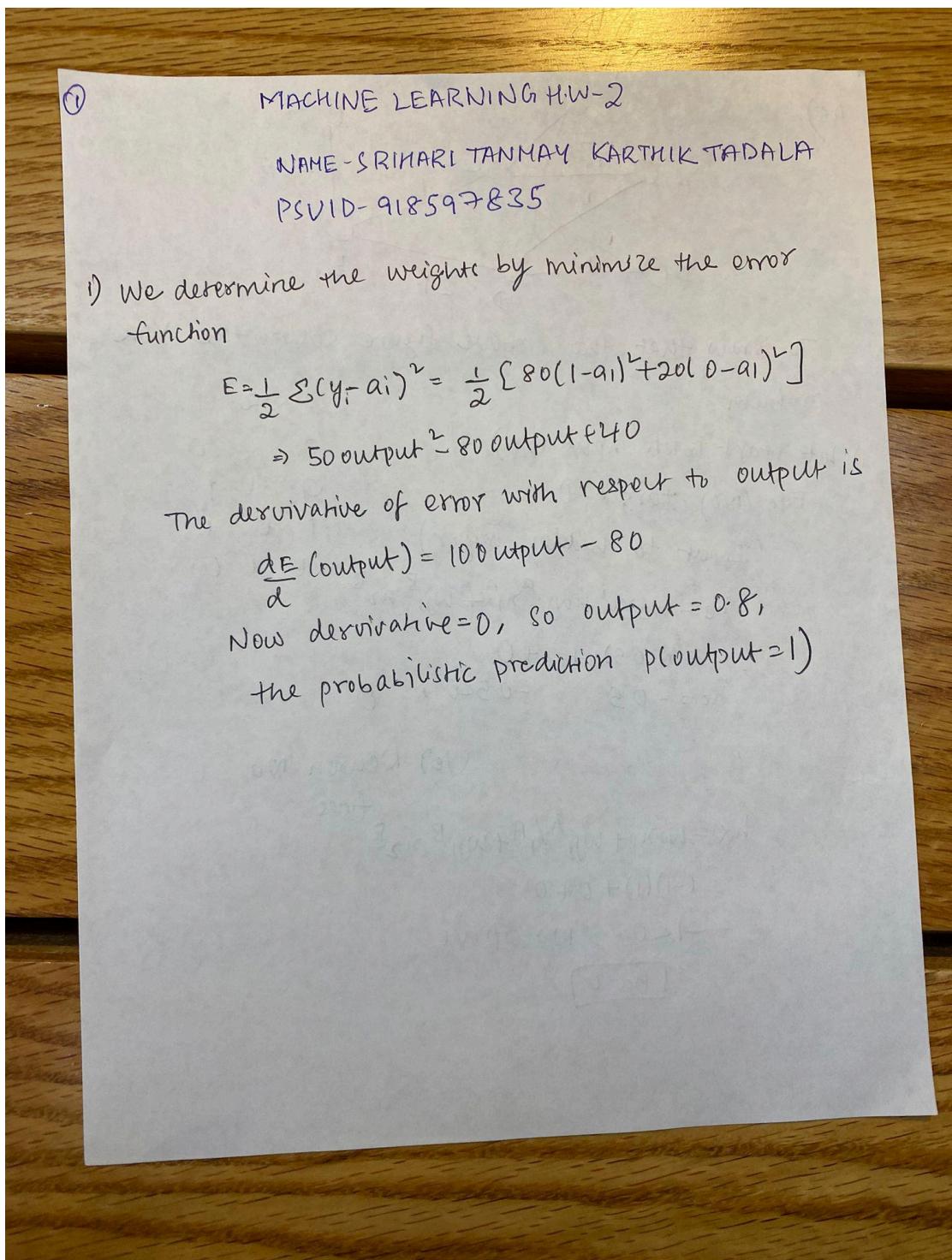


**Machine Learning,
Spring 2024 Homework #2
SRIHARI TANMAY KARTHIK TADALA**

Solution 1 -



Solution 2-

The order in which training samples are presented to a neural network during its training phase is pivotal, as it greatly impacts the learning process. When the network is subjected to the complete set of training examples repeatedly, for example, 1000 times, it gains an advantage by acquiring knowledge from each example during every cycle. This approach guarantees a well-rounded learning experience by incorporating data from the entire dataset.

