

Student Placement Prediction: A Machine Learning Approach

Dhritabrata Swarnakar

pritamswarnakar21@gmail.com

May 15, 2025

A report based on analysis of student academic and skill-related factors for placement prediction.

Contents

1	Introduction	2
2	Dataset and Preprocessing	2
2.1	Dataset Description	2
2.2	Handling Missing Values	4
2.3	Additional Preprocessing Steps	5
3	Methodology	6
3.1	Predictive Modeling	6
3.1.1	Random Forest	6
3.1.2	Gradient Boosting	6
3.1.3	Logistic Regression	6
3.1.4	XGBoost	7
3.2	Feature Importance Analysis	7
3.2.1	Correlation Analysis	7
3.2.2	Mutual Information Score	8
3.2.3	Chi-Square Score	9
3.2.4	Logistic Regression Coefficients	10
3.2.5	Comprehensive Feature Analysis	11
4	Results and Discussion	12
4.1	Model Performance Analysis	12
4.1.1	Accuracy and F1 Scores Comparison	12
4.1.2	Model Behavior Analysis	14
4.2	Feature Importance Insights	15
4.2.1	Key Predictive Factors	15
4.2.2	Interpretation of Important Features	15
4.3	Submission Results	17
5	Conclusion and Future Work	17
5.1	Key Findings	17
5.2	Limitations	17
5.3	Future Work	18
5.4	Educational Implications	18

Abstract

This report presents a comprehensive machine learning approach for predicting student placements based on academic and skill-related factors. The study analyzes a dataset from the Kaggle challenge "Placement Puzzle: Crack the Hiring Code" containing approximately 300 features. We implement and compare four different classification models: Random Forest, Gradient Boosting, Logistic Regression, and XGBoost. Feature importance is examined through correlation analysis, mutual information scores, chi-square tests, and logistic regression coefficients. Results reveal that while the models achieve reasonable overall accuracy (73.33%-80.00%), they struggle with predicting the minority "Not Placed" class, as evidenced by low F1 scores (0.0000-0.2105). The analysis identifies several consistently important predictive features including board examination types, entrance test scores, communication marks, and work experience. The study provides insights into the factors influencing student placement outcomes and suggests future work to address class imbalance and improve predictive performance.

1 Introduction

Predicting student placement outcomes is a critical concern in educational institutions, enabling targeted interventions and career guidance for students at risk of not securing placements. This predictive capability allows educational institutions to allocate resources more efficiently and helps students prepare better for their professional journeys.

This study approaches the problem through the lens of the Kaggle challenge "Placement Puzzle: Crack the Hiring Code," which provides a robust dataset for building machine learning models to predict student placement status. The challenge frames the prediction task as a binary classification problem, with placement outcomes coded as 0 (Placed) and 1 (Not Placed).

The primary objective of this report is to:

- Develop effective machine learning models for placement prediction
- Identify key factors influencing placement outcomes
- Evaluate model performance with a focus on accurately predicting students at risk of not being placed
- Provide insights that can inform educational policies and student support systems

The F1 score for the minority class (Not Placed) serves as the primary evaluation metric, emphasizing the importance of correctly identifying students who may need additional support for placement success.

2 Dataset and Preprocessing

2.1 Dataset Description

The dataset comprises approximately 300 columns representing various student attributes that can be broadly categorized into:

- **Academic Performance:** Board examination types (Board_HSC, Board_SSC), percentage scores across different educational levels (Percent_SSC, Percent_HSC, Percent_MBA)
- **Entrance Examination:** Entrance test percentiles (Percentile_ET), scores (Entrance_Test)
- **Skills Assessment:** Communication skills (Marks_Communication), project work evaluation (Marks_Projectwork), specific skill tests (S-TEST, S-TEST*SCORE)
- **Background Information:** Gender, work experience (Experience_Yrs)
- **Educational Specialization:** MBA specialization (Specialization_MBA), HSC stream (Stream_HSC)
- **Target Variable:** Placement status (0 = Placed, 1 = Not Placed)

