

### 1.1

Nearest Neighbor Classifier is using labeled data as reference to classify a data without label by choosing the nearest labeled cluster. K-means clustering is a process to cluster unlabeled datas into cluster by iterations of choosing random centroid as a clustering reference. So Nearest Neighbor Classifier can be viewed as K means clustering after choosing centroids to cluster to. Yes, finding a cluster from the existing ones for a new data point in the test set K-Means clustering can be viewed as Nearest Neighbor Classifier because the K means centroids have been chosen.

### 1.2

Yes it is possible. We can at first preprocess the data by taking the natural logarithm of all  $y$ , so the  $N$  datapoint will become  $N$  data points  $\{(x_{1-n}, \ln(y_{1-n}))\}$ , then we can use LinearRegression model from ~~sklearn~~ to fit this processed data. If the score of the model (which models the loss) is within a reasonable range, then we can this data set has a exponential relationship.

2.

$$\textcircled{1} \quad L = \frac{1}{n} \sum_i^n [y^i - (\theta_0 + \theta_1 x^i)]^2 + \beta (e^{\theta_0} + e^{\theta_1})$$

$$\textcircled{2} \quad \frac{\partial L}{\partial \theta_0} = -\frac{2}{n} \sum_i^n (y^i - \theta_0 - \theta_1 x^i) + \beta e^{\theta_0}$$

$$\textcircled{3} \quad \frac{\partial L}{\partial \theta_1} = -\frac{2}{n} \sum_i^n x_i (y^i - \theta_0 - \theta_1 x^i) + \beta e^{\theta_1}$$

$$\textcircled{4} \quad \theta_0 \leftarrow \theta_0 - \alpha \frac{\partial L}{\partial \theta_0}$$

$\textcircled{5}$  I think it won't work well because if  $\beta$  is large, which means a big regularization penalty,

for  $\beta e^{\theta_0} / \beta e^{\theta_1}$  to minimize,  $\theta_0 / \theta_1 \rightarrow -\infty$ , which against the first purpose of trying to reduce  $\theta_0 / \theta_1$  to 0 to simplify the model.

3.

(1) The 3 nearest neighbours are

$x_7$ : Euclidean distance 1

$x_4, x_5$ , Euclidean distance  $\sqrt{2}$ ,

Label should be  $x_4 x_5$ 's label,  $y=+1$

(2) The 4<sup>th</sup> and 5<sup>th</sup> nearest neighbour is  $x^2, x^6$ , distance = 2

Label should be  $y=-1$

(3) The five nearest points are:

$x_7$ ,  $d=1$

$x_4$   $d=2$

$x_5$   $d=2$

$x_2$   $d=4$

$x_6$   $d=4$ , label should be  $y=-1$

$$4.1 \quad A = \begin{bmatrix} 1, 2 \\ 3, 4 \end{bmatrix} \quad A^T = \begin{bmatrix} 2, 4 \\ 1, 3 \end{bmatrix}$$

$$\begin{bmatrix} 1, 2 \\ 3, 4 \end{bmatrix} - \begin{bmatrix} 2, 4 \\ 1, 3 \end{bmatrix} = \begin{bmatrix} -1, -2 \\ 2, 1 \end{bmatrix} \neq \begin{bmatrix} -2, 1 \\ -1, 2 \end{bmatrix}$$

$$4.2 \quad \left( \|X\|_2 \right)^2 = X^T X$$

$$X^T X \cdot X^T X \cdot X^T X \cdot X^T X = 2 \times 2 \times 2 \times 2 = 16$$

## 5 (10 points) Clustering

Assume there is a dataset  $S = \{x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}, x^{(5)}\}$  where:

$$\begin{aligned}x^{(1)} &= -8, & x^{(2)} &= -10, & x^{(3)} &= 0, \\x^{(4)} &= 11, & x^{(5)} &= 12.\end{aligned}$$

- If we have initial centroids  $c_1 = 0$ ,  $c_2 = 11$ ,  $c_3 = 12$  and perform  **$K$ -means** clustering algorithm on dataset  $S$  using Euclidean distance. Fill in the table until the algorithm converges. To show the algorithm has converged, it is necessary to fill in one more iteration that duplicates the previous one. You may not use all the rows. Some example entries are filled. For each of the iteration, find the positions of the centroids, assign the data points to the centroids, and compute the loss.

5.2

Centroids positions : -9, 0, 11.5

$$SSE = 1 + 1 + 0.25 + 0.25 = 2.5$$

5.3

As mentioned above, -9, 0, 11.5

(6)

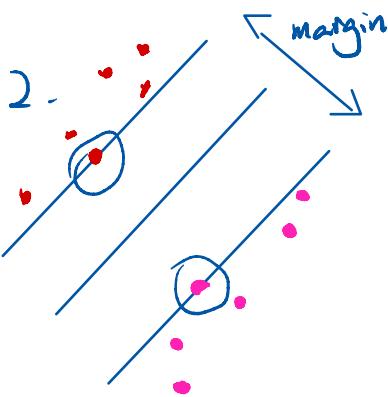
1.

$$\frac{\partial f}{\partial w} = 2\gamma w + \frac{1}{m} \sum_{i=1}^m -y_i x_i, \quad 1 - y_i(w^T x_i - b) > 0$$

note, only sum the  $-y_i x_i$

if that  $x_i, y_i$  satisfies

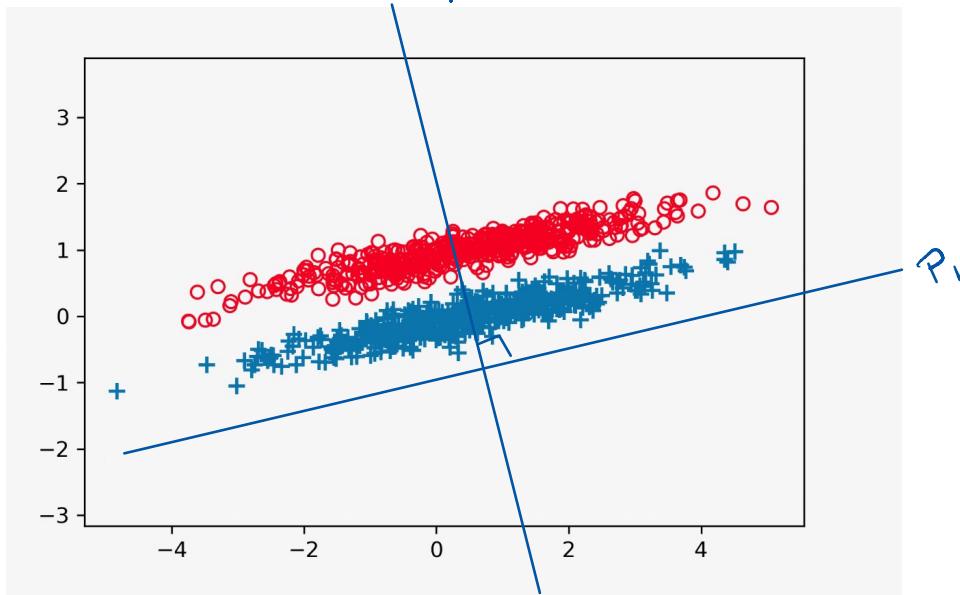
$1 - y_i(w^T x_i - b) > 0$ , in total  
there is  $m$  points satisfied  
the condition,  $0 \leq m \leq n$



For a binary-label dataset  
that has minimized support  
vectors, the minimized vector  
is 2, as showed in figure

7.

(1)



(2)  $P_2$

8.

$$\text{if } \hat{y} - y = [1, 1, 1, \dots, 1]$$

$\underbrace{\hspace{10em}}$   
 $n$

that is, each prediction  $\hat{y}_i$  is 1 more than  $\hat{y}$ .

then :

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| = 1$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = 1$$

$$ME = \max_{i=1}^n |y_i - \hat{y}_i| = 1$$

$$9. \textcircled{1} k(x, y) = (x \cdot y + 2)^2$$

$$\Phi(x_i) = \langle x_1^2, x_2^2, \dots, x_d^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \dots, \sqrt{2}x_1x_d, \sqrt{2}x_2x_3, \dots, \sqrt{2}x_2x_d, \dots, \sqrt{2}x_{d-1}x_d, 2x_1, \dots, 2x_d, 2 \rangle$$

$$\Phi(y_i) = \langle y_1^2, y_2^2, \dots, y_d^2, \sqrt{2}y_1y_2, \sqrt{2}y_1y_3, \dots, \sqrt{2}y_1y_d, \dots, \sqrt{2}y_{d-1}y_d, 2y_1, \dots, 2y_d, 2 \rangle$$

$$\textcircled{2} \quad k(x, y) = (xy - 1)^3$$

$$d=2 \quad x = [1, 0] \quad y = [0, 1]$$

$$K = \begin{bmatrix} k(x, x), & k(x, y) \\ k(y, x), & k(y, y) \end{bmatrix} = \begin{bmatrix} 0, & -1 \\ -1, & 0 \end{bmatrix}$$

$$a = [1, 1] \quad \text{aka } a^T = [1, 1] \begin{bmatrix} 0, -1 \\ -1, 0 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix} = -2 < 0$$

$K$  is not positive semidefinite matrix, so not a valid kernel function