

HOMWORK 3

Saikumar Yadugiri
9083079468, saikumar@cs.wisc.edu

Instructions: Use this latex file as a template to develop your homework. Submit your homework on time as a single pdf file to Canvas. Late submissions may not be accepted. Please wrap your code and upload to a public GitHub repo, then attach the link below the instructions so that we can access it. You can choose any programming language (i.e. python, R, or MATLAB). Please check Piazza for updates about the homework.

GitHub Link: https://github.com/saikumarysk/cs760_hw3

1 Questions (50 pts)

1. (9 pts) Explain whether each scenario is a classification or regression problem. And, provide the number of data points (n) and the number of features (p).
 - (a) (3 pts) We collect a set of data on the top 500 firms in the US. For each firm we record profit, number of employees, industry and the CEO salary. We are interested in predicting CEO salary with given factors.
As CEO's salary can be continuous real value, this is a regression problem. Here, $n = 500$ and $p = 3$.
 - (b) (3 pts) We are considering launching a new product and wish to know whether it will be a success or a failure. We collect data on 20 similar products that were previously launched. For each product we have recorded whether it was a success or failure, price charged for the product, marketing budget, competition price, and ten other variables.
This is a binary classification problem with $n = 20$ and $p = 13$.
 - (c) (3 pts) We are interested in predicting the % change in the US dollar in relation to the weekly changes in the world stock markets. Hence we collect weekly data for all of 2012. For each week we record the % change in the dollar, the % change in the US market, the % change in the British market, and the % change in the German market.
As % change is continuous real value, this can be thought of as a regression problem. Here, $n = 52$ (for each week of the year) and $p = 3$.
2. (6 pts) The table below provides a training data set containing six observations, three predictors, and one qualitative response variable.

X_1	X_2	X_3	Y
0	3	0	Red
2	0	0	Red
0	1	3	Red
0	1	2	Green
-1	0	1	Green
1	1	1	Red

Suppose we wish to use this data set to make a prediction for Y when $X_1 = X_2 = X_3 = 0$ using K-nearest neighbors.

- (a) (2 pts) Compute the Euclidean distance between each observation and the test point, $X_1 = X_2 = X_3 = 0$.
The euclidean distance for (0, 0, 0) for each row from top to bottom are 3, 2, $\sqrt{10}$, $\sqrt{5}$, $\sqrt{2}$, $\sqrt{3}$ respectively.

