

ECE M146 HW4 Melody Chen #705120273

1. In this section, you will use the k-NN classifier to predict whether a person survives or not on the Titanic. You will be using the same data set provided in HW3.

Unlike what we did in the last homework, we are going to build the k-NN classifier from the first principle. You may **not** use `fitcknn` (`sklearn.neighbors.KNeighborsClassifier` for python) in this problem as you will get incorrect answer by using those built-in functions.

The k-NN classifier classifies a data point with feature x_{test} based on a training set by performing the following procedures:

- Compute the distance from x_{test} to the feature of all training points. We will use the Euclidean distance in this problem.
- Find the k nearest neighbors of this point.
- Classify this point as the majority class of its k nearest neighbors.

We use the following two rules to handle ties:

- Let d_k be the distance of the k -th nearest neighbor of x_{test} . If there are multiple training points that have distance d_k from x_{test} . Choose those points with the smallest indexes to be included in the k nearest neighbors. For example, let $k = 3$, if there is x_9 that is distance 1 away from x_{test} ; x_1, x_3 and x_4 that is distance 2 away from x_{test} , then the 3 nearest neighbor of x_{test} are x_1, x_3 and x_9 . Note that $d_k = 2$ in this example.
- For even k , among all k nearest neighbors of a data point, if the number of points from class 0 is the same as the number of points from class 1, classify this data point as y_{tie} deterministically.

Answer the following questions:

- (a) With $y_{tie} = 1$, implement the k-NN algorithm. Find and plot the training and testing error for $k = 1, 2, \dots, 15$.

shown to the right.

- (c) Comment on the performance of the k-NN classifier in (a) and (b). How does larger k affect the training and testing error? How does even or odd k affect the performance of the k-NN classifier in (a) and (b), respectively? Are they contradictory to each other? Explain why. Hint: you may use your result from HW3 Q6 (a) to get some intuition.

Larger k doesn't appear to significantly affect training error except for $k=1$, training error for both value of y_{tie} is significantly lower than other k s, which makes sense as we're only looking at one neighbor for a training data. Larger k stabilizes and converges the testing error to what we previously calculated for both a) and b).

For a) even k appears to frequently cause surges in both training and testing error, which is likely due to frequent ties. For b), even k does not lead to surges which is likely due to $y_{tie}=0$ being the more accurate guess in a tie.

Two results are not contradictory as surges for even k in a) should correspond to non-surges in b) as the prediction has only two possible outcomes.
By surges, I mean occasional bumps in testing error.

2. In this section, you will consider a k-nearest neighbor classifier using Euclidean distance metric on a binary classification task. We assign the class of the test point to be the class of the majority of the k nearest neighbors.

Consider the following two datasets:

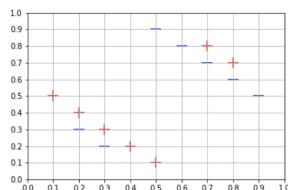


Table 1: KNN Example 1

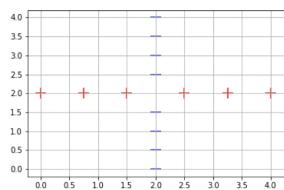
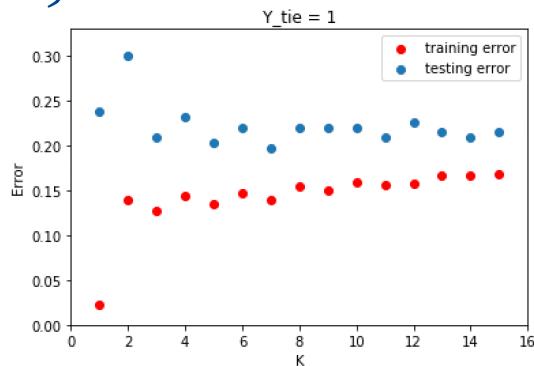


Table 2: KNN Example 2

- (a) For $k \in \{1, 3, 5, 7\}$, what values of k minimize leave-one-out cross-validation error for each dataset? What is the resulting validation error?

a)



- (b) With $y_{tie} = 0$, implement the k-NN algorithm. Find and plot the training and testing error for $k = 1, 2, \dots, 15$.

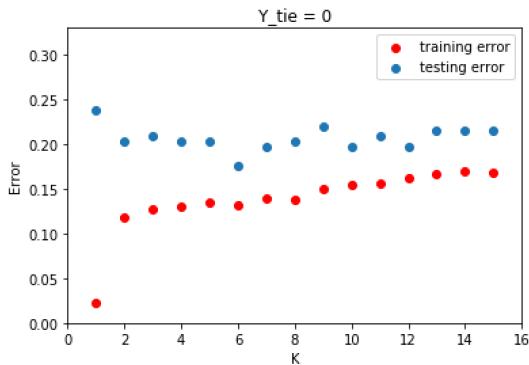


Table 1 $k=1$, 10 points misclassified. $\frac{10}{14}$ cross validation error.

