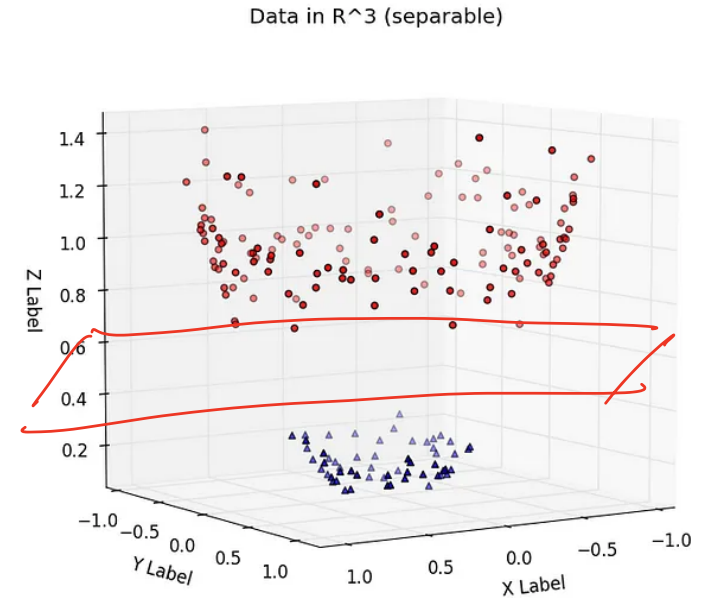
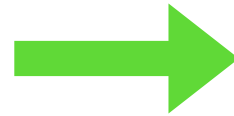
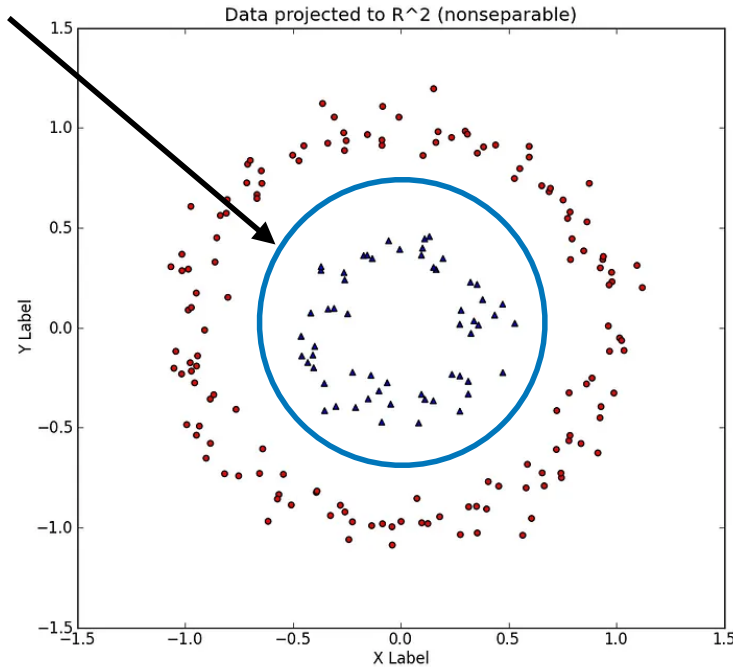


**Kernel**

# Objective today

Use kernels to design nonlinear regression & classification models

Goal: Non-linear decision boundary



# Outline

1. Kernel

2. Kernel trick and Kernel regression

3. Kernel SVM

# Common Kernels

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

# Common Kernels

Linear kernel:  $k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$

Polynomial kernel:  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^p$

$\rightarrow \phi(x) =$

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \\ x_1^2 \\ \vdots \\ x_d^2 \\ x_1 x_2 \\ \vdots \\ x_1^3 \\ \vdots \\ x_1^p \end{bmatrix}$$

$$\begin{aligned} & \phi(x)^\top \phi(z) \\ &= (\mathbf{x}^\top \mathbf{z} + 1)^p \end{aligned}$$

# Common Kernels

Linear kernel:  $k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$

Polynomial kernel:  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^p$

Gaussian kernel (aka RBF):  
 $k(\mathbf{x}, \mathbf{z}) = \exp\left(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2\right)$

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

$x = z \quad 1$

$\|\mathbf{x} - \mathbf{z}\|_2^2 \rightarrow +\infty$

$\exp(-+\infty) \rightarrow 0$

# Well-defined Kernels

Given any symmetric function  $k(\mathbf{x}, \mathbf{z})$ , can it be used as a kernel?