- (b) (2 pts) What is our prediction with $K = 1$? Why?
 The nearest neighbor with euclidean distance $\sqrt{2}$ is $(-1, 0, 1)$ whose label is Green. So, we predict Green.
- (c) (2 pts) What is our prediction with $K = 3$? Why?
 The 3 nearest neighbors with euclidean distances $\sqrt{2}$, $\sqrt{3}$, and 2 are $(-1, 0, 1)$, $(1, 1, 1)$, and $(2, 0, 0)$ respectively. Their predictions are Green, Red, and Red respectively. So, we predict **Red** as it has the highest frequency.
3. (12 pts) When the number of features p is large, there tends to be a deterioration in the performance of KNN and other local approaches that perform prediction using only observations that are near the test observation for which a prediction must be made. This phenomenon is known as the curse of dimensionality, and it ties into the fact that non-parametric approaches often perform poorly when p is large.
- (a) (2pts) Suppose that we have a set of observations, each with measurements on $p = 1$ feature, X . We assume that X is uniformly (evenly) distributed on $[0, 1]$. Associated with each observation is a response value. Suppose that we wish to predict a test observation's response using only observations that are within 10% of the range of X closest to that test observation. For instance, in order to predict the response for a test observation with $X = 0.6$, we will use observations in the range $[0.55, 0.65]$. On average, what fraction of the available observations will we use to make the prediction?
 Since the values of X come from a uniform distribution, for every $x \in [0.05, 0.95]$, on an average we will look at 10% of the available observations to make a prediction.
- (b) (2pts) Now suppose that we have a set of observations, each with measurements on $p = 2$ features, X_1 and X_2 . We assume that predict a test observation's response using only observations that (X_1, X_2) are uniformly distributed on $[0, 1] \times [0, 1]$. We wish to be within 10% of the range of X_1 and within 10% of the range of X_2 closest to that test observation. For instance, in order to predict the response for a test observation with $X_1 = 0.6$ and $X_2 = 0.35$, we will use observations in the range $[0.55, 0.65]$ for X_1 and in the range $[0.3, 0.4]$ for X_2 . On average, what fraction of the available observations will we use to make the prediction?
 Since X_1 and X_2 are modelled as i.i.d from uniform distribution, we will see $0.1 \times 0.1 = 1\%$ of all observations from $[0, 1] \times [0, 1]$ to make a prediction for every $(x_1, x_2) \in [0.05, 0.95] \times [0.05, 0.95]$.
- (c) (2pts) Now suppose that we have a set of observations on $p = 100$ features. Again the observations are uniformly distributed on each feature, and again each feature ranges in value from 0 to 1. We wish to predict a test observation's response using observations within the 10% of each feature's range that is closest to that test observation. What fraction of the available observations will we use to make the prediction?
 We will see $(0.1)^{100}$ fraction of available observations from $[0, 1]^{100}$ to make a prediction for every $(x_1, \dots, x_{100}) \in [0.05, 0.95]^{100}$.
- (d) (3pts) Using your answers to parts (a)–(c), argue that a drawback of KNN when p is large is that there are very few training observations “near” any given test observation.
 As p increases, the fraction of neighborhood we observe to make a prediction also decreases exponentially by 10^{-p} (for appropriate inputs). So, we will use very few observations near the test observation.
- (e) (3pts) Now suppose that we wish to make a prediction for a test observation by creating a p -dimensional hypercube centered around the test observation that contains, on average, 10% of the training observations. For $p = 1, 2$, and 100, what is the length of each side of the hypercube? Comment what happens to the length of the sides as $\lim_{p \rightarrow \infty}$.
 Assuming that all the training data points are still uniformly distributed from $[0, 1]^p$, we want the volume of the hypercube to $0.1 \times$ volume of the hypercube formed by $[0, 1]^p = 0.1$. Let the side of the p -dimensional hypercube be s . In this case, the volume of is s^p . So, we need $s^p = 0.1$. That is, $s = 0.1^{\frac{1}{p}}$. For $p = 1, 2$, and 100, the length of the side of the p -dimensional hypercube is 0.1, 0.31623, and $10^{-\frac{1}{100}}$ respectively. Moreover as $p \rightarrow \infty$, $\lim_{p \rightarrow \infty} s = \lim_{p \rightarrow \infty} (0.1)^{\frac{1}{p}} = 1$.
4. (6 pts) Suppose you trained a classifier for a spam detection system. The prediction result on the test set is summarized in the following table.

		Predicted class	
		Spam	not Spam
Actual class	Spam	8	2
	not Spam	16	974

Calculate

- (a) (2 pts) Accuracy **0.982**.
 - (b) (2 pts) Precision **0.333**.
 - (c) (2 pts) Recall **0.8**.
5. (9pts) Again, suppose you trained a classifier for a spam filter. The prediction result on the test set is summarized in the following table. Here, "+" represents spam, and "-" means not spam.

Confidence positive	Correct class
0.95	+
0.85	+
0.8	-
0.7	+
0.55	+
0.45	-
0.4	+
0.3	+
0.2	-
0.1	-

- (a) (6pts) Draw a ROC curve based on the above table.

