

CS540 Introduction to Artificial Intelligence

Lecture 5

Young Wu

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

July 5, 2022

Adverse Selection

Quiz

Q1

- Suppose the last two digits of your 10-digit student ID is the expected grade (out of 100) you will get in a course. Choose between the two courses:
- A : a course in which you get your expected grade.
- B : a course in which you get the average expected grade of everyone taking this course.

Course Rhythm

Admin

- ① (*Q*) In-class Quizzes, 0.5 points ($T - F$)
- ② (*D*) Group discussion (reply to the Discussion post, also make it resolved), 0.5 points (*M*)
- ③ (*D*) Sharing solutions (create a note, not question, and tag *m2, m3, d1*), 0.5 points each (*M*)
- ④ (*M*) Math homework, 1 point (*M*)
- ⑤ (*P*) Programming homework, 8 points (*M*) ←
- ⑥ (*X*) Exams, see Midterm page for past exams (same format this year).

Course Grades

Admin

- Final grade = $0.3 \cdot X + 0.1 \max(X, Q) + 0.1 \max(X, D) + 0.1 \max(X, M) + 0.4 \cdot P$
- Additional discussion points used in borderline grades (for example 89 to A).

Sharing Solutions

Admin

- 1 Use LaTeX (Word, Maple, MyScript etc).

$\text{sqrt}((a_1^2) / (2 \text{ pi}))$ is difficult to read compared to $\sqrt{\frac{a_1^2}{2\pi}}$.

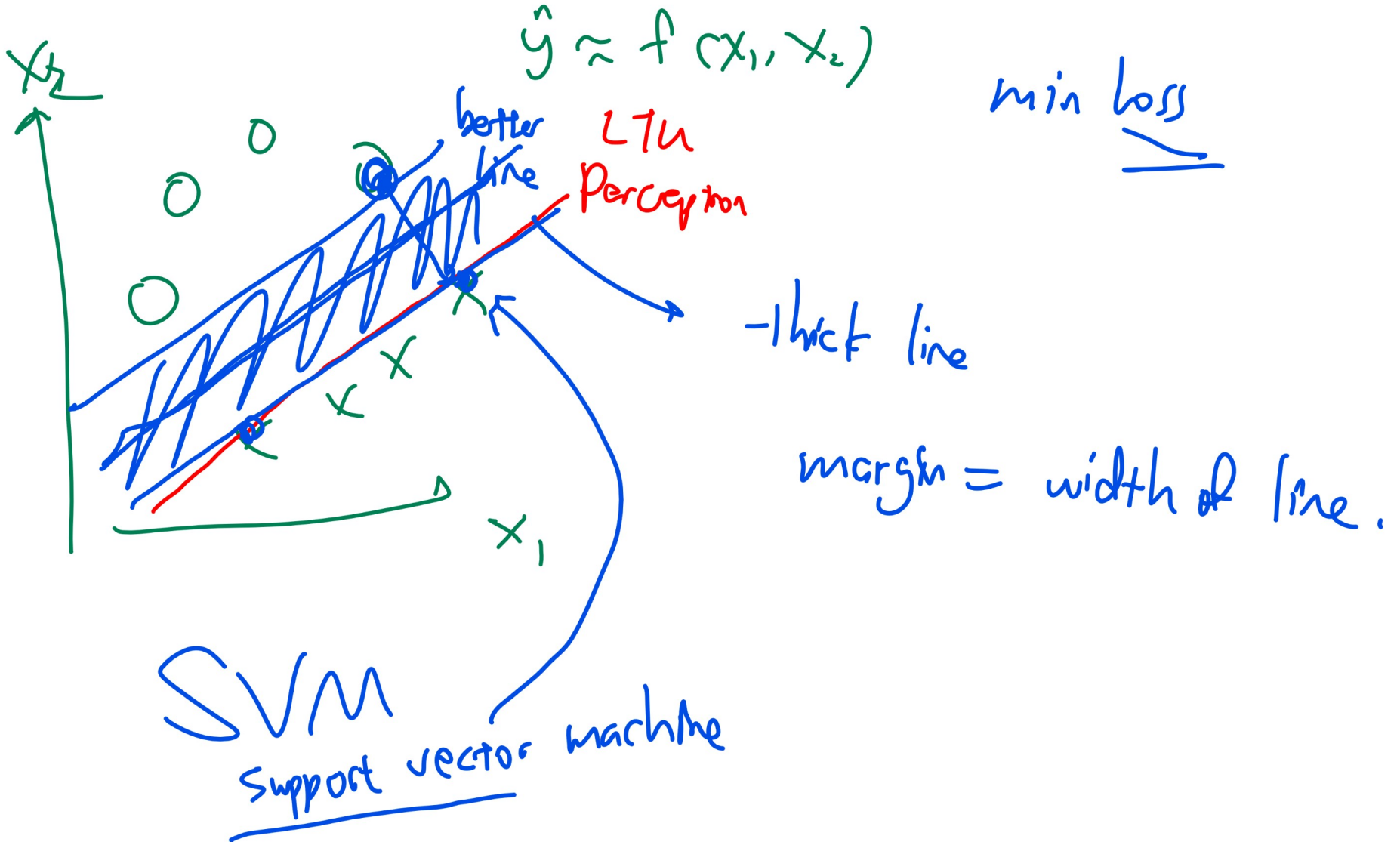
- 2 Handwritten on tablet or on paper and photo or scan (Office Lens).
- 3 Other suggestions?