In contrast, when the network is trained by repeatedly presenting each example 1000 times before moving on to the next, it has a tendency to focus on learning from the current example and neglect the previous ones. The network's concentrated approach can result in the loss of knowledge acquired from previous examples, which is referred to as catastrophic forgetting. Consequently, the network's weights are primarily modified according to the most recent examples, which may compromise its capacity to precisely classify earlier samples.

Hence, the sequence and presentation of training examples significantly impact the network's capacity to acquire and preserve information throughout the dataset.

Solution 3-

- 1) Networks of perceptrons with linear activation functions are unsatisfactory because a single perceptron can replicate the functionality of an entire network with linear activations. The redundancy of the network arises from the fact that both a single perceptron and a network have the ability to compute a weighted sum of inputs, which is the fundamental operation of linear activation.
- 2) When every layer in a network carries out a linear transformation, the overall result, regardless of the number of layers, is the same as applying a single linear transformation to the input. This suggests that increasing the number of layers does not improve the network's capacity to represent intricate patterns that go beyond simple linear relationships.
- 3) The utilization of linear activation functions hinders networks from effectively solving intricate problems that display non-linear connections, such as the XOR problem, image

recognition, or language processing. This constraint arises from their fundamental linear characteristic, which restricts them to linear approximations and makes them inadequate for tasks that require the representation of non-linearities.

4) The utilization of linear activation functions limits the network's capacity to learn, limiting its ability to adapt to complex data patterns and diverse datasets. This greatly limits the range of issues that these networks can effectively address.

Perceptron networks with linear activation functions are considered simplistic because they are naturally unable to capture nonlinear dynamics in data. They have limited utility in tasks that are straightforward and linearly separable, which are relatively rare in practical situations. In order to effectively model complex phenomena and data that display non-linear relationships, it is essential to incorporate non-linear activation functions. This allows neural networks to fully utilize their potential and adaptability.

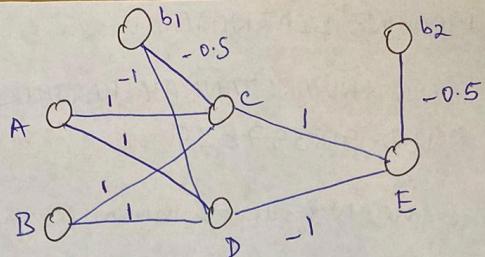
Solution 4-

Overfitting is prone to happen under specific circumstances, including:

- 1) Insufficient Training Data: The probability of overfitting rises as the size of the training dataset decreases. With a reduced number of data points, the model may prioritize memorizing individual instances rather than comprehending general patterns. This can result in poor performance when faced with unfamiliar data, as the model learns the specific noise or irregularities present in the training set.
- 2) Increased Model Complexity: A model that has a significant number of parameters is more susceptible to overfitting. The reason for this is that a complex model possesses the capacity to accurately conform to the training data, capturing not only the fundamental patterns but also the incidental variations and unpredictability. Consequently, although it may perform exceptionally well on the training data, its capacity to apply that knowledge to new data is compromised, resulting in overfitting.

Solution 5-

A5)



To show that the above figure solves the XOR

problem

1) let start with input (0,0)

For (0,0), target = 0

(Input to hidden layer)

$$h_C = w_0 x_0 + w_{j1}^A x_1 + w_{j1}^B x_2$$

$$= (-0.5)1 + 0 + 0$$

$$h_C = -0.5 \quad -0.5 < 0, \quad 0$$

XOR

0	0	0
0	1	1
1	0	1
1	1	0

(ie) Neuron No

fines:

$$h_D = w_0 x_0 + w_{j1}^A x_1 + w_{j1}^B x_2$$

$$= (-1)(1) + 0 + 0$$

$$= -1 < 0 \quad \text{No fines}$$

$$\boxed{h_D = 0}$$

(2)

Output $y(E)$ to hidden layer

$$w_{0n_0} + w_{k_1 j} n_1 c + w_{k_2 j} n_2 D \\ = (-0.5)1 + 0 + 0 = -0.5 < 0 \Rightarrow 0$$