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

# Well-defined Kernels

Given any symmetric function  $k(\mathbf{x}, \mathbf{z})$ , can it be used as a kernel?

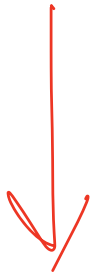
$$\exists \phi, \text{ s.t.}, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$$



# Well-defined Kernels

Given any symmetric function  $k(\mathbf{x}, \mathbf{z})$ , can it be used as a kernel?

$$\exists \phi, \text{ s.t.}, k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z}), \forall \mathbf{x}, \mathbf{z}$$



$$K_{ij} = k(x_i, x_j)$$

$\mathbb{R}^{m \times m}$

~~$\exists \phi, \text{ s.t.},$~~   $\forall \mathbf{x}_1, \dots, \mathbf{x}_m$ , the kernel matrix  $K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_m) \\ k(\mathbf{x}_2, \mathbf{x}_1) & \dots & k(\mathbf{x}_2, \mathbf{x}_m) \\ \dots & \dots & \dots \\ k(\mathbf{x}_m, \mathbf{x}_1) & \dots & k(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$  is PSD

# Construction of well-defined kernels

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$   $k_i = x^T z$

1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$

$$k_1(x, z) = \phi_1^T(x) \phi_1(z)$$

$$c \cdot k_1(x, z) = \underbrace{\left( \sqrt{c} \phi_1(x) \right)^T}_{\phi'} \left( \sqrt{c} \phi_1(z) \right)$$

# Construction of well-defined kernels

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$

1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$

2.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$

# Construction of well-defined kernels

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$

1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$
2.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$
3.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) \cdot k_2(\mathbf{x}, \mathbf{z})$

# Construction of well-defined kernels

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$

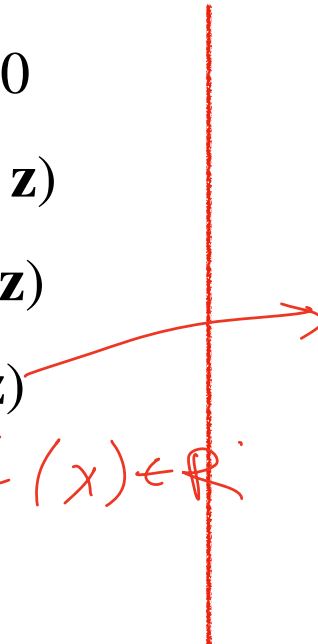
1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$

2.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$

3.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) \cdot k_2(\mathbf{x}, \mathbf{z})$

4.  $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{z})f(\mathbf{z})$

$f(x) \in \mathbb{R}$


$$\underbrace{\left( f(x)\phi_1(x) \right)^T}_{\phi'(x)} \left( f(z)\phi_1(z) \right)$$

# Construction of well-defined kernels

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$

1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$
2.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$
3.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) \cdot k_2(\mathbf{x}, \mathbf{z})$
4.  $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{z})f(\mathbf{z})$
5.  $k(\mathbf{x}, \mathbf{z}) = \exp(k_1(\mathbf{x}, \mathbf{z}))$

# Construction of well-defined kernels

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$

1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$  ✓

2.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$

3.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) \cdot k_2(\mathbf{x}, \mathbf{z})$

4.  $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{z})f(\mathbf{z})$  ✓

5.  $k(\mathbf{x}, \mathbf{z}) = \exp(k_1(\mathbf{x}, \mathbf{z}))$

... (see lecture note)

