

COSC522 HOMEWORK 4

STUDENT NAME: JOSEPH OCTHE AGADA

NET ID: JAGADA

①

PROBLEM 1

(1) I will present details of calculation here while implementation will be in the notebook attached to this submission. I will start from Minimum Euclidean Distance classifier (MD) to Support Vector Machine (SVM) and finally to perceptron.

MD

The steps involve in plotting the decision boundary for MD are:

- calculate the centroids
- determine the midpoint of the centroids
- calculate the boundary line

I will calculate the centroids as follow:

AND Gate has the following 4 points

1. (0, 0) with class label of 0

2. (0, 1) with class label of 0

3. (1, 0) with class label of 0

4. (1, 1) with class label of 1

Centroid for class 0

$$x\text{-coordinate} = \frac{0+0+1}{3} = \frac{1}{3}$$

$$y\text{-coordinate} = \frac{0+1+0}{3} = \frac{1}{3}$$

The centroid for class 0 is $(\frac{1}{3}, \frac{1}{3})$

Class 1 has only one point (1, 1) hence the centroid for class 1 is (1, 1).

Next step is to determine the midpoint. Midpoint for the 2 centroid is

$$M = \left(\frac{\frac{1}{3} + 1}{2}, \frac{\frac{1}{3} + 1}{2} \right) = \left(\frac{2}{3}, \frac{2}{3} \right)$$

The last step is to calculate the decision boundary line.

The decision boundary is a line perpendicular to the line joining the centroids. If the centroids are C_0 and C_1 then we can use the vector $(C_0 C_1)$ to determine the slope and direction of the decision boundary.

If the vector

$$\vec{J} = \left(1 - \frac{1}{3}, 1 - \frac{1}{3} \right) = \left(\frac{2}{3}, \frac{2}{3} \right)$$

Since the decision boundary is perpendicular to \vec{J} , and slope of $\vec{J} = 1$, then the slope of the decision boundary is -1 .

Hence the equation of the decision boundary is

$$y - \frac{2}{3} = -1(x - \frac{2}{3})$$

The decision boundary is

$$y = -x + \frac{4}{3}$$

The plot for the decision boundary is in the note book attached

SVM

To calculate the decision boundary of SVM, we determine the equation of the hyperplane that maximizes the distance between the

hyperplane that maximizes the distance between(3)
the hyperplane and the closest support vector.
As already known,

for AND gate

class 0 has points $(0, 0), (0, 1), (1, 0)$

class 1 has point $(1, 1)$

The objective is to find a line that maximizes
the distance between the line and the
closest points (Support Vector(s)). Let the line
be

$$w_1x_1 + w_2x_2 + b = 0$$

Where $w = [w_1, w_2]$ is the weight and b is the
bias of the hyperplane. For SVM, the classes are
labeled as -1 and 1 hence we set the constraint
as follow to ensure that each datapoint is on
the correct side of the boundary. the constraint
are:

$$\text{for class 0 } w_i x_i + b \leq -1$$

$$\text{for class 1 } w_i x_i + b \geq 1$$

Since the samples are linearly separable, the
objective is to minimize $\|w\|^2$ subject to

$$y_i(w_i x_i + b) \geq 1 + \epsilon$$

where y_i is the class label for each sample
($y_i = -1$ for class 0 and 1 for class 1)

We solve this optimization problem using
Lagrangian for optimization as follow:

- Identify the Support Vectors
- calculate w and b based on the constraint

Suppose we have a solution $w = [1, 1]$, we
solve for b point $(1, 1)$ with label +1

$$wx+b \geq 1 \Rightarrow (1)(1) + (1)(1) + b \geq 1 \\ \Rightarrow b \geq -1$$

Support vector machine decision boundary in this case becomes

$$x_1 + x_2 - 1 \geq 0$$

With this boundary, the margin boundaries are

$$x_1 + x_2 = 0 \text{ for class } 0 (-1 \text{ side of the boundary})$$

$$x_1 + x_2 = 2 \text{ for class } 1 (+1 \text{ side of the boundary})$$

These are parallel lines equidistant from the decision boundary illustrating the maximum boundaries between classes.