$y(E) = 0$
 Since the output $y(E) = 0$ and target = 0 matches and
 it satisfies for input (0, 0)

2) For input (0, 1) and target = 1

$$h_C = w_{0n_0} b_1 + w_{j_1 i} A_{n_1 A} + w_{j_2 i} B_{n_2 B} \\ = (-0.5)(1) + 0 + 1(1) \\ = -0.5 + 1 \Rightarrow 0.5 \quad \underline{0.5 > 0}$$

<u>$h_C = 1$</u>	<u>Neuron fires</u>
-----------------------------	---------------------

$$h_D = w_{0n_0} b_1 + w_{j_1 i} A_{n_1 A} + w_{j_2 i} B_{n_2 B} \\ = (-1)(1) + 0 + 1 \Rightarrow -1 + 1 \Rightarrow 0 \geq 0$$

<u>$h_D = 0$</u>	<u>Neuron fires</u>
-----------------------------	---------------------

hidden to output layer

$$y(E) = (-0.5) + 1 + (-1)10 \\ = -0.5 + 1 + 0 = +0.5$$

Neuron fires

 $y(E) = 1$ satisfied target

3. For input $(1,0)$ & target 1

$$h_C = (-0.5)(1) + 1(1) + 0$$

$$= -0.5 + 1 \Rightarrow 0.5 > 0 \quad \text{Neuron fires}$$

$$h_C = 1$$

$$h_C = (-1)(1) + (1)(1) + 0 = 0 \leq 0 \quad \text{No fire}$$

$$\boxed{h_d = 0}$$

$$y(E) = (-0.5)(1) + (1)(1) + 0$$

$$= 0.5 \Rightarrow 0.5 > 0$$

$$y(E) = 1$$

hence the output got from input $(1,0)$ is 1

hence the matches the target $= 1$

4. For input $(1,1)$ and target $= 0$

$$h_C = (-0.5)(1) + 1 + 1 \Rightarrow -0.5 + 1 + 1$$

$$\Rightarrow 1.5$$

Neuron fires.