# Construction of well-defined kernels

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$

1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$
2.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$
3.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) \cdot k_2(\mathbf{x}, \mathbf{z})$
4.  $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{z})f(\mathbf{z})$
5.  $k(\mathbf{x}, \mathbf{z}) = \exp(k_1(\mathbf{x}, \mathbf{z}))$

... (see lecture note)

In class exercise:

Given  $k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$  being well defined,

Prove Gaussian kernel  
 $\exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$  is well defined



# Construction of well-defined kernels $\exp(-\mathbf{x}^T \mathbf{z})$

Kernels built by recursively applying the following one or more rules are well-defined kernels

Given well-defined  $k_1, k_2$

1.  $k(\mathbf{x}, \mathbf{z}) = ck_1(\mathbf{x}, \mathbf{z}), c > 0$
2.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$
3.  $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) \cdot k_2(\mathbf{x}, \mathbf{z})$
4.  $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{z})f(\mathbf{z})$
5.  $k(\mathbf{x}, \mathbf{z}) = \exp(k_1(\mathbf{x}, \mathbf{z}))$

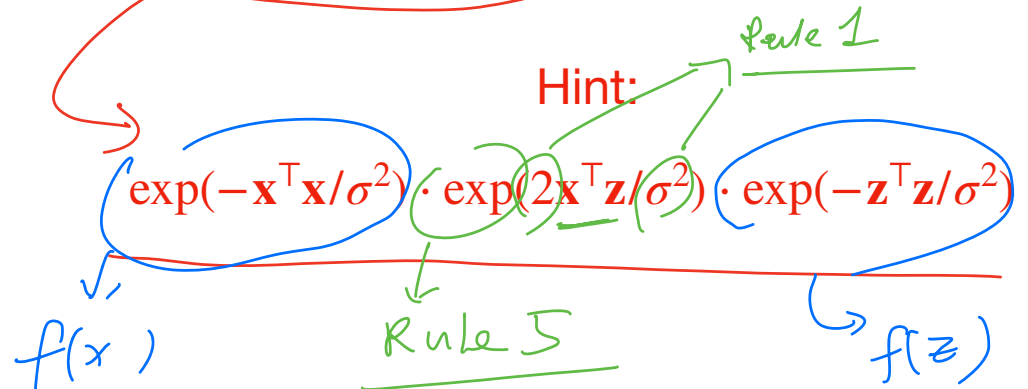
... (see lecture note)

In class exercise:

Given  $k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$  being well defined,

Prove Gaussian kernel

$\exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$  is well defined



# Outline

1. Kernel

2. Kernel trick and Kernel regression

3. Kernel SVM

# Kernel Trick

We wanted to do linear regression in the new features  $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)$ ,

# Kernel Trick

We wanted to do linear regression in the new features  $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)$ ,

**BUT**,  $\phi(\mathbf{x})$  can be very high-dim or even infinite-dim....



# Kernel Trick

We wanted to do linear regression in the new features  $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)$ ,

**BUT**,  $\phi(\mathbf{x})$  can be very high-dim or even infinite-dim....

Solution: recall linear regression can be done by  
just using inner product of two features!



# The kernel trick

A recipe:

# The kernel trick

A recipe:

1. Write the learning algorithm in terms of  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$

# The kernel trick

A recipe:

1. Write the learning algorithm in terms of  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$
2. Define a kernel  $k(\mathbf{x}, \mathbf{z})$  (e.g., Gaussian kernel, poly kernel)



# The kernel trick

A recipe:

1. Write the learning algorithm in terms of  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$
2. Define a kernel  $k(\mathbf{x}, \mathbf{z})$  (e.g., Gaussian kernel, poly kernel)
3. Replace all  $\langle \mathbf{x}, \mathbf{z} \rangle$  operation in the Alg by  $k(\mathbf{x}, \mathbf{z})$