This is the explanation for SVM.

Implementation is in the code notebook attached.

Perceptron

Given some optimal weights w_1, w_2, w_0 , the decision boundary for perceptron is gotten by setting

$$w_1 x_1 + w_2 x_2 + w_0 = 0$$

To calculate the decision boundary, I will determine the optimal weights for the AND gate using online learning with initial weights of

$$w_1 = 0.1$$

$$w_2 = 0.2$$

$$w_0 = 0.5$$

and learning rate

$$\eta = 1$$

Other initialization parameters may give a different decision boundary. This is reason why decision boundary for perceptron is not unique.

Epoch 1

Sample 1: $(0, 0)$ target $= 1$

$$Z = W_1x_1 + W_2x_2 + W_0 = 0.1x_0 + 0.2x_0 + 0.5 = 0.5$$

$$Z = 0.5 > 0 \text{ hence } \hat{Y} = 1$$

but $Y = 0$, therefore we update weight as follow:

$$\begin{pmatrix} W_1 \\ W_2 \\ W_0 \end{pmatrix} = \begin{pmatrix} 0.1 \\ 0.2 \\ 0.5 \end{pmatrix} + 1(0-1) \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0.1 \\ 0.2 \\ -0.5 \end{pmatrix}$$

Sample 2: $(0, 1)$; target $= 0$

$$Z = W_1x_1 + W_2x_2 + Z = 0.1x_0 + 0.2x_1 - 0.5 = -0.3$$

$$Z = -0.3 < 0 \text{ hence } \hat{Y} = 0 \text{ and } Y = 0$$

Therefore, no weight update

Sample 3: $(1, 0)$ target $= 0$

$$Z = W_1x_1 + W_2x_2 + W_0 = 0.1x_1 + 0.2x_0 - 0.5 = -0.4$$

$Z = -0.4 < 0$ hence $\hat{Y} = 0$ and $Y = 0$, therefore no weight update

Sample 4 $(1, 1)$ target $= 0$

$$Z = W_1x_1 + W_2x_2 + W_0 = 0.1x_1 + 0.2x_1 - 0.5 = -0.2$$

$Z = -0.2 < 0$ hence $\hat{Y} = 0$ and $Y = 1$ hence we update weight as follow:

$$\begin{pmatrix} W_1 \\ W_2 \\ W_0 \end{pmatrix} = \begin{pmatrix} 0.1 \\ 0.2 \\ -0.5 \end{pmatrix} + (1-0) \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1.1 \\ 1.2 \\ 0.5 \end{pmatrix}$$

Epoch 2 will not result in any weight change hence final weights are $W_1 = 1.1$, $W_2 = 1.2$, $W_0 = 0.5$

Hence the decision boundary is

$$1.1x_1 + 1.2x_2 + 0.5 = 0$$

Reason Why Decision Boundaries for perceptron are Not Unique

Decision boundary for perceptron is not unique for the following reasons

- Dependency on Initial Weights and Order of data presentation — the final weight of a perceptron depend on the Initial Weights and for different set of Weights, there is a different decision boundary. Hence decision boundary largely depend on Initialization strategy/initial weights used.
- Non-Unique Solution in linearly Separable Data When sample is linearly separable, there are infinitely many possible hyperplanes that separate the classes. The perceptron algorithm stop as soon as it finds a solution, but other hyperplanes could also correctly classify the sample. Thus, depending on the Initialization strategy/Weights and the specific order in which data points are presented, perceptron could converge to different solutions.
- Non Convex Solution Space - For perceptron, the optimization problem is non-convex, which means there are multiple solutions that could satisfy same objective. Since perceptron algorithm does not seek to minimize an objective

function globally but only locally, it can find any feasible solution leading to non-unique boundaries.

