

# CS446 HW1

Junsheng Huang

January 25, 2024

## 1

### 1.1

For each test image, we need to calculate the distance for every dataset node in worst case. We have  $M$  test images and  $N$  training images and each distance calculate need  $d$  dimensions, so the answer is  $O(M*N*d)$

### 1.2

$k = 10$  should have a "smoother" decision boundary. Because for bigger  $k$ , we would have a simpler model and more training nodes to "vote" for single test node. For example, when the  $k$  is too big to contain all the training set, the decision boundary is the whole set since we only have one decision. When  $k = 1$ , the model is over-fit and decision boundary is made up of many single nodes.

### 1.3

Consider  $\mathbf{x} = \begin{bmatrix} x1 \\ x2 \end{bmatrix}$ ,  $\mathbf{w} = [1, 1]$ ,  $b = 0$ , so  $\mathbf{w} \cdot \mathbf{x} + b = x1 + x2$  and this matches ground truth  $g$ .

### 1.4

The relationship between the largest singular value of  $A$  (named it as  $x$ ) and the largest eigenvalue of  $A^T A$  (named it as  $\lambda$ ) is  $x = \sqrt{\lambda}$ .

This is because: (remember  $V$  is orthogonal matrix)

$$A^T \cdot A = (U\Sigma V^T)^T (U\Sigma V^T) \quad (1)$$

$$= V\Sigma^T U^T \cdot U\Sigma V^T \quad (2)$$

$$= V\Sigma^T \Sigma V^T \quad (3)$$

$$= V\Sigma^T \Sigma V^{-1} \quad (4)$$

That means  $\Sigma^2$  is the eigenvalue matrix of  $A^T A$ , so the relationship is  $x = \sqrt{\lambda}$ .

### 1.5

The Naive Bayes assumption is violated at text world. For example, the sequence "of course" is more common than it would be if the words "of" and "course" were conditionally independent. That is,  $Pr(X = of\ course|Y) > Pr(x1 = of|Y)Pr(x2 = course|Y)$ .

## 2

### 2.1

We can write  $X \in R^{n \times d}$  to  $A \in R^{d \times d}$ ,  $w' \in R^d$ ,  $B \in R^d$ , all new rows are filled with 0, so the rank of  $A$  is  $n$ , which is not full-rank, thus  $w'$  has at least one solution, and if there exist  $w'$  that satisfy  $Aw' = B$ , there must exist  $w$  that satisfy  $Xw = y$ , and the empirical risk with squared loss is zero. (Or you can consider that the dimension of  $X$ 's column space is  $n$ , and  $\mathbf{y} \in R^n$ , so  $\mathbf{y}$  is in the column space of  $X$ , thus we can say that there always exist a  $\mathbf{w}$  that  $X\mathbf{w} = \mathbf{y}$ .)

## 2.2

The rank of  $x$  is  $n$ , so we have  $n$  eigenvalue and singular value. Since  $\Sigma$  should have the singular value at its diagonal, the rank of  $\Sigma$  is  $n$ .

## 2.3

$$X \cdot X^T = (U\Sigma V^T)(U\Sigma V^T)^T \quad (5)$$

$$= U\Sigma V^T \cdot V\Sigma^T U^T \quad (6)$$

$$= U\Sigma\Sigma^T U^T \quad (7)$$

since  $U$ ,  $\Sigma^2$  and  $U^T$  are all full-rank matrices,  $XX^T$  is all full-rank, thus is invertible.

## 3

### 3.1

The smallest number of support vectors is 2. This happens when exactly 2 support vectors on the hard margin of the SVM.

### 3.2

since the optimal solution to the dual only has 3 values that is not zero, there are at least 3 support vectors. So the smallest number of support vectors is 3, the largest number of support vectors is the total number of  $\alpha$  (or 10000 if using first part's data number).

