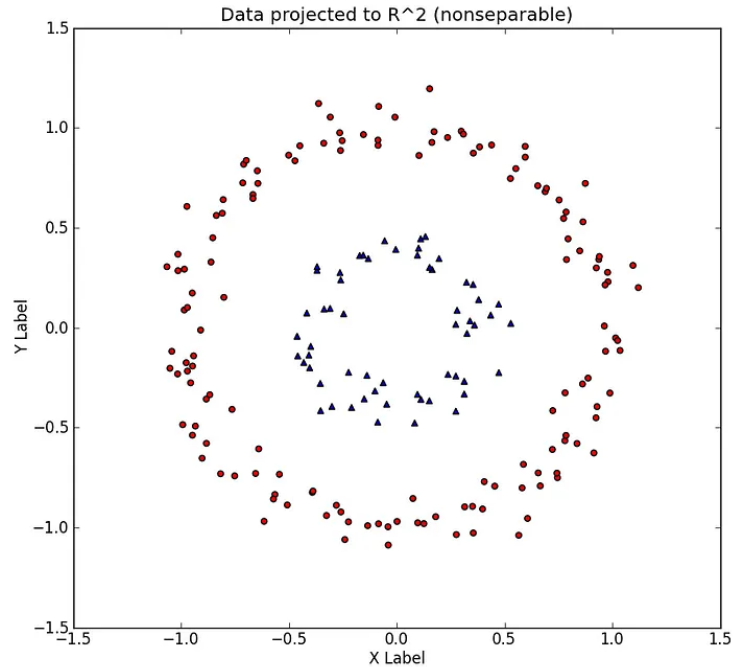
HW5 and P5 are released (due in one week)

Objective today (and next Tuesday)

Use kernels to design nonlinear ML models (regression & classification)

Objective today (and next Tuesday)

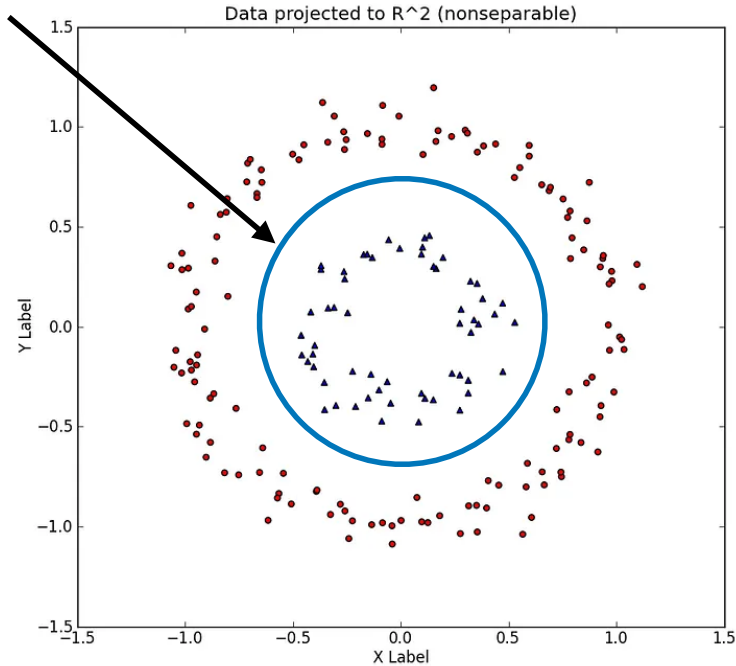
Use kernels to design nonlinear ML models (regression & classification)



Objective today (and next Tuesday)

Use kernels to design nonlinear ML models (regression & classification)

