

# *Neural Network Implementation on Medical Appointment No-Show Dataset*

**COPS Summer of Code 2025**

Intelligence Guild



*Club of Programmers, IIT (BHU) Varanasi*

*By*  
***Tanishq Singh***

24045164

[hello@tanishqsingh.com](mailto:hello@tanishqsingh.com)

Official IG Website: <https://cops-iitbhu.github.io/IG-website>

# Contents

<b>1</b>	<b>Implementation Overview</b>	<b>3</b>
<b>2</b>	<b>Data Analysis</b>	<b>3</b>
2.1	Import Data . . . . .	3
2.2	Missing Data . . . . .	4
<b>3</b>	<b>Exploratory Data Analysis</b>	<b>4</b>
3.1	Creating New Features . . . . .	4
3.1.1	Waiting Period . . . . .	4
3.1.2	Appointment Weekday . . . . .	4
3.1.3	Patients missed appointment before . . . . .	4
3.1.4	Track Neighbour-hood with Maximum Missed Appointment . . . . .	5
3.2	Features Relation . . . . .	6
<b>4</b>	<b>Data Preprocessing</b>	<b>7</b>
4.1	LabelEncoding . . . . .	7
4.2	Removing False Data . . . . .	7
4.3	Train, Test Split . . . . .	7
4.4	Standardizing Data Values . . . . .	7
4.5	Processed Data . . . . .	8
<b>5</b>	<b>Model Implementation</b>	<b>8</b>
5.1	Mathematics Behind Neural Network . . . . .	8
5.1.1	What are nodes? . . . . .	8
5.1.2	Artificial Neural Network . . . . .	9
5.1.3	Mathematics of Complex ANN . . . . .	9
5.1.4	Forward Propagation in our Implementation . . . . .	10
5.1.5	Backpropagation: Theory and Detailed Steps . . . . .	10
5.2	Neural Network Implementation from Scratch . . . . .	11
5.2.1	Assumptions In Our Neural Network . . . . .	11
5.2.2	Defining Layer Class . . . . .	12
5.2.3	Initializing Layer Class . . . . .	12
5.2.4	Defining Model Class . . . . .	13
5.2.5	Initializing Model Class . . . . .	16
5.2.6	Defining Metrics Function . . . . .	16
5.2.7	Training Custom ANN . . . . .	16
5.3	Neural Network Implementation using PyTorch . . . . .	17
5.3.1	Assumptions . . . . .	17
5.3.2	Importing Libraries and Converting Numpy to tensor . . . . .	17
5.3.3	Defining Model and Layers . . . . .	17
5.3.4	Initializing Model . . . . .	17
5.3.5	Training PyTorch Model . . . . .	18
5.3.6	Model Evaluation . . . . .	18

<b>6</b>	<b>Evaluation and Analysis</b>	<b>18</b>
6.1	Convergence Time	18
6.1.1	Convergence of PyTorch is faster than core?	18
6.1.2	Comparision of training time and cost	19
6.2	Performance Metrics	20
6.2.1	Accuracy	20
6.2.2	F1 Score	20
6.2.3	Precision-Recall Area Under Curve (PR AUC)	21
6.3	Memory Usage	21
6.4	Confusion Matrix and Inference	21
6.5	Analysis and Discussion	23
6.5.1	Convergence Speed, Performance, and Memory Usage	23
6.5.2	Some Other Factors	23

# 1 Implementation Overview

The project involves building two distinct neural network implementations (given below) to analyze their performance differences on a real-world healthcare dataset.

1. One from scratch using basic `Python` libraries, and
2. Another using `PyTorch` framework.

we compare two neural networks trained on the same dataset, with different performance metrics such as `F1 Score`, `PR AUC`, `Accuracy`, and the `Confusion Matrix`. We also take a closer look at how efficient each approach is, both in terms of prediction quality and computational effort.

## 2 Data Analysis

We use the [joniarroba/noshowappointments](#) dataset to analyze the different implementations of neural network mentioned above.

### 2.1 Import Data

From the kaggle Data Dictionary, there are 14 columns, which are :-

Table 1: Data Properties

Index	Column Name	Description
0	PatientId	Identification of a patient.
1	AppointmentID	Identification of each appointment.
2	Gender	Male or Female.
3	ScheduledDay	The day someone called or registered the appointment, this is before appointment of course.
4	AppointmentDay	The day of the actual appointment, when they have to visit the doctor.
5	Age	How old is the patient.
6	Neighbourhood	Where the appointment takes place.
7	Scholarship	True or False.
8	Hypertension	True or False.
9	Diabetes	True or False.
10	Alcoholism	True or False.
11	Handicap	True or False.
12	SMS_received	True or False.
13	No-Show	Yes or No.

As from the table we can see than `AppointmentDay` and `ScheduledDay` are date, and there are some spelling mistakes also, which we take care of.

```
1 import pandas as pd
2
3 # Import Data from csv file.
4 df = pd.read_csv("medical-no-show.csv",
5                  parse_dates=["ScheduledDay", "AppointmentDay"])
6
7 # Correct spelling mistakes.
8 df.rename(columns={"Hipertension": "Hypertension",
9                  "Handcap": "Handicap",
10                 "No-show": "no_show"}, inplace=True)
```

## 2.2 Missing Data

As there are only 20 rows where data is missing from more than 1 lakh rows, we drop those rows. As it'll not make a significant difference in our prediction.

```
1 print(len(df)) # output: 110527
2 print(df.isna().sum().sum()) # output: 20
3
4 df.dropna(inplace=True) # Rows dropped with missing data
```

## 3 Exploratory Data Analysis

### 3.1 Creating New Features

#### 3.1.1 Waiting Period

Waiting period is referred here as number of days a patient waited from booking day to appointment day.

```
1 df["days_waited"] = (df["AppointmentDay"] - df["ScheduledDay"]).dt.days
```

#### 3.1.2 Appointment Weekday

The weekday on which the appointment is made.

```
1 # sunday(0), monday(1), tuesday(2), and so on.
2 df["appointment_weekday"] = df["AppointmentDay"].dt.weekday
```

#### 3.1.3 Patients missed appointment before

Number of patients who missed appointment before, and if they did then number of times they missed appointment before.

```
1 # Dict with key as PatientId and value as no. of missed appointments.
2 patient_missed_before = df.groupby("PatientId")["no_show"].sum()
3 patient_missed_before = patient_missed_before.to_dict()
4
5 df["missed_appointment_before"] = df["PatientId"].map(lambda x:
6                                                         patient_missed_before.get(x) > 0)
7 df["times_missed_before"] = df["PatientId"].map(lambda x:
8                                                         patient_missed_before.get(x))
```

### 3.1.4 Track Neighbour-hood with Maximum Missed Appointment

Here we sort neighbourhood based on maximum number of missed appointments and then assign each patients with it's index.