### 3.3

#### 3.3.1

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^2 \quad (8)$$

$$= (x_1 z_1 + x_2 z_2 + 1)^2 \quad (9)$$

$$= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 x_2 z_2 + 2x_1 z_1 + 2x_2 z_2 + 1 \quad (10)$$

$$= \phi(\mathbf{x})^T \phi(\mathbf{z}) \quad (11)$$

so for  $\mathbf{x} = (x_1, x_2)$ ,  $\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$

#### 3.3.2

since we only consider the sign of  $x_1$  and  $x_2$  (same for 1 and different for -1): we can only consider the  $\sqrt{2}x_1 x_2$  in  $\phi(\mathbf{x})$ , which means  $w = (0, 0, 1, 0, 0, 0)$  can work.

## 4

### 4.1

$$P(y = +1|x) = \frac{P(x, y = +1)}{P(x)} \quad (12)$$

$$= \frac{P(x|y = +1) \cdot P(y = +1)}{P(x|y = +1) \cdot P(y = +1) + P(x|y = -1) \cdot P(y = -1)} \quad (13)$$

$$= \frac{P(x|y = +1) \cdot p}{P(x|y = +1) \cdot p + P(x|y = -1) \cdot (1 - p)} \quad (14)$$

$$= \frac{1}{1 + \frac{P(x|y=-1) \cdot (1-p)}{P(x|y=+1) \cdot p}} \quad (15)$$

$$= \frac{1}{1 + \exp(\log \frac{A}{B})} \quad (16)$$

$$(17)$$

where  $A = P(x|y = -1) \cdot (1 - p)$ ,  $B = P(x|y = +1) \cdot p$

## 4.2

assumption: each  $x_j$  is conditional independent,  $x = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}$

Thus, we have:

$$\log\left(\frac{A}{B}\right) = \log\left(\frac{P(x|y = -1) \cdot (1-p)}{P(x|y = +1) \cdot p}\right) = \log\left(\frac{1-p}{p}\right) + \log\left(\frac{P(x|y = -1)}{P(x|y = +1)}\right) \quad (18)$$

$$\log\left(\frac{P(x|y = -1)}{P(x|y = +1)}\right) = \log\left(\exp\left(-\frac{1}{2} \cdot \sum (x_j - \mu_{-,j})^2 + \frac{1}{2} \cdot \sum (x_j - \mu_{+,j})^2\right)\right) \quad (19)$$

$$= -\frac{1}{2} \sum (\mu_{-,j})^2 + \frac{1}{2} \sum (\mu_{+,j})^2 + \sum x_j (\mu_{-,j} - \mu_{+,j}) \quad (20)$$

so  $\log\left(\frac{A}{B}\right)$  can be written in the form  $\mathbf{w}^T \mathbf{x} + b$ ,  $\mathbf{w} = \boldsymbol{\mu}_- - \boldsymbol{\mu}_+$ ,  $b = \log\left(\frac{1-p}{p}\right) + \frac{1}{2} \boldsymbol{\mu}_+^T \boldsymbol{\mu}_+ - \frac{1}{2} \boldsymbol{\mu}_-^T \boldsymbol{\mu}_-$

## 4.3

$$P(y|\mathbf{x}) = \frac{1}{1 + \exp(y \cdot (\mathbf{w}^T \mathbf{x} + b))} \quad (21)$$

## 5

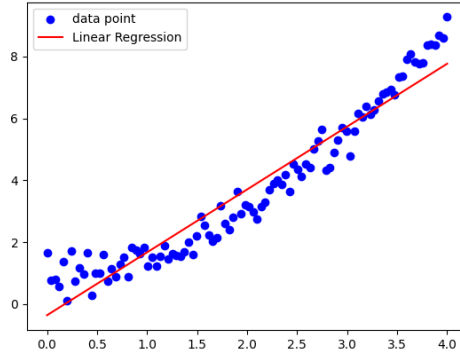


Figure 1: Linear Regression