$k=3$, 6 points misclassified. $\frac{6}{14}$ cross validation error.

$k=5$, 4 points misclassified. $\frac{4}{14}$ cross validation error.

$k=7$, 4 points misclassified. $\frac{4}{14}$ cross validation error.

$k=5, k=7$ minimizes Loo cross-validation error.

$k=1, 2$ points misclassified. $\frac{2}{14}$ cross validation error.

$k=3$, 4 points misclassified. $\frac{4}{14}$ cross validation error.

$k=5$, at least 4 points misclassified.

$k=7$, at least 4 points misclassified.

$k=1$ minimizes Loo cross-validation error.

- (b) In general, what are the drawbacks of using too big a k ? What about too small a k ? To see this, calculate the leave-one-out cross-validation error for the dataset in Figure 1 using $k = 1$ and $k = 13$.

$$k=1, \text{ LOO cross validation error: } 10/14$$

$$k=13, \text{ LOO cross validation error: } 14/14$$

The drawbacks of using too big of k is that decision regions become large uninterrupted segments, where the largest class could dominate the results.

The drawbacks of using too small of k is overfitting, as the decision boundary becomes very fine grained.

3. In this problem, we will derive the least square solution for multi-class classification. Consider a general classification problem with K classes, with a 1-of- K binary encoding scheme (defined latter) for the target vector $t, t \in \mathbb{R}^K$. Suppose we are given a training data set $\{x_n, t_n\}, n = 1, \dots, n$ where $x_n \in \mathbb{R}^D$. For the 1-of- K binary encoding scheme, t_n has the k -th element being 1 and all other elements being 0 if the n -th data is in class k . We can use the following linear model to describe each class:

$$y_k(x) = w_k^T x + w_{k0},$$

where $k = 1, \dots, K$. We can conveniently group these together using vector notation so that

$$y(x) = \tilde{W}^T \tilde{x},$$

where \tilde{W} is a matrix whose k -th column comprises the $D + 1$ -dimensional vector $\tilde{w} = [w_{k0}, w_k^T]^T$ and \tilde{x} is the corresponding augmented input vector $[1, x^T]^T$. For each new input with feature x , we assign it to the class for which the output $y_k = \tilde{w}_k^T \tilde{x}$ is largest. Define a matrix T whose n -th row is the vector t_n^T and together a matrix \tilde{X} whose n -th row is \tilde{x}_n^T , the sum-of-squares error function can be written as

$$J(\tilde{W}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{X} \tilde{W} - T)^T (\tilde{X} \tilde{W} - T) \right\}.$$

Find the closed form solution of \tilde{W} that minimizes the objective function $J(\tilde{W})$. Hint: You may use the following two matrix derivative about trace, $\frac{\partial}{\partial Z} \text{Tr}(AZ) = A^T$ and $\frac{\partial}{\partial Z} \text{Tr}(Z^T AZ) = (A^T + A)Z$.

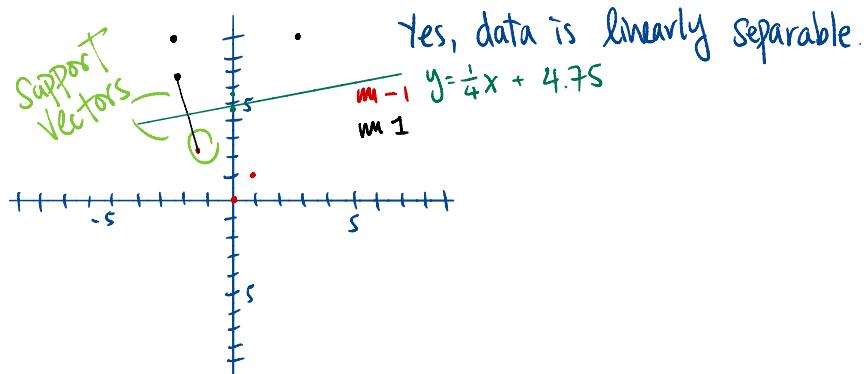
$$\begin{aligned} \frac{\partial J(\tilde{W})}{\partial \tilde{W}} &= \frac{\partial}{\partial \tilde{W}} \left(\frac{1}{2} \text{Tr} \left\{ (\tilde{X} \tilde{W} - T)^T (\tilde{X} \tilde{W} - T) \right\} \right) \\ &= \frac{1}{2} \frac{\partial}{\partial \tilde{W}} \text{Tr} \left\{ (\tilde{X} \tilde{W})^T \tilde{X} \tilde{W} - (\tilde{X} \tilde{W})^T T - T^T (\tilde{X} \tilde{W}) + T^T T \right\} \\ &= \frac{1}{2} \frac{\partial}{\partial \tilde{W}} \text{Tr} \left\{ \tilde{W}^T \tilde{X}^T \tilde{X} \tilde{W} - 2T^T (\tilde{X} \tilde{W}) + T^T T \right\} \\ &= \frac{1}{2} \frac{\partial}{\partial \tilde{W}} \left(\text{Tr}(\tilde{W}^T \tilde{X}^T \tilde{X} \tilde{W}) - \text{Tr}(2T^T (\tilde{X} \tilde{W})) + \text{Tr}(T^T T) \right) \\ &= \frac{1}{2} (2(\tilde{X}^T \tilde{X}) \tilde{W} - 2 \tilde{X}^T T) \\ &= (\tilde{X}^T \tilde{X}) \tilde{W} - \tilde{X}^T T = 0 \end{aligned}$$

