

# CS540 Introduction to Artificial Intelligence

## Lecture 5

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

June 9, 2020

# Survey Question

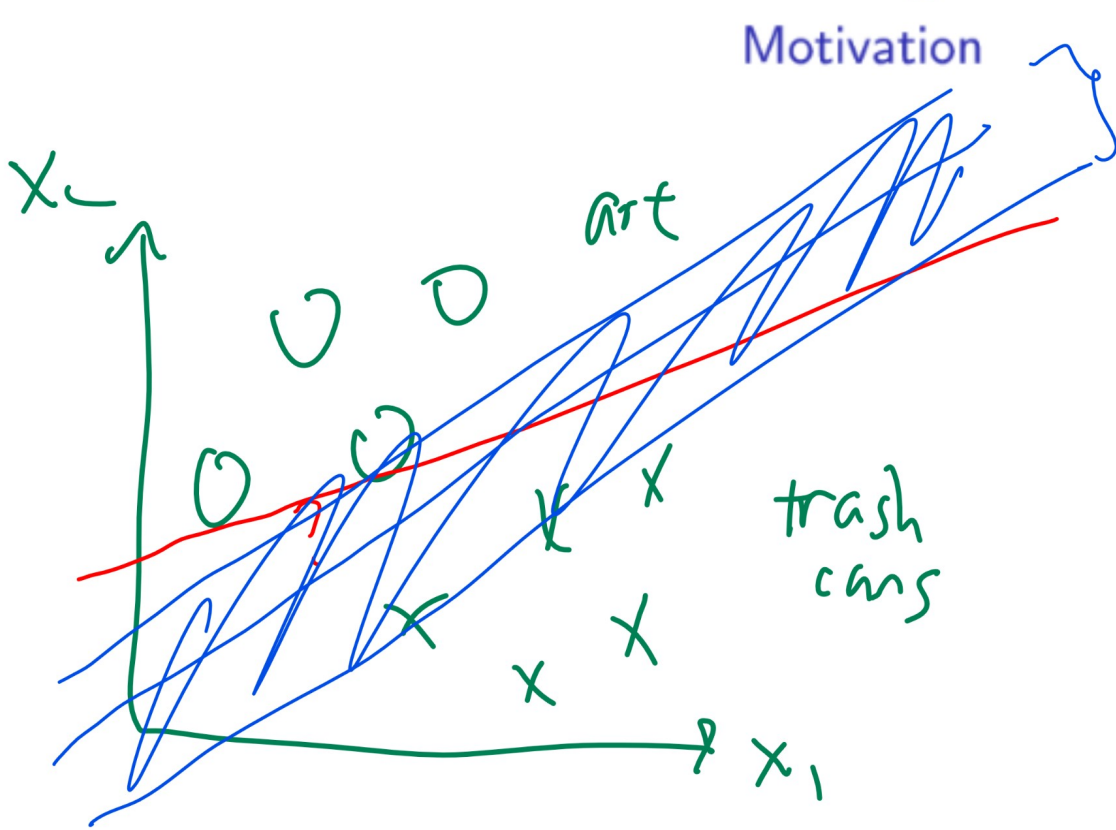
Admin

Socratic. App Room: CS540C

- Which prerecorded lecture videos have you watched?
- A: Yes
- B: Lectures 1, 2, 3, 4, 5, 6
- C: Lectures 1, 2, 3, 4
- D: Lectures 1, 2
- E: No