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 no_shows_df = sorted_neighbourhoods.reset_index()
5
6 plt.figure(figsize=(12, 6))
7 sns.barplot(
8     data=no_shows_df,
9     x="Neighbourhood",
10    y="no_show",
11    hue="no_show",
12    palette="viridis"
13 )
14
15 plt.xticks(rotation=90)
16 plt.title("Total No-Shows by Neighbourhood")
17 plt.xlabel("Neighbourhood")
18 plt.ylabel("Number of No-Shows")
19 plt.tight_layout()
20 plt.show()

```

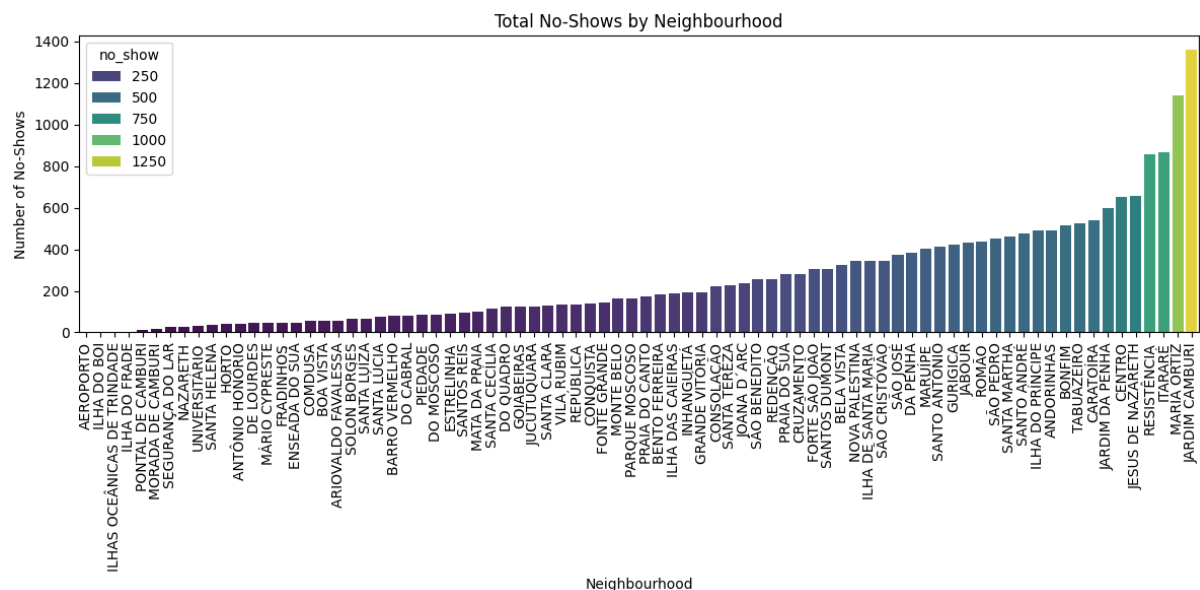


Figure 1: Total No-Shows by each Neighbourhood

```

1 sorted_neighbourhoods = df.groupby("Neighbourhood")["no_show"]
2                           .sum()
3                           .sort_values()
4
5 df["neighbourhood_weight"] = df["Neighbourhood"].map(lambda x:
6                                                         sorted_neighbourhoods.index
7                                                         .to_list().index(x))

```

### 3.2 Features Relation

We can plot a simple heatmap of features correlations with No-Show, to find out which features are highly correlated with our needs.

```

1 correlation = df.drop(columns=["PatientId", "AppointmentID",
2                               "ScheduledDay", "AppointmentDay",
3                               "Neighbourhood"])
4               .corrwith(df["no_show"])
5               .drop("no_show")
6
7 plt.figure(figsize=(10, 6))
8 sns.barplot(
9     x=correlation.values,
10    y=correlation.index,
11    hue=correlation.index,
12    palette="coolwarm"
13 )
14
15 plt.title("Correlation with no_show")
16 plt.ylabel("Features")
17 plt.xlabel("Correlation Coefficient")
18
19 plt.grid(axis="x", linestyle="--", alpha=0.7)
20 plt.tight_layout()
21 plt.show()

```

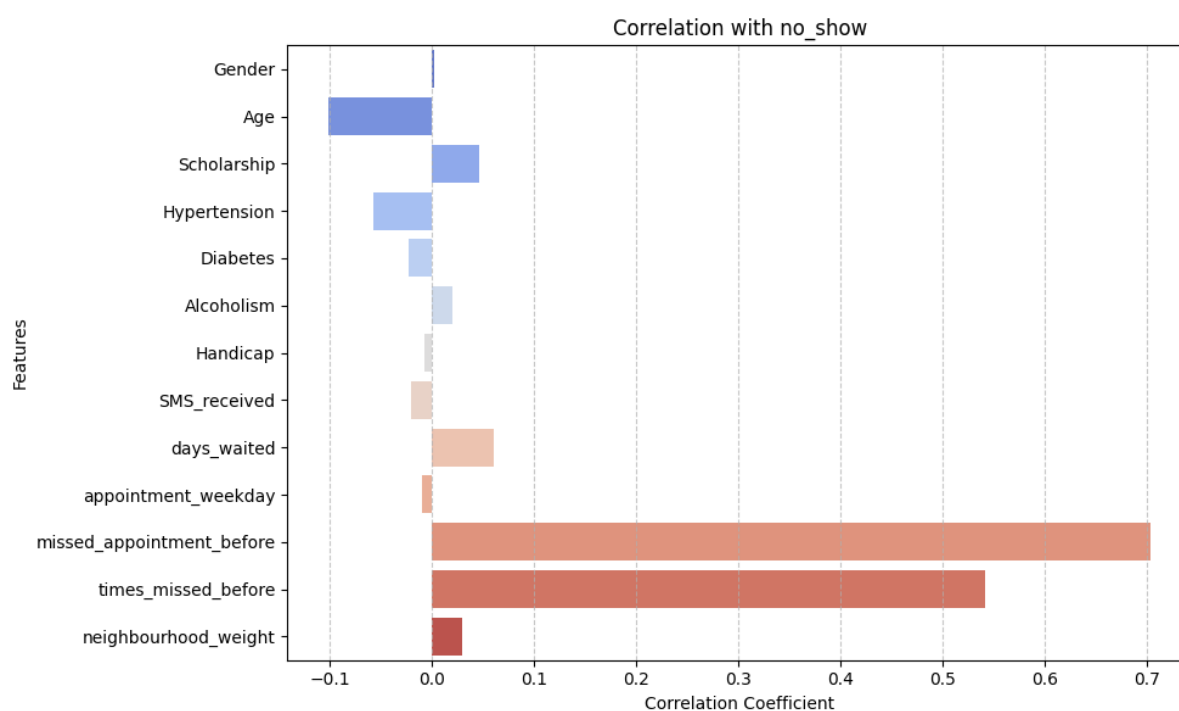


Figure 2: Correlation with No-Show

## 4 Data Preprocessing

Here we'll process the data to make fit for the model to train, like encoding 'yes' / 'no' to 1 and 0, splitting train, test set, and related things.

### 4.1 LabelEncoding

Rows like 'no\_show', and 'Gender' contains value 'Yes' / 'No', and 'Female', 'Male' respectively. Which should be encoded to 0 and 1.