$\exp\left(-\frac{\|\mathbf{x}-\mathbf{z}\|^2}{\sigma^2}\right)$

$\swarrow O(d)$

# Kernel ridge regression

1. Recall linear regression can be done via just using inner product:

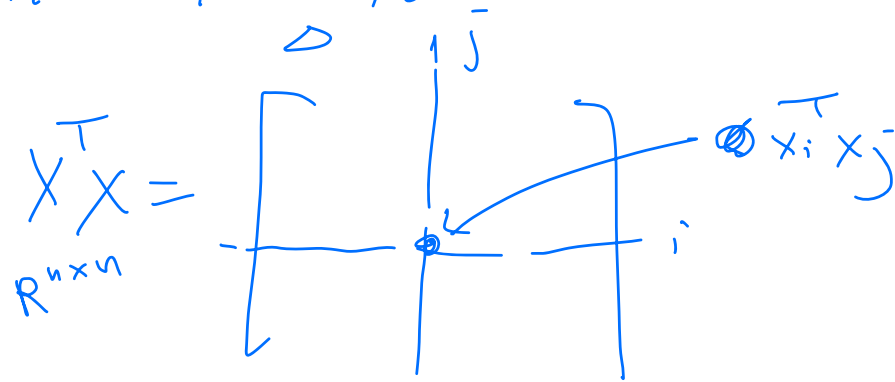
$$\alpha = (X^T X + \lambda I)^{-1} Y \in \mathbb{R}^n$$

$$W = \sum_{i=1}^n \alpha_i x_i = X \alpha$$

$$\|X^T W - Y\|_2^2 + \lambda \|W\|_2^2$$

$$\Leftrightarrow \|X^T X \alpha - Y\|_2^2 + \lambda \|X \alpha\|_2^2$$

$$X = \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \\ \vdots & & \vdots \end{bmatrix} \in \mathbb{R}^{d \times n}$$



# Kernel ridge regression

1. Recall linear regression can be done via just using inner product:

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

2. Define a kernel, e.g.,  $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$

# Kernel ridge regression

1. Recall linear regression can be done via just using inner product:

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

2. Define a kernel, e.g.,  $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$

3. Replace  $X^T X$  by a **kernel matrix K**

Handwritten diagram illustrating the matrix product  $X^T X$ . The matrix is shown with a dot at the intersection of row  $i$  and column  $j$ . An arrow points from this dot to the expression  $x_i^T x_j$ .

Handwritten diagram illustrating the kernel matrix  $K$ . The matrix is shown with a dot at the intersection of row  $i$  and column  $j$ . An arrow points from this dot to the expression  $K(x_i, x_j)$ .

# Kernel ridge regression

1. Recall linear regression can be done via just using inner product:

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

2. Define a kernel, e.g.,  $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|_2^2 / \sigma^2)$

3. Replace  $X^T X$  by a **kernel matrix K**

$$K \in \mathbb{R}^{n \times n}, K_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\alpha = (K + \lambda I)^{-1} Y$$

# Kernel ridge regression

In test time, recall linear regression makes prediction at  $\mathbf{x}$ :

$$\hat{y} = \sum_{i=1}^n \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle$$

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

$w^T x$

# Kernel ridge regression

In test time, recall linear regression makes prediction at  $\mathbf{x}$ :

$$\hat{y} = \sum_{i=1}^n \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle$$

Replace it w/  $k(\mathbf{x}_i, \mathbf{x})$ :

# Kernel ridge regression

In test time, recall linear regression makes prediction at  $\mathbf{x}$ :

$$\hat{y} = \sum_{i=1}^n \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle$$