# Maximum Margin Diagram

Motivation



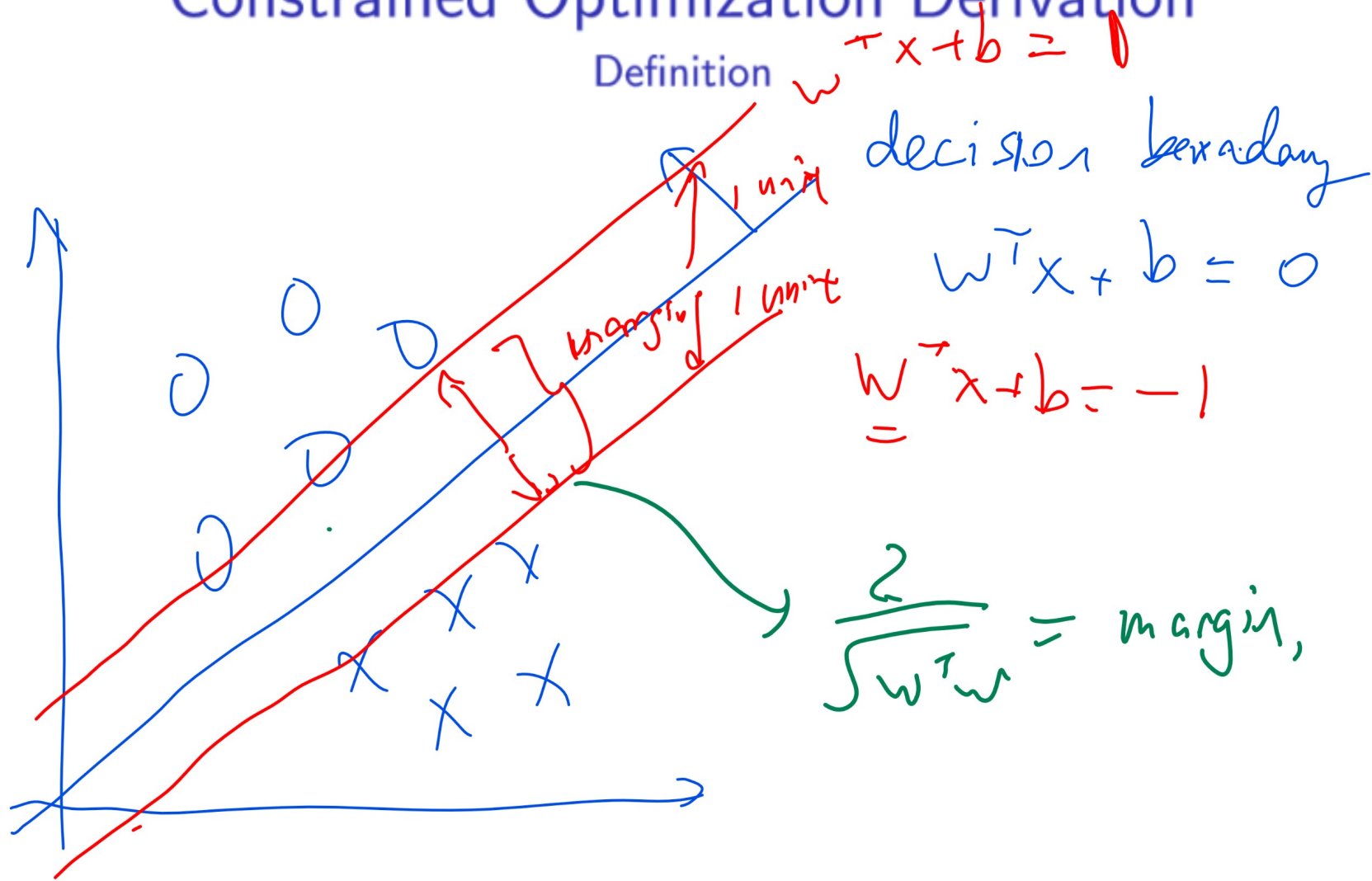
margin

Perceptron

min mistake (loss)  
(Cost)

max margin

# Constrained Optimization Derivation



# Constrained Optimization

## Definition

- The goal is to maximize the margin subject to the constraint that the plus plane and the minus plane separates the instances with  $y_i = 0$  and  $y_i = 1$ .

$$\max_w \frac{2}{\sqrt{w^T w}} \text{ such that } \begin{cases} (w^T x_i + b) \leq -1 & \text{if } y_i = 0 \\ (w^T x_i + b) \geq 1 & \text{if } y_i = 1 \end{cases}, i = 1, 2, \dots, n$$

- The two constraints can be combined.

$$\max_w \frac{2}{\sqrt{w^T w}} \text{ such that } (2y_i - 1)(w^T x_i + b) \geq 1, i = 1, 2, \dots, n$$

# Hard Margin SVM

## Definition

$$\max_w \frac{2}{\sqrt{w^T w}} \text{ such that } (2y_i - 1)(w^T x_i + b) \geq 1, i = 1, 2, \dots, n$$

- This is equivalent to the following minimization problem, called hard margin SVM.

$$\min_w \frac{1}{2} w^T w \text{ such that } (2y_i - 1)(w^T x_i + b) \geq 1, i = 1, 2, \dots, n$$

instance  $i$   
in features.

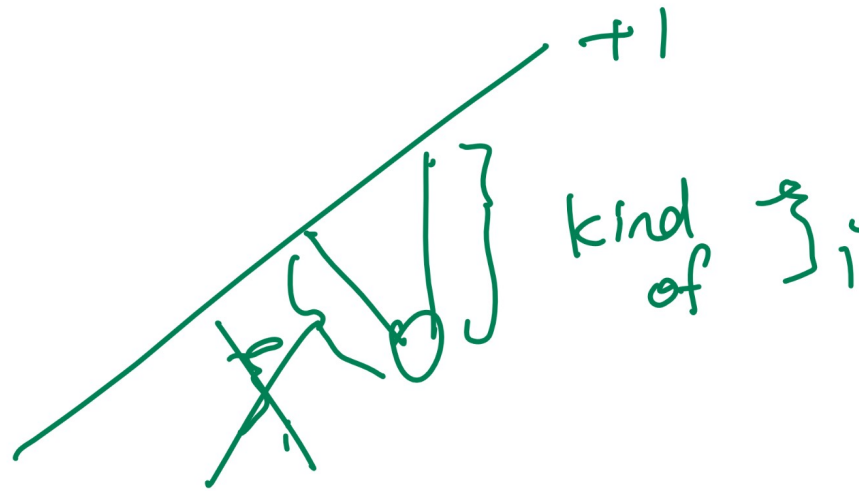
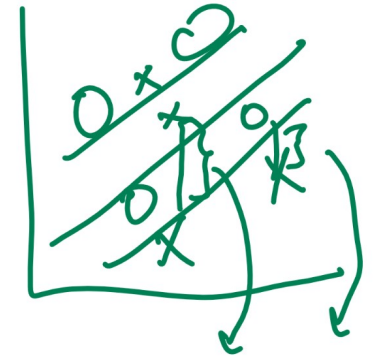
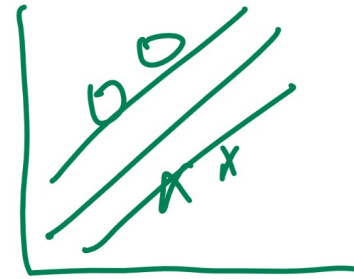
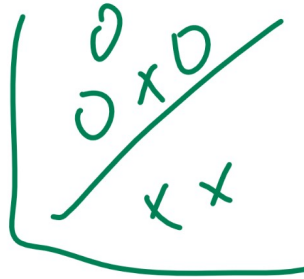
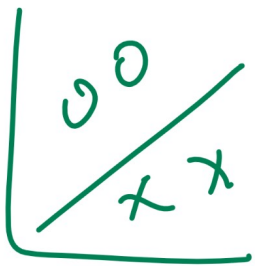
# Soft Margin Diagram

## Definition

LTU  $\rightarrow$  Logistic

Hard-margin

Soft



$x_i$   $\xi_i > 0$   $\xi_i = 0$

# Soft Margin SVM

## Definition

$$\min_w \frac{1}{2} w^T w + \frac{1}{\lambda n} \sum_{i=1}^n \xi_i$$

such that  $(2y_i - 1)(w^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0, i = 1, 2, \dots, n$

*Handwritten annotations:*  
 -  $\frac{1}{2} w^T w$  is circled in green and labeled "margin".  
 -  $\frac{1}{\lambda n} \sum_{i=1}^n \xi_i$  is circled in green and labeled "average mistake".  
 - The constraint term  $1 - \xi_i$  is circled in green and labeled "relative importance".