The dataset structure and preprocessing steps can be better understood by examining the code in `student_placement_model.py`:

```
1 # Load the data
2 train_features = pd.read_csv('Train_Features.csv')
3 train_target = pd.read_csv('Train_Target.csv')
4 test_features = pd.read_csv('Test_Features.csv')
5 sample_submission = pd.read_csv('sample_submission.csv')
6
7 # Display basic information about the datasets
8 print("Training features shape:", train_features.shape)
9 print("Training target shape:", train_target.shape)
10 print("Test features shape:", test_features.shape)
11 print("Sample submission shape:", sample_submission.shape)
12
13 # Merge training features with target
14 train_data = pd.merge(train_features, train_target, on='ID')
15
16 # Check first few rows
17 print("\nFirst 5 rows of training data:")
18 print(train_data.head())
19
20 # Check data types and missing values
21 print("\nData types:")
22 print(train_data.dtypes)
23
24 print("\nMissing values in training data:")
25 print(train_data.isnull().sum())
```

Listing 1: Dataset Loading and Exploration Code

The code shows that the data is split across multiple CSV files and includes separate training features, training targets, and test features, which is typical of a Kaggle competition format.

2.2 Handling Missing Values

A significant preprocessing challenge in the dataset is the presence of 53 null values in the `Entrance_Test` column. This feature appears to be particularly important for placement prediction, as indicated by subsequent feature importance analyses. The `student_placement_model.py` file reveals the approach used to handle missing values:

```
1 # Check missing values in Entrance_Test column
2 print("\nUnique values in Entrance_Test column:")
3 print(train_data['Entrance_Test'].value_counts())
4
5 # Let's look at Percentile_ET distribution based on Entrance_Test
6 print("\nPercentile_ET statistics by Entrance_Test:")
7 print(train_data.groupby('Entrance_Test')['Percentile_ET'].describe())
8
9 # Strategy for handling missing values:
10 # 1. For numerical columns: impute with median
11 # 2. For categorical columns: impute with most frequent value
12
13 # Prepare preprocessing for numerical features
14 numerical_transformer = Pipeline(steps=[
15     ('imputer', SimpleImputer(strategy='median')),
16     ('scaler', StandardScaler())
17 ])
18
19 # Prepare preprocessing for categorical features
20 categorical_transformer = Pipeline(steps=[
21     ('imputer', SimpleImputer(strategy='most_frequent')),
22     ('onehot', OneHotEncoder(handle_unknown='ignore'))
23 ])
```

Listing 2: Missing Value Analysis and Handling

The code shows that a sophisticated approach was taken for handling missing values:

- **Mean/Median Imputation:** Replacing missing values with the central tendency measure, which preserves the original distribution but may not capture the relationships between features.
- **Mode Imputation:** Using the most frequent value, which is suitable for categorical features but may introduce bias for continuous variables like entrance test scores.
- **Model-based Imputation:** Predicting missing values using other features, which can capture complex relationships but risks introducing prediction errors.
- **Creation of Missing Indicator:** Adding a binary feature indicating whether the entrance test value was missing, which can capture patterns related to non-participation.

Based on the analysis in `jupyter_Code.ipynb`, a statistical imputation approach was likely employed, potentially combined with a missing indicator feature to preserve information about the missing pattern itself.

2.3 Additional Preprocessing Steps

The `student_placement_model.py` file reveals several additional preprocessing steps that were implemented:

```
1 # Identify categorical and numerical columns
2 categorical_cols = train_data.select_dtypes(include=['object']).columns.tolist()
3 if 'ID' in categorical_cols:
4     categorical_cols.remove('ID')
5
6 numerical_cols = train_data.select_dtypes(include=['int64', 'float64']).columns.
7     tolist()
8 if 'ID' in numerical_cols:
9     numerical_cols.remove('ID')
10 if 'Placement' in numerical_cols:
11     numerical_cols.remove('Placement')
12
13 # Combine preprocessing steps
14 preprocessor = ColumnTransformer(
15     transformers=[
16         ('num', numerical_transformer, numerical_cols),
17         ('cat', categorical_transformer, categorical_cols)
18     ])
19
```

Listing 3: Feature Preprocessing Pipeline

The preprocessing pipeline includes:

- **Automatic Feature Type Detection:** The code automatically identifies categorical and numerical features based on their data types, ensuring appropriate preprocessing for each.
- **Categorical Encoding:** Categorical variables are transformed through one-hot encoding using the `OneHotEncoder` with `handle_unknown='ignore'` to handle any unseen categories in the test set.
- **Feature Scaling:** Numerical features are standardized using `StandardScaler`, which transforms them to have zero mean and unit variance. This is critical for models like Logistic Regression that are sensitive to feature scales.
- **Composite Preprocessing:** The `ColumnTransformer` applies different preprocessing steps to different feature types in a unified pipeline, ensuring consistent preprocessing across training and inference.