PROBLEM 1

- (2) Here, I calculated the projections for FLD and PCA, and applied Minimum Euclidean to the projections to arrive a new set of decision boundaries for FLD + MD and PCA + RID (Recall that MD = Minimum Euclidean Distance). Details of the calculations are as follow:

for class 0

$$M_0 = \left(\frac{0+0+\beta}{3}, \frac{0+1+0}{3} \right) = \left(\frac{1}{3}, \frac{1}{3} \right)$$

for class 1

$$M_1 = (1, 1)$$

overall mean

$$M = \left(\frac{0+0+1+1}{4}, \frac{0+1+0+1}{4} \right) = \left(\frac{1}{2}, \frac{1}{2} \right)$$

Within class Scatter matrix (S_W)

for class 0

$$S_W^{(0)} = \sum_{i=1}^3 (\mathbf{x}_i - M_0)(\mathbf{x}_i - M_0)^T$$

for class 1

$$S_W^{(1)} = \sum_{i=1}^3 (\mathbf{x}_i - M_1)(\mathbf{x}_i - M_1)^T$$

Between class scatter matrix

$$S_B = n_0 (M_0 - M)(M_0 - M)^T + n_1 (M_1 - M)(M_1 - M)^T$$

Next step is eigen decomposition by solving for ⑧ the eigen values and eigen vectors using $S^{-1}S_B$ to find the direction that maximize separation (For FLD). The eigen vector corresponding to the highest eigen value becomes the weight projection weight (W_{FLD})

PCA Projection

For PCA projection, we define

$$X_{\text{centered}} = X - \bar{M}$$

We compute the covariance matrix

$$\text{Cov}(X) = \frac{1}{n-1} X_{\text{centered}}^T X_{\text{centered}}$$

Next step is to calculate the eigen values and eigen vectors of $\text{Cov}(X)$ so that the eigen vector corresponding to the the highest eigen value become the projection weight for PCA (W_{PCA})

Projections for FLD and PCA?

$$X_{FLD} = X \cdot W_{FLD}$$

$$X_{PCA} = X \cdot W_{PCA}$$

The projection was implemented in the R notebook attached.

PROBLEM 2

(9)

a. Causes of High precision Low Recall

$$\text{First, precision} = \frac{\text{TP}}{\text{TP+FN}} \quad \text{Recall} = \frac{\text{TP}}{\text{TP+FN}}$$

Here, high precision itself is not a problem but it becomes a problem if it is at the expense of low precision because it implies that the model is missing or failing to identify many positive classes. possible causes are

→ Imbalanced Training Sample - When training sample is imbalanced with the positive class being the minority class, recall is expected to ~~be~~ be low because the model, even though may perform well generally, it will tend to misclassify the positive class leading to low recall.

→ Classification threshold - the model may be using a high ~~th~~ threshold for classifying instance as positive which reduces false positive, hence high precision but also leads to missing many ~~positive~~ actual positive cases implying low recall.

→ overfitting - the model might have been over fine tune to perform well on positive class leading to predicting positive only when it is very ~~sure~~ thereby certain thereby missing many actual positive implying low recall.

b. Low precision and high Recall

(10)

Low precision means the model has high rate of false positive meaning many samples predicted positive are actually negative, while high recall indicates that it identifies most of the positive cases correctly, possible causes are:

- Lax classification threshold - the model might employ a low threshold for classifying instances as positive leading to too many predictions being marked as positive which increase false positive and lower precision.
- the goal may be to prioritize identifying the positive class (for example in the case of medical diagnosis), this will lead to high recall while precision will be intentionally sacrificed.
- Complex or Noisy Data - if the data set is noisy or have overlapping classes, the model may find it easier to classify a larger number of instances as positive.