Replace it w/  $k(\mathbf{x}_i, \mathbf{x})$ :

$$\hat{y} = \sum_{i=1}^n \alpha_i \cdot k(\mathbf{x}_i, \mathbf{x})$$



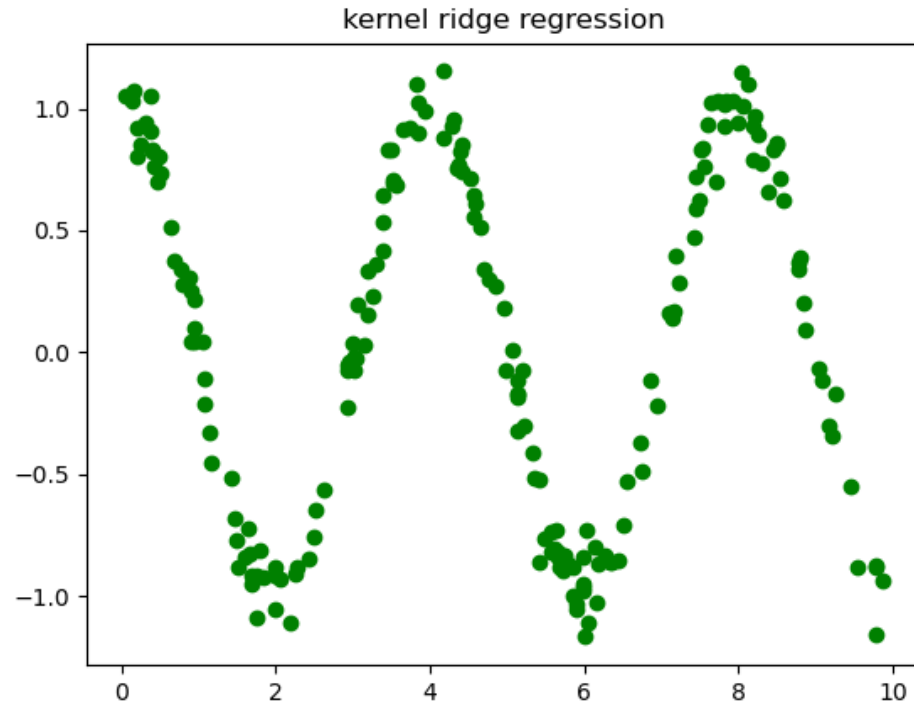
# Demo

Training data is generated as follows:  $x \sim \text{uniform}[0,10]$ ,

$$y = \sin(x\pi/2) + \epsilon, \epsilon \sim \mathcal{N}(0,0.1)$$

# Demo

Training data is generated as follows:  $x \sim \text{uniform}[0,10]$ ,  
 $y = \sin(x\pi/2) + \epsilon, \epsilon \sim \mathcal{N}(0,0.1)$



# Outline

1. Kernel

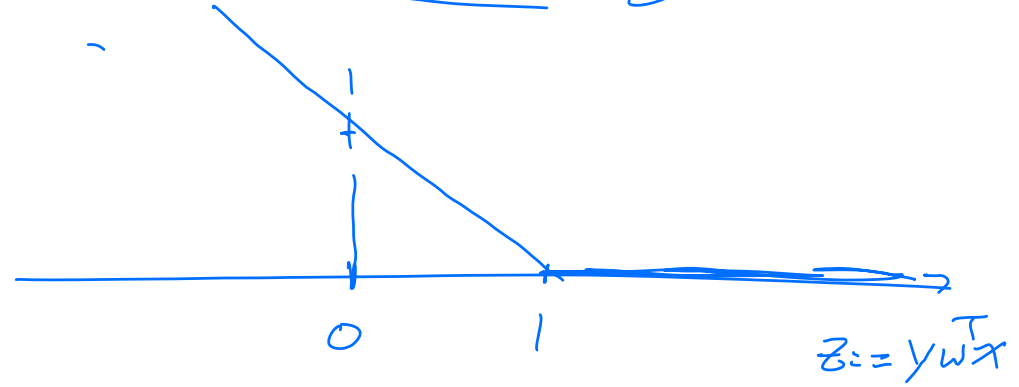
2. Kernel trick and Kernel regression

3. Kernel SVM

# Recall the soft-margin SVM formulation

$$\min_w \|w\|_2^2/2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^T \mathbf{x}_i)\}$$

hinge loss



## Recall the soft-margin SVM formulation

$$\min_w \|w\|_2^2/2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^\top \mathbf{x}_i)\}$$