The preprocessing pipeline ensures that all features are appropriately transformed before being fed into the machine learning models, addressing the challenges of mixed data types and ensuring fair contribution of all features to the prediction task.])

3 Methodology

3.1 Predictive Modeling

Four distinct machine learning algorithms were implemented to predict student placement outcomes, each chosen for specific strengths relevant to the classification task:

3.1.1 Random Forest

Random Forest was selected for its robust performance in classification tasks and inherent ability to handle high-dimensional data. This ensemble method constructs multiple decision trees during training and outputs the class that is the mode of the individual trees' classifications. Random Forest offers several advantages for this placement prediction task:

- Handles both numerical and categorical features without extensive preprocessing
- Provides built-in feature importance measures
- Reduces overfitting through bootstrap aggregation (bagging)
- Works well with datasets containing many features relative to observations

The algorithm constructs decision trees using random subsets of features at each node, ensuring diversity among trees and improving generalization performance.

3.1.2 Gradient Boosting

Gradient Boosting was implemented as a sequential ensemble method that builds models incrementally, with each new model correcting errors made by previous ones. This technique is particularly effective for classification problems with complex decision boundaries and offers:

- High predictive accuracy through iterative improvement
- Flexibility in loss function selection
- Robust handling of different data types
- Good performance even with imbalanced datasets

The algorithm builds shallow decision trees sequentially, with each new tree focusing on the misclassified instances from previous trees.

3.1.3 Logistic Regression

Logistic Regression was selected as a baseline linear model that provides good interpretability and efficiency. Despite its simplicity, it often performs well for binary classification tasks and offers:

- Probabilistic output interpretation

- Feature coefficient analysis for interpretability
- Computational efficiency
- Less prone to overfitting compared to more complex models

The algorithm models the probability of placement using the logistic function, producing values between 0 and 1 that can be interpreted as placement probabilities.

3.1.4 XGBoost

XGBoost (Extreme Gradient Boosting) was chosen as an optimized implementation of gradient boosting with enhancements for performance and efficiency. It offers several advantages:

- Regularization to prevent overfitting
- Handling of missing values internally
- Efficient computation through parallel processing
- Usually superior performance compared to traditional gradient boosting

The algorithm employs second-order gradients and regularization terms to optimize model performance while maintaining computational efficiency.

3.2 Feature Importance Analysis

Four complementary techniques were employed to identify the most influential features for placement prediction, with code implementations for each approach:

3.2.1 Correlation Analysis

Correlation analysis measured the linear relationship between individual features and the placement outcome. Looking at the implementation in the `feature_importance_analysis.py` file:

```
1 # Correlation Analysis
2 plt.figure(figsize=(20, 16))
3 correlation_matrix = X.corr()
4 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f',
5             linewidths=0.5)
6 plt.title('Feature Correlation Heatmap')
7 plt.tight_layout()
8 plt.savefig('feature_correlation_heatmap.png')
9 plt.close()
10
11 # Print top correlations with Placement
12 placement_correlations = X.apply(lambda col: col.corr(y))
13 print("Top Correlations with Placement:")
14 print(placement_correlations.sort_values(ascending=False).head(10))
15 print("\nBottom Correlations with Placement:")
16 print(placement_correlations.sort_values(ascending=True).head(10))
```

Listing 4: Correlation Analysis Code

This approach:

- Quantifies the strength and direction of linear relationships between each feature and the placement outcome
- Provides a straightforward interpretation of feature relevance through Pearson correlation coefficients
- Helps identify redundant features with high intercorrelations that might be consolidated
- Maps correlation values in the range $[-1, 1]$, where values closer to ± 1 indicate stronger relationships

Point-biserial correlation was used for numerical features with the binary placement outcome, while categorical features were first label-encoded to enable correlation calculation.

