



**Indian Institute of Technology Bombay**  
**CS 419 End Semester Examination**

**Date: November 21, 2024**

**Max Marks: 50**

**Duration: 180 minutes**

---

**Instructions:**

- Answer all questions.
- Show all work clearly, be as clear and precise as possible.
- The exam is open notes.

# PART I : Perceptron Convergence Theorem

The **Perceptron Convergence Theorem** states that if the dataset is linearly separable, the Perceptron algorithm will find a separating hyperplane in a finite number of steps. In this question, we will see the proof of this theorem. The problem setup and notations are defined first.

## Problem Setup

### Assumptions

- The dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  is linearly separable
- $\mathbf{x}_i \in \mathbb{R}^d$  is the feature vector, and  $y_i \in \{-1, 1\}$  is the label.
- There exists a weight vector  $\mathbf{w}^*$  such that  $y_i(\mathbf{w}^* \cdot \mathbf{x}_i) > 0 \quad \forall i$

### Perceptron Update Rule

- Initialize  $\mathbf{w} = \mathbf{0}$ .
- For each misclassified point  $(\mathbf{x}_i, y_i)$  when  $y_i(\mathbf{w}_t \cdot \mathbf{x}_i) < 0$ :

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_i \mathbf{x}_i$$

### Key Quantities

- Define  $\gamma$  as the **margin** of the dataset. Concretely,  $\gamma$  is given as

$$\gamma = \min_i \frac{y_i(\mathbf{w}^* \cdot \mathbf{x}_i)}{\|\mathbf{w}^*\|}$$

- The norm of feature vectors is bounded by  $M$

$$M = \max_i \|\mathbf{x}_i\|$$

### The Cauchy-Schwarz inequality

The Cauchy-Schwarz inequality states that for any vectors  $\mathbf{u}$  and  $\mathbf{v}$  in  $\mathbb{R}^n$ :

$$(\mathbf{u} \cdot \mathbf{v})^2 \leq \|\mathbf{u}\|^2 \|\mathbf{v}\|^2,$$

where,  $\mathbf{u} \cdot \mathbf{v} = \sum_{i=1}^n u_i v_i$  is the dot product, and  $\|\mathbf{u}\|$  is the euclidean norm.

## Proof

### Bounding the Increase in $\|\mathbf{w}\|$ :

(i) When the Perceptron makes an update, we have:

$$\|\mathbf{w}_{t+1}\|^2 = \|\mathbf{w}_t + y_i \mathbf{x}_i\|^2 = \underbrace{\| \quad (a) \quad \|^2}_{\text{}} + 2y_i \underbrace{(\quad (b) \quad)}_{\text{}} + \underbrace{\| \quad (c) \quad \|^2}_{\text{}}$$

(ii) Since  $\|y_i \mathbf{x}_i\|^2 = \|\mathbf{x}_i\|^2 \leq M^2$ , we get:

$$\|\mathbf{w}_{t+1}\|^2 \leq \|\mathbf{w}_t\|^2 + \underbrace{\quad (d) \quad}_{\text{}}$$

### Increase in Alignment with $\mathbf{w}^*$ :

(iii) The dot product with the target vector  $\mathbf{w}^*$  increases by at least  $\gamma$ :

$$\mathbf{w}_{t+1} \cdot \mathbf{w}^* = (\mathbf{w}_t + y_i \mathbf{x}_i) \cdot \mathbf{w}^* = \mathbf{w}_t \cdot \mathbf{w}^* + y_i (\mathbf{x}_i \cdot \mathbf{w}^*) \geq \mathbf{w}_t \cdot \mathbf{w}^* + \gamma$$

$$\mathbf{w}_{t+1} \cdot \mathbf{w}^* \geq \mathbf{w}_t \cdot \mathbf{w}^* + \underbrace{\quad (e) \quad}_{\text{}}$$

### Bounding the Number of Updates:

(iv) Let  $N$  be the number of updates made. Then:

$$\|\mathbf{w}_N\|^2 \leq \underbrace{\quad (f) \quad}_{\text{}}, \quad \mathbf{w}_N \cdot \mathbf{w}^* \geq \underbrace{\quad (g) \quad}_{\text{}}$$

(v) Prove  $N \leq \left( \frac{M\|\mathbf{w}^*\|}{\gamma} \right)^2$  (Apply Cauchy-Schwartz inequality and some more steps)

[15 Marks: (i): 1+1+1, (ii): 2, (iii): 3, (iv): 2+2, (v): 3 ]

## PART II : Kernels

(A) Prove that the following are valid kernels by finding a transformation  $\phi$  such that  $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$  .

1. Let  $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$ . [2.5 marks]

2. Let  $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2 + 3 \cdot (\mathbf{x}^T \mathbf{x}') + 2$ . [2.5 marks]

## PART III : Neural Networks

### Implementing boolean functions with neural networks

Consider a binary logic function  $f : \{0, 1\}^2 \rightarrow \{0, 1\}$ . A neural network with a single hidden layer of two neurons and ReLU activation is used to approximate this function. Specifically, the function represented by the neural network is  $f(x) = w_2(\sigma(w_1x + b_1)) + b_2$ , where  $x = [x_1, x_2] \in \mathbb{R}^2$ ,  $w_1 \in \mathbb{R}^{2 \times 2}$ ,  $w_2 \in \mathbb{R}^{1 \times 2}$ . The function  $\sigma : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  and is applied pointwise, meaning  $\sigma([a, b]) = [\sigma(a), \sigma(b)]$ . Please note that  $\sigma(\cdot)$  is the ReLU activation itself. Design  $w_1, w_2$  and  $\sigma$  (the activation function) to compute **AND**, **XOR**, **OR** and **XNOR** functions. [10 marks]

Refer to the following truth tables for the exact definition of boolean functions: AND, OR, XOR and XNOR.

#### AND

$x_1$	$x_2$	AND( $x_1, x_2$ )
0	0	0
0	1	0
1	0	0
1	1	1

#### OR

$x_1$	$x_2$	OR( $x_1, x_2$ )
0	0	0
0	1	1
1	0	1
1	1	1

#### XOR

$x_1$	$x_2$	XOR( $x_1, x_2$ )
0	0	0
0	1	1
1	0	1
1	1	0

#### XNOR

$x_1$	$x_2$	XNOR( $x_1, x_2$ )
0	0	1
0	1	0
1	0	0
1	1	1

## PART IV : Logistic Regression

In this problem, we are tasked with performing multi-class classification with three possible classes, denoted as  $C_1$ ,  $C_2$ , and  $C_3$ . There is always an option to use the softmax approach but some genius student in CS 419 course came up with the idea of extending the logistic regression to three classes using two sigmoids. Here is how they propose to model probabilities of the first two classes,  $C_1$  and  $C_2$ , and derive the probability of the third class,  $C_3$ , based on these values.

### Modeling Probabilities

The probabilities  $p_1$  and  $p_2$  are modeled using logistic sigmoid functions. Let  $x$  be the input feature vector,  $w_1$  and  $w_2$  be the weight vectors for classes  $C_1$  and  $C_2$ , and  $b_1$  and  $b_2$  be the bias terms for these classes, respectively. The logistic sigmoid function is defined as:

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

We can express  $p_1$  and  $p_2$  as follows:

$$p_1 = \sigma(w_1^T x + b_1) = \frac{1}{1 + e^{-(w_1^T x + b_1)}}$$
$$p_2 = \sigma(w_2^T x + b_2) = \frac{1}{1 + e^{-(w_2^T x + b_2)}}$$

The third class probability  $p_3$  is derived as:

$$p_3 = 1 - p_1 - p_2$$

**Assume that weight vectors  $w_1$  and  $w_2$  can be any general vectors in the  $R^d$ .** Thus, the sum of the probabilities is always equal to 1, i.e.,  $p_1 + p_2 + p_3 = 1$ .

1. This, in theory, seems like a great idea but there is an immediate problem. Point out the problem? [2 mark]
2. The values for  $p_1, p_2$  and  $p_3$  should satisfy an additional constraint if the above formulation is valid. What is that? [3 mark]
3. Given a dataset  $\{(x_i, y_i) : i \in \{1, \dots, 100\}\}$ , design a loss function to estimate  $w_1, w_2$ , which also encodes the constraint introduced in item 2. [5 marks]

## PART V : Coding and PyTorch

Read the following part carefully before reading the code.

In the context of autoencoders, the model learns a transformation from the input  $x$  to a lower-dimensional representation  $z$ , and then reconstructs the input  $x'$  from  $z$ . The process can be described as:

$$x \xrightarrow{\text{Encoder}} z \xrightarrow{\text{Decoder}} x'$$

Here,  $x \in \mathbb{R}^d$  is the original input,  $z \in \mathbb{R}^{d'}$  is the compressed representation, where  $d' < d$ , and  $x'$  is the reconstructed output, ideally close to  $x$ . The network is trained to minimize the distance between  $x$  and  $x'$ , ensuring that the learned representation captures the most important features of the input data.

In the code below, we aim to train an autoencoder using the distance minimization scheme given above. We provide an autoencoder implementation and the corresponding training and early stopping code. However, there are a number of bugs currently - due to dimension mismatch, error in early stopping and learning rate. Identify and report where the bugs are, and how to fix them. **[10 marks]**

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from torch.utils.data import DataLoader, TensorDataset,
    ↪ random_split
5
6 # Define the Autoencoder class
7 class Autoencoder(nn.Module):
8     def __init__(self, input_dim, hidden_dims, latent_dim):
9         super(Autoencoder, self).__init__()
10
11         # Encoder
12         self.encoder = nn.Sequential(
13             nn.Linear(input_dim, hidden_dims[0]),
14             nn.ReLU(),
15             nn.Dropout(0.2),
16             nn.Linear(hidden_dims[0], hidden_dims[1]),
17             nn.ReLU(),
18             nn.Linear(hidden_dims[1], latent_dim)
19         )
20
21         # Decoder
22         self.decoder = nn.Sequential(
23             nn.Linear(latent_dim, hidden_dims[1]),
24             nn.ReLU(),
25             nn.Linear(hidden_dims[1], hidden_dims[0]),
26             nn.ReLU(),
27             nn.Linear(hidden_dims[0], input_dim),
28             nn.Sigmoid()
29         )
30
31     def forward(self, x):
32         # Forward pass
33         encoded = self.encoder(x)
34         decoded = self.decoder(encoded)
35         return encoded, decoded
36
37 # Generate some dummy data # NO ERROR IN NEXT 2 LINES
38 torch.manual_seed(42)
39 data = get_training_data() # 1000 samples, 20 features
40 train_size = int(0.8 * len(data))
41 es_size = len(data) - train_size
42 train_data, es_data = random_split(data, [train_size, es_size
    ↪ ])
43 # Create DataLoader
44 batch_size = 32
45 train_loader = DataLoader(TensorDataset(train_data),
    ↪ batch_size=batch_size, shuffle=True)
46 es_loader = DataLoader(TensorDataset(es_data), batch_size=
    ↪ batch_size, shuffle=False)

```



```

45 # Initialize the model, loss function, and optimizer
46 input_dim = 32
47
48 # Two hidden layers with 128 and 64 units
49 hidden_dims = [128, 64]
50
51 # Latent representation size
52 latent_dim = 10
53
54 model = Autoencoder(input_dim=input_dim,
55                     hidden_dims=hidden_dims,
56                     latent_dim=latent_dim)
57
58 criterion = nn.MSELoss() # Reconstruction loss
59 optimizer = torch.optim.Adam(model.parameters(), lr=1e-2)
60
61 # Training loop
62 epochs = 20
63 patience = 5
64 best_es_loss = float('inf')
65 best_model_state = None
66 no_improve_epochs = 0
67
68 for epoch in range(epochs):
69     total_loss = 0
70     for batch in train_loader:
71         inputs = batch[0] # NO ERROR IN THIS LINE
72
73         # Forward pass
74         output1, output2 = model(inputs)
75         y_label = inputs # NO ERROR IN THIS LINE
76         y_pred = output2
77         loss = criterion(y_label, y_pred)
78
79         # Backward pass
80         optimizer.zero_grad()
81         loss.backward()
82         optimizer.step()
83
84         total_loss += loss.item()
85
86     print(f"Epoch [{epoch+1}/{epochs}]")
87     print(f"Loss: {total_loss/len(train_loader):.4f}")
88
89     model.eval()
90     es_loss = 0
91     with torch.no_grad():
92         for es_batch in train_loader:
93             es_inputs = es_batch[0]
94             _, es_outputs = model(es_inputs)

```

```

95         es_loss += criterion(es_inputs, es_outputs).item()
96     es_loss /= len(es_loader)
97
98     print(f"Epoch [{epoch+1}/{epochs}]")
99     print(f"Train Loss: {total_loss/len(train_loader):.4f}")
100    print(f"Early Stopping Loss: {es_loss:.4f}")
101
102    # Early stopping logic
103    if es_loss < best_es_loss:
104        best_es_loss = es_loss
105        best_model_state = model.state_dict()
106        no_improve_epochs = 0
107
108    else:
109        no_improve_epochs += 1
110        if no_improve_epochs >= patience:
111            print("Early stopping triggered.")
112            break
113
114    # Load the best model state
115    if best_model_state:
116        model.load_state_dict(best_model_state)
117
118    # Example: Get latent representations for the data
119    model.eval()
120    with torch.no_grad():
121        latent, reconstructed = model(data)
122
123    print("Latent representation shape:", latent.shape)

```