Claim: the optimal solution  $\hat{w}$  is also in  $\text{span}(X)$

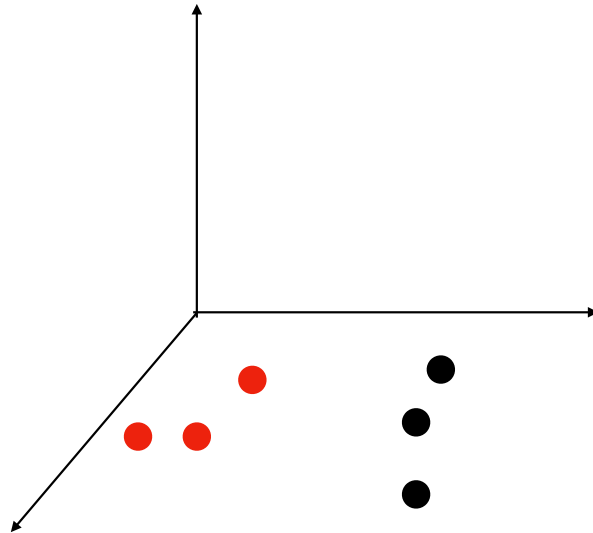
$$\hat{w} = \sum_{i=1}^n \alpha_i X_i$$

# Recall the soft-margin SVM formulation

$$\min_w \|w\|_2^2/2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^\top \mathbf{x}_i)\}$$

Claim: the optimal solution  $\hat{w}$  is also in  $\text{span}(X)$

Intuitive proof:

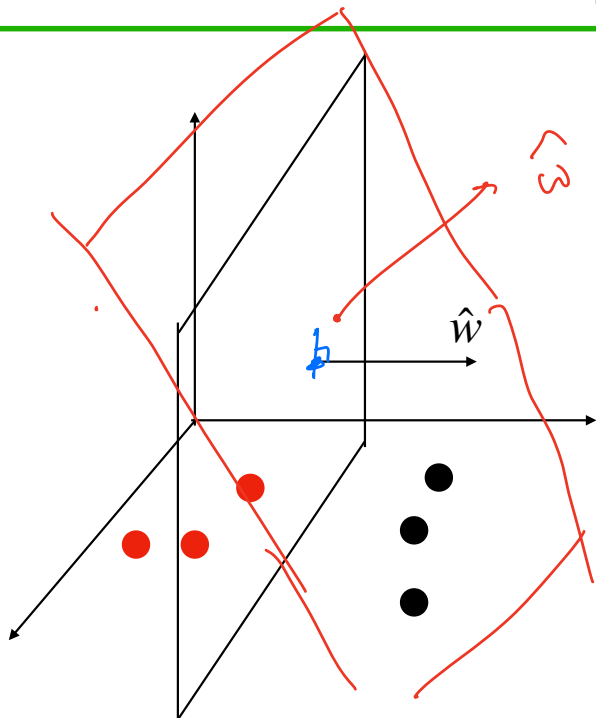


# Recall the soft-margin SVM formulation

$$\min_w \|w\|_2^2/2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^\top \mathbf{x}_i)\}$$

Claim: the optimal solution  $\hat{w}$  is also in  $\text{span}(X)$

Intuitive proof:



# A new formulation of soft-margin SVM formulation

Re-parameterize  $w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha$  ← learn  $\alpha \in \mathbb{R}^n$

---

Replace  $w$  by  $X\alpha$  in soft-margin SVM objective



# A new formulation of soft-margin SVM formulation

$$\text{Re-parameterize } w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha$$

$$\min_{\alpha} \|X\alpha\|_2^2/2 + C \sum_{i=1}^n \max \{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