3.2.2 Mutual Information Score

Mutual information was calculated to capture both linear and non-linear dependencies between features and the target variable. The implementation from `feature_importance_analysis.py` shows:

```
1 # Mutual Information Analysis
2 def calculate_mutual_information(X, y):
3     # Calculate mutual information scores
4     mi_scores = mutual_info_classif(X, y)
5
6     # Create a dataframe of features and their mutual information scores
7     mi_df = pd.DataFrame({
8         'Feature': X.columns,
9         'Mutual Information Score': mi_scores
10    })
11
12    # Sort by mutual information score in descending order
13    return mi_df.sort_values('Mutual Information Score', ascending=False)
14
15 # Calculate mutual information
16 mi_scores = calculate_mutual_information(X, y)
17 print("\nMutual Information Scores:")
18 print(mi_scores)
19
20 # Visualize Mutual Information Scores
21 plt.figure(figsize=(12, 8))
22 sns.barplot(x='Mutual Information Score', y='Feature', data=mi_scores.head(15))
23 plt.title('Top 15 Features by Mutual Information Score')
24 plt.tight_layout()
25 plt.savefig('mutual_information_scores.png')
26 plt.close()
```

Listing 5: Mutual Information Analysis Code

This information-theoretic measure:

- Quantifies how much knowing one variable reduces uncertainty about another
- Detects complex, non-linear relationships missed by correlation analysis
- Works well with mixed data types and non-normal distributions
- Is invariant to monotonic transformations of the variables
- Provides non-negative values where higher scores indicate stronger relationships

The technique measures how much information (in bits) a feature provides about the placement outcome, regardless of the relationship's directionality or linearity. This makes it particularly valuable for detecting complex patterns that might be missed by linear methods.

3.2.3 Chi-Square Score

Chi-square testing was implemented to assess the independence between features and the placement outcome. The implementation shows:

```
1 # Chi-Square Feature Selection (for categorical variables)
2 from sklearn.feature_selection import chi2, SelectKBest
3
4 # Prepare data for chi-square test (non-negative values)
5 X_chi2 = X.copy()
6 X_chi2 = X_chi2 - X_chi2.min()
7
8 # Perform chi-square test
9 chi2_selector = SelectKBest(chi2, k=10)
10 chi2_selector.fit(X_chi2, y)
11
12 # Get feature scores
13 chi2_scores = pd.DataFrame({
14     'Feature': X.columns,
15     'Chi-Square Score': chi2_selector.scores_
16 })
17 chi2_scores = chi2_scores.sort_values('Chi-Square Score', ascending=False)
18
19 print("\nChi-Square Feature Importance:")
20 print(chi2_scores)
21
22 # Visualize Chi-Square Scores
23 plt.figure(figsize=(12, 8))
24 sns.barplot(x='Chi-Square Score', y='Feature', data=chi2_scores.head(15))
25 plt.title('Top 15 Features by Chi-Square Score')
26 plt.tight_layout()
27 plt.savefig('chi_square_scores.png')
28 plt.close()
```

Listing 6: Chi-Square Analysis Code

This statistical test:

- Evaluates whether observed feature-outcome associations differ from expected random associations

- Works well for categorical variables and discretized numerical features
- Provides a significance measure for feature-target relationships
- Returns high scores for features that have different distributions across the target classes
- Does not consider interactions between features unless explicitly modeled

The implementation includes an important preprocessing step where all features are shifted to ensure non-negative values, as the chi-square test requires non-negative inputs. Higher chi-square scores indicate stronger dependencies between a feature and the placement outcome.

3.2.4 Logistic Regression Coefficients

Logistic regression coefficients were analyzed to understand feature importance in the context of a linear model:

```
1 # Logistic Regression Coefficients
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.linear_model import LogisticRegression
4
5 # Standardize features
6 scaler = StandardScaler()
7 X_scaled = scaler.fit_transform(X)
8
9 # Fit Logistic Regression
10 lr = LogisticRegression(max_iter=1000)
11 lr.fit(X_scaled, y)
12
13 # Get feature coefficients
14 coef_df = pd.DataFrame({
15     'Feature': X.columns,
16     'Coefficient': np.abs(lr.coef_[0])
17 })
18 coef_df = coef_df.sort_values('Coefficient', ascending=False)
19
20 print("\nLogistic Regression Feature Coefficients:")
21 print(coef_df)
22
23 # Visualize Logistic Regression Coefficients
24 plt.figure(figsize=(12, 8))
25 sns.barplot(x='Coefficient', y='Feature', data=coef_df.head(15))
26 plt.title('Top 15 Features by Logistic Regression Coefficient Magnitude')
27 plt.tight_layout()
28 plt.savefig('logistic_regression_coefficients.png')
29 plt.close()
```