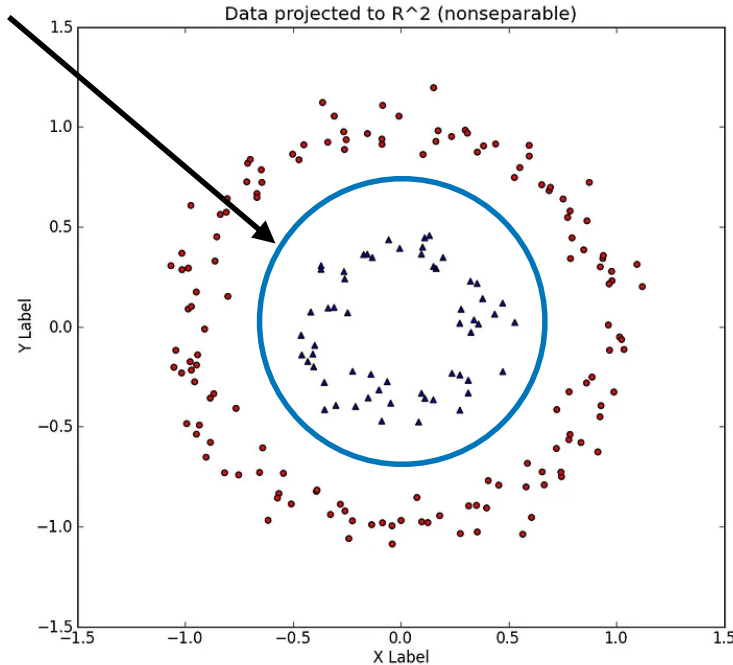
Goal: Non-linear decision boundary



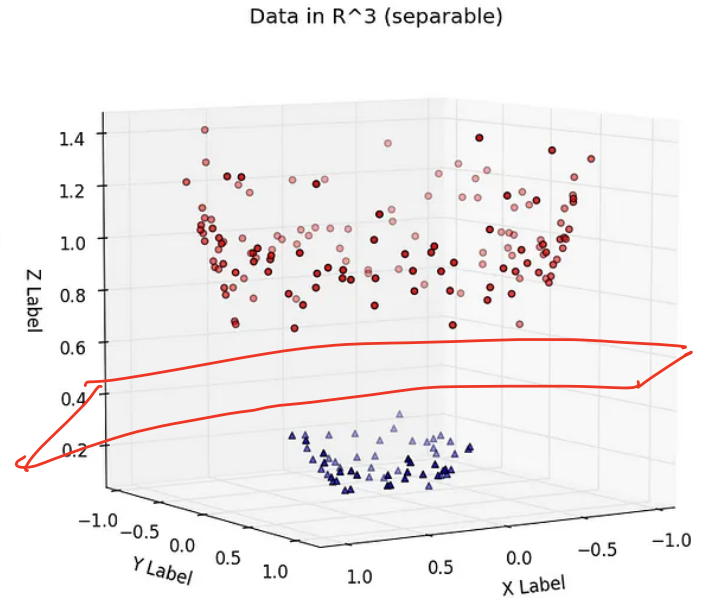
Objective today (and next Tuesday)

Use kernels to design nonlinear ML models (regression & classification)

Goal: Non-linear decision boundary



Our approach



Outline

1. A new perspective on ridge linear regression
2. Feature mapping and Kernel
3. Kernel trick and demo of kernel regression

Linear regression revisited

Dataset $\mathcal{D} = \{\mathbf{x}_i, y_i\}, \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

Linear regression revisited

Dataset $\mathcal{D} = \{\mathbf{x}_i, y_i\}$, $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

Ridge Linear regression solves the following problem:

$$\arg \min_w \sum_{i=1}^n (w^\top \mathbf{x}_i - y_i)^2 + \lambda \|w\|_2^2$$

Linear regression revisited

Dataset $\mathcal{D} = \{\mathbf{x}_i, y_i\}, \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

Ridge Linear regression solves the following problem:

$$\arg \min_w \sum_{i=1}^n (w^\top \mathbf{x}_i - y_i)^2 + \lambda \|w\|_2^2$$

Closed-form solution exists, i.e.,

$$\hat{w} = (XX^\top + \lambda I)^{-1}XY$$

$$X = \begin{bmatrix} | & & | \\ x_1 & \dots & x_n \\ | & & | \end{bmatrix}$$
$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Linear regression revisited

Claim: $\hat{w} = (XX^T + \lambda I)^{-1}XY \in \text{Span}(X)$

$$\hat{w} = \sum_{i=1}^n \alpha_i X_i$$

Linear regression revisited

Claim: $\hat{w} = (XX^T + \lambda I)^{-1}XY \in \text{Span}(X)$

An intuitive proof: GD (or SGD)