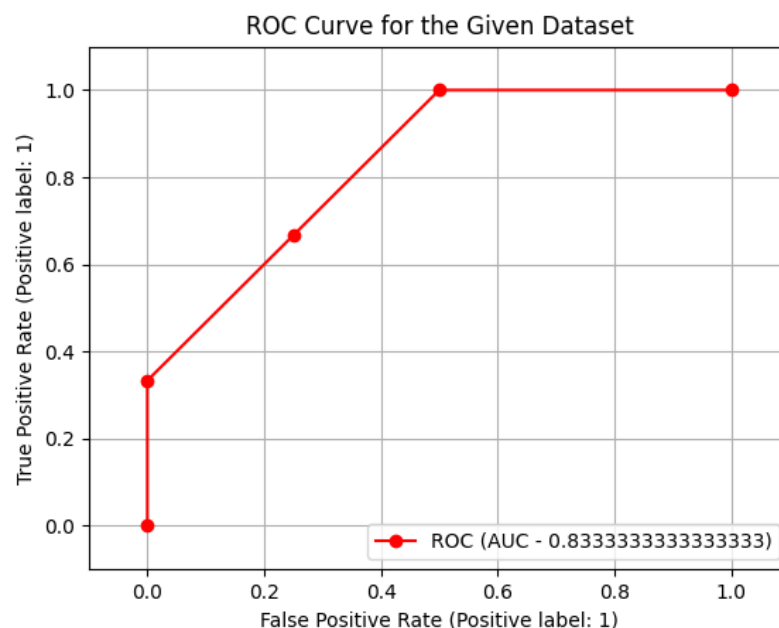


Figure 1: ROC Curve for the Given Data

- (b) (3pts) (Real-world open question) Suppose you want to choose a threshold parameter so that mails with confidence positives above the threshold can be classified as spam. Which value will you choose? Justify your answer based on the ROC curve.

For spam classification, I would choose the threshold $c = 0.4$. This is because using the the above curve, (TPR = 0.6, FPR = 0.2) looks like a really good operating point. It has low false positive rate and high enough recall. Also, it is okay to have a few false positives for spam classification as it is not a critical task such as mushroom edibility classification.

6. (8 pts) In this problem, we will walk through a single step of the gradient descent algorithm for logistic regression. As a reminder,

$$\hat{y} = f(x, \theta)$$

$$f(x; \theta) = \sigma(\theta^\top x)$$

$$\text{Cross entropy loss } L(\hat{y}, y) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})]$$

$$\text{The single update step } \theta^{t+1} = \theta^t - \eta \nabla_{\theta} L(f(x; \theta), y)$$

- (a) (4 pts) Compute the first gradient $\nabla_{\theta} L(f(x; \theta), y)$.

Given $L(f(x; \theta), y) = -[y \log(f(x; \theta)) + (1 - y) \log(1 - f(x; \theta))] = y \log(1 + e^{-\theta^\top x}) + (1 - y) \log(1 + e^{\theta^\top x})$. So, consider for some $j \in \{0, \dots, d\}$,

$$\begin{aligned} \frac{\partial L(f(x; \theta), y)}{\partial \theta_j} &= \frac{y}{1 + e^{-\theta^\top x}} \cdot e^{-\theta^\top x} \cdot -x_j + \frac{1 - y}{1 + e^{\theta^\top x}} \cdot e^{\theta^\top x} \cdot x_j \\ &= -y \cdot (1 - \hat{y}) \cdot x_j + (1 - y) \cdot \hat{y} \cdot x_j \\ &= x_j \cdot (-y + y \cdot \hat{y} + \hat{y} - y \cdot \hat{y}) \\ &= (\hat{y} - y) \cdot x_j \end{aligned}$$

Hence, $\nabla_{\theta} L(f(x; \theta), y) = (\hat{y} - y)x$

- (b) (4 pts) Now assume a two dimensional input. After including a bias parameter for the first dimension, we will have $\theta \in \mathbb{R}^3$.

$$\text{Initial parameters : } \theta^0 = [0, 0, 0]$$

$$\text{Learning rate } \eta = 0.1$$

$$\text{data example : } x = [1, 3, 2], y = 1$$

Compute the updated parameter vector θ^1 from the single update step.

$$\begin{aligned} \hat{y} = f(x; \theta) &= f\left(\begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}\right) = \sigma(0) = 0.5. \text{ So, the gradient } \nabla_{\theta} L(f(x; \theta), y) = (\hat{y} - y)x = \\ &= \begin{bmatrix} -0.5 \\ -1.5 \\ -1 \end{bmatrix}. \text{ Hence, } \theta^1 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} -0.5 \\ -1.5 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.05 \\ 0.15 \\ 0.1 \end{bmatrix}. \end{aligned}$$