$$h_d = (-1)(1) + 1 + 1$$

$$\Rightarrow -1 + 1 + 1 \Rightarrow 1$$

$$y(E) = (-0.5)(1) + (-1)$$

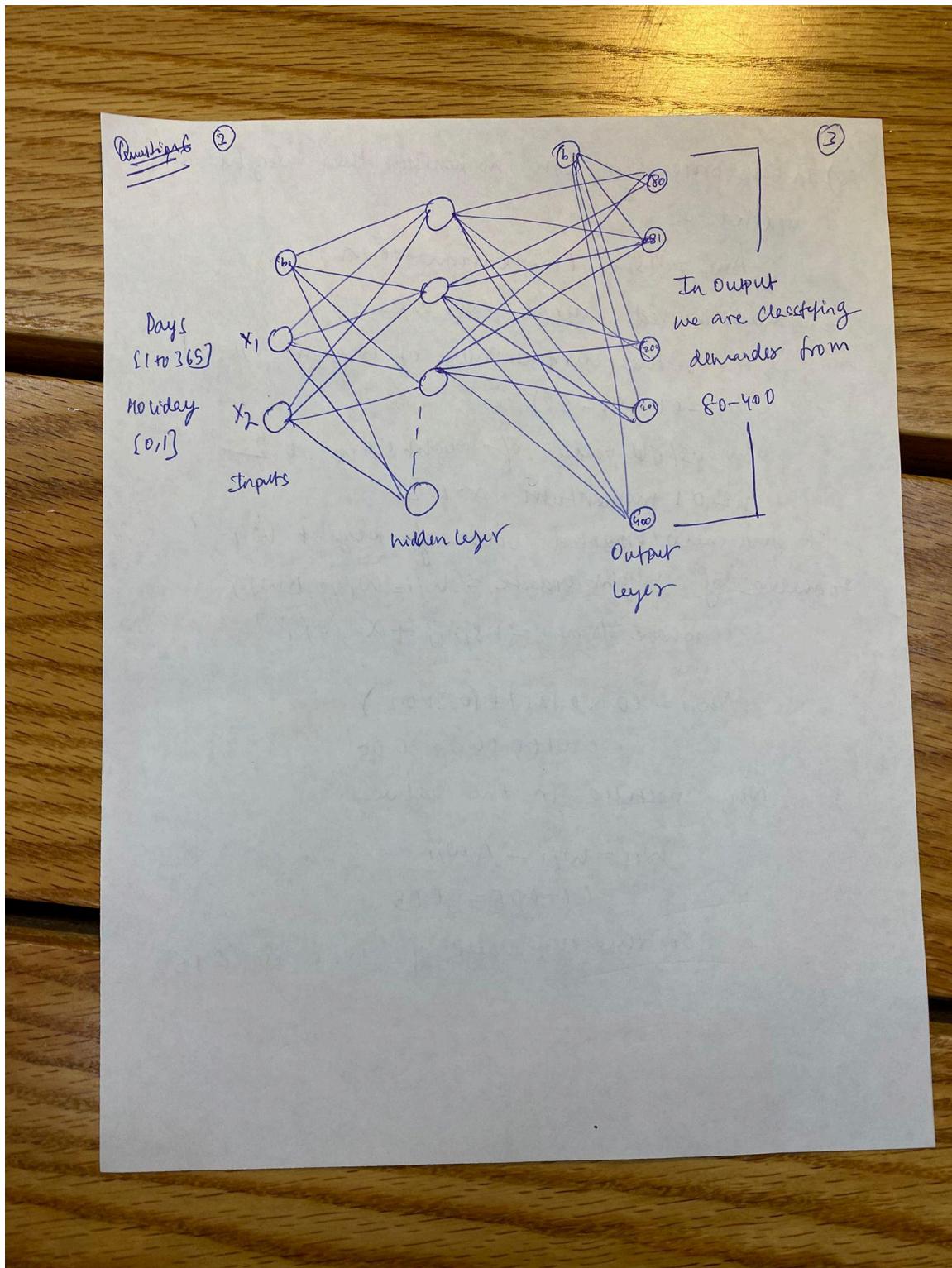
$$= -0.5$$

$$= -0.5 < 0$$

$$\underline{\underline{y(E) = 0}}$$

Hence the output for input $(1,1)$ the output is 0 & target given was also 0 so the target given input satisfies & matches. \therefore Hence the given figure & weight satisfies or solves the problem of XOR

Solution 6-



1. Data Gathering: Collect historical data: Gather the daily demand statistics from the past 5 years, including daily values and any relevant attributes.

2. Data Processing and Normalization: Normalization is essential to streamline the data for effective model training.

- Input Features: Allocate a range of features to the input, including:
 - Days (1 to 365) to encompass a full year, and then duplicate for the following years,

Holiday Indicator: Use a binary input to represent holidays, with 0 indicating no holiday and 1 indicating a holiday.

Additional functionalities can be incorporated as required.

- Hidden Neurons: The ideal number of hidden neurons is determined by experimenting and testing different counts. Initially, the model is configured with 5 concealed layers, and subsequent modifications can be made depending on its performance.

- The number of output neurons ranges from 80 to 400, and they are used to classify electricity demand patterns over the last 5 years. For instance, when the 200th neuron is activated, it indicates the categorization of electricity demand for the upcoming 5 days.

3. Data Division: Divide the dataset into separate sets for training, testing, and validation.

4. Training and Evaluating the Model: Begin by training the model using data from the initial 4-year period. Afterwards, assess the model's accuracy by using the data from the fifth year to evaluate its performance.

5. Optimizing the Learning Rate: It is recommended to use a learning rate within the range of 0.001 to 0.01. Commence with a learning rate of 0.001, modifying it as necessary depending on the training results of the model.

6 2) To optimize your MLP-based system for electricity demand prediction by incorporating weather forecast data, adhere to the following procedure:

- 1) Increase the number of input neurons by incorporating daytime and nighttime temperatures as additional input features, with a temperature range of -50 to 50 degrees Celsius.
 - 2) Standardize Weather Data: Ensure that weather inputs are rescaled to match the existing features in order to achieve optimal performance of the model.
 - 3) Modify the input layer of the model to incorporate these supplementary weather characteristics.
 - 4) Enhance Data Sets: Integrate weather-related inputs into both the training and testing datasets.
 - 5) Enhance Model Training: Utilize the enriched dataset to train the MLP model, incorporating the effects of weather on electricity demand.
6. 3) The success or limitations of a system in predicting electricity consumption are influenced by various factors that affect its effectiveness.