Linear regression revisited

$$\text{Claim: } \hat{w} = (XX^T + \lambda I)^{-1}XY \in \text{Span}(X)$$

An intuitive proof: GD (or SGD)

$$w_0 = \mathbf{0}, w^{t+1} = w^t - \eta \left[\sum_{i=1}^n (\mathbf{x}_i^T w^t - y_i) \mathbf{x}_i + \lambda w^t \right]$$

$$w^t \in \text{Span}(X)$$

$$(\cdot) \in \text{Span}(X)$$

A new perspective of linear regression

Since we know optimal solution lives in $\text{span}(X)$, we can re-parameterize

$$w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha, \quad \alpha_i \in \mathbb{R}, \forall i$$

We will learn α instead

A new perspective of linear regression

Since we know optimal solution lives in $\text{span}(X)$, we can re-parameterize

$$w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha, \quad \alpha_i \in \mathbb{R}, \forall i$$

$$\arg \min_w \sum_{i=1}^n \left\| X^T w - Y \right\|_2^2 + \lambda \|w\|_2^2$$

Handwritten notes: $w = \sum_{i=1}^n \alpha_i x_i$ with arrows pointing from the w in the norm terms to the $\alpha_i x_i$ terms.

Original formulation

A new perspective of linear regression

Since we know optimal solution lives in $\text{span}(X)$, we can re-parameterize

$$w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha, \quad \alpha_i \in \mathbb{R}, \forall i$$

$$\arg \min_w \sum_{i=1}^n \|X^\top w - Y\|_2^2 + \lambda \|w\|_2^2 \quad \longrightarrow \quad \arg \min_\alpha \| \underbrace{X^\top X}_\triangle \alpha - Y \|_2^2 + \underbrace{\lambda}_{\triangle} \| \underbrace{X\alpha}_\triangle \|_2^2$$

Original formulation

New formulation w/ α as our variables

A new perspective of linear regression

*n # of data
points*

$$\arg \min_{\alpha} \left\| \underbrace{X^T X}_{\delta} \alpha - Y \right\|_2^2 + \lambda \left\| \underbrace{X \alpha}_{\delta} \right\|_2^2$$

$$\underline{\underline{\alpha \in \mathbb{R}^n}}$$

$$w = \lambda \alpha$$

A new perspective of linear regression

$$\arg \min_{\alpha} \left\| X^T X \alpha - Y \right\|_2^2 + \lambda \|X \alpha\|_2^2$$

Solution:

$$\alpha = \underbrace{(X^T X + \lambda I)}_{n \times n}^{-1} \underbrace{Y}_{\mathbb{R}^n} \in \mathbb{R}^n$$

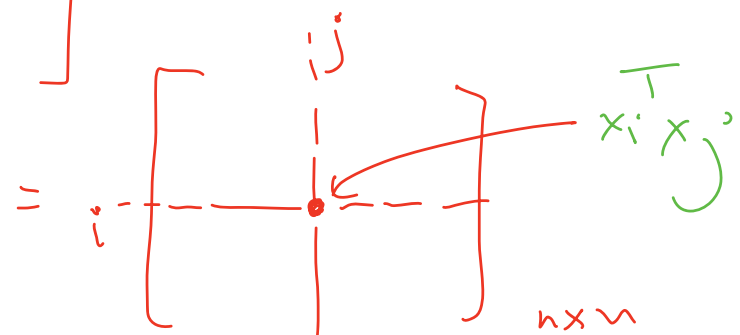
① $\alpha \in \mathbb{R}^n$

② $X = \begin{bmatrix} | & & | \\ x_1 & \dots & x_n \\ | & & | \end{bmatrix} \in \mathbb{R}^{d \times n}$

$$X^T X = \begin{bmatrix} - & x_1^T & - \\ - & x_2^T & - \\ \dots & \dots & \dots \\ - & x_i^T & - \\ \dots & \dots & \dots \\ - & x_n^T & - \end{bmatrix}$$

$$\begin{bmatrix} | & & | \\ x_1 & x_2 & \dots & x_n \\ | & & & | \end{bmatrix}$$

$$X^T X \in \mathbb{R}^{n \times n}$$



A new perspective of linear regression

$$\arg \min_{\alpha} \left\| X^T X \alpha - Y \right\|_2^2 + \lambda \|X \alpha\|_2^2$$

Solution:

$$\alpha = (X^T X + \lambda I)^{-1} Y \in \mathbb{R}^n$$

$$X^T X \in \mathbb{R}^{n \times n}, (X^T X)_{i,j} = \mathbf{x}_i^T \mathbf{x}_j = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$