2 Programming (50 pts)

1. (10 pts) Use the whole D2z.txt as training set. Use Euclidean distance (i.e. $A = I$). Visualize the predictions of 1NN on a 2D grid $[-2 : 0.1 : 2]^2$. That is, you should produce test points whose first feature goes over $-2, -1.9, -1.8, \dots, 1.9, 2$, so does the second feature independent of the first feature. You should overlay the training set in the plot, just make sure we can tell which points are training, which are grid.

The expected figure looks like this.

The required plot is provided in figure 2.

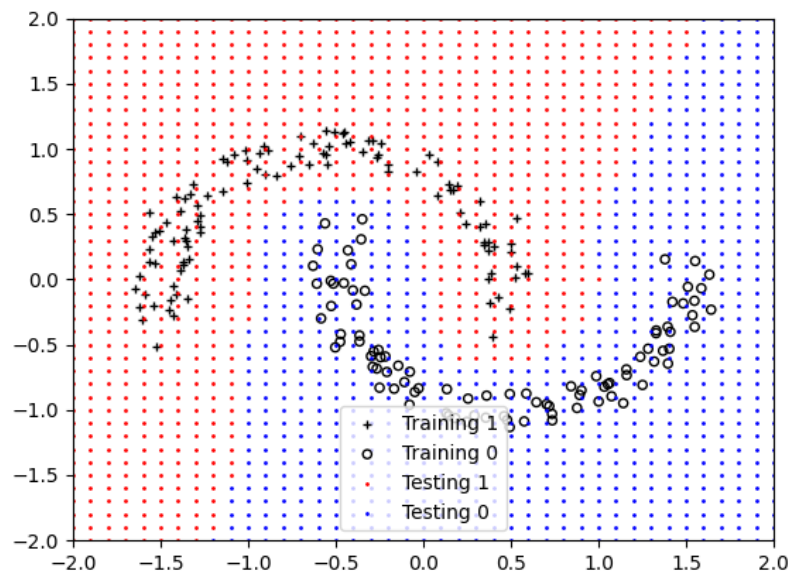


Figure 2: Scatter Plot and Predictions for D2z.txt

Spam filter Now, we will use 'emails.csv' as our dataset. The description is as follows.

- Task: spam detection
 - The number of rows: 5000
 - The number of features: 3000 (Word frequency in each email)
 - The label (y) column name: 'Predictor'
 - For a single training/test set split, use Email 1-4000 as the training set, Email 4001-5000 as the test set.
 - For 5-fold cross validation, split dataset in the following way.
 - Fold 1, test set: Email 1-1000, training set: the rest (Email 1001-5000)
 - Fold 2, test set: Email 1000-2000, training set: the rest
 - Fold 3, test set: Email 2000-3000, training set: the rest
 - Fold 4, test set: Email 3000-4000, training set: the rest
 - Fold 5, test set: Email 4000-5000, training set: the rest
2. (8 pts) Implement 1NN, Run 5-fold cross validation. Report accuracy, precision, and recall in each fold.
- The required details are provided in table 1.

Fold Index	Accuracy	Precision	Recall
1	0.825	0.65449	0.81754
2	0.853	0.68571	0.86643
3	0.862	0.72121	0.83803
4	0.851	0.71642	0.81633
5	0.775	0.60574	0.75817

Table 1: Accuracy, Precision, and Recall for 5-Fold Cross Validation for 1NN learning method

3. (12 pts) Implement logistic regression (from scratch). Use gradient descent (refer to question 6 from part 1) to find the optimal parameters. You may need to tune your learning rate to find a good optimum. Run 5-fold cross validation. Report accuracy, precision, and recall in each fold.

The required details are provided in table 2.