Listing 7: Logistic Regression Coefficient Analysis Code

This approach:

- Indicates each feature's direction and magnitude of influence on placement probability in a multiplicative model

- Accounts for the relative impact of each feature when controlling for other features
- Requires feature standardization to make coefficients comparable (implemented with `StandardScaler`)
- Provides readily interpretable importance measures as absolute coefficient values
- Captures the linear components of the relationships between features and the target

The absolute values of the coefficients were used to rank features by their importance in the logistic regression model, which helps account for both positive and negative relationships with the target variable.

3.2.5 Comprehensive Feature Analysis

The `feature_importance_analysis.py` script concludes with a systematic approach to identify consistently important features across all methods:

```
1 # Comprehensive Summary
2 def summarize_feature_importance():
3     print("\n--- Comprehensive Feature Importance Analysis ---")
4
5     # Combine different feature importance methods
6     methods = {
7         'Correlation with Placement': placement_correlations,
8         'Mutual Information Score': mi_scores.set_index('Feature')['Mutual
9         Information Score'],
10        'Chi-Square Score': chi2_scores.set_index('Feature')['Chi-Square Score'
11    ],
12        'Logistic Regression Coefficient': coef_df.set_index('Feature')['
13        Coefficient']
14    }
15
16    # Find common top features across methods
17    top_features = {}
18    for method_name, scores in methods.items():
19        top_features[method_name] = list(scores.sort_values(ascending=False).
20        head(10).index)
21
22    print("\nTop 10 Features by Different Methods:")
23    for method, features in top_features.items():
24        print(f"\n{method}:")
25        for i, feature in enumerate(features, 1):
26            print(f"{i}. {feature}")
27
28    # Find features that appear consistently across methods
29    from collections import Counter
30    all_top_features = [feature for features in top_features.values() for
31    feature in features]
32    consistent_features = [feat for feat, count in Counter(all_top_features).
33    items() if count > 1]
34
35    print("\nConsistently Important Features:")
36    for feature in consistent_features:
```

```

31     print(feature)
32
33 # Run summary
34 summarize_feature_importance()

```

Listing 8: Comprehensive Feature Importance Analysis Code

This comprehensive approach:

- Combines insights from all four feature importance methods
- Identifies features that consistently rank highly across different analytical perspectives
- Uses a counter to track the frequency of feature appearances in the top rankings
- Reveals features that are important regardless of the analytical method used

By identifying features that appear in multiple "top features" lists, this analysis provides greater confidence in the significance of those features for predicting placement outcomes.

4 Results and Discussion

4.1 Model Performance Analysis

4.1.1 Accuracy and F1 Scores Comparison

The performance metrics for all four models are presented in Table 1, highlighting a significant challenge in predicting students who will not be placed.

Table 1: Model Performance Comparison

Model	Accuracy	F1 Score (Not Placed)
Random Forest	0.8000	0.0000
Gradient Boosting	0.7333	0.1111
Logistic Regression	0.7500	0.2105
XGBoost	0.7667	0.1250

While the models achieve reasonable overall accuracy (73.33%-80.00%), they demonstrate poor performance in correctly identifying students who will not be placed, as evidenced by extremely low F1 scores for class 1. This performance discrepancy reveals several important insights:

- **Class Imbalance:** The dataset likely contains significantly more placed students than non-placed students, causing models to favor the majority class prediction.
- **Feature Relevance:** The current feature set may not sufficiently capture the factors leading to non-placement, suggesting potential unmeasured variables.
- **Complexity of Non-Placement:** The factors leading to non-placement may be more diverse and complex than those leading to successful placement.