```
1 from sklearn.preprocessing import LabelEncoder
2
3 # Label Encode 'no_show' to 0's and 1's
4 le_noshow = LabelEncoder()
5 df["no_show"] = le_noshow.fit_transform(df["no_show"])
6
7 # Label Encode 'Gender' to 0's and 1's
8 le_gender = LabelEncoder()
9 df["Gender"] = le_gender.fit_transform(df["Gender"])
```

### 4.2 Removing False Data

There may be some data which are not true like age to be negative, or having appointment date to be before than schedule date. So we'll remove those rows.

```
1 df = df[df["Age"] >= 0]
2 df = df[df["days_waited"] >= 0]
```

### 4.3 Train, Test Split

Splitting all the data we've to **train** (90%) and **test** (10%).

```
1 from sklearn.model_selection import train_test_split
2
3 x = df.drop(columns=["PatientId", "AppointmentID", "ScheduledDay",
4                     "AppointmentDay", "Neighbourhood", "no_show"])
5 y = df["no_show"]
6
7 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=.1)
8
9 # Output: ((64763, 13), (64763,)), (7196, 13), (7196,))
10 x_train.shape, y_train.shape, x_test.shape, y_test.shape
```

### 4.4 Standardizing Data Values

While training a neural network we should make sure all our data are standardized before. Because of this required time to reach convergence is reduced.

```
1 from sklearn.preprocessing import StandardScaler
2
3 scaler = StandardScaler()
4
5 x_train = scaler.fit_transform(x_train)
6 x_test = scaler.transform(x_test)
```



## 4.5 Processed Data

After performing EDA, and pre-processing, our final data looks something like this.