Fold Index	Accuracy	Precision	Recall
1	0.94	0.90036	0.88772
2	0.968	0.92683	0.96029
3	0.937	0.94378	0.82746
4	0.949	0.92042	0.90476
5	0.927	0.85196	0.92157

Table 2: Accuracy, Precision, and Recall for 5-Fold Cross Validation for Logistic Regression learning method

4. (10 pts) Run 5-fold cross validation with kNN varying k ($k=1, 3, 5, 7, 10$). Plot the average accuracy versus k , and list the average accuracy of each case.

Expected figure looks like this.

The required plot is provided in figure 3.

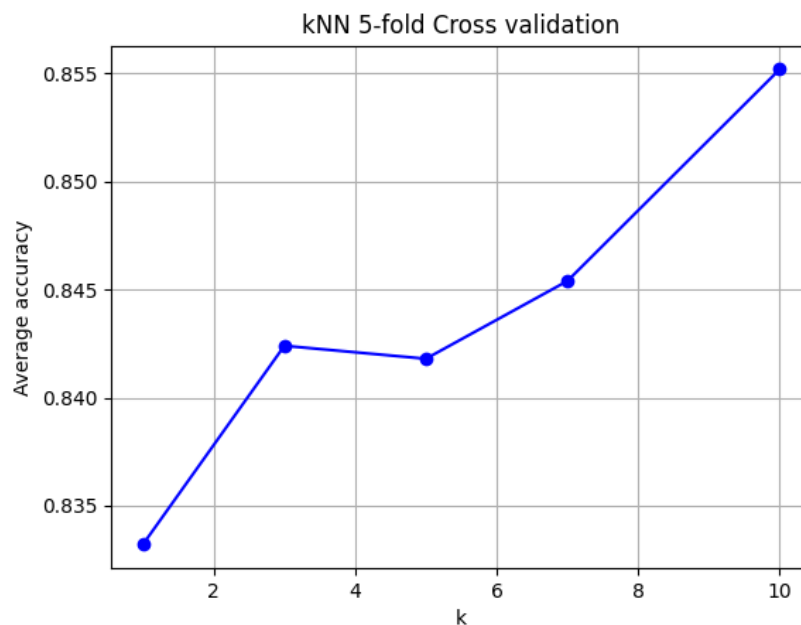


Figure 3: Average Accuracy vs. K for 5-Fold Cross Validation for k-NN learning method

5. (10 pts) Use a single training/test setting. Train kNN ($k=5$) and logistic regression on the training set, and draw ROC curves based on the test set.

Expected figure looks like this.

The required plot is provided in figure 4.

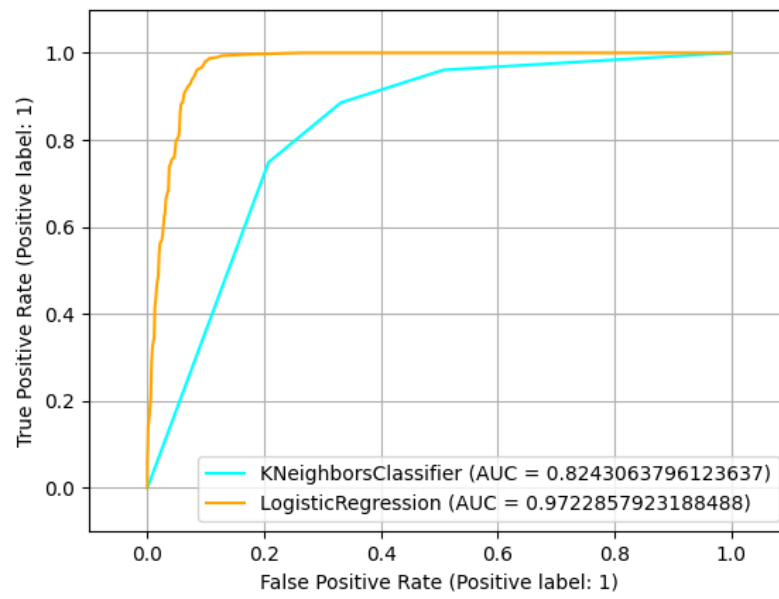


Figure 4: ROC for 5-NN and Logistic Regression Learning Methods