A new perspective of linear regression

When we make prediction on a test example $\mathbf{x} \in \mathbb{R}^d$, we have:

$$\hat{w}^\top \mathbf{x} = \left(\sum_{i=1}^n \alpha_i \mathbf{x}_i \right)^\top \mathbf{x} = \sum_{i=1}^n \alpha_i \cdot \langle \mathbf{x}_i, \mathbf{x} \rangle$$

$\hat{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$

A new perspective of linear regression

When we make prediction on a test example $\mathbf{x} \in \mathbb{R}^d$, we have:

$$\hat{w}^\top \mathbf{x} = \left(\sum_{i=1}^n \alpha_i \mathbf{x}_i \right)^\top \mathbf{x} = \sum_{i=1}^n \alpha_i \cdot \langle \mathbf{x}_i, \mathbf{x} \rangle$$

$$\hat{w} = \left(\underbrace{\mathbf{x} \mathbf{x}^\top}_{d \times d} + \lambda \mathbf{I} \right)^{-1} \mathbf{x} \mathbf{y}$$

Notice a theme here:

Linear regression can be done by just using inner product of features
 $\langle \mathbf{x}, \mathbf{z} \rangle$, $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{z} \in \mathbb{R}^d$

Outline

1. A new perspective on ridge linear regression
2. Feature mapping and Kernel
3. Kernel trick and demo of kernel regression

Feature mapping

Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

Feature mapping

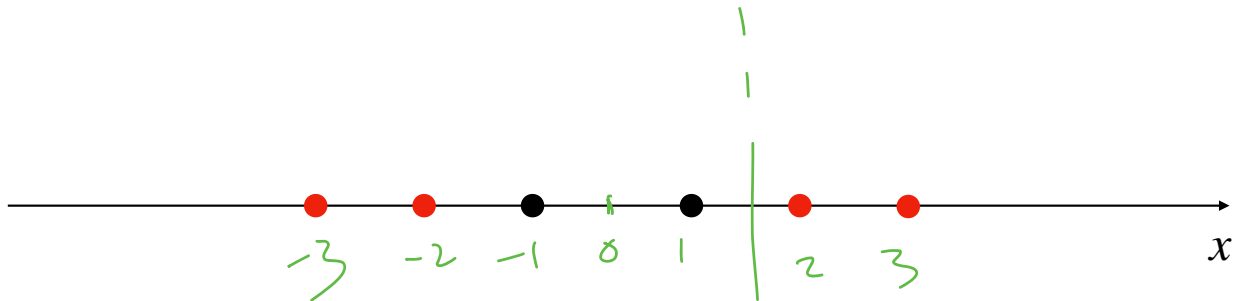
Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

$$\text{Ex 1: } \mathbf{x} \in \mathbb{R}, \phi(\mathbf{x}) = [x, x^2]^\top \in \mathbb{R}^2 = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

Feature mapping

Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

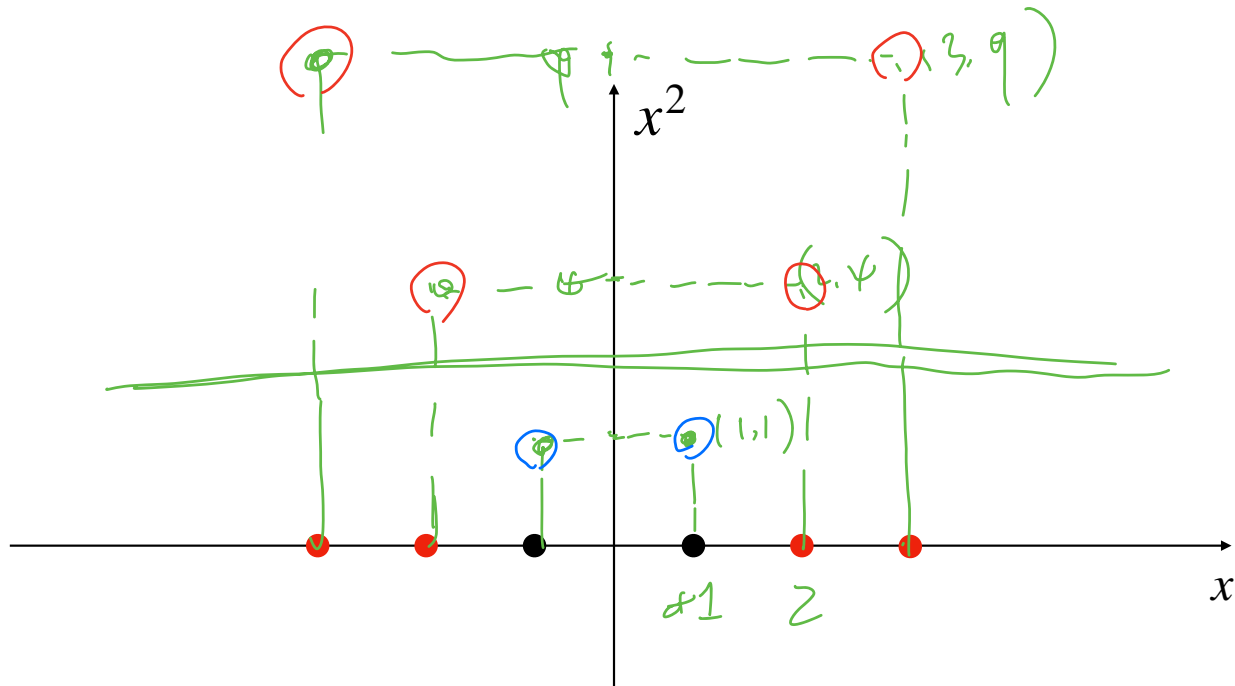
$$\text{Ex 1: } \mathbf{x} \in \mathbb{R}, \phi(\mathbf{x}) = [x, x^2]^\top \in \mathbb{R}^2$$