```
1 print(pd.DataFrame(x_train, columns=x.columns).sample(2))
```

	Gender	Age	Scholarship	Hypertension	Diabetes	Alcoholism	Handicap	SMS_received	days_waited
13467	1.419619	1.722408	-0.319686	1.951761	3.522963	-0.161941	6.369641	1.014010	-0.098341
56435	1.419619	0.761867	-0.319686	-0.512358	-0.283852	-0.161941	-0.129147	-0.986183	-0.886512
	appointment_weekday	missed_appointment_before	times_missed_before	neighbourhood_weight					
	-0.634270		-0.898783				-0.672176		-1.009461
	0.095551		1.112616				2.291211		0.460754

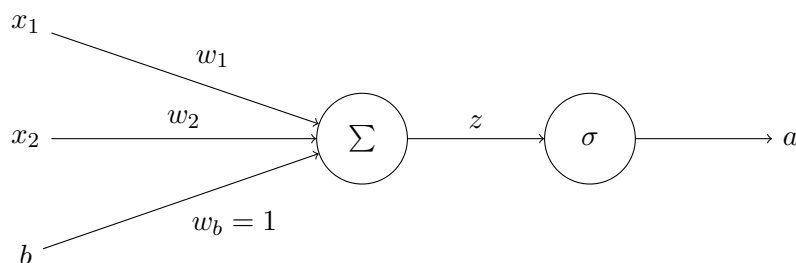
Figure 3: Final Processed Data

## 5 Model Implementation

### 5.1 Mathematics Behind Neural Network

#### 5.1.1 What are nodes?

A node takes in inputs, processes them (usually by applying a weighted sum and an activation function), and passes the result to the next layer. It's responsible for learning patterns in data.



$$z = w_1x_1 + w_2x_2 + 1 \cdot b$$

$$a = \sigma(z)$$

Figure 4: A single ANN node with two inputs, bias, and activation function

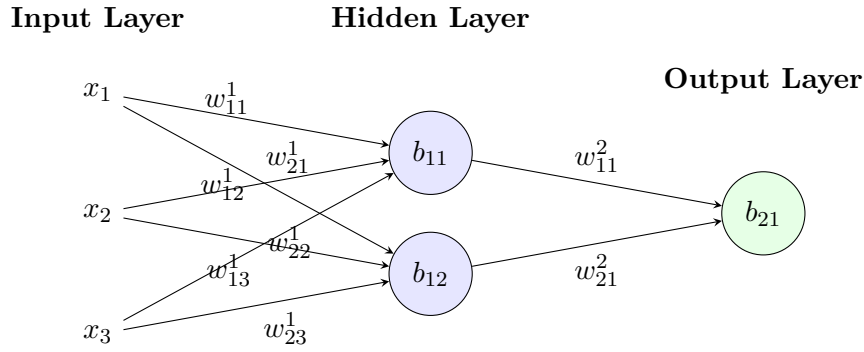
There can be n-number of **inputs** (features).

$$z = w_1x_1 + w_2x_2 + \dots + W_nx_n + 1 \cdot b$$

$$a = \sigma(z)$$

### 5.1.2 Artificial Neural Network

An ANN is a network of these nodes connected to each other which are organized in layers. An **Artificial Neural Network** is a computational model inspired by the way biological brains process information. It's made up of layers of simple processing units called **neurons** or **nodes**, connected in a network.



$$\begin{aligned}
 o_{11} &= \sigma(w_{11}^1 x_1 + w_{21}^1 x_2 + w_{31}^1 x_3 + b_{11}) \\
 o_{12} &= \sigma(w_{12}^1 x_1 + w_{22}^1 x_2 + w_{32}^1 x_3 + b_{12}) \\
 \hat{y} = o_{21} &= \sigma(w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_{21})
 \end{aligned}$$

Figure 5: A ANN three inputs, one hidden, and single output layer

### 5.1.3 Mathematics of Complex ANN

This contains multiple inputs (features), multiple hidden layers, with different number of nodes in each layer, and there can be single output, or even multiple.

$$\mathbf{A}^{[0]} = \mathbf{X} = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{bmatrix} \quad \mathbf{A}^{[l]} = \begin{bmatrix} o_{l1} \\ o_{l2} \\ \vdots \\ o_{lm} \end{bmatrix} \quad \mathbf{W}^{[l]} = \begin{bmatrix} w_{11}^1 & w_{12}^1 & \cdots & w_{1m}^1 \\ w_{21}^1 & w_{22}^1 & \cdots & w_{2m}^1 \\ \vdots & \vdots & \ddots & \vdots \\ w_{l1}^l & w_{l2}^l & \cdots & w_{lm}^l \end{bmatrix} \quad \mathbf{B}^{[l]} = \begin{bmatrix} b_{l1} \\ b_{l2} \\ \vdots \\ b_{lm} \end{bmatrix}$$

Then for each layer  $l = 1, 2, \dots, L$ :

$$\mathbf{Z}^{[l]} = \mathbf{W}^{[l]} \mathbf{A}^{[l-1]} + \mathbf{b}^{[l]}$$

$$\mathbf{A}^{[l]} = f^{[l]}(\mathbf{Z}^{[l]})$$

$f^{[l]}(Z)$  (activation function) can be different for each layer.

$$\text{ReLU: } f(x) = \max(0, x) \quad \text{Sigmoid: } \sigma(x) = \frac{1}{1 + e^{-x}} \quad \text{Tanh: } \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

### Some of the Activation Functions

### 5.1.4 Forward Propagation in our Implementation

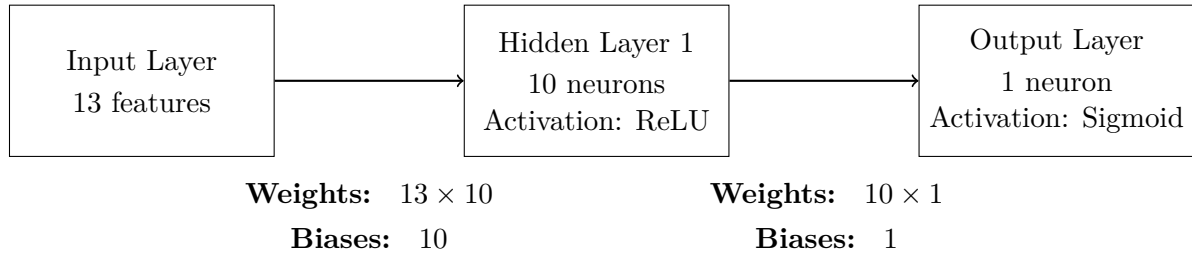


Figure 6: Fully Connected Feedforward Neural Network Diagram

$$A^{[1]} = \text{ReLU}(W_{10 \times 13}^{[1]} X_{13 \times 1} + B_{10 \times 1}^{[1]})$$

$$\hat{Y} = A^{[2]} = \sigma(W_{1 \times 10}^{[2]} A_{10 \times 1}^{[1]} + B_{1 \times 1}^{[2]})$$

As we are performing binary classification, the `loss function` will be `binary classentropy`

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

Values of `weights`, and `biases` are found using a method called `BackPropagation`.

### 5.1.5 Backpropagation: Theory and Detailed Steps

Backpropagation is the key algorithm used to train artificial neural networks. It computes the gradient of the loss function with respect to each weight by the chain rule, enabling gradient descent optimization.

**Gradient Descent:** is an iterative optimization algorithm. In each iteration, we update the weights and bias using the partial derivatives of the loss function with respect to the weights and bias.

$$w_{new,i} := w_{old,i} - \alpha \cdot \frac{\partial L}{\partial w_i}$$

$$b_{new,i} := b_{old,i} - \alpha \cdot \frac{\partial L}{\partial b_i}$$

Where  $\alpha$  is the `learning rate`, and  $L$  is `loss function`. It controls the step-size of each update. Repeated application of these updates gradually make the model parameters (i.e weights and bias) converge. Since, the loss function used has only one minima, for this application the convergence will always be at the global minima.

### ANN Equations for Our Model

$$\begin{aligned} A^{[1]} &= \text{ReLU}(Z^{[1]}) \quad , \quad Z^{[1]} = W_{10 \times 13}^{[1]} X_{13 \times 1} + B_{10 \times 1}^{[1]} \\ \hat{Y} = A^{[2]} &= \sigma(Z^{[2]}) \quad , \quad Z^{[2]} = W_{1 \times 10}^{[2]} A_{10 \times 1}^{[1]} + B_{1 \times 1}^{[2]} \end{aligned}$$

#### Step 1: Compute Gradient for Output Layer ( $2^{nd} - layer$ )

$$\begin{aligned} \frac{\partial L}{\partial W^{[2]}} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial W^{[2]}} \\ \frac{\partial L}{\partial B^{[2]}} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial B^{[2]}} \\ \frac{\partial L}{\partial \hat{y}} &= \left( \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \quad , \quad \frac{\partial \hat{y}}{\partial Z^{[2]}} = \hat{y} \cdot (1 - \hat{y}) \quad , \quad \frac{\partial Z^{[2]}}{\partial W^{[2]}} = A^{[1]} \quad , \quad \frac{\partial Z^{[2]}}{\partial B^{[2]}} = 1 \end{aligned}$$

#### Step 2: Compute Gradient for Hidden Layer ( $1^{st} - layer$ )

$$\begin{aligned} \frac{\partial L}{\partial W^{[1]}} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial W^{[1]}} \\ \frac{\partial L}{\partial B^{[1]}} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial B^{[1]}} \\ \frac{\partial Z^{[2]}}{\partial B^{[1]}} &= \begin{cases} W^{[2]} & \text{if } Z^{[1]} > 0 \\ 0 & \text{otherwise} \end{cases} \quad , \quad \frac{\partial Z^{[2]}}{\partial W^{[1]}} = W^{[2]} \cdot \frac{\partial A^{[1]}}{\partial W^{[1]}} \cdot \frac{\partial Z^{[1]}}{\partial W^{[1]}} \\ \frac{\partial A^{[1]}}{\partial W^{[1]}} &= \begin{cases} 1 & \text{if } Z^{[1]} > 0 \\ 0 & \text{otherwise} \end{cases} \quad , \quad \frac{\partial Z^{[1]}}{\partial W^{[1]}} = X_{[i]} \end{aligned}$$

#### Step 3: Updates the weights, and biases

$$\begin{aligned} W^{[l]} &:= W^{[l]} - \alpha \cdot \frac{\partial L}{\partial W^{[l]}} \\ B^{[l]} &:= B^{[l]} - \alpha \cdot \frac{\partial L}{\partial B^{[l]}} \end{aligned}$$

Then for each layer  $l = 1, 2$ .

Where  $\alpha$  is learning rate.

## 5.2 Neural Network Implementation from Scratch

### 5.2.1 Assumptions In Our Neural Network

We'll implement **Stochastic Gradient Descent** which updates parameters after each row, because of this in **PyTorch NN** we'll not use any **DataLoaders** as then the performance will differ.

### 5.2.2 Defining Layer Class

This `Layer` class defines a single neural network layer with support for **ReLU** and **Sigmoid** activations. It stores weights, biases, pre-activation (`z`), and activated outputs (`a`), and includes methods for the forward pass and activation function derivatives which are essential for backpropagation.