- This is equivalent to the following minimization problem, called soft margin SVM.

$$\min_w \frac{\lambda}{2} w^T w + \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - (2y_i - 1)(w^T x_i + b) \right\}$$



# SVM Weights

## Quiz

- Fall 2005 Final Q15 and Fall 2006 Final Q15
- Find the weights  $w_1, w_2$  for the SVM classifier

$\mathbb{1}_{\{w_1 x_{i1} + w_2 x_{i2} + 1 \geq 0\}}$  given the training data  $x_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  and

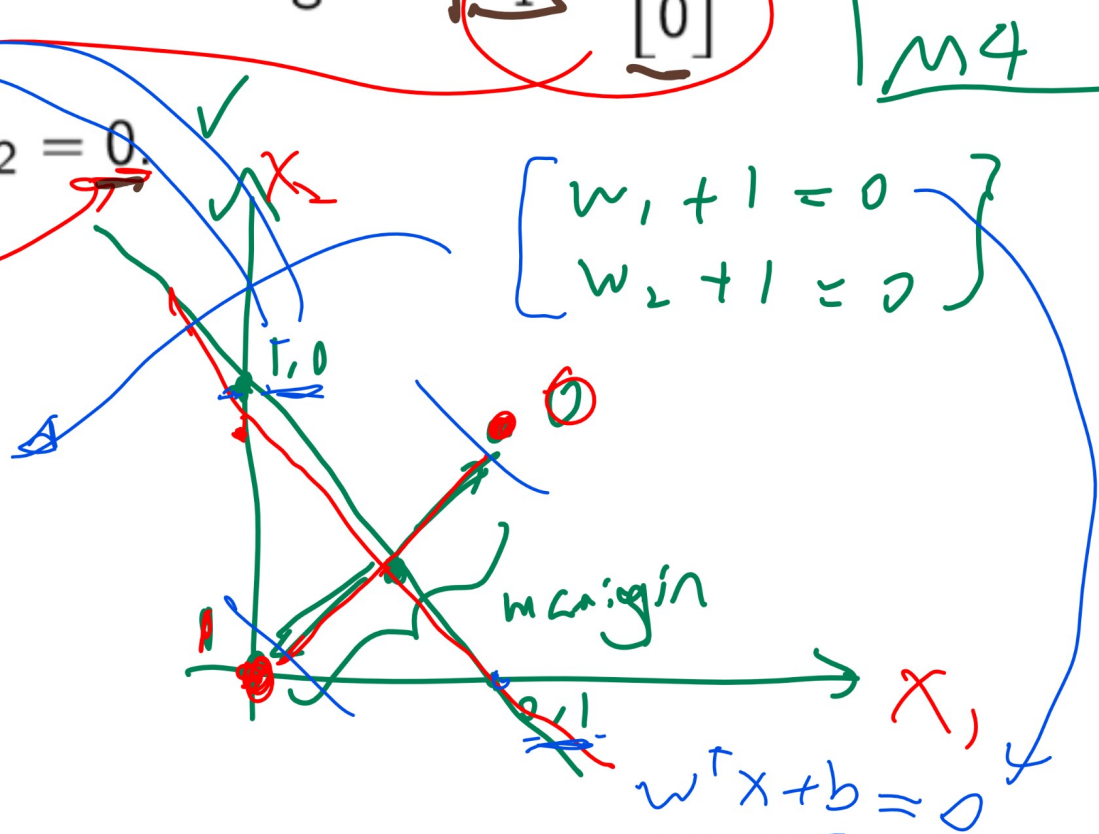
$x_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  with  $y_1 = 1, y_2 = 0$ .

- A:  $w_1 = 0, w_2 = -2$
- B:  $w_1 = -2, w_2 = 0$
- C:  $w_1 = -1, w_2 = -1$
- D:  $w_1 = -2, w_2 = -2$
- E: none of the above

hand-margin.

on midterm M4

$$\begin{cases} w_1 + 1 = 0 \\ w_2 + 1 = 0 \end{cases}$$



# SVM Weights Diagram

## Quiz

# SVM Weights 2

## Quiz

Q2

- Find the weights  $w_1, w_2$  for the SVM classifier *trainingly*  
 $\mathbb{1}_{\{w_1 x_{i1} + w_2 x_{i2} + 2 \geq 0\}}$  given the training data  
 $y = \neg(x_1 \vee x_2), x_1, x_2, y \in \{0, 1\}$ .

*min*  
~~*w1w*~~

- A:  $w_1 = -3, w_2 = -3$
- B:  $w_1 = -4, w_2 = -3$
- C:  $w_1 = -3, w_2 = -4$
- D:  $w_1 = -4, w_2 = -4$
- E:  $w_1 = -8, w_2 = -8$

*w1w*  
 constraints.