*(Handwritten red annotations: a circle around  $\alpha$ , a circle around  $X\alpha$  with a  $w$  below it, and a circle around  $X\alpha$  in the max term with a  $w$  below it.)*

# A new formulation of soft-margin SVM formulation

Re-parameterize  $w = \sum_{i=1}^n \alpha_i x_i = X\alpha$

$$\min_{\alpha} \|X\alpha\|_2^2/2 + C \sum_{i=1}^n \max\{0, 1 - y_i(\mathbf{x}_i^T X\alpha)\}$$

$l(\alpha)$

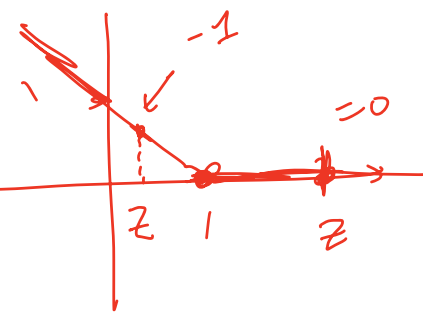
$\nabla_{\alpha} l(\alpha)$

Alg: gradient descent to optimize  $\alpha \in \mathbb{R}^n$

$x^T x \alpha$

$X^T X = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}$   
 (with a dot at the intersection of row  $i$  and column  $j$ )  
 $x_i^T x_j$

$C \cdot \sum_{j=1}^n -1 \left( \underbrace{y_i x_i^T X \alpha}_{=z_i} \leq 1 \right) (y_i x_i^T)$



# A new formulation of soft-margin SVM formulation

$$\text{Re-parameterize } w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha$$

$$\min_{\alpha} \|X\alpha\|_2^2/2 + C \sum_{i=1}^n \max \{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

Alg: gradient descent to optimize  $\alpha \in \mathbb{R}^n$

$$\nabla_{\alpha} \ell(\alpha) = \cancel{X}^\top X\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^\top X\alpha) \leq 1\} (-y_i X^\top x_i)$$

# A new formulation of soft-margin SVM formulation

$$\text{Re-parameterize } w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha$$

$$\min_{\alpha} \|X\alpha\|_2^2/2 + C \sum_{i=1}^n \max \{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

Alg: gradient descent to optimize  $\alpha \in \mathbb{R}^n$

$$\nabla_{\alpha} \ell(\alpha) = 2X^\top X\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^\top X\alpha) \leq 1\} (-y_i X^\top x_i)$$

$$\alpha' = \alpha - \eta \nabla_{\alpha} \ell(\alpha)$$

# A new formulation of soft-margin SVM formulation

Re-parameterize  $w = \sum_{i=1}^n \alpha_i \mathbf{x}_i = X\alpha$

$$\min_{\alpha} \|X\alpha\|_2^2/2 + C \sum_{i=1}^n \max\{0, 1 - y_i(\mathbf{x}_i^T X\alpha)\} = \begin{bmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{bmatrix} x_i$$

Alg: gradient descent to optimize  $\alpha \in \mathbb{R}^n$

$$\nabla_{\alpha} \ell(\alpha) = 2X^T X\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^T X\alpha) \leq 1\} (-y_i X^T x_i) = \begin{bmatrix} x_1^T x_1 \\ \vdots \\ x_n^T x_i \end{bmatrix}$$

$$\alpha' = \alpha - \eta \nabla_{\alpha} \ell(\alpha)$$

Q: Can we apply kernel trick??