```

1 class Layer:
2     def __init__(self, activation, units):
3         self.activation = activation
4         self.units = units # Number of nodes
5
6         self.z = None # Pre-Activation
7         self.a = None # Activated
8
9         self.weights = None
10        self.bias = None
11
12    def activation_function(self, xi):
13        if self.activation == "relu":
14            return np.maximum(0, xi)
15
16        elif self.activation == "sigmoid":
17            return 1 / (1 + np.exp(-xi))
18
19    def activation_derivative(self, xi):
20        if self.activation == "relu":
21            return (xi > 0).astype(float)
22
23        elif self.activation == "sigmoid":
24            sig = self.activation_function(xi)
25            return sig * (1 - sig)
26
27    def forward(self, xi):
28        self.z = np.dot(self.weights, xi) + self.bias
29        self.a = self.activation_function(self.z)
30
31        return self.a

```

### 5.2.3 Initializing Layer Class

We are defining two layers:

1. One hidden with ten nodes and relu as activation function, and
2. Another one output layer with one node and sigmoid as activation function

```

1 layer_1 = Layer(units=10, activation="relu") # Hidden Layer
2 layer_2 = Layer(units=1, activation="sigmoid") # Output Layer

```

### 5.2.4 Defining Model Class

The `Model` class defines a simple feedforward neural network framework with customizable layers. It includes support for weight initialization (He/Xavier), forward propagation, binary cross-entropy loss, and backpropagation with gradient descent. It also provides methods to compile, train (fit), and predict using the model.

```

1 class Model:
2     def __init__(self, layers, input_features):
3         self.layers = layers # List of layers in the model
4         self.input_features = input_features # Number of inputs
5
6         self.compiled = False
7         self.metrics = False # To show metrics after each epochs
8
9         self.learning_rate = 0.01 # Learning rate of the model
10
11     """
12     We'll define these functions later
13     """
14     def compile(self, metrics=False): pass
15
16     def summary(self): pass
17
18     def predict(self, x_test): pass
19
20     def loss_calculate(self, yi_true, yi_pred): pass
21
22     def fit(self, x_train, y_train, val_data, epochs=1, learning_rate=0.01):
23         pass

```

#### Summary Function Initialize

This function tells about the units, activation, trainable parameters in each layer, and also the total trainable parameters in the model.

```

1 def summary(self):
2     total_trainable_parameters = 0
3
4     for layer_i in range(len(self.layers)):
5         layer = self.layers[layer_i]
6         layer_trainable_parameters = self.input_features * layer.units
7         total_trainable_parameters += layer_trainable_parameters
8
9         print(" " * 16, f"Layer: {layer_i + 1}")
10        print(f"Units: {layer.units}")
11        print(f"Activation: {layer.activation}")
12        print(f"Trainable Parameters: {layer_trainable_parameters}")
13        print("=" * 40)
14
15    print(f"Total Trainable Parameters: {total_trainable_parameters}")
16    print("=" * 40, "\n")

```

#### Loss Function Initialize

Binary cross-entropy loss measures the difference between predicted probabilities and actual binary labels. It penalizes incorrect predictions by calculating the -ve log-likelihood of the true class, encouraging the model to output probabilities close to 0 or 1.

```

1 def loss_calculate(self, yi_true, yi_pred):
2     yi_pred = yi_pred + (1e-8) # Handling 0 case
3
4     return -1 * np.mean(
5         (yi_true * np.log(yi_pred)) +
6         ((1 - yi_true) * np.log(1 - yi_pred))
7     )

```

## Weights, and Biases Initialize

1. Weights are initialized differently based on the activation function to maintain stable gradients and avoid problems like vanishing or exploding gradients during training.
2. For layers with ReLU activation, weights are initialized using He initialization, which scales the weights by  $\sqrt{\frac{2}{\text{fan\_in}}}$  to keep variance consistent through the network.
3. For sigmoid activations, Xavier (Glorot) initialization is used, scaling weights by  $\sqrt{\frac{1}{\text{fan\_in}}}$  to help activations avoid saturation and maintain gradient flow.
4. Bias terms for all layers are initialized as zero vectors to provide a baseline offset without introducing bias in the initial output.
5. Once weights and biases are set, the network is marked as compiled and ready for training, with an option to track additional metrics during the process.

```

1 def compile(self, metrics=False):
2     for layer_i in range(len(self.layers)):
3         layer = self.layers[layer_i]
4
5         fan_in = self.input_features if layer_i == 0 else self.layers[
6             layer_i - 1].units
7
8         # He initialization
9         if layer.activation == "relu":
10             layer.weights = np.random.randn(layer.units, fan_in) *
11                 np.sqrt(2. / fan_in)
12
13         # Xavier initialization
14         elif layer.activation == "sigmoid":
15             layer.weights = np.random.randn(layer.units, fan_in) *
16                 np.sqrt(1. / fan_in)
17
18         # Default small random values
19         else:
20             layer.weights = np.random.randn(layer.units, fan_in) * 0.01
21
22         # Biases initialized to zero
23         layer.bias = np.zeros((layer.units, 1))
24
25     self.compiled = True
26     self.metrics = metrics

```

### Fit Function Initialize

The fit function in an ANN trains the model by iteratively updating its weights and biases using different inputs. It performs forward propagation, calculates loss, backpropagates errors, and optimizes parameters to minimize the loss over multiple epochs.