Sharing Solutions

- For solution sharing, please make sure it is Piazza note, not a Piazza question.
- For actual questions, please use a different name, e.g. "M2Q1 Question" or "Question about M2Q1".
- Make sure you tag the post correctly: $m2$, $m3$, or $d1$ in order to get the points.
- Please sign up before making the post and please do not sign up for more than 4 questions per week.
- I will either "good note" the post or leave a comment: if I leave a comment, please update your answers, reply to my comment, and remember to make the reply "unresolved" so I can see.

Maximum Margin Diagram

Motivation



SVM Weights

Quiz

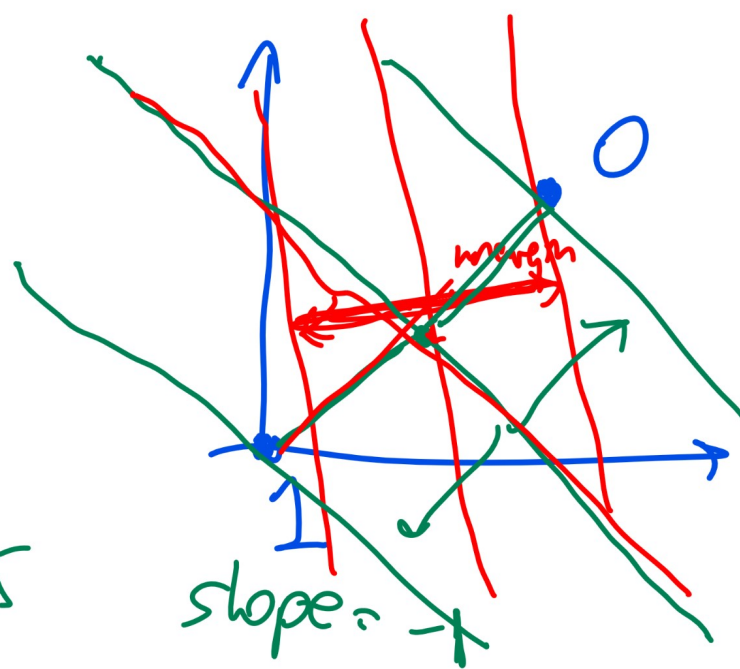
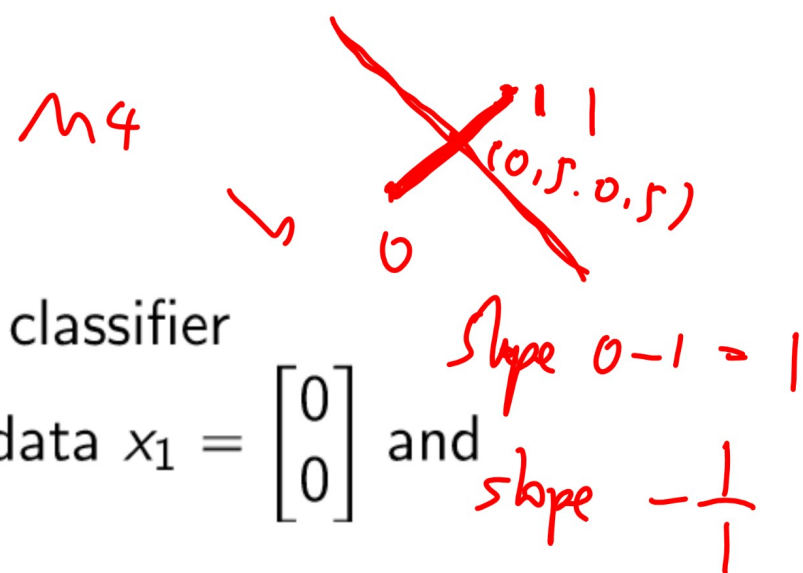
- Find the weights w_1, w_2 for the SVM classifier