$x_1$	$x_2$	$y$
0	0	1
0	1	0
1	0	0
1	1	0

# SVM Weights 2 Diagram

Quiz

$$-3x_1 - 3x_2 + 2 = 0$$

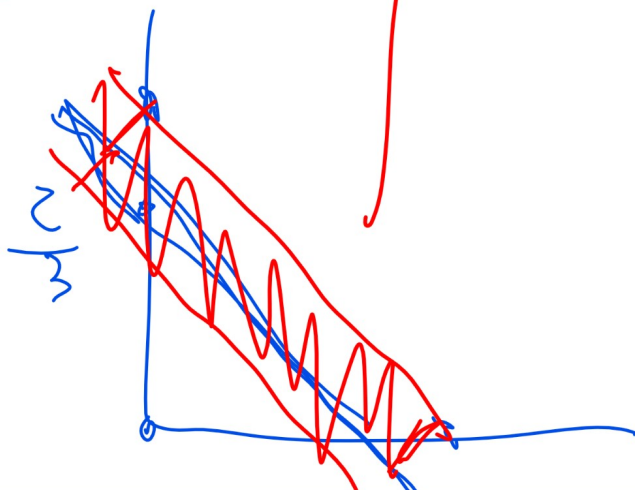
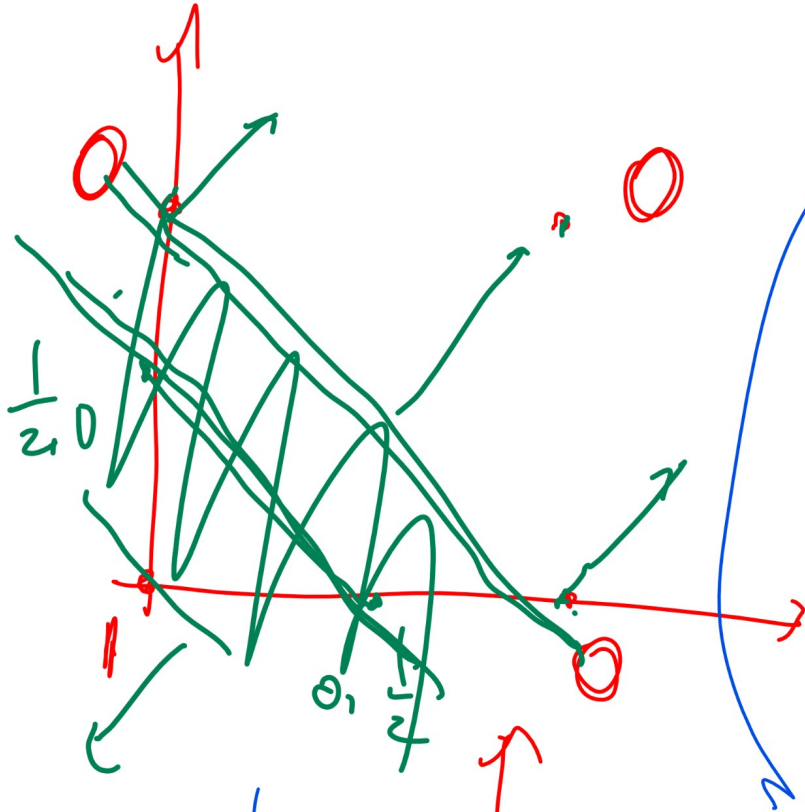


$$w_1 \frac{1}{2} + w_2 0 + 2 = 0$$

$$w_1 0 + w_2 \frac{1}{2} + 2 = 0$$

$$w_1 = -4, w_2 = -4$$

$$x_2 = \frac{2}{3} - x_1$$



# Soft Margin

## Quiz

- Fall 2011 Midterm Q8 and Fall 2009 Final Q1
- Let  $w = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$  and  $b = 3$ . For the point  $x = \begin{bmatrix} 4 \\ 5 \end{bmatrix}$ ,  $y = 0$ , what is the smallest slack variable  $\xi$  for it to satisfy the margin constraint?

$$(2y_i - 1)(w^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0$$

$$\min \xi_i = 18$$

$$\rightarrow (1, 2) \begin{pmatrix} 4 \\ 5 \end{pmatrix} + 3 \geq 1 - \xi_i, \xi_i \geq 0$$

$$\xi_i \geq 18, \xi_i \geq 0$$

# Soft Margin 2

## Quiz

Q.3

• Let  $w = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$  and  $b = 3$ . For the point  $x = \begin{bmatrix} -1 \\ -2 \end{bmatrix}$ ,  $y = 0$ , what is the smallest slack variable  $\xi$  for it to satisfy the margin constraint?

