

# Artificial Neural Network Project Report

Ammar Saleem(i210327) , Wajeeh ul Hassan(i210337)

May 18, 2024

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Data Preprocessing</b>	<b>3</b>
2.1	Loading the Dataset . . . . .	3
2.2	Data Cleaning and Transformation . . . . .	3
<b>3</b>	<b>Model Training</b>	<b>5</b>
<b>4</b>	<b>Model Evaluation</b>	<b>7</b>
<b>5</b>	<b>Visual Diagrams</b>	<b>8</b>
5.1	Correlation Heatmap . . . . .	9
5.2	Feature Distributions . . . . .	10
<b>6</b>	<b>Challenges Faced</b>	<b>11</b>
6.1	Unbalanced Data . . . . .	11
6.2	Model Overfitting . . . . .	11
<b>7</b>	<b>Conclusion</b>	<b>12</b>

# Chapter 1

## Introduction

The purpose of this project is to develop an Artificial Neural Network (ANN) model to predict loan approval status based on various financial and personal attributes. The project involves preprocessing the dataset, defining the ANN architecture, training the model, and evaluating its performance. Additionally, the report discusses the challenges faced during the project and provides visual diagrams to better understand the data distribution and model performance.

# Chapter 2

## Data Preprocessing

Data preprocessing is a crucial step in building a machine learning model. It involves cleaning the data, handling missing values, encoding categorical variables, and scaling numerical features.

### 2.1 Loading the Dataset

The dataset used in this project is the loan approval dataset, which contains various features such as income, age, experience, marital status, house ownership, car ownership, profession, city, and state. The target variable is the ‘*RiskFlag*’, indicating whether the loan was approved (1) or not (0).

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3
4 df = pd.read_csv('loan_approval_dataset_preprocessed.csv')
5 X_train, X_test, y_train, y_test = train_test_split(
6     df.drop("RiskFlag", axis=1), df["RiskFlag"], test_size
    =0.2, stratify=df["RiskFlag"], random_state=42)
```

Listing 2.1: Loading the Dataset

### 2.2 Data Cleaning and Transformation

The data cleaning process includes handling missing values, encoding categorical variables, and scaling numerical features. Below is the code snippet for this preprocessing step.

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 df = pd.read_json('/content/loan_approval_dataset.json')
6 df = df.drop(['Id'], axis=1)
7
8 categorical_columns = ['Risk_Flag', 'Married/Single', '
    House_Ownership', 'Car_Ownership', 'Profession', 'CITY', '
    STATE']
9
10 # Encoding categorical variables
11 for column in categorical_columns:
12     frequency_map = df[column].value_counts(normalize=True)
13     df[f'{column}'] = df[column].map(frequency_map)
14
15 # Scaling numerical features
16 from sklearn.preprocessing import MinMaxScaler
17 numeric_columns = ['Income', 'Age', 'Experience', '
    CURRENT_JOB_YRS', 'CURRENT_HOUSE_YRS']
18 scaler = MinMaxScaler()
19 df[numeric_columns] = scaler.fit_transform(df[numeric_columns])
20
21 df.to_csv('loan_approval_dataset_preprocessed.csv', index=False)

```

Listing 2.2: Data Cleaning and Transformation

# Chapter 3

## Model Training

The neural network model is defined using PyTorch, and it consists of three fully connected layers. The model is trained using the Binary Cross-Entropy Loss function and the Adam optimizer.

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 from torch.utils.data import DataLoader, TensorDataset
5
6 # Convert data to PyTorch tensors
7 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
8 X_train_tensor = torch.tensor(X_train.values, dtype=torch.float32).to(device)
9 y_train_tensor = torch.tensor(y_train.values, dtype=torch.float32).to(device)
10 X_test_tensor = torch.tensor(X_test.values, dtype=torch.float32).to(device)
11 y_test_tensor = torch.tensor(y_test.values, dtype=torch.float32).to(device)
12
13 # Create TensorDatasets
14 train_dataset = TensorDataset(X_train_tensor, y_train_tensor)
15 test_dataset = TensorDataset(X_test_tensor, y_test_tensor)
16
17 # Create DataLoaders
18 batch_size = 64
19 train_loader = DataLoader(dataset=train_dataset, batch_size=batch_size, shuffle=True)
```