Feature mapping

Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

Ex 1: $\mathbf{x} \in \mathbb{R}$, $\phi(\mathbf{x}) = [x, x^2]^\top \in \mathbb{R}^2$



Feature mapping

Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

Ex 2: quadratic
feature mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top, \in \mathbb{R}^2$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2]^\top$$

Feature mapping

Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

Ex 2: cubic feature
mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top,$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2]^\top$$

Feature mapping

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{pmatrix}$$

$p=3$

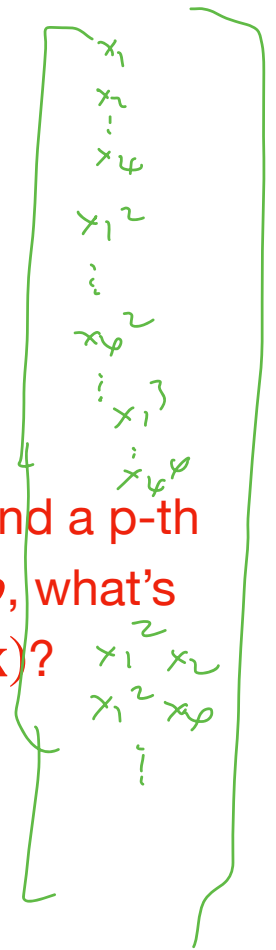
Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

Ex 2: cubic feature mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top,$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2]^\top$$

Q: in general, for $\mathbf{x} \in \mathbb{R}^d$, and a p -th order polynomial feature ϕ , what's the dimension of $\phi(\mathbf{x})$?



Feature mapping

Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

Ex 2: cubic feature mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top,$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2]^\top$$

Q: in general, for $\mathbf{x} \in \mathbb{R}^d$, and a p -th order polynomial feature ϕ , what's the dimension of $\phi(\mathbf{x})$?

at least $\binom{d}{p} \approx d^p$

$$\begin{aligned} d &= 100 \\ p &= 10 \end{aligned}$$

Feature mapping

Define $\phi(\mathbf{x}) \in \mathbb{R}^m$ as a feature mapping (often $m > d$)

Ex 2: cubic feature
mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top,$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2]^\top$$

Q: in general, for $\mathbf{x} \in \mathbb{R}^d$, and a p -th
order polynomial feature ϕ , what's
the dimension of $\phi(\mathbf{x})$?

at least $\binom{d}{p}$

Dim of $\phi(\mathbf{x})$ can be very large!