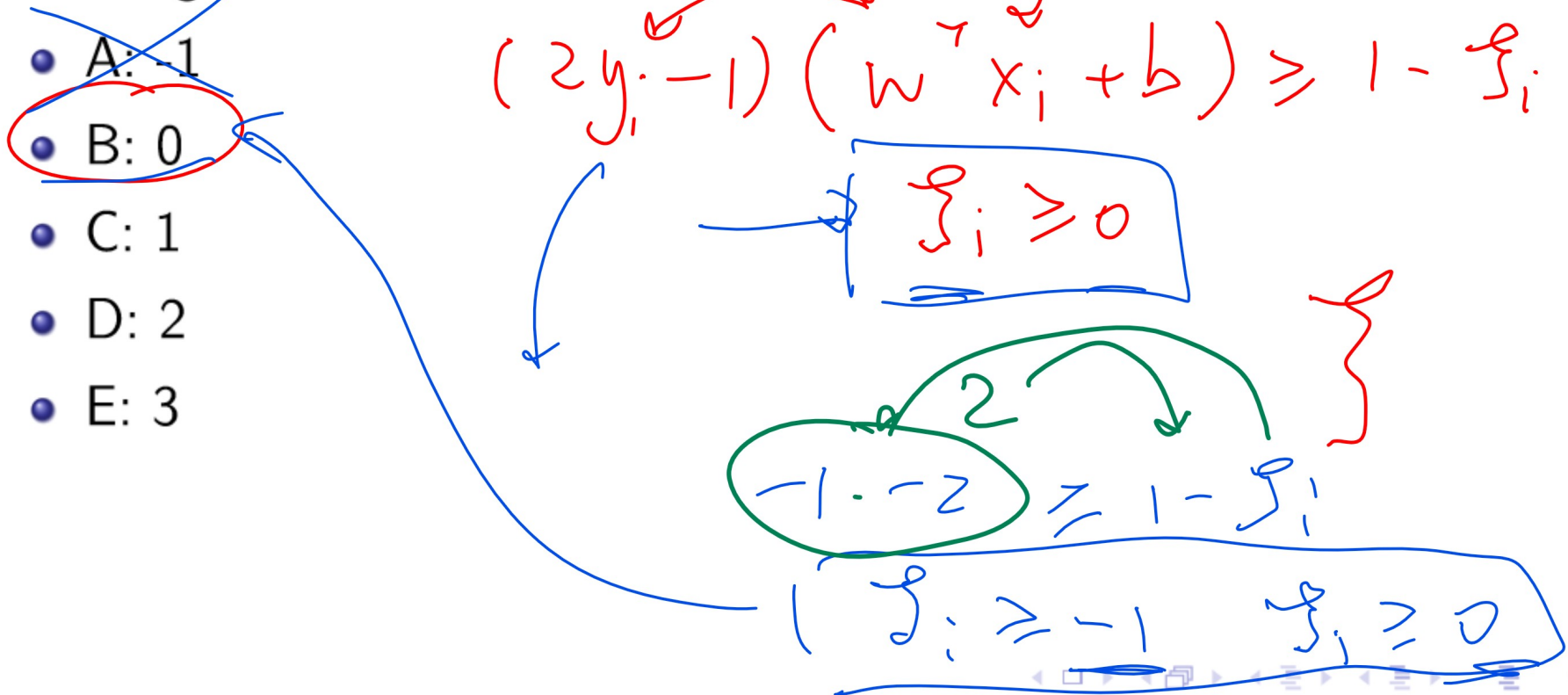
- A: -1
- B: 0
- C: 1
- D: 2
- E: 3

$$(2y_i - 1)(w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

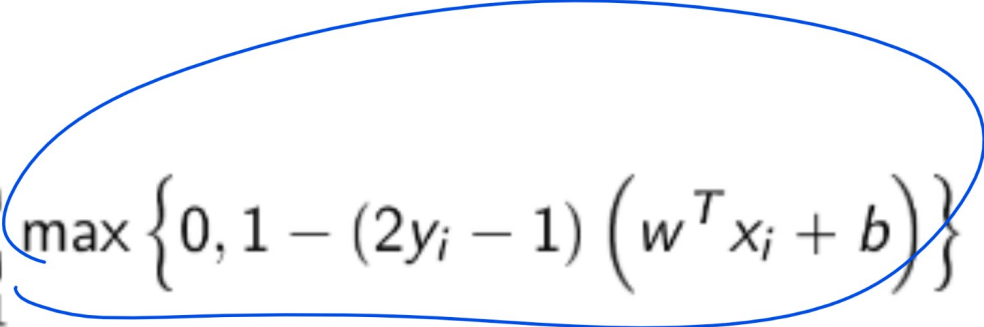
$$-1 \cdot -2 \geq 1 - \xi_i$$

$$\xi_i \geq -1 \quad \xi_i \geq 0$$



# Subgradient Descent

## Definition

$$\min_w \frac{\lambda}{2} w^T w + \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - (2y_i - 1) (w^T x_i + b) \right\}$$


- The gradient for the above expression is not defined at points with  $1 - (2y_i - 1) (w^T x_i + b) = 0$ .
- Subgradient can be used instead of gradient.

# Subgradient

- The subderivative at a point of a convex function in one dimension is the set of slopes of the lines that are tangent to the function at that point.
- The subgradient is the version for higher dimensions.
- The subgradient  $\partial f(x)$  is formally defined as the following set.

$$\partial f(x) = \left\{ v : f(x') \geq f(x) + v^T (x' - x) \quad \forall x' \right\}$$

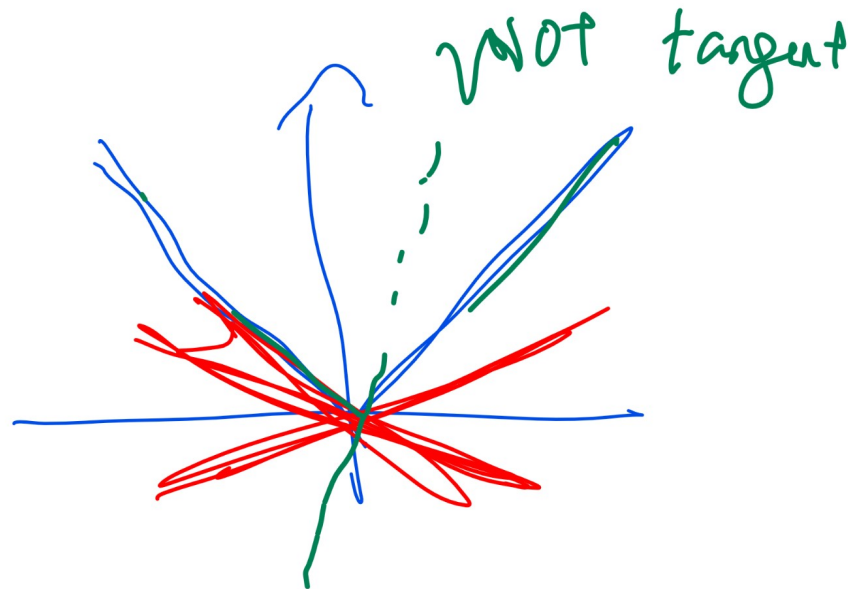


# Subgradient 1

## Quiz

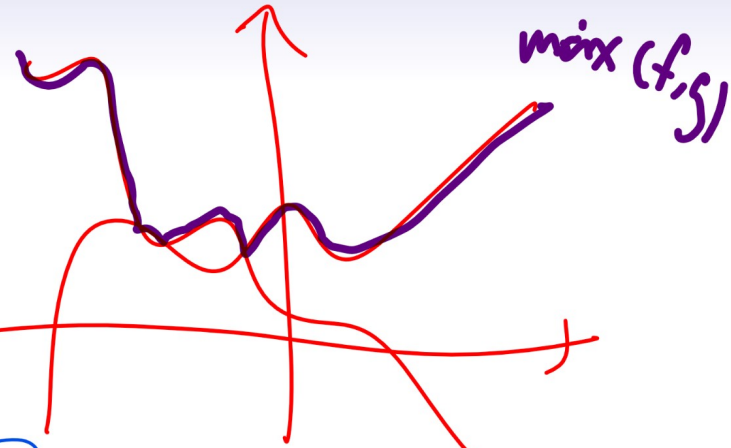
• Which ones (multiple) are subderivatives of  $|x|$  at  $x = 0$ ?