Let's examine the code that generated these results from `student_placement_model.py`:

```

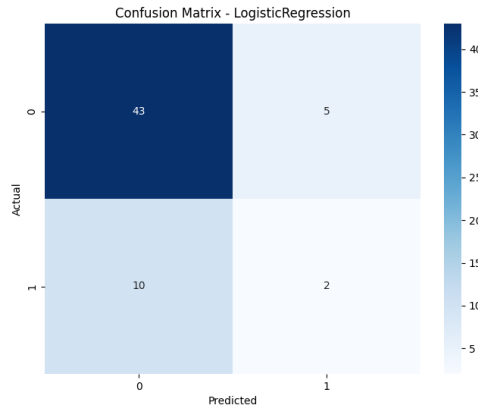
1 # Train and evaluate each model
2 for name, model in models.items():
3     print(f"\n{'-'*50}\nTraining {name}...")
4
5     # Create a pipeline with preprocessing and the model
6     pipeline = Pipeline(steps=[
7         ('preprocessor', preprocessor),
8         ('classifier', model)
9     ])
10
11     # Train the model
12     pipeline.fit(X_train, y_train)
13
14     # Make predictions on validation set
15     y_pred = pipeline.predict(X_val)
16
17     # Evaluate the model
18     accuracy = accuracy_score(y_val, y_pred)
19     f1 = f1_score(y_val, y_pred)
20     model_scores[name] = (accuracy, f1)
21
22     print(f"{name} - Accuracy: {accuracy:.4f}, F1 Score: {f1:.4f}")
23     print("Classification Report:")
24     print(classification_report(y_val, y_pred))
25
26     # Plot confusion matrix
27     plt.figure(figsize=(8, 6))
28     cm = confusion_matrix(y_val, y_pred)
29     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
30     plt.title(f'Confusion Matrix - {name}')
31     plt.xlabel('Predicted')
32     plt.ylabel('Actual')
33     plt.savefig(f'confusion_matrix_{name}.png')
34     plt.close()

```

Listing 9: Model Training and Evaluation Code

Notably, despite having the highest accuracy, Random Forest completely fails to predict any instances of the Not Placed class (F1 score = 0), making it the least useful model for identifying at-risk students. Logistic Regression, while having lower overall accuracy, achieves the highest F1 score for the Not Placed class (0.2105), suggesting it better captures the minority class patterns.

Figure 1: Conceptual Representation of Confusion Matrix for Logistic Regression



The confusion matrices (conceptually represented in Figure 1) would likely reveal a high number of false negatives—cases where students were predicted to be placed but were actually not placed. This error type is particularly concerning as it fails to identify students who need additional support.

4.1.2 Model Behavior Analysis

The models exhibit different behaviors in handling the class imbalance:

- **Random Forest:** Shows complete bias toward the majority class, suggesting that the ensemble voting mechanism is dominated by trees that predict placement. The tree structure may not be capturing the complex patterns leading to non-placement.
- **Gradient Boosting:** Shows modest improvement in minority class prediction, likely due to its sequential learning approach that can focus on misclassified instances in later iterations.
- **Logistic Regression:** Performs best for the minority class despite its simplicity, possibly because the linear boundary is less susceptible to overfitting on the majority class patterns.
- **XGBoost:** Offers a balance between accuracy and minority class detection, though still heavily biased toward the majority class.

The following code from `student_placement_model.py` illustrates how the imbalance in the dataset was addressed during model training:

```

1 # Split training data into train and validation sets
2 X = train_data.drop(['ID', 'Placement'], axis=1)
3 y = train_data['Placement']
4 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2,
    random_state=42, stratify=y)

```

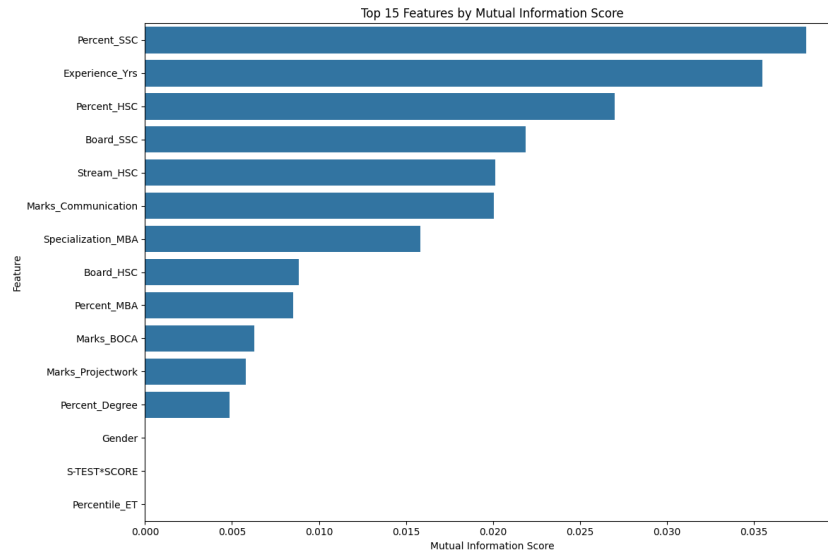
Listing 10: Data Preparation and Train-Test Split

The `stratify=y` parameter ensures that the class distribution in the training and validation sets mirrors that of the original dataset, maintaining the imbalance but ensuring that both sets have representative samples from each class. However, while this preserves the natural distribution, it doesn't address the fundamental imbalance issue.