Fit linear functions in the high-dim feature space

The feature mapping $\phi(\mathbf{x}) \in \mathbb{R}^m$ allows us to perform linear regression in the ϕ space

Fit linear functions in the high-dim feature space

The feature mapping $\phi(\mathbf{x}) \in \mathbb{R}^m$ allows us to perform linear regression in the ϕ space

Ex: cubic feature mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top, \quad \phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2]^\top$$

$$\begin{aligned} \omega^\top \phi(\mathbf{x}) &= \omega_0 + \omega_1 x_1 + \dots + \omega x_1^2 + \omega' x_2^2 \\ &+ \dots + \omega'' x_1^3 \\ &+ \dots + \omega x_1^2 x_2 \end{aligned}$$

Fit linear functions in the high-dim feature space

The feature mapping $\phi(\mathbf{x}) \in \mathbb{R}^m$ allows us to perform linear regression in the ϕ space

Ex: cubic feature mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top, \quad \phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2]^\top$$

$w^\top \phi(\mathbf{x})$ now can represent a 3-order polynomials!

Fit linear functions in the high-dim feature space

The feature mapping $\phi(\mathbf{x}) \in \mathbb{R}^m$ allows us to perform linear regression in the ϕ space

Ex: cubic feature mapping ϕ

$$\mathbf{x} = [x_1, x_2]^\top, \quad \phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1x_2^2, x_1^2x_2]^\top$$

$w^\top \phi(\mathbf{x})$ now can represent a 3-order polynomials!

To fit a 3-order polynomial in \mathbf{x} , we can instead do linear regression in $\phi(\mathbf{x})$

$$\begin{pmatrix} x_1 & \dots & x_n \end{pmatrix} \rightarrow \begin{pmatrix} \phi(x_1) & \dots & \phi(x_n) \\ 1 & & 1 \end{pmatrix}$$

Fit linear functions in the high-dim feature space

Perform linear regression in ϕ space, i.e.,

$$\min_w \sum_{i=1}^n (w^T \phi(\mathbf{x}_i) - y_i)^2 + \lambda \|w\|_2^2$$

Fit linear functions in the high-dim feature space

Perform linear regression in ϕ space, i.e.,

$$\min_w \sum_{i=1}^n (w^\top \phi(\mathbf{x}_i) - y_i)^2 + \lambda \|w\|_2^2$$

Linear in ϕ , but high-order poly in \mathbf{x}

Fit linear functions in the high-dim feature space

Perform linear regression in ϕ space, i.e.,

$$\min_w \sum_{i=1}^n (w^\top \phi(\mathbf{x}_i) - y_i)^2 + \lambda \|w\|_2^2$$

Linear in ϕ , but high-order poly in \mathbf{x}

What is the potential problem of doing this?

$\phi(x)$
—————
 d^p

$d = \dim \text{ of } \mathbf{x}$
 $p = \text{order poly}$

Fit linear functions in the high-dim feature space

Perform linear regression in ϕ space, i.e.,

$$\min_w \sum_{i=1}^n (w^\top \phi(\mathbf{x}_i) - y_i)^2 + \lambda \|w\|_2^2$$

Linear in ϕ , but high-order poly in \mathbf{x}

What is the potential problem of doing this?

This is where the new perspective of linear regression and kernels come to rescue!

Kernel

Kernel $k(\mathbf{x}, \mathbf{z}) \in \mathcal{R}$

$\mathbf{x}^T \mathbf{z}$

Kernel

Kernel $k(\mathbf{x}, \mathbf{z})$

A valid kernel is a kernel such that $\exists \phi, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$

$$k(x, z) = x^\top z$$

$$= \langle x, z \rangle$$

ϕ : identity

Kernel

Kernel $k(\mathbf{x}, \mathbf{z})$

A valid kernel is a kernel such that $\exists \phi, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$

Ex: quadratic kernel

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

$$\phi(\mathbf{x}) \quad ??$$

$$\phi(\mathbf{x})^\top \phi(\mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

Kernel

Kernel $k(\mathbf{x}, \mathbf{z})$

A valid kernel is a kernel such that $\exists \phi, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$

Ex: quadratic kernel

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2]^\top$$

$$\phi(\mathbf{x})^\top \phi(\mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

Kernel

Kernel $k(\mathbf{x}, \mathbf{z})$

A valid kernel is a kernel such that $\exists \phi, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$

Ex: quadratic kernel

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

$x \in \mathbb{R}^d$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2]^\top$$

$(\mathbf{x}^\top \mathbf{z} + 1) \leftarrow \mathcal{O}(d)$

Q: what's the computation of $k(\mathbf{x}, \mathbf{z})$?