PROBLEM 3

Hrs (x_k)	1.0	2.0	3.0	4.0	5.0
Prob (P_k)	0.07	0.26	0.61	0.87	0.97
Pass (y_k)	0	0	0	1	1

Q. To build a model that will predict probability of pass, I will model the relationship between x_k and P_k using linear basis function where P_k is treated as the dependent variable while x_k is the independent variable. The general regression equation is

$$P_k(w, x_k) = \sum_{j=0}^{n-1} w_j \phi_j(x_k)$$

where $\phi_j(x_k)$ is a linear basis function, for this particular problem, the regression equation can be reduced to

$$P_k = w_0 + w_1 x_k$$

Where w_0 and w_1 can be derived to be

$$\hat{w}_1 = \frac{\sum_k (x_k - \bar{x})(P_k - \bar{P})}{\sum (x_k - \bar{x})^2}$$

$$\hat{w}_0 = \bar{P} - \hat{w}_1 \bar{x}$$

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

(12)

$$\hat{P} = \frac{0.07 + 0.26 + 0.61 + 0.87 + 0.97}{5} = 0.556$$

To calculate w_1 , we first calculate (x_{k-2}) and $(P_k - \hat{P})$ for each pair

for $x_1 = 1$ $P_1 = 0.07$ (ie $k=1$)

$$(x_1 - \bar{x}) = (1 - 3) = -2$$

$$(P_1 - \hat{P}) = (0.07 - 0.556) = -0.486$$

for $K = 2$

$$(x_2 - \bar{x}) = (2 - 3) = -1$$

$$(P_2 - \hat{P}) = (0.26 - 0.556) = -0.296$$

for $K = 3$

$$(x_3 - \bar{x}) = (3 - 3) = 0$$

$$(P_3 - \hat{P}) = (0.61 - 0.556) = 0.054$$

for $K = 4$

$$(x_4 - \bar{x}) = (4 - 3) = 1$$

$$(P_4 - \hat{P}) = (0.87 - 0.556) = 0.314$$

for $K = 5$

$$(x_5 - \bar{x}) = (5 - 3) = 2$$

$$(P_5 - \hat{P}) = (0.97 - 0.556) = 0.414$$

$$\begin{aligned} \text{Hence } \sum_K (x_k - \bar{x})(P_k - \hat{P}) &= (-2)(-0.486) + (-1)(-0.296) \\ &\quad + (0 \times 0.054) + (1)(0.314) + (2)(0.414) \\ &= 0.972 + 0.296 + 0 + 0.314 + 0.828 \\ &= \underline{\underline{-2}} = 2.41 \end{aligned}$$

(13)

hence $\sum_k (x_k - \bar{x})^2 = (-2)^2 + (-1)^2 + (0)^2 + 1^2 + 2^2$
 $= 1^2 + 2^2 + 1^2 + 2^2 = 10$

Hence

$$W_1 = \frac{2 \cdot 41}{10} = \underline{\underline{0.241}}$$

$$W_0 = \bar{P} - W_1 \bar{x} = 0.556 - 0.241 \times 3 = -0.167$$

the linear regression model becomes

$$\boxed{P_k = -0.167 + 0.241 x_k}$$

When $x = 3.27$

$$P_k = -0.167 + 0.241(3.27)$$

$$P_k = \underline{\underline{0.62}}$$

b: the general regression equation is

$$y(W, x) = \sum_{j=0}^{m-1} W_j \phi_j(x)$$

In this particular case

$\phi_j(x)$ is a Sigmoid basis function given as

$$\sigma(x_k) = \frac{1}{1 + e^{-x_k}}$$

This regression problem is equivalent to

$$P(y_k | x_k) = \frac{1}{1 + \exp(-(B_0 + B_1 x_k))}$$

Where B_0 is intercept and

B_1 is slope

β_0 and β_1 are estimated using a combination of maximum likelihood estimation and iterative approach such as Newton-Raphson as (14)

$$\hat{\beta}_0 = -3.748$$

$$\beta_1 = 1.047$$

Hence

$$P(Y_k | x_k) = \frac{1}{1 + \exp[-(-3.748) + 1.047x_k]}$$

A student who studies for 3.25 hours has the following estimated probability of passing.

$$P(Y | x_k = 3.25) = \frac{1}{1 + \exp[-(-3.748) + 1.047(3.25)]}$$

$$= \underline{\underline{0.415}}$$

A student who studies for 3.25 has probability of passing as 0.415