$\mathbb{1}_{\{w_1 x_{i1} + w_2 x_{i2} + 1 \geq 0\}}$ 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 : $w_1 = -2, w_2 = -2$

$-x_1 - x_2 + 1 \geq 0$ 0.5, 0.5
 $x_2 = -x_1 + 1$



SVM Weights Diagram

Quiz

SVM Weights

Quiz

- Find the weights w_1, w_2 for the SVM classifier

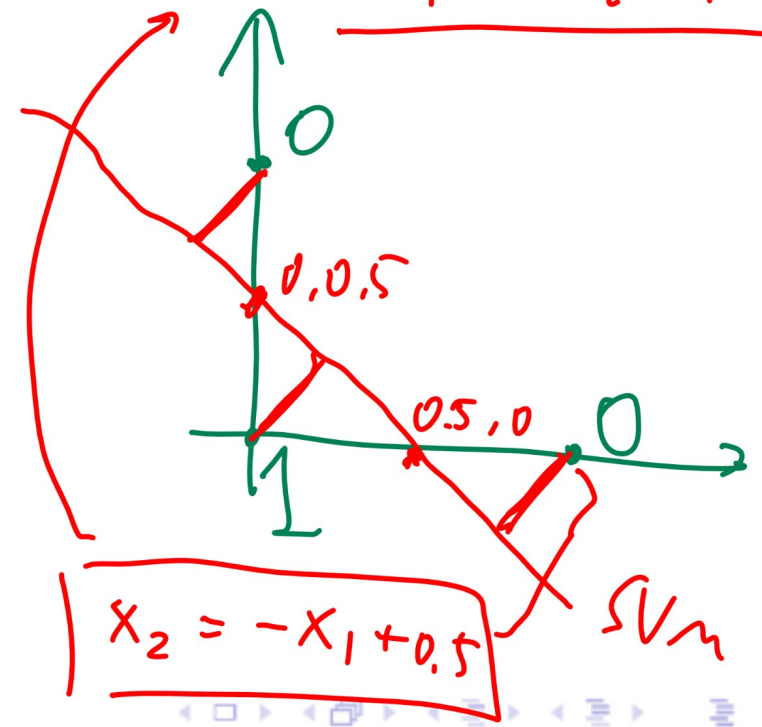
$\mathbb{I}\{w_1x_{i1} + w_2x_{i2} + 1 \geq 0\}$ given the training data

$x_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, x_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, x_3 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ with $y_1 = 1, y_2 = y_3 = 0$.

- A: $w_1 = -1.5, w_2 = -1.5$
- B: $w_1 = -2, w_2 = -1.5$
- C: $w_1 = -1.5, w_2 = -2$
- D: $w_1 = -2, w_2 = -2$**
- E: I don't understand SVM