- A: -1
- B: -0.5
- C: 0
- D: 0.5
- E: 1



# Subgradient 2

## Quiz

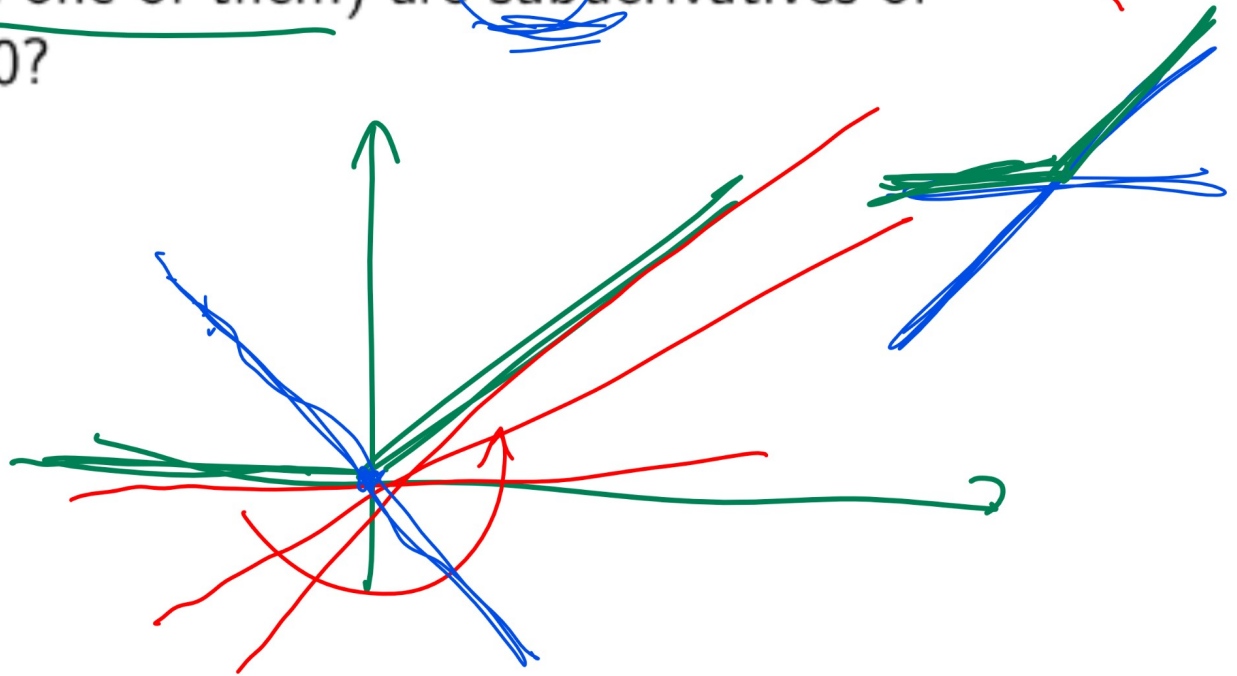


Q4 (last)

Which ones (select one of them) are subderivatives of

$\max\{x, 0\}$  at  $x = 0$ ?

- A: -1
- B: -0.5
- C: 0
- D: 0.5
- E: 1



# Subgradient Descent Step

## Definition

- One possible set of subgradients with respect to  $w$  and  $b$  are the following.

$$\partial_w C \ni \lambda w - \sum_{i=1}^n (2y_i - 1) x_i \mathbb{1}_{\{(2y_i - 1)(w^T x_i + b) \geq 1\}}$$

$$\partial_b C \ni - \sum_{i=1}^n (2y_i - 1) \mathbb{1}_{\{(2y_i - 1)(w^T x_i + b) \geq 1\}}$$

- The gradient descent step is the same as usual, using one of the subgradients in place of the gradient.

# PEGASOS Algorithm

## Algorithm

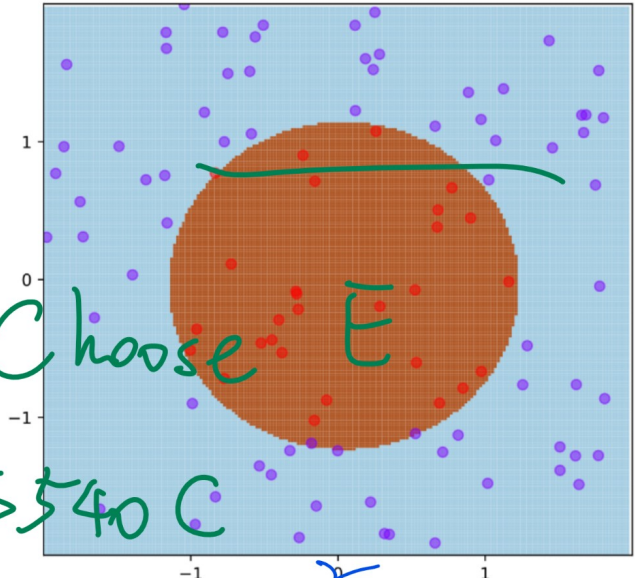
- Inputs: instances:  $\{x_i\}_{i=1}^n$  and  $\{z_i = 2y_i - 1\}_{i=1}^n$
- Outputs: weights:  $\{w_j\}_{j=1}^m$
- Initialize the weights.

$$w_j \sim \text{Unif} [0, 1]$$

- Randomly permute (shuffle) the training set and perform subgradient descent for each instance  $i$ .

$$w = (1 - \lambda) w - \alpha z_i \mathbb{1}_{\{z_i w^T x_i \geq 1\}} x_i$$

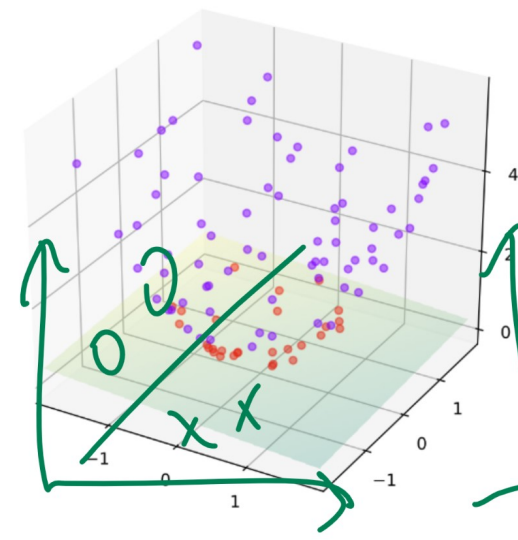
- Repeat for a fixed number of iterations.