Key Factors for Success:

- Data Integrity: Ensuring the accuracy of historical demand data and weather forecasts, as well as including relevant features, enhances the effectiveness of predictions.
- Model Design: An appropriately configured MLP architecture is essential for detecting complex patterns in electricity demand and the corresponding weather conditions.
- Feature Management: The careful selection and processing of important features, particularly weather information, improves the accuracy of the model.
- Overfitting Prevention: Utilizing regularization techniques aids in mitigating overfitting, thereby preserving the model's ability to generalize.

Possible constraints:

- Unforeseen external factors: Demand changes caused by unexpected events such as natural disasters, policy alterations, or economic fluctuations may be difficult to predict, which can impact the accuracy of the model.
- Timeframe for Forecasting: Multilayer Perceptrons (MLPs) perform exceptionally well in predicting outcomes in the short term, but their accuracy in making long-term projections may be less dependable.
- Aberrant Weather Conditions: The model may struggle to accurately forecast demand in the presence of extreme weather scenarios that deviate significantly from historical data.
- Data Quality Concerns: Insufficient data quality, which includes the absence or inaccuracy of historical and weather data, can have a negative impact on the accuracy of forecasts.
- Overfitting can occur when a model is excessively complex or improperly tuned, leading to inaccurate forecasts.

Solution 7-

1) Data Preprocessing and Normalization: Standardize all inputs by transforming them to have a mean of zero and a standard deviation of one. This process helps in preparing the data for neural network training.

- Input Neuron Configuration: The inputs will consist of a diverse range of features, as specified:

The term "Week" is used to represent a numerical value ranging from 1 to 52, which indicates the specific week of the year. The daytime temperature ranges from -50 to 50 degrees Celsius.

Seasons are divided into four categories: Spring, Summer, Autumn, and Winter, and they need to be properly normalized.

The presence of an epidemic is indicated by a binary input: 1 for presence and 0 for absence.

- Optimal number of hidden neurons: The most suitable number of hidden neurons is determined through experimentation and fine-tuning. Initially, there are 5 concealed layers, but more layers can be investigated if needed.

- Output Neurons: The number of output neurons is determined by the capacity of the geriatric ward, where each neuron may represent a group of 10 beds.

2) Data Division: Following the process of normalization, the dataset should be partitioned into separate groups for Training, Testing, and Validation.

3) Model Training and Evaluation: Commence with data from the initial four years for preliminary training. Afterwards, use the data from the last year to evaluate the accuracy and performance of the model.

4) Learning Rate Adjustment: Begin with initial learning rates ranging from 0.001 to 0.01, and make necessary adjustments to improve training results. The effectiveness of the system depends on various factors, such as the data's quality and quantity, the choice of features, the model's design, and the complexity of the patterns in bed demand.

Solution 8-

(4)

A8) In back propagation with momentum the weights are updated as

$$\Delta w_{ji} = n s_j \alpha_{ji} + \text{momentum term}$$

Q) In multi layered neural network

given us the current value of weight $w_{ji} = 0.1$
 $\alpha_{ji} = 0.1, s_j = 0.2$

Previous/old values of weight $s \Delta w_i = 0.2$

$n = 0.1$ momentum, $\alpha = 0.2$

To find new/updated value of weight w_{ji}' ?