Q2

$-2x_1 - 2x_2 + 1 = 0$



Slope = -1
 $\begin{bmatrix} 0 & 0.5 \\ 0.5 & 0 \end{bmatrix}$

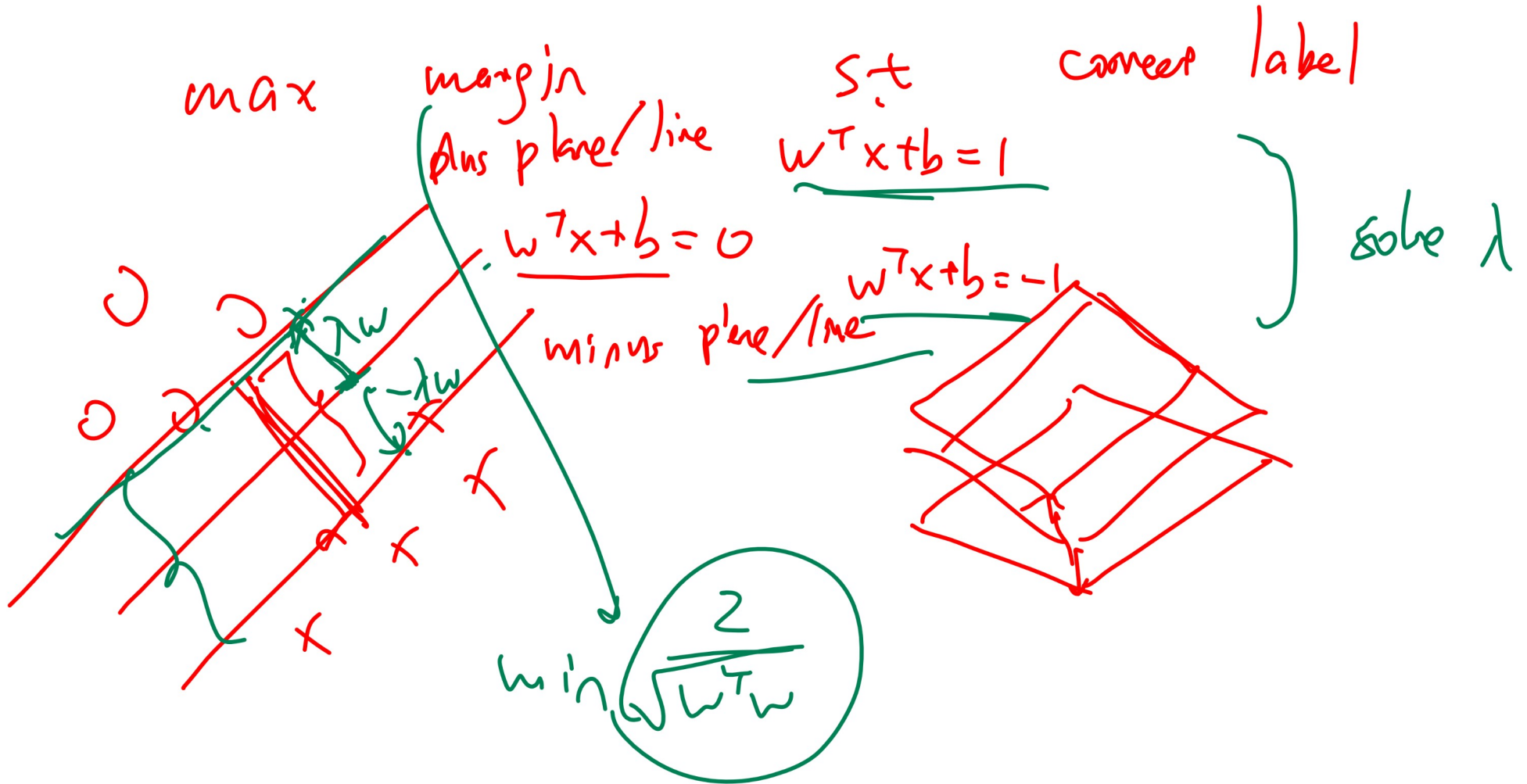


SVM Weights Diagram

Quiz

Constrained Optimization Diagram

Definition



Constrained Optimization Derivation

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 $y_i = 1$.

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

margin

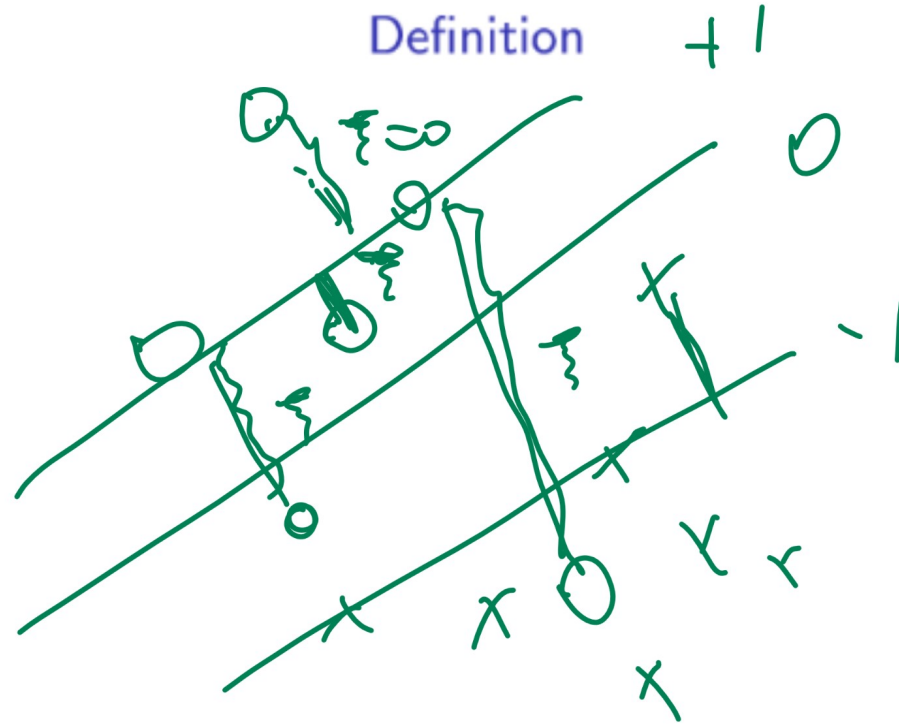
- This is equivalent to the following minimization problem, called hard margin SVM.