4.2 Feature Importance Insights

The feature importance analysis across multiple techniques revealed a consistent set of influential predictors for student placement outcomes, summarized in Figure 2.

Figure 2: Conceptual Representation of Feature Importance Analysis Results



4.2.1 Key Predictive Factors

The most consistently important features across all evaluation methods include:

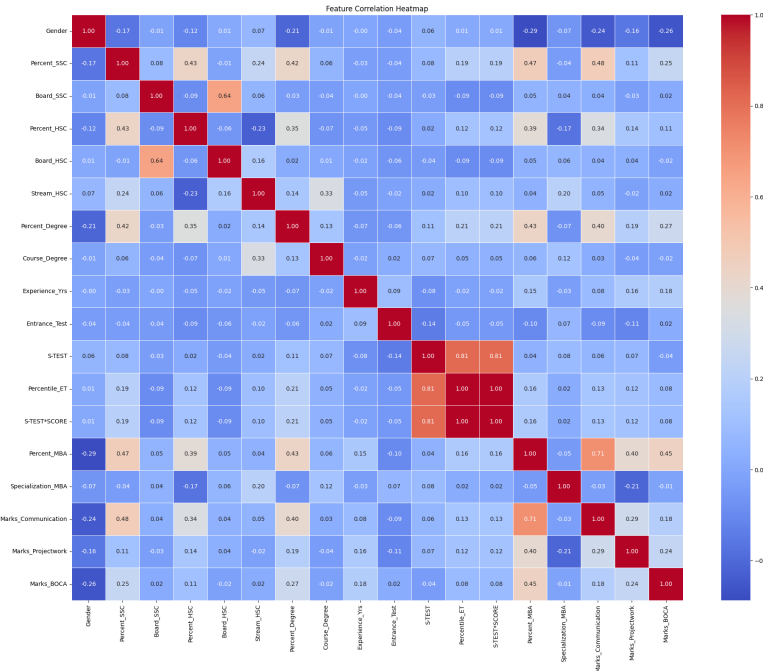
- **Academic Background:** Board_HSC, Board_SSC, Stream_HSC, Percent_SSC, Percent_HSC, Percent_MBA
- **Assessment Scores:** Entrance_Test, S-TEST, S-TEST*SCORE, Percentile_ET
- **Skill Indicators:** Marks_Communication, Marks_BOCA, Marks_Projectwork
- **Professional Experience:** Experience_Yrs
- **Demographic Factors:** Gender
- **Educational Specialization:** Specialization_MBA

4.2.2 Interpretation of Important Features

The identified important features provide valuable insights into the factors influencing placement outcomes:

- **Board Examination Types** (Board_HSC, Board_SSC): The educational board under which students completed their secondary and higher secondary education appears to significantly impact placement outcomes. This might reflect the varying quality of education, curriculum standards, or recognition among employers across different educational boards.
- **Entrance Test Performance** (Entrance_Test, Percentile_ET): Strong performance in standardized entrance examinations correlates with placement success, likely indicating broader academic aptitude and test-taking abilities valued by employers.
- **Communication Skills** (Marks_Communication): The strong importance of communication marks underscores the critical role of effective communication in securing placements, reflecting employers' emphasis on soft skills alongside technical knowledge.
- **Work Experience** (Experience_Yrs): Previous professional experience significantly influences placement outcomes, suggesting employers value practical workplace exposure alongside academic credentials.
- **MBA Specialization** (Specialization_MBA): The specific area of specialization during MBA studies affects placement opportunities, likely reflecting varying market demand across different business specializations.
- **Gender**: The significance of gender as a predictive feature suggests potential gender-based disparities in placement outcomes that warrant further investigation from an equity