Kernel

Kernel $k(\mathbf{x}, \mathbf{z})$

A valid kernel is a kernel such that $\exists \phi, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$

Ex: quadratic kernel

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2]^\top$$

Q: what's the computation of $k(\mathbf{x}, \mathbf{z})$?

Q: what's the computation of $\phi(\mathbf{x})^\top \phi(\mathbf{z})$?

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= \phi(\mathbf{x})^\top \phi(\mathbf{z}) \\ &= \mathcal{O}(d^2) \\ &\phi(\mathbf{x}) \in \mathbb{R}^d \end{aligned}$$

Kernel

Kernel $k(\mathbf{x}, \mathbf{z})$

$\phi(x) \in \mathbb{R}^{d^3}$
 Computation $\left(\phi(x)^\top \phi(z) \right) \leftarrow \mathcal{O}(d^3) \quad ??$

A valid kernel is a kernel such that $\exists \phi, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$

Ex: quadratic kernel

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2]^\top$$

Q: what's the computation of $k(\mathbf{x}, \mathbf{z})$?

Q: what's the computation of $\phi(\mathbf{x})^\top \phi(\mathbf{z})$?

Ex: cubic feature mapping ϕ

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^3$$

Computation $\mathcal{O}(d)$

$$\phi(x) = \begin{bmatrix} 1 \\ \sqrt{3}x_1 \\ \sqrt{3}x_2 \\ x_1^3 \\ x_2^3 \\ \sqrt{3}x_1x_2 \\ \vdots \\ x_1^3 \\ \vdots \\ x_d^3 \end{bmatrix}$$

$$\phi(x)^\top \phi(z) = k(x, z)$$

Kernel

Kernel $k(\mathbf{x}, \mathbf{z})$

A valid kernel is a kernel such that $\exists \phi, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$

Ex: quadratic kernel

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^2$$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2]^\top$$

Q: what's the computation of $k(\mathbf{x}, \mathbf{z})$?

Q: what's the computation of $\phi(\mathbf{x})^\top \phi(\mathbf{z})$?

Ex: cubic feature mapping ϕ

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^3 \quad p=20$$

Generalizing to p -th order polynomials:

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^p \quad \mathcal{O}(d)$$

$$\phi(x) \in \mathbb{R}^{(d^p)}$$

$$\phi(x)^\top \phi(z) \approx \mathcal{O}(d^p)$$

$\phi(x)^T \phi(z)$ **Kernel**

Gaussian Kernel: $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$

part of
Gaussian

The mapping $\phi(\mathbf{x})$ is infinite-dimensional

$\phi(x) \in \mathbb{R}^{\infty}$

Kernel

Gaussian Kernel: $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$

The mapping $\phi(\mathbf{x})$ is infinite-dimensional

Ex: $\mathbf{x} \in \mathbb{R}$, the mapping $\phi(\mathbf{x})$: *$d=1$* *i -th polynomial*

$$\phi(\mathbf{x}) = \left[\dots, \frac{1}{\sqrt{i!}} \exp\left(-\frac{x^2}{2\sigma^2}\right) x^i, \dots \right]^T \in \mathbb{R}^\infty$$

i -th element

Kernel

Gaussian Kernel: $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$

2. Linear function $w^T \phi(\mathbf{x})$ can model any indefinitely differentiable function f

exp(x)

sin(x)

Kernel

Gaussian Kernel: $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$

2. Linear function $w^T \phi(\mathbf{x})$ can model any indefinitely differentiable function f

Why? ϕ contains all polynomials, and f can be written as an infinite Taylor series..

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$
$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

Summary so far

1. Feature mapping $\phi(\mathbf{x})$ lifts \mathbf{x} into high-dimensional space (e.g., high-order polynomials)

Summary so far

1. Feature mapping $\phi(\mathbf{x})$ lifts \mathbf{x} into high-dimensional space (e.g., high-order polynomials)

2. A kernel $k(\mathbf{x}, \mathbf{z})$ is a symmetric function, such that there exists a ϕ , so that

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

Summary so far

1. Feature mapping $\phi(\mathbf{x})$ lifts \mathbf{x} into high-dimensional space (e.g., high-order polynomials)

2. A kernel $k(\mathbf{x}, \mathbf{z})$ is a symmetric function, such that there exists a ϕ , so that

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

3. Kernel allows us to compute $\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$ without ever explicitly computing ϕ

($k(\mathbf{x}, \mathbf{z})$ is easy to compute but $\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$ is hard to compute)