```

1  def fit(self, x_train, y_train, val_data, epochs=1, learning_rate=0.01):
2      self.learning_rate = learning_rate
3      history_core = []
4      counter_core = 0
5      time_start = time.time()
6
7      for epoch in range(epochs):
8          for xi, yi in zip(x_train, y_train):
9              if counter_core % 200 == 0:
10                 cost = self.loss_calculate(val_data[1],
11                                             self.predict(val_data[0]))
12                 history_core.append([counter_core / 200, cost,
13                                     time.time()-time_start])
14                 counter_core = counter_core + 1
15
16                 yi_pred = (self.predict([xi]))[0]
17                 loss = self.loss_calculate(yi, yi_pred)
18
19                 z1, a1 = self.layers[0].z, self.layers[0].a
20                 yi_pred = yi_pred + (1e-8) # Handling 0 case
21
22                 # Calculating derivatives (discussed in 5.1.5)
23                 dl_dy_hat = -1 * (yi/yi_pred - (1 - yi)/(1 - yi_pred))
24                 dy_hat_dz2 = yi_pred * (1 - yi_pred)
25                 dz2_dw2 = a1.T
26                 dz2_db2 = 1
27                 dz2_dw1 = self.layers[1].weights.T * xi *
28                             self.layers[0].activation_derivative(z1)
29                 dz2_db1 = self.layers[0].activation_derivative(z1)
30
31                 """
32                 Calculating partial derivatives of loss wrt weight/bias
33                 Layer 1, and 2
34                 """
35                 dl_dw1 = np.array(dl_dy_hat * dy_hat_dz2 * dz2_dw1)
36                 dl_db1 = np.array(dl_dy_hat * dy_hat_dz2 * dz2_db1)
37                 dl_dw2 = np.array(dl_dy_hat * dy_hat_dz2 * dz2_dw2)
38                 dl_db2 = np.array(dl_dy_hat * dy_hat_dz2 * dz2_db2)
39
40                 # Updating weights and biases
41                 self.layers[1].weights -= self.learning_rate * dl_dw2
42                 self.layers[1].bias -= self.learning_rate * dl_db2
43                 self.layers[0].weights -= self.learning_rate * dl_dw1
44                 self.layers[0].bias -= self.learning_rate * dl_db1
45
46                 metrics_core = metrics(val_data[1], self.predict(val_data[0]))
47
48                 if self.metrics: # Print metrics for each epoch
49                     print(f"Epoch: {epoch+1}, metrics: {metrics_core}")
50
51             return (metrics_core, history_core)

```



### Predict Function Initialize

The predict function in an ANN models predicts the output for any given input, by passing the input through all the layers in the neural network called feed forward.

```

1 def predict(self, x_test):
2     if self.compiled:
3         y_pred = []
4
5         for xi in x_test: # different cases
6             for i in range(len(self.layers)):
7                 xi = self.layers[i].forward(xi.reshape(-1, 1))
8
9             y_pred.append(xi[0][0])
10        return np.array(y_pred)
11    else: print("Model is not yet compiled.")

```

### 5.2.5 Initializing Model Class

Define custom model with layers, and input features.

```

1 custom_model = Model(layers=[layer_1, layer_2],
2                       input_features=x_train.shape[1])
3
4 custom_model.summary() # show trainable parameters, layers, nodes, etc
5 custom_model.compile(metrics=True) # show metrics after each epoch

```

### 5.2.6 Defining Metrics Function

```

1 from sklearn.metrics import accuracy_score, f1_score, confusion_matrix
2 from sklearn.metrics import precision_recall_curve, auc
3
4 def metrics(y_true, y_hat):
5     # Convert y_hat (probability) to 0(<0.5) to 1(>0.5)
6     y_pred = np.where(y_hat > .5, 1, 0)
7
8     accuracy = accuracy_score(y_true, y_pred) # Accuracy Score
9     f1 = f1_score(y_true, y_pred) # F1 Score
10    cm = confusion_matrix(y_true, y_pred) # Confusion Matrix
11
12    # Precision-Recall AUC
13    precision, recall, _ = precision_recall_curve(y_true, y_hat)
14    pr_auc = auc(recall, precision)
15
16    return {
17        "accuracy": accuracy, "pr_auc": pr_auc,
18        "f1_score": f1, "confusion_matrix": cm
19    }

```

### 5.2.7 Training Custom ANN

```

1 metrics_core, history_core = custom_model.fit(x_train, y_train,
2                                              val_data=(x_test, y_test))

```

## 5.3 Neural Network Implementation using PyTorch

### 5.3.1 Assumptions

Already covered (please refer [5.2.1](#)).

### 5.3.2 Importing Libraries and Converting Numpy to tensor

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 from sklearn.metrics import accuracy_score
5
6 # Convert NumPy arrays to PyTorch tensors
7 x_train_tensor = torch.tensor(x_train, dtype=torch.float32)
8 y_train_tensor = torch.tensor(y_train.to_numpy(), dtype=torch.float32)
9                     .view(-1, 1)
10
11 x_test_tensor = torch.tensor(x_test, dtype=torch.float32)
12 y_test_tensor = torch.tensor(y_test.to_numpy(), dtype=torch.float32)
13                     .view(-1, 1)
```

### 5.3.3 Defining Model and Layers

Here we are making the exact same model we made using core python, this consist of sequential layers:

1. One hidden with ten nodes and relu as activation function, and
2. Another one output layer with one node and sigmoid as activation function

```
1 class BinaryClassifier(nn.Module):
2     def __init__(self, input_dim):
3         super(BinaryClassifier, self).__init__()
4
5         self.model = nn.Sequential(
6             nn.Linear(input_dim, 10), # 10 nodes in first layer
7             nn.ReLU(),
8
9             nn.Linear(10, 1), # one node in last (output) layer
10            nn.Sigmoid()
11        )
12
13    def forward(self, x):
14        return self.model(x)
```