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

$$y_i = 1 \Rightarrow 1$$

Soft Margin Diagram

Definition



ξ_i

Soft Margin Derivation

Definition

SVM Formulations

Definition

- Hard margin:

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

margin

- Soft margin:

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

margin *loss (mistake)*

$$\min \frac{\lambda}{2} w^T w + \frac{1}{n} \sum \xi_i$$

Soft Margin

Quiz

- 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 ξ for it to satisfy the margin constraint? $\xi = 18$

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

$$(2 \cdot 0 - 1) \left(\underbrace{\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 4 \\ 5 \end{bmatrix} + 3}_{17} \right) \geq 1 - \xi_i$$



$$-17 \geq 1 - \xi_i$$

$\xi_i \geq 18$ AND $\xi_i \geq 0$

Soft Margin 2

Quiz

• Let $w = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $b = 3$. For the point $x = \begin{bmatrix} -4 \\ -5 \end{bmatrix}$, $y = 0$, Q3
 what is the smallest slack variable ξ for it to satisfy the margin constraint?

- A : -10
- B : 0
- C : 10
- D : None of the above
- → E : I don't understand what is ξ

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

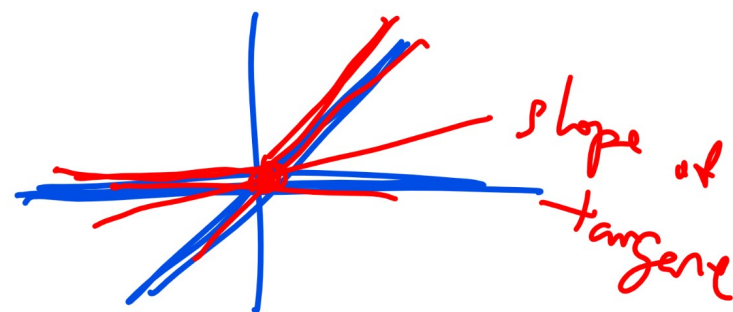
$$\xi_i \geq 0$$

$$\xi_i \geq -10$$

max ξ

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 a gradient.

Subgradient 1

Quiz

• Which ones are subderivatives of $\max\{x, 0\}$ at $x = 0$?

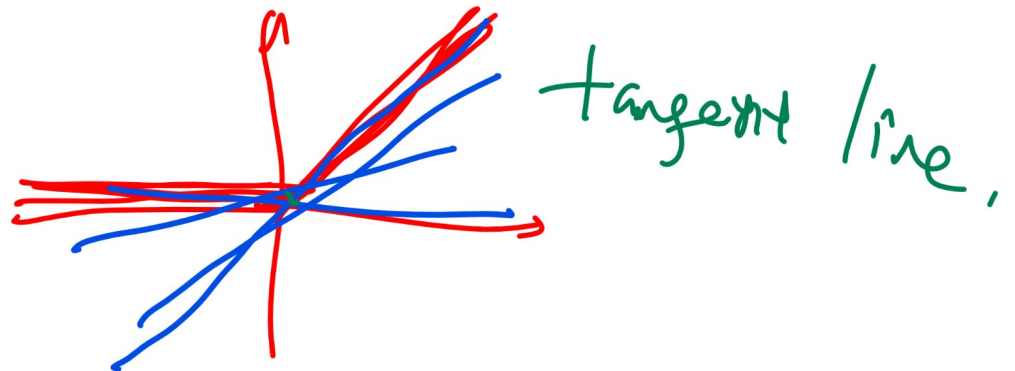
~~• A: -1~~

~~• B: -0.5~~

• C: 0

• D: 0.5

• E: 1

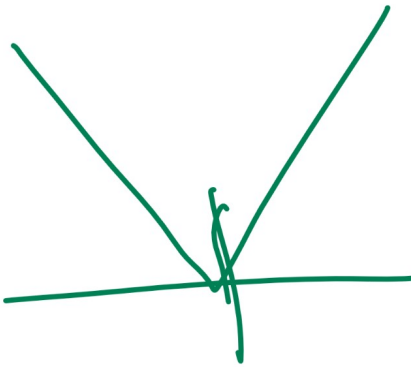


Subgradient 2

Quiz

• Which ones are subderivatives of $|x|$ 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.