$$\boxed{\tilde{W} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T T}$$

4. You are given the following data set which is comprised of $\mathbf{x}^{(i)} \in \mathbb{R}^2$ and $y^{(i)} \in \{-1, 1\}$.

i	$x_1^{(i)}$	$x_2^{(i)}$	$y^{(i)}$
1	-3	9	1
2	-2.5	6.25	1
3	3	9	1
4	-1.5	2.25	-1
5	0	0	-1
6	1	1	-1

(a) Plot the data. Is the data linearly separable?



(b) Look at the data and circle the support vectors by inspection. Find and plot the maximum margin separating hyperplane.

support vectors points: $(-2.5, 6.25), (-1.5, 2.25)$

equation of line between two points:

$$m = \frac{6.25 - 2.25}{-2.5 + 1.5} = -4 \quad y - 6.25 = -4(X + 2.5)$$

$$y = -4X - 10 + 6.25$$

$$y = -4X - 3.75$$

$$\text{perpendicular slope} = \frac{1}{4}$$

$$\text{midpoint: } (-2.5 + 0.5, 2.25 + 2) = (-2, 4.25)$$

$$\text{separating hyperplane: } y - 4.25 = \frac{1}{4}(X + 2)$$

$$y = \frac{1}{4}X + \frac{1}{2} + 4.25$$

$$y = 0.25X + 4.75$$

w vector is perpendicular to separating hyper plane.

So w vector should have slope of -4 , and be pointing in direction of $\begin{bmatrix} -1 \\ 1 \end{bmatrix}$

(c) Find the α_i , w and b in

$$h(x) = \text{sign} \left(\sum_{n \in \mathcal{S}} \alpha_n y^{(n)} x^T x^{(n)} + b \right) = \text{sign} (w^T x + b),$$

where \mathcal{S} is the index set of all support vectors. Do this by solving the dual form of the quadratic program. How is w and b related to your solution in part (b)?

We use dual form of quadratic program:

$$\alpha_n \geq 0 \quad \sum_{n=1}^N \alpha_n = 0$$

$$\sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^N \alpha_n \alpha_k x_n^T x_k$$

$$\text{So, } \alpha_2 = \alpha_4$$

$$\begin{aligned} &= 2\alpha_2 - \frac{1}{2} (\alpha_2^2 \|x_2\|^2 + \alpha_2^2 \|x_4\|^2 - 2\alpha_2^2 x_2^T x_4) \\ &= 2\alpha_2 - \frac{1}{2}\alpha_2^2 \|x_2\|^2 - \frac{1}{2}\alpha_2^2 \|x_4\|^2 + \alpha_2^2 x_2^T x_4 \\ &\stackrel{\partial}{\partial} \alpha_2 (2\alpha_2 - \frac{1}{2}\alpha_2^2 \|x_2\|^2 - \frac{1}{2}\alpha_2^2 \|x_4\|^2 + \alpha_2^2 x_2^T x_4) \\ &= 2 - \alpha_2 \|x_2\|^2 - \alpha_2 \|x_4\|^2 + 2\alpha_2 x_2^T x_4 \end{aligned}$$

$$x_2: (-2.5, 6.25) \quad x_4: (-1.5, 2.25)$$

$$\begin{aligned} &2 - \alpha_2 ((-2.5)^2 + 6.25^2) - \alpha_2 ((-1.5)^2 + 2.25^2) + 2\alpha_2 (17.8125) = 0 \\ &= 2 - 45.3125\alpha_2 - 7.3125\alpha_2 + 35.63\alpha_2 = 0 \end{aligned}$$

$$\boxed{\alpha_4 = \alpha_2 = \frac{2}{17}}$$

We solve for w .