```

20 test_loader = DataLoader(dataset=test_dataset, batch_size=
    batch_size, shuffle=False)
21
22 # Define the neural network architecture
23 class NeuralNetwork(nn.Module):
24     def __init__(self, input_size):
25         super(NeuralNetwork, self).__init__()
26         self.fc1 = nn.Linear(input_size, 64)
27         self.fc2 = nn.Linear(64, 32)
28         self.fc3 = nn.Linear(32, 1)
29         self.relu = nn.ReLU()
30         self.sigmoid = nn.Sigmoid()
31
32     def forward(self, x):
33         x = self.relu(self.fc1(x))
34         x = self.relu(self.fc2(x))
35         x = self.sigmoid(self.fc3(x))
36         return x
37
38 # Initialize the model and move it to GPU
39 input_size = X_train.shape[1]
40 model = NeuralNetwork(input_size).to(device)
41
42 # Define loss function and optimizer
43 criterion = nn.BCELoss()
44 optimizer = optim.Adam(model.parameters(), lr=0.001)
45
46 # Train the model
47 num_epochs = 100
48 for epoch in range(num_epochs):
49     model.train()
50     running_loss = 0.0
51     for inputs, labels in train_loader:
52         optimizer.zero_grad()
53         outputs = model(inputs)
54         loss = criterion(outputs, labels.view(-1, 1))
55         loss.backward()
56         optimizer.step()
57         running_loss += loss.item() * inputs.size(0)
58     epoch_loss = running_loss / len(train_loader.dataset)
59     print(f"Epoch {epoch+1}/{num_epochs}, Loss: {epoch_loss:.4f}")

```

Listing 3.1: Model Training

## Chapter 4

# Model Evaluation

After training the model, it is essential to evaluate its performance on the test set. The accuracy of the model is calculated as follows:

```
1 # Evaluate the model
2 model.eval()
3 correct = 0
4 total = 0
5 with torch.no_grad():
6     for inputs, labels in test_loader:
7         outputs = model(inputs)
8         predicted = torch.round(outputs)
9         total += labels.size(0)
10        correct += (predicted == labels.view(-1, 1)).sum().item
11    ()
12 accuracy = correct / total
13 print(f"Test Accuracy: {accuracy:.4f}")
```

Listing 4.1: Model Evaluation



# Chapter 5

## Visual Diagrams

Visual diagrams are useful to understand the data distribution and the correlations between different features. Below are some visualizations created using Matplotlib and Seaborn.

## 5.1 Correlation Heatmap

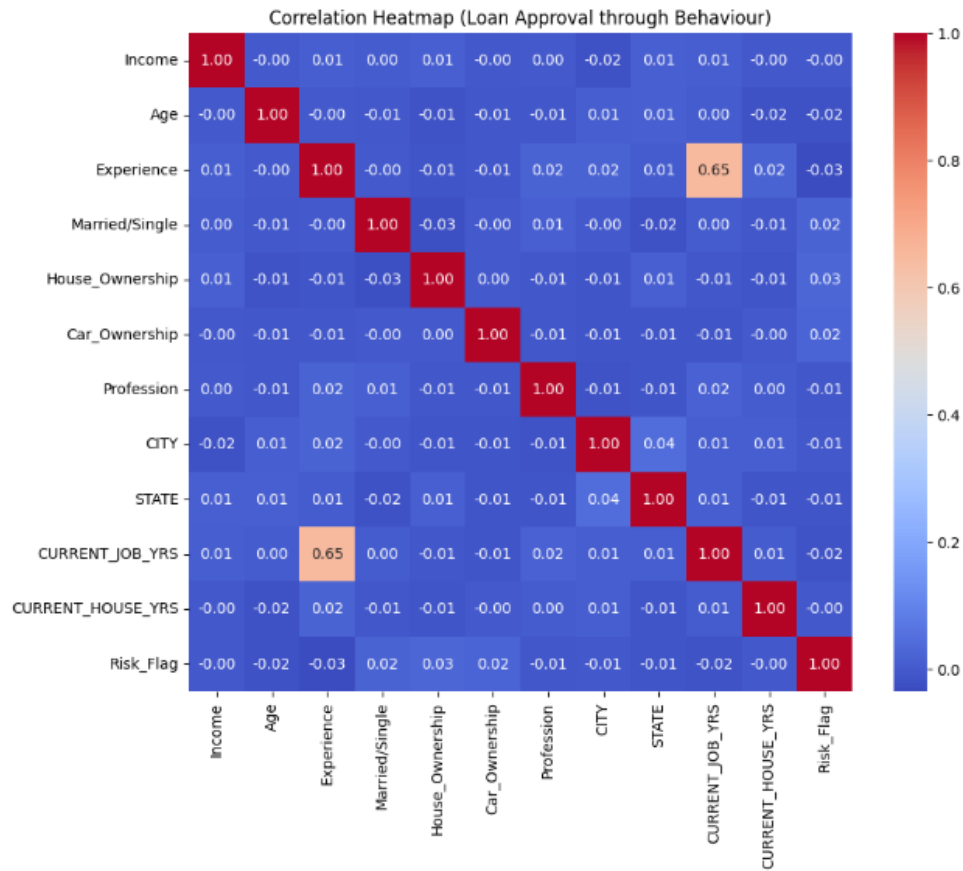


Figure 5.1: Correlation Heatmap (Loan Approval through Behaviour)

## 5.2 Feature Distributions

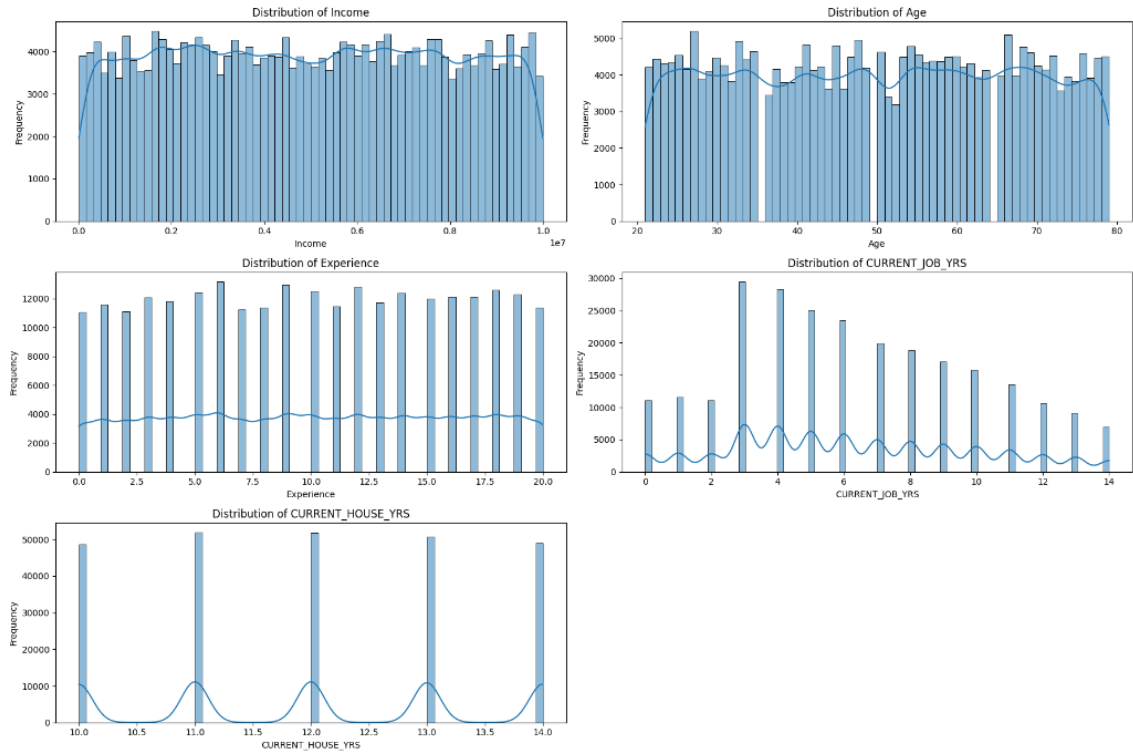


Figure 5.2: Distributions of Numeric Features

# Chapter 6

## Challenges Faced

### 6.1 Unbalanced Data

One of the significant challenges faced during this project was the unbalanced dataset. The target variable ' $Risk_{Flag}$ ' had a disproportionate number of approved and unapproved loans, which could lead to a biased model.

### 6.2 Model Overfitting

Another common challenge was the potential for model overfitting. Given the complexity of neural networks, there was a risk of the model learning the training data too well and not generalizing effectively to new data.

## Chapter 7

### Conclusion

This project demonstrated the application of an Artificial Neural Network for predicting loan approval status. Despite challenges such as unbalanced data and potential overfitting, the model achieved a satisfactory level of accuracy. Future work could involve exploring more sophisticated models, addressing class imbalance through techniques such as SMOTE, and performing hyperparameter tuning to further enhance model performance.