Regularization Parameter

Definition

$$w = w - \alpha \sum_{i=1}^n z_i \mathbb{1}_{\{z_i w^T x_i \geq 1\}} x_i - \lambda w$$

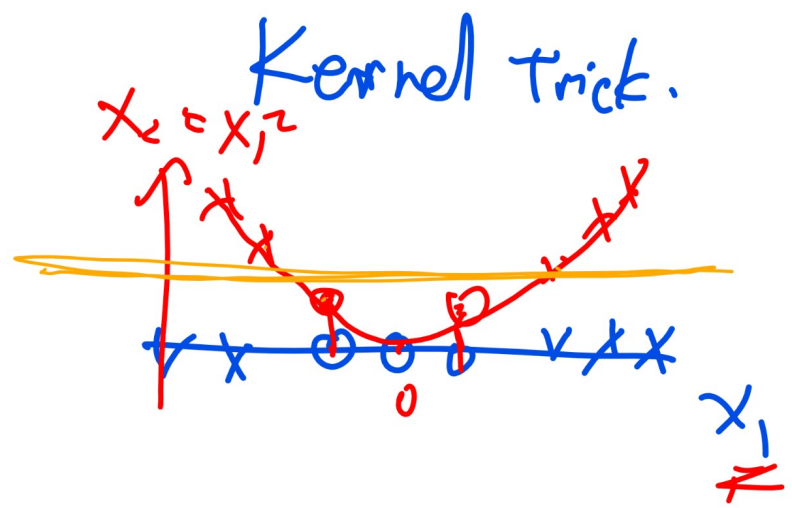
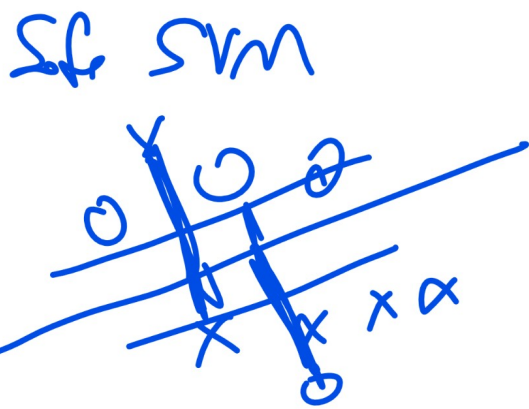
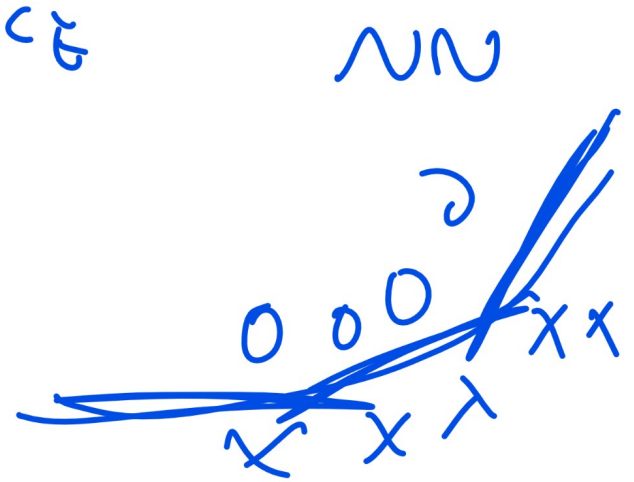
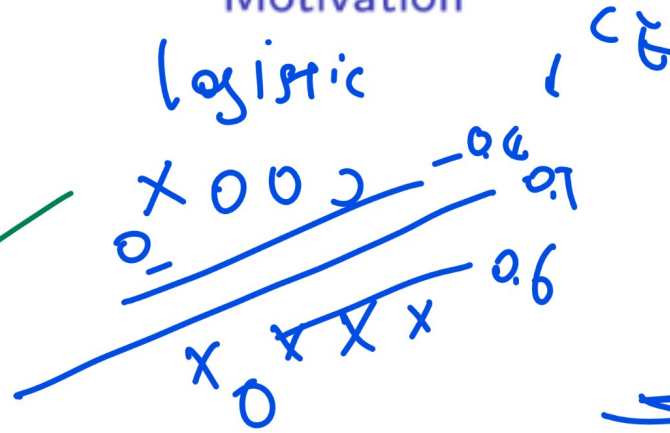
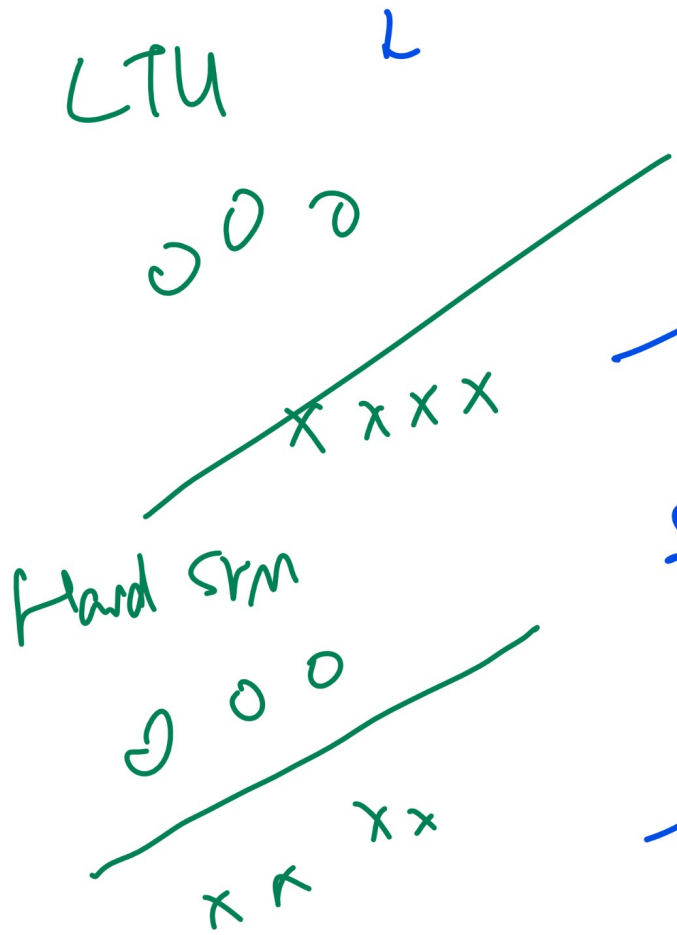
$$z_i = 2y_i - 1, i = 1, 2, \dots, n$$