### 5.3.4 Initializing Model

```
1 # Initialize model
2 model = BinaryClassifier(input_dim=x_train.shape[1])
3
4 criterion = nn.BCELoss() # Binary Cross Entropy as Loss Function
5 optimizer = optim.SGD(model.parameters(), lr=0.01) # Gradient Descent
```

### 5.3.5 Training PyTorch Model

```
1 history_torch = []
2 counter_torch = 0
3 time_start = time.time()
4
5 # Training loop
6 for epoch in range(1): # Single Epochs
7     model.train()
8     for xi, yi in zip(x_train_tensor, y_train_tensor):
9         if counter_torch % 200 == 0:
10             cost = criterion(model(x_test_tensor), y_test_tensor).item()
11             history_torch.append([counter_torch / 200, cost,
12                                   time.time()-time_start])
13             counter_torch = counter_torch + 1
14
15             optimizer.zero_grad() # Making gradient zero
16
17             loss = criterion(model(xi), yi)
18             loss.backward() # Calculate gradient
19
20             optimizer.step()
```

### 5.3.6 Model Evaluation

```
1 model.eval()
2 with torch.no_grad():
3     y_preds = model(x_test_tensor).detach().numpy()
4
5 metrics_pytorch = metrics(y_test, y_preds)
```

## 6 Evaluation and Analysis

### 6.1 Convergence Time

#### 6.1.1 Convergence of PyTorch is faster than core?

The time it takes for the model to learn enough to predict well is called **convergence time**. After this time the loss reaches minimum, and the graph flattens.

The PyTorch model converged slightly faster than the **core Python** implementation, both in terms of iterations and training time. This difference wasn't very large (check figure 7 mentioned below), some of the possible reasons can be:

1. I didn't use DataLoader in PyTorch.
2. Updated weight/biases after each row, that means in each epoch parameters are updated more than 60k times.
3. Both models were trained on CPU.

### 6.1.2 Comparision of training time and cost

```

1  sns.set(style="whitegrid") # style to show grid
2
3  # convert to numpy array
4  history_core = np.array(history_core)
5  history_torch = np.array(history_torch)
6
7  # set figure size
8  plt.figure(figsize=(10, 4))
9
10 # Red Line for Core Implementation
11 sns.lineplot(x=history_core[:117, 2], y=history_core[:117, 1],
12              color="red", label="Core Implementation")
13
14 # Blue Line for PyTorch Implementation
15 sns.lineplot(x=history_torch[:100, 2], y=history_torch[:100, 1],
16              color="blue", label="PyTorch Implementation")
17
18 plt.xlabel("Time (seconds)")
19 plt.ylabel("Cost")
20 plt.title("Time(second) v/s Cost: for PyTorch and Core Implementation")
21
22 plt.tight_layout()
23 plt.show()

```

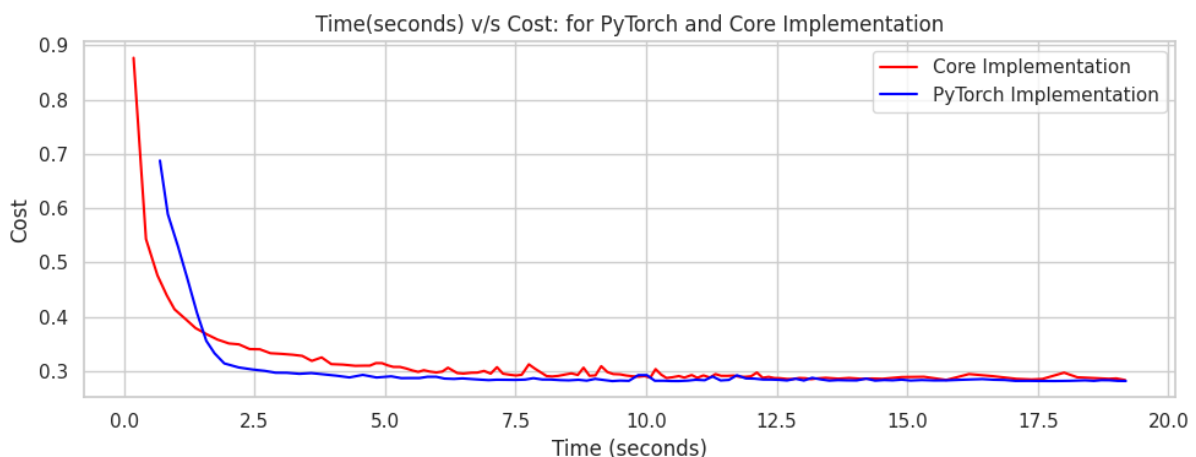


Figure 7: Time (seconds) v/s Loss: for PyTorch and Core Implementation

#### Conclusion:

1. Here you can see that the convergence time for PyTorch was comparatively less than Core Implementation.
2. The cost for PyTorch got more and more stable with time,
3. But with the other one i.e Core Implementation it was not the case, we can still see some sharp pitches in the graph.

## 6.2 Performance Metrics

We compare the two models built above using the following metrics:

### 6.2.1 Accuracy

The ratio of correct predictions out of all the predictions is called **accuracy** for the given model. Higher the accuracy more better out predictions.

$$\text{Accuracy} = \frac{\text{total correct predictions}}{\text{total predictions}}$$

Implementation	Accuracy
PyTorch Implementations	83.94%
Core Implementations	83.86%

Table 2: Accuracy comparison between **Core** and **PyTorch Implementations**

here you can see, the difference in accuracy is not quite large between both the implementations. Both the models are predicting around  $\approx 84\%$  *correctly*.

### 6.2.2 F1 Score

The harmonic mean of the precision and recall is called **F1 Score** for the given model. Where as **precision** is the ratio of true predicted positives by total predicted positives, and **recall** is the ratio of true predicted positives by total true positives.

Higher the **F1 Score** means both **precision**, and **recall** is high, which identify that our model is correctly classifying right positives and not misclassifying too many negatives as positives.

$$F1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Precision} = \frac{\text{true predicted positives}}{\text{total predicted positives}}, \quad \text{Recall} = \frac{\text{true predicted positives}}{\text{total positives}}$$

Implementation	F1 Score
PyTorch Implementations	76.71%
Core Implementations	75.65%

Table 3: F1-Score comparison between **Core** and **PyTorch Implementations**

The difference in the **F1 Score** is nearly  $\approx 1\%$ , which signifies that **PyTorch** model is more nicely predicting compared to **Core Implementation** model.

