

CS540 Introduction to Artificial Intelligence

Lecture 5

Young Wu

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

June 6, 2022

Guess the Percentage

Admin

- Guess what percentage of the students (who are here) started $P1$?
- A : 0 to 20 percent.
- B : 20 to 40 percent.
- C : 40 to 60 percent.
- D : 60 to 80 percent.
- E : 80 to 100 percent.

The Percentage

Admin

- Did you start $P1$?
- A :
- B : Yes.
- C :
- D : No.
- E :

Guess the Percentage

Admin

- Guess what percentage of the students (who are here) submitted $P1$?
- A : 0 to 20 percent.
- B : 20 to 40 percent.
- C : 40 to 60 percent.
- D : 60 to 80 percent.
- E : 80 to 100 percent.

Sharing Solutions

Admin

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

sqrt((a_1^2) / (2 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?
 - Remember to make it a public Piazza Note (not a Question).
 - 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).

Shared Solution List and Feedback

Admin

- The shared solutions are listed in the Main post on Piazza.
- Thank you for the feedback! I posted the responses to those in the Feedback post on Piazza.
- Question: can you see the labels in the 3D diagrams?
- A : Yes.
- B : Cannot see label.
- C : Cannot see anything.

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 : $w_1 = -2, w_2 = -2$

SVM Weights Diagram

Quiz

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

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

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

SVM Weights Diagram

Quiz

Constrained Optimization Diagram

Definition

Constrained Optimization Derivation

Definition

Soft Margin Diagram

Definition

Soft Margin Derivation

Definition

SVM Formulations

Definition

- Hard margin:

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

- Soft margin:

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

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?

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

Soft Margin 2

Quiz

- Let $w = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $b = 3$. For the point $x = \begin{bmatrix} -4 \\ -5 \end{bmatrix}$, $y = 0$, what is the smallest slack variable ξ for it to satisfy the margin constraint?
- A : -12
- B : -10
- C : 0
- D : 10
- E : 12

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

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

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

- λ is usually called the regularization parameter because it reduces the magnitude of w the same way as the parameter λ in $L2$ 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

- 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.

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

Kernel Trick for XOR

Quiz

- 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

Kernel Trick for XOR

Quiz

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

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

- A : True.
- B : False.

Kernel Matrix

Definition

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

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

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.

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

Examples of Kernel Matrix Derivation

Definition

Popular Kernels

Discussion

- Other popular kernels include the following.

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

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

- ③ Radial Basis Function (Gaussian) kernel:

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

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

Kernel Matrix

Quiz

- What is the feature vector $\varphi(x)$ induced by the kernel $K_{ii'} = \exp(x_i + x_{i'}) + \sqrt{x_i x_{i'}} + 3$?

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

Kernel Matrix Math 2

Quiz