Handwritten notes: "learning rate" points to α ; "min loss + $\frac{\lambda}{2} w^2$ " and "L2 regularization" are written above the update equation; " $-\lambda w$ " is written below the update equation.

- λ is usually called the regularization parameter because it reduces the magnitude of w the same way as the parameter λ in L_2 regularization.
- The stochastic subgradient descent algorithm for SVM is called PEGASOS: Primal Estimated sub-GrAdient SOLver for Svm.

Kernel Trick 1D Diagram

Motivation



Kernelized SVM

Definition

x_{11} x_{12} x_{13} x_{1n}

- With a feature map φ , 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.

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

$$\varphi\left(\begin{pmatrix} x_1 \\ y_2 \end{pmatrix}\right) = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 \end{pmatrix}$$

$$\varphi'\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{pmatrix}$$

Kernel Trick for XOR

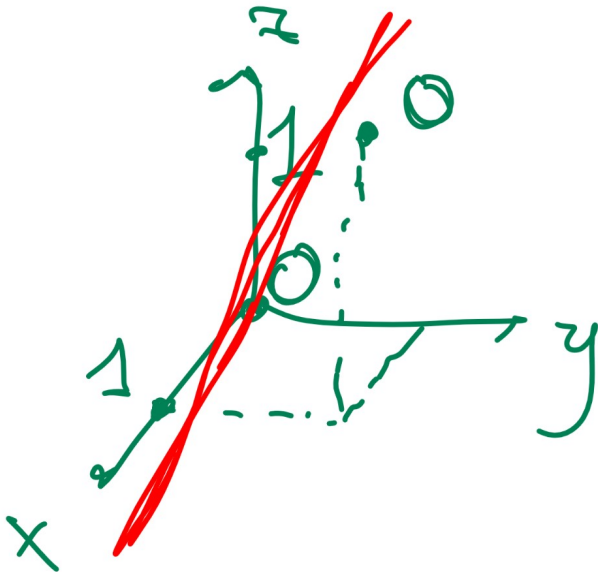
Quiz

XOR

- SVM with quadratic kernel $\varphi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$ can correctly classify the following training set?