Formulae of weight update $= w_{ji}' = w_{ji} + \Delta w_{ji}'$

where $\Delta w_{ji}' = n s_j \alpha_{ji} + \alpha \Delta w_{ji}^{(t-1)}$

$$\Delta w_{ji} = (0.1 \times 0.1 \times 1) + (0.2 \times 0.2)$$

$$= 0.01 + 0.04 = 0.05.$$

Now substitute in the value

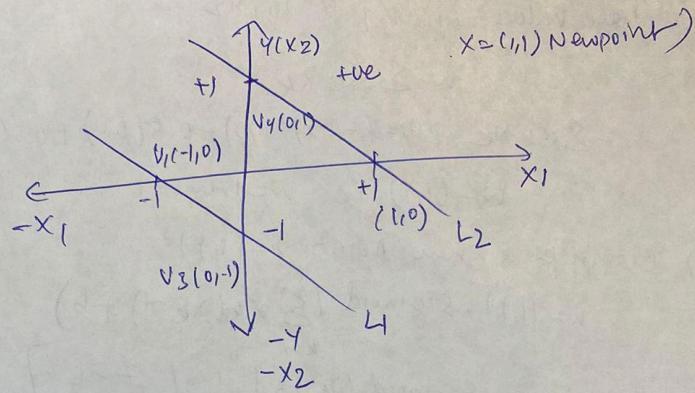
$$w_{ji}' = w_{ji} + \Delta w_{ji}'$$

$$= 0.1 - 0.05 = 0.05$$

So the new value of w_{ji}' is 0.05.

Solution 9 -

A9-a)



Given 4 support values

	x_1	x_2	Class	α
v_1	-1	0	-1	$\alpha_1 = -0.5$
v_2	1	0	1	$\alpha_2 = 0.5$
v_3	0	-1	-1	$\alpha_3 = -0.5$
v_4	0	1	1	$\alpha_4 = 0.5$

$$\text{bias} = 0$$

The equation of a blue line $w_1x_1 + w_2x_2 + b = 0$
Now to find w_1 & w_2 , the formulae to find

vector weights is $w = \sum_{k \in \text{Training examples}} (\alpha_k x_k)$

(4)

We have values of x_1 & n_1, n_2

$$w = \sum x_k n_k$$

$$w = \sum x_k n_k ; w = (-0.5)(-1,0) + 0.5(1,0) + 0.5(0,1)$$

$$w = (1,1) \text{ i.e. } w_1 = 1, w_2 = 1$$

a) Classifying a new point $n = (1,1)$

$$n = (1,1), \text{ sigmoid } \left(\sum_{k=1}^m x_k (n_k n) + b \right)$$

$$\text{where sigmoid}(z) = \begin{cases} 1 & \text{if } z > 0 \\ -1 & \text{if } z \leq 0 \end{cases}$$

$$n(1,1) \Rightarrow \text{sigmoid} \left[[-0.5(-1,0)(1,1)] + [0.5(1,0)(1,1)] + [-0.5(0,1)(1,1)] + [0.5(1,1)(1,1)] \right] + 0$$

$$n(1,1) \Rightarrow \text{sigmoid}[?]$$

$$= 1$$

$$\boxed{\text{sgmd}(z) = 1 \text{ if } z > 0 \\ z \geq 0}$$

This means that new input $n = (1,1)$ is on the positive side of the boundary

$$n(1,1) = 1$$

Q3)
a) b) The equation of separation line is

$$w_1x_1 + w_2x_2 + b = 0$$

putting values of w_1 & w_2

$$1(x_1) + 1(x_2) + 0 = 0$$

$$x_1 + x_2 = 0 \quad | \quad x_2 = -x_1$$

finding the value of x_2 when $x_1 = 0$

$$x_2 = -x_1$$

$$(x_1 = 0, x_2 = 0)$$

putting $x_1 = 0, x_2 = 0$

$$x_1 = 1, x_2 = -1$$

$$x_1 = 2, x_2 = -2$$

$$x_1 = -1, x_2 = 1$$

$$x_1 = -2, x_2 = 2$$

We get the decision boundary on separating line as

so from the above is the decision boundary separation line crossing

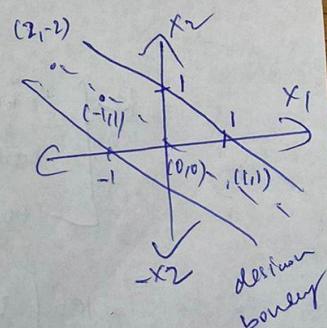
from center when $x_1 = 0, x_2 = 0$

when $x_1 = 1, x_2 = 1$

$$x_2 = 1, x_1 = -1$$

$$x_1 = 2, x_2 = -2$$

$$x_1 = -2, x_2 = 2$$



Solution 10-

Hence the new point $X = (1, 1)$ lie to the positive side of the decision boundary