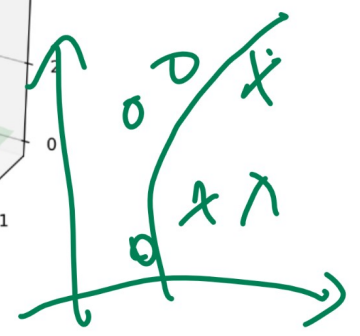
Q1 Choose E

Room: CS540C

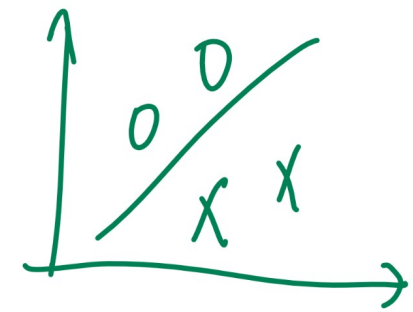
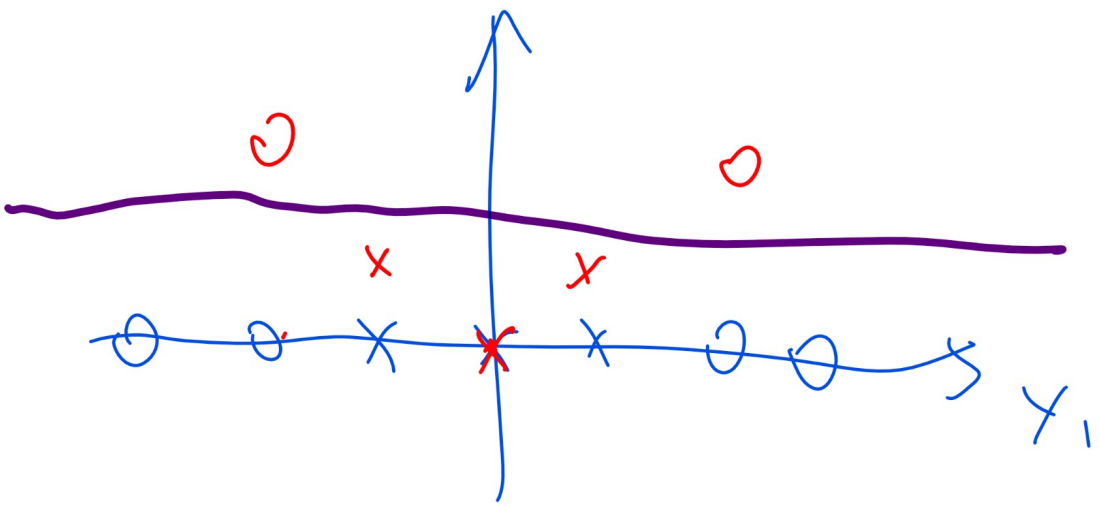
0  $x_2 = x_1$  0



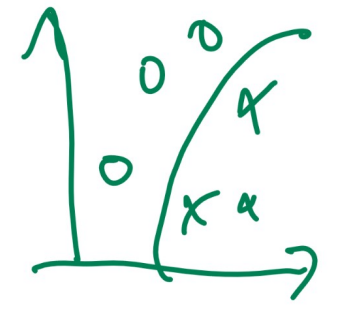
Perceptron



NN



SVM



kernelized SVM

# Kernelized SVM

## Definition

- With a feature map  $\varphi$ , the SVM can be trained on new data points  $\{(\varphi(x_1), y_1), (\varphi(x_2), y_2), \dots, (\varphi(x_n), y_n)\}$ .
- The weights  $w$  correspond to the new features  $\varphi(x_i)$ .
- Therefore, test instances are transformed to have the same new features.

$$\hat{y}_i = \mathbb{1}_{\{w^T \varphi(x_i) \geq 0\}}$$

# Kernel Matrix

## Definition

- The feature map is usually represented by a  $n \times n$  matrix  $K$  called the Gram matrix (or kernel matrix).

$$\underline{K_{ii'}} = \underline{\varphi(x_i)}^T \underline{\varphi(x_{i'})}$$

# Examples of Kernel Matrix

## Definition

- For example, if  $\varphi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$ , then the kernel matrix can be simplified.

$$K_{ii'} = (x_i^T x_{i'})^2$$

- Another example is the quadratic kernel  $K_{ii'} = (x_i^T x_{i'} + 1)^2$ . It can be factored to have the following feature representations.

$\varphi(x_i)^T \varphi(x_{i'})$

$\varphi(x) = (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$  ✓



# Examples of Kernel Matrix Derivation

$x_{11}, x_{21}$  2 different images  
Definition

$$k_{ij} = (x_i^T x_j + 1)^2 = \left( \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} (x_{j1} \ x_{j2}) + 1 \right)^2$$

$$= (x_{i1}x_{j1} + x_{i2}x_{j2} + 1)^2$$

$$= \boxed{x_{i1}^2} \boxed{x_{j1}^2} + \boxed{x_{i2}^2} \boxed{x_{j2}^2} + \boxed{1} \boxed{1} + \boxed{\frac{\sqrt{2}x_{i1}x_{j2}}{\sqrt{2}}} \boxed{\frac{\sqrt{2}x_{j1}x_{i2}}{\sqrt{2}}}$$

$$+ \boxed{\frac{\sqrt{2}x_{i1}x_{j1}}{\sqrt{2}}} \boxed{\frac{\sqrt{2}x_{j1}x_{i1}}{\sqrt{2}}} + \boxed{\frac{\sqrt{2}x_{i2}x_{j2}}{\sqrt{2}}} \boxed{\frac{\sqrt{2}x_{j2}x_{i2}}{\sqrt{2}}}$$

$$\Rightarrow \varphi(x) = (x_1^2, x_2^2, 1, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2)$$

# Popular Kernels

## Discussion

- Other popular kernels include the following.

- ① Linear kernel:  $K_{ii'} = x_i^T x_{i'}$  ← SVM

- ② Polynomial kernel:  $K_{ii'} = (x_i^T x_{i'} + 1)^d$  ←

- ③ Radial Basis Function (Gaussian) kernel:

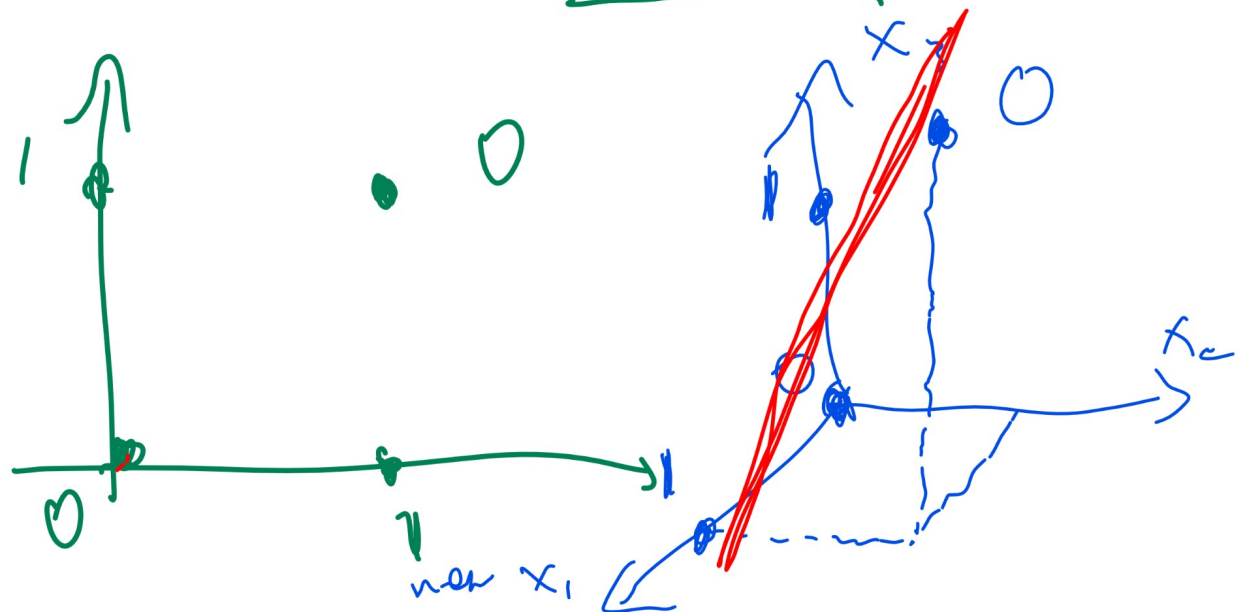
$K_{ii'} = \exp\left(-\frac{1}{\sigma^2} (x_i - x_{i'})^T (x_i - x_{i'})\right)$  Kernel SVM  
n x n

- Gaussian kernel has infinite dimensional feature representations. There are dual optimization techniques to find  $w$  and  $b$  for these kernels.

# Kernel Trick for XOR

Quiz y	$x_1$	$x_2$	new $x_1$	$x_2$	$x_3$
0	0	0	0	0	0
1	0	1	0	0	1
1	1	0	1	0	0
0	1	1	1	$\sqrt{2}$	1

- March 2018 Final Q17
- SVM with quadratic kernel  $\varphi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$  can correctly classify the training set for  $y = \overline{x_1} \text{ XOR } x_2$ .
- A: True.
- B: False.



# Kernel Trick for XOR 2

## Quiz

Q2

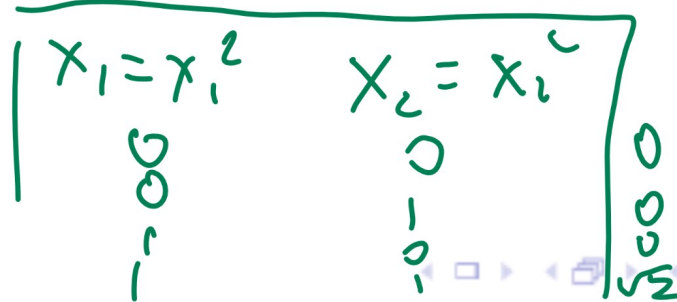
$$\phi(x_1, x_2) = (x_1^3, x_1^2, x_1)$$

- SVM with quadratic kernel  $\phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$  can correctly classify the training set for  $y = x_1 \text{ NAND } x_2$ . NAND is just "not and".

- A: True
- B: False.

feature mapping

$x_1$	$x_2$	NAND
0	0	1
0	1	1
1	0	1
1	1	0



# Kernel Matrix

## Quiz

- Fall 2009 Final Q2
- What is the feature vector  $\varphi(x)$  induced by the kernel

$K_{ij} = \exp(x_i + x_j) + \sqrt{x_i x_j} + 3$  in 1D

- A:  $(\exp(x), \sqrt{x}, 3)$
- B:  $(\exp(x), \sqrt{x}, \sqrt{3})$
- C:  $(\sqrt{\exp(x)}, \sqrt{x}, 3)$
- D:  $(\sqrt{\exp(x)}, \sqrt{x}, \sqrt{3})$
- E: None of the above

$x_i = x_{i1}$  one feature

$$\boxed{\exp(x_i)} \cdot \boxed{\exp(x_{j1})} + \boxed{\sqrt{x_i}} \cdot \boxed{\sqrt{x_{j1}}} + \boxed{\sqrt{3}} \cdot \boxed{\sqrt{3}}$$

# Kernel Matrix Math

## Quiz

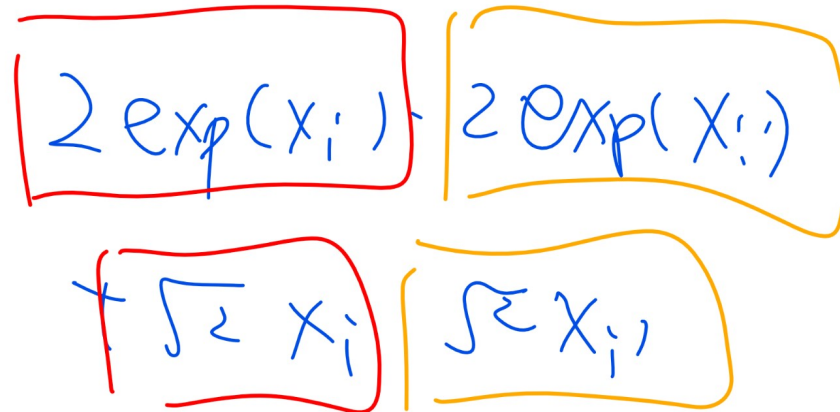
# Kernel Matrix 2

## Quiz

~~Q3~~ Q4

• What is the feature vector  $\varphi(x)$  induced by the kernel  $K_{ii'} = 4 \exp(x_i + x_{i'}) + \underline{2x_i x_{i'}}$ ?

- A:  $(4 \exp(x), 2\sqrt{x})$
- ~~• B:  $(2 \exp(x), \sqrt{2}\sqrt{x})$~~
- C:  $(4 \exp(x), 2x)$
- D:  $(2 \exp(x), \sqrt{2}x)$
- E: None of the above



# Kernel Matrix Math 2

## Quiz