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

x_1^2	$\sqrt{2}x_1x_2$	x_2^2
0	0	0
0	0	1
1	0	0
1	$\sqrt{2}$	1



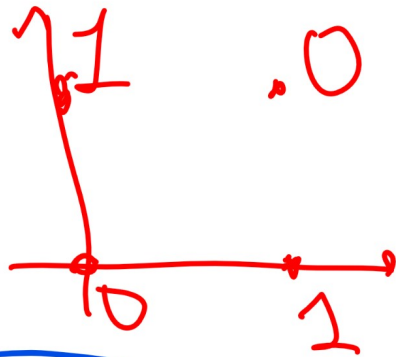
Kernel Trick for XOR

Quiz

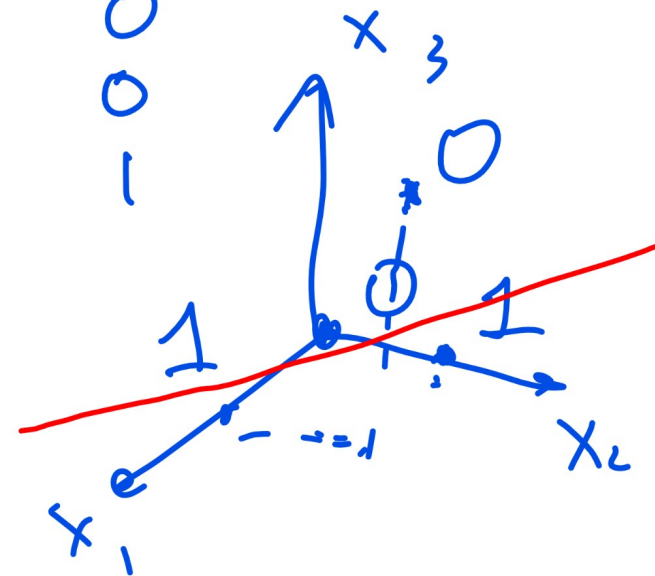
- SVM with kernel $\varphi(x) = (x_1, x_1x_2, x_2)$ can correctly classify the following training set.

Q2

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



$\underbrace{x_1}_{x_1}$
 $x_1 x_2$
 x_2
 0
 0
 0
 1
 1
 0



- A : True.
- B : False.

C: what is kernel?

Kernel Matrix

Definition

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

$$K_{ij'} = \varphi(x_i)^T \varphi(x_{j'})$$

Symmetric

p. s. d.
eigenvalue ≥ 0

$K \geq 0$

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.

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

Examples of Kernel Matrix Derivation

Definition

Popular Kernels

Discussion

- Other popular kernels include the following.

① Linear kernel: $K_{ij'} = x_i^T x_{i'}$

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

- ③ Radial Basis Function (Gaussian) kernel:

$$K_{ij'} = \exp\left(-\frac{1}{\sigma^2} (x_i - x_{i'})^T (x_i - x_{i'})\right)$$

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

Kernel Matrix

Quiz

1d-feature ith. jth item/instance

- What is the feature vector $\phi(x)$ induced by the kernel

$$K_{ij} = \exp(x_i + x_j) + \sqrt{x_i x_j} + 3?$$

$$\phi(x) = \begin{pmatrix} \exp(x) \\ \sqrt{x} \\ \sqrt{3} \end{pmatrix}$$

$$K_{ij} = \phi^T(x_i) \phi(x_j)$$

$$= \begin{bmatrix} \exp(x_i) \\ \sqrt{x_i} \\ \sqrt{3} \end{bmatrix} \cdot \begin{bmatrix} \exp(x_j) \\ \sqrt{x_j} \\ \sqrt{3} \end{bmatrix} + \begin{bmatrix} \sqrt{x_i} \\ \sqrt{x_j} \\ \sqrt{3} \end{bmatrix} \cdot \begin{bmatrix} \sqrt{x_i} \\ \sqrt{x_j} \\ \sqrt{3} \end{bmatrix} + \sqrt{3} \cdot \sqrt{3}$$

$$\phi(x_i) = \begin{pmatrix} \exp(x_i) \\ \sqrt{x_i} \\ \sqrt{3} \end{pmatrix}^T \begin{pmatrix} \exp(x_j) \\ \sqrt{x_j} \\ \sqrt{3} \end{pmatrix} = \phi(x_j)$$

Kernel Matrix Math

Quiz

Kernel Matrix 2

Quiz

- What is the feature vector $\varphi(x)$ induced by the kernel

$$K_{ii'} = 4 \exp(x_i + x_{i'}) + 2x_i x_{i'}$$

- 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

Q3

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

$$\begin{pmatrix} 2 \exp(x_i) \\ \sqrt{2} x_i \end{pmatrix}^T \begin{pmatrix} 2 \exp(x_{i'}) \\ \sqrt{2} x_{i'} \end{pmatrix}$$

Kernel Matrix Math 2

Quiz