Q1.8:

Ans(0) To define the kernel function as

$$K(x_i, y_i) = \sum \text{OR}(x_i, y_i) \quad \text{OR copied}$$

$$\text{Given } x_1 = (0, 0, 0, 0)$$

$$x_2 = (1, 1, 1, 1)$$

$$x_3 = (1, 0, 0, 1)$$

For kernel function $K(x_i, y_i) = \sum \text{OR}(x_i, y_i)$

$$\text{Kernel function} = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & K(x_1, x_3) \\ K(x_2, x_1) & K(x_2, x_2) & K(x_2, x_3) \\ K(x_3, x_1) & K(x_3, x_2) & K(x_3, x_3) \end{bmatrix}$$

Here $k = \sum \text{OR}$

$$K = \begin{bmatrix} 0+0+0+0 & 1+1+1+1 & 1+0+0+1 \\ 1+1+1+1 & 1+1+1+1 & 1+1+1+1 \\ 1+0+0+1 & 1+1+1+1 & 1+0+0+1 \end{bmatrix}$$

$$K = \begin{bmatrix} 0 & 4 & 2 \\ 4 & 4 & 4 \\ 2 & 4 & 2 \end{bmatrix}$$

~~Ans~~
this is the final answer

Solution 11 & 12 -

⑥ Ans II)

Let support vectors s_1, s_2, s_3 be

$$s_1 = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}, s_2 = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}, s_3 = \begin{bmatrix} 1 \\ 1.5 \end{bmatrix}$$

Adding bias (b) to it would be

$$s_1 = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}, s_2 = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}, s_3 = \begin{bmatrix} 1 \\ 1.5 \end{bmatrix}$$

We know

$$\text{-ve class } x_1 \cdot \bar{s}_1 \bar{s}_1 + x_2 \bar{s}_2 \bar{s}_1 + x_3 \bar{s}_3 \bar{s}_1 = -1 \quad ①$$

$$\text{-ve class } x_1 \cdot \bar{s}_1 \bar{s}_2 + x_2 \cdot \bar{s}_1 \bar{s}_2 + x_3 \bar{s}_3 \bar{s}_2 = -1 \quad ②$$

$$\text{+ve class } x_1 \cdot \bar{s}_1 \bar{s}_3 + x_2 \cdot \bar{s}_1 \bar{s}_3 + x_3 \bar{s}_1 \bar{s}_3 = +1 \quad ③$$

Solving the first equation

$$x_1(\bar{s}_1 \cdot \bar{s}_1) + x_2(\bar{s}_2 \cdot \bar{s}_1) + x_3(\bar{s}_3 \cdot \bar{s}_1)$$

$$x_1 \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} + x_2 \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} + x_3 \begin{bmatrix} 1 \\ 1.5 \end{bmatrix} \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \quad ④$$

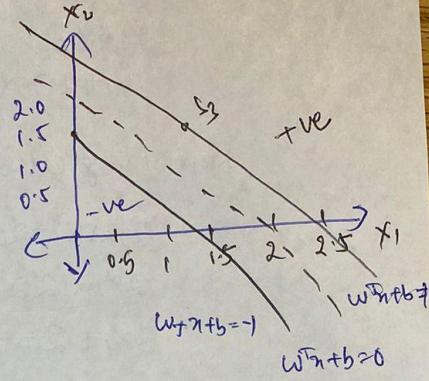
$$2.25x_1 + 2x_2 + 3x_3 = -1 \quad ③$$

equation ② :-

$$2x_1 + 2 - 2.5x_2 + 2.25x_3 = -1 \quad ⑤$$

equation ③ :-

$$x_1 \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1.5 \end{bmatrix} + x_2 \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} \begin{bmatrix} 1 \\ 1.5 \end{bmatrix} + x_3 \begin{bmatrix} 1 \\ 1.5 \end{bmatrix} \begin{bmatrix} 1 \\ 1.5 \end{bmatrix}$$



⑦

Equation of line separation is $n_1 + n_2 = 2$ | here 2 is
 $n_1 = 1 - n_2$ intercept & slope $m=1$

Now putting the value of n_1 & we get n_2

$$\begin{array}{l} n_1 = 2 - n_2 = 0 \\ n_1 = 1 - n_2 = 1 \end{array} \quad \left| \begin{array}{l} n_1 = 1.5, n_2 = 0.5 \\ n_1 = 0, n_2 = 2 \end{array} \right.$$

Also, the line of equation with maximum margin

$$n_1 + n_2 = 2$$

Ans(2) For the polynomial kernel function

$$k(x, z) = [(x \cdot z) + 1]^d$$

To show that:

$$k(x, z) = \phi(x) \cdot \phi(z)$$

$$\text{where } \phi(z) = [1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, z_2^2, \sqrt{2}z_1 z_2]$$

$$\text{Given } k(x, z) = [(x \cdot z) + 1]^d$$

$$= [x_1 z_1 + x_2 z_2 + 1]^d$$

$$k(x, z) = x_1 z_1 + x_2 z_2 + 1 + 2x_1 z_1 x_2 z_2 + 2x_1 z_1 -$$

$$2x_2 z_2 \quad \text{①}$$

$$\text{Now we need to show that } k(x, z) = [\phi(x) \cdot \phi(z)]$$

$$\text{where } \phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2)$$

$$\text{similarly } \phi(z) = (1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, z_2^2, \sqrt{2}z_1 z_2)$$

Putting $\phi(z) \cdot \phi(x)$ on the RHS of the equation we get

$$\begin{aligned}
 [\phi(n) \cdot \phi(z)] &= [(1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2)(1, \sqrt{2}z_1, \\
 &\quad \sqrt{2}z_2, z_1^2, z_2^2, \sqrt{2}z_1 z_2)] \\
 &= 1 + \sqrt{2}x_1 \sqrt{2}z_1 + \sqrt{2}x_2 \sqrt{2}z_2 + x_1^2 z_1^2 + x_2^2 z_2^2 + \sqrt{2}x_1 x_2 \sqrt{2} \\
 &\quad z_1 z_2 \\
 &= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 x_2 z_2 \quad (2)
 \end{aligned}$$

Since eq(1) & eq(2) are same we have proved that

$$K(x_1 z) = [\phi(n) \cdot \phi(z)]$$

when showing both left hand side $[K(x_1 z)]$ & right hand right side $[\phi(n) \cdot \phi(z)]$ separately, we got same answer for both, hence it is proved.

Solution 13-

Observation and comment: _

The optimal learning rate, n=0.1, achieved the best result with 11 steps to reach the global minimum. _ The optimal learning rate, n=0.01, achieved the best result with 129 steps to reach the global minimum.

The observation for the learning rate n=0.001 is that, despite taking 500 steps, it did not reach the global minimum. However, it came very close, with a value of $f(x,y)=2.31$. This suggests that a higher learning rate leads to a faster convergence to the global minimum.

```
import random

def calculate_f(x, y):
    return 52 - 4*x + 2*x**2 + 24*y + 3*y**2 # Derivative of f(x,y) with respect to x
def derivative_x(x):
    return 4*x - 4
# Derivative of f(x,y) with respect to y

def derivative_y(y):
    return 6*y + 24

def gradient_descent(x_init, y_init, num_steps, learning_rate):
    x, y = x_init, y_init
    for step in range(num_steps):
        x -= learning_rate * derivative_x(x)
        y -= learning_rate * derivative_y(y)
        z = calculate_f(x, y)
        print(f"STEP[{step}]: f(x,y) => f({{round(x,2)},{{round(y,2)}}}) ==> {round(z,2)}")
    learning_rate = 0.001
    num_steps = 500
    x_initial = random.randint(-10, 10)
    y_initial = random.randint(-10, 10)
    print(f"\nInitial random points: \n{x_initial},{y_initial}") gradient_descent(x_initial, y_initial, num_steps, learning_rate) #Main Program Ends
```

```
learning_rate = 0.001
num_steps = 500
x_initial = random.randint(-10, 10)
y_initial = random.randint(-10, 10)
print(f"\nInitial random points: \n({x_initial},{y_initial})") gradient_descent(x_initial,
y_initial, num_steps, learning_rate) #Main Program Ends
```