# Kernelized GD for SVM

$$\min_{\alpha} \|X\alpha\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

While not converged:

$$g = X^\top X\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^\top X\alpha) \leq 1\} (-y_i X^\top x_i)$$

$$\alpha' = \alpha - \eta g$$

# Kernelized GD for SVM

$$\min_{\alpha} \|X\alpha\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

Pick a well-defined kernel  $k$ ;

While not converged:

$$g = X^\top X\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^\top X\alpha) \leq 1\} (-y_i X^\top x_i)$$

$$\alpha' = \alpha - \eta g$$

# Kernelized GD for SVM

$$\min_{\alpha} \|X\alpha\|_2^2 + C \sum_{i=1}^n \max\{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

Pick a well-defined kernel  $k$ ;

Replace  $X^\top X$  by kernel matrix  $K$

While not converged:

$K \in \mathbb{R}^{n \times n}$

$$g = X^\top X \alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^\top X\alpha) \leq 1\} (-y_i X^\top x_i)$$

$$\alpha' = \alpha - \eta g$$



# Kernelized GD for SVM

$$\min_{\alpha} \|X\alpha\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

Pick a well-defined kernel  $k$ ;

Replace  $X^\top X$  by kernel matrix  $K$

While not converged:

$$g = X^\top X\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^\top X\alpha) \leq 1\} (-y_i X^\top x_i)$$

$$\alpha' = \alpha - \eta g$$

Replace  $X^\top \mathbf{x}_i$  by  $\mathbf{k}_i = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_i) \\ k(\mathbf{x}_2, \mathbf{x}_i) \\ \dots \\ k(\mathbf{x}_n, \mathbf{x}_i) \end{bmatrix} \in \mathbb{R}^n$

$\rightarrow k(x_j, x_i)$

# Kernelized GD for SVM

$$\min_{\alpha} \|X\alpha\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(\mathbf{x}_i^\top X\alpha)\}$$

Pick a well-defined kernel  $k$ ;

Replace  $X^\top X$  by kernel matrix  $K$

Replace  $X^\top \mathbf{x}_i$  by  $\mathbf{k}_i = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_i) \\ k(\mathbf{x}_2, \mathbf{x}_i) \\ \dots \\ k(\mathbf{x}_n, \mathbf{x}_i) \end{bmatrix}$

While not converged:

$$g = X^\top X\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(x_i^\top X\alpha) \leq 1\} (-y_i X^\top x_i)$$

$$\alpha' = \alpha - \eta g$$

Replace

$$g = K\alpha + C \sum_{i=1}^n \mathbf{1}\{y_i(\mathbf{k}_i^\top \alpha) \leq 1\} (-y_i \mathbf{k}_i)$$

# Summary for kernel SVM so far

1. Ideally, want to do the SVM in the lifted high-dim feature space, i.e.,

$$\min_{\alpha} \|w\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^\top \phi(x_i))\}$$

# Summary for kernel SVM so far

1. Ideally, want to do the SVM in the lifted high-dim feature space, i.e.,

$$\min_{\alpha} \|w\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^\top \phi(x_i))\}$$

But  $\phi$  can be high-dim (e.g., infinite-dim in Gaussian kernel case)..

# Summary for kernel SVM so far

1. Ideally, want to do the SVM in the lifted high-dim feature space, i.e.,

$$\min_{\alpha} \|w\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^\top \phi(x_i))\}$$

But  $\phi$  can be high-dim (e.g., infinite-dim in Gaussian kernel case)..

2. Via the re-parameterization step, we see GD can be implemented **via just using**  $\langle \mathbf{x}, \mathbf{z} \rangle$

# Summary for kernel SVM so far

1. Ideally, want to do the SVM in the lifted high-dim feature space, i.e.,

$$\min_{\alpha} \|w\|_2^2 + C \sum_{i=1}^n \max \{0, 1 - y_i(w^\top \phi(x_i))\}$$

But  $\phi$  can be high-dim (e.g., infinite-dim in Gaussian kernel case)..

2. Via the re-parameterization step, we see GD can be implemented **via just using**  $\langle \mathbf{x}, \mathbf{z} \rangle$

3. We apply kernel trick, i.e., replace all  $\langle \mathbf{x}, \mathbf{z} \rangle$  by  $k(\mathbf{x}, \mathbf{z})$

## Take-home message today

Kernel trick allows us to do regression / classification in  $\phi(\mathbf{x})$  space  
(possibly infinite dim) **without ever explicitly computing  $\phi(\mathbf{x})$ !**