# Campus Placement Prediction Report

**Student:** Devansh Patel (C0928483)

**Course:** AML-3104 Neural Networks and Deep Learning

**Submitted to:** Ishant Gupta

## Introduction

Predicting a student's likelihood of being placed on campus based on a variety of academic and demographic criteria is the aim of this project. Employability test scores, job experience, specialisation, degree %, secondary education percentage, and higher secondary education percentage are among the features included in the dataset. This report outlines the dataset, preprocessing steps, model selection, evaluation, and results.

---

## DataSet Description

**The dataset contains the following features:**

- Numerical Features:

    - **ssc_p:** Secondary education percentage (10th grade).

    - **hsc_p:** Higher secondary education percentage (12th grade).

    - **degree_p:** Degree percentage (undergraduate).

    - **etest_p:** Employability test percentage.

    - **mba_p:** MBA percentage.

- **Categorical Features:**

    - **gender:** Gender of the student (Male/Female).

    - **ssc_b:** Board of education for secondary education (Central/Others).

    - **hsc_b:** Board of education for higher secondary education (Central/Others).

    - **hsc_s:** Stream in higher secondary education (Commerce/Science/Arts).

    - **degree_t:** Type of degree (Sci&Tech/Comm&Mgmt).

- o **workex:** Work experience (Yes/No).

- o **specialisation:** Postgraduate specialization (Mkt&HR/Mkt&Fin).

- **Target Variable:**

  - o **status:** Placement status (Placed/Not Placed).

---

## Preprocessing Steps

The following preprocessing steps were applied to prepare the dataset for modeling:

1. **Handling Missing Values:**

   - o Since it was irrelevant for predicting placement status, the pay column— which included missing values for students who were not placed—was removed.

2. **Encoding Categorical Variables:**

   - o To transform them into numerical format, categorical features including gender, ssc_b, hsc_b, hsc_s, degree_t, workex, and specialisation were one-hot encoded.

3. **Splitting the Dataset:**

   - o To assess model performance, the dataset was divided into training (70%) and testing (30%) sets.

4. **Feature Scaling:**

   - o To guarantee that every feature was on the same scale, numerical features were standardised.

---

## Model Selection

Three models were chosen for this project:

1. **The Logistic Regression Model:**

   - o An easy-to-understand paradigm that works well for binary classification tasks.

o It helps comprehend the connection between attributes and the goal variable and offers a baseline performance.

**2. Random Forest:**

   o An ensemble approach that manages feature interactions and non-linear relationships.
   o It performs well with both numerical and categorical data and is resistant to overfitting.

**3. XGBoost:**

   o a strong gradient boosting technique with a good accuracy rate.
   o It offers feature importance and manages unbalanced datasets effectively.

**4. Voting Classifier:**

   o XGBoost, Random Forest, and Logistic Regression combined with soft voting.
   o It enhances overall performance by utilising the advantages of all three models.

---

# Model Evaluation

The following measures were used to assess the models:

   o **Accuracy:** The percentage of placements that were accurately anticipated.

   o **Precision:** The percentage of accurately predicted placements.

   o **Recall:** The percentage of actual placements that were accurately anticipated is known as recall.

   o **F1-Score:** The precision and recall harmonic mean.

**Results:**

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 0.85 | 0.86 | 0.90 | 0.88 |
| Random Forest | 0.88 | 0.89 | 0.92 | 0.90 |
| XGBoost | 0.89 | 0.90 | 0.93 | 0.91 |
| Voting Classifier | 0.90 | 0.91 | 0.94 | 0.92 |

**Confusion Matrices:**

- o To visualise true positives, true negatives, false positives, and false negatives, confusion matrices were plotted for every model.

**ROC Curves:**

- o The true positive rate (TPR) and false positive rate (FPR) for every model were compared using ROC curves. The Voting Classifier's AUC score was the highest.

---

# Conclusion

- o With an F1-score of 0.92 and an accuracy of 90%, the Voting Classifier outperformed the rest.

- o Specialisation, etest_p, and degree_p are important variables that affect placement forecasts.

- o Based on student information, the model can be used as a Streamlit app to forecast placement status.