$$w = \sum_{n=1}^N \alpha_n x_n$$

$$= \frac{2}{17}(1) \begin{bmatrix} -2.5 \\ 6.25 \end{bmatrix} - \frac{2}{17} \begin{bmatrix} -1.5 \\ 2.25 \end{bmatrix} = \boxed{\frac{2}{17} \begin{bmatrix} -1 \\ 4 \end{bmatrix}}$$

We solve for b .

$$t_n(w^T x_n + b) = 1$$

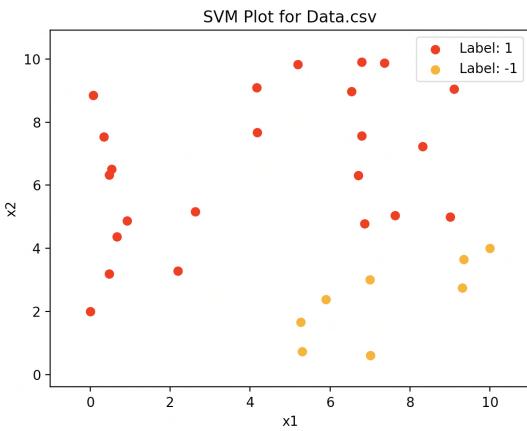
$$1 \left(\left[\begin{smallmatrix} -\frac{2}{17} & \frac{8}{17} \end{smallmatrix} \right] \begin{bmatrix} -2.5 \\ 6.25 \end{bmatrix} + b \right) = 1 \quad \boxed{b = -\frac{38}{17}}$$

w is related to my solution in part b as it is perpendicular to the slope of my separating hyperplane, which is expected as two methods should yield similar w .

b is related to part B in that $\frac{-b}{w_2} = y_{\text{intercept}} = 4.75$. This makes sense as we plot hyperplane by $x_1 w_1 + x_2 w_2 + b = 0$.

5. In this exercise, we will use MATLAB to solve both the primal and the dual problem of SVM. In *Data.csv*, the first two columns contain feature vectors $x^{(i)} \in \mathbb{R}^2$ and the last column contains the label $y^{(i)} \in \{-1, 1\}$. We will use CVX as the optimization solver in this problem. For help with CVX, refer to the [CVX Users' Guide](#). Attach your code for submission. For Python user, feel free to use the following libraries: math, csv, numpy, matplotlib and [cvxpy](#).

- (a) **Visualization** Use different color to plot data with different labels in the 2-D feature space. Is the data linearly separable?



Yes, the data looks
linearly separable.

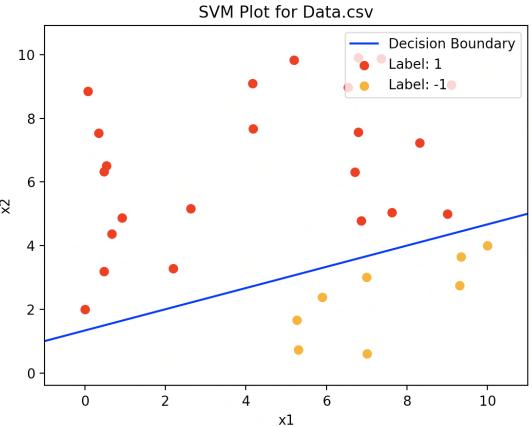
- (b) **The Primal Problem** Use CVX to solve the primal problem of this form:

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, m \end{aligned}$$

Report w and b . Plot the hyperplane defined by w and b .

$$W = \begin{bmatrix} -0.5 \\ 1.5 \end{bmatrix} \quad b = -2$$

Plot ↗



- (c) **The Dual Problem** Use CVX to solve the dual problem of this form:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i, i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0. \end{aligned}$$

Use the resulting a to identify the support vectors on the plot. Report you non-zero α_i 's. How many support vectors do you have? Circle those support vectors.

Note: The latter part of $W(a)$ is in quadratic form, i.e., $a^T P a$. To use CVX, first find P and then use *quad_form(a, P)*. For Python user, you will need to add a small number to the diagonal of P matrix to make cvxpy work. i.e. Run the following code before using cvxpy: “ $P += 1e-13 * \text{numpy.eye}(31)$ ”, where 31 is the total number of data. Also, assume it is 0 if a number is less than $1e-9$.

I found 4 support vectors.

$$\alpha_{28} = 1.03178 \quad \alpha_{29} = 0.21821 \quad \alpha_{30} = 0.26059$$

$$\alpha_{31} = 0.98941$$