### 6.2.3 Precision-Recall Area Under Curve (PR AUC)

It measures the model's ability to distinguish between categories, especially when dealing with imbalanced datasets (e.g., 90% negatives, 10% positives).

Higher the PR AUC means higher the ability of the model to correctly classify between positives and negatives.

Implementation	PR AUC
PyTorch Implementations	73.75%
Core Implementations	73.51%

Table 4: PR-AUC comparison between `Core` and `PyTorch` Implementations

Although the margin is not quite large, but `PyTorch` model is classifying more correctly positives and negatives as that compared to `Core Implementation` model.

## 6.3 Memory Usage

1. The `PyTorch` implementation consumed more memory than the `Core Implementation` model
2. As `PyTorch` includes additional overhead for managing tensors, autograd (automatic differentiation), and internal buffers.
3. On the other hand, the `Core Implementation` model only stores the bare minimum: weights, inputs, and gradients manually.

## 6.4 Confusion Matrix and Inference

A `Confusion Matrix` is a performance measurement tool for classification models. It tells how well a model is performing by showing the correct and incorrect predictions made by the model compared to the actual outcomes.

Actual / Predicted	Positive	Negative
Positive	True Positive (TP)	False Negative (FN)
Negative	False Positive (FP)	True Negative (TN)

Table 5: Typical Confusion Matrix for Binary Classification

With the help of this, we can easily calculate other performance metrics such as:

1. Accuracy
2. F1 Score
3. Precision-Recall Area Under Curve

```

1 # Plot 2 graphs on same axis, with required size
2 fig, axs = plt.subplots(1, 2, figsize=(12, 5))
3
4 # Confusion Matrix for both implementation
5 cm_core = metrics_core["confusion_matrix"]
6 cm_torch = metrics_pytorch["confusion_matrix"]
7
8 # Plotting graph for Core Implementation
9 sns.heatmap(cm_core, annot=True, fmt="d", cmap="Reds", ax=axs[0])
10 axs[0].set_title("Confusion Matrix: Core Implementation")
11 axs[0].set_xlabel("Predicted")
12 axs[0].set_ylabel("Actual")
13
14 # Plotting graph for PyTorch Implementation
15 sns.heatmap(cm_torch, annot=True, fmt="d", cmap="Blues", ax=axs[1])
16 axs[1].set_title("Confusion Matrix: PyTorch")
17 axs[1].set_xlabel("Predicted")
18 axs[1].set_ylabel("Actual")
19
20 plt.tight_layout()
21 plt.show()

```

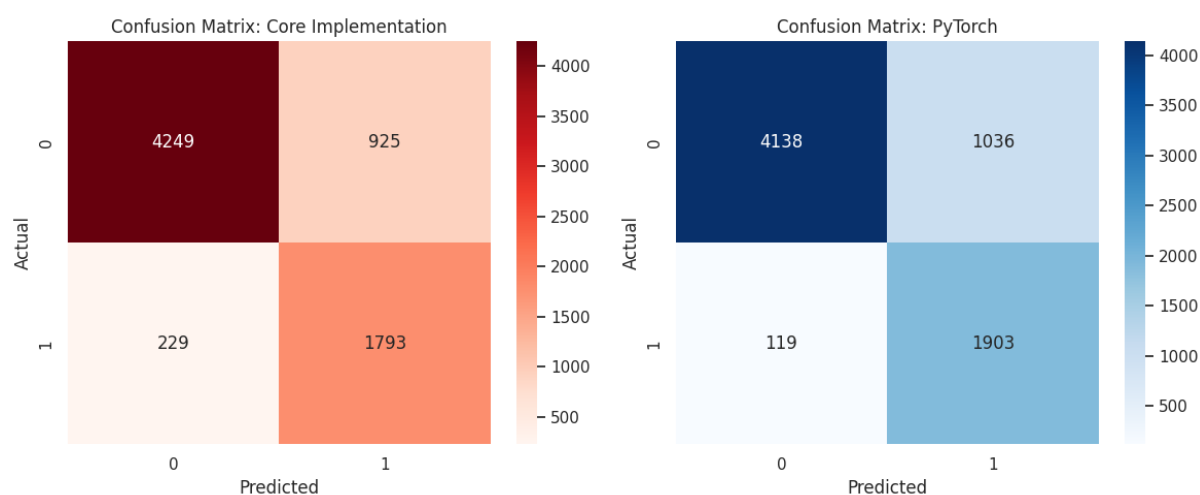


Figure 8: Confusion Matrix: for Both Implementation

## Conclusions

1. Overall accuracy is almost same in both implementations.
2. PyTorch model is better at identifying true positives compared to Core Implementation model.
3. Core Implementation model is better at identifying true negatives compared to PyTorch model.
4. PyTorch model makes more mistake identifying false positives, and Core Implementation model makes more mistakes identifying false negatives compared to each other.
5. If you want higher true positives then go with PyTorch, else Core Implementation model will be better for higher true negatives.

## 6.5 Analysis and Discussion

### 6.5.1 Convergence Speed, Performance, and Memory Usage

Discussion for these topics have been already done above in their respective domains:

1. Convergence Speed [6.1](#)
2. Performance Metrics [6.2](#)
3. Memory Usage [6.3](#)

### 6.5.2 Some Other Factors

#### Optimizations

1. You can use `DataLoaders` in PyTorch model which significantly improves the performance.
2. Instead of updating `weights/biases` after each rows which is like if a data set contains 60 k rows, then in two epochs parameters are updates nearly 120 k times. Instead this you can use batch gradient descent.
3. Change in `learning rate` can help reaching convergence quickly.

#### Hardware Acceleration

1. Here we used `CPU` in both the model which is slow for these type of tasks, In PyTorch model we could upgrade the tensors to use `GPU` which are significantly faster.
2. Distributing load to multiple `GPU`, or if possible multiples machines can reduce training time.

#### Numerical Stability

Standardizing all the values in the dataset, can help reduce the training time and help reach convergence fast. Example in a dataset containing fields `IQ`, `CPI` and `salary`, predicting salary can be not efficient as it can be in lakhs and other features will be in hundreds.

#### Code Efficiency

1. Instead of writing all code in single go, better approach will be to define `function` of some general tasks.
2. If some functions share related task, then define `class` for it.
3. Efficiently writing loops, Correctly naming `variables/functions/classes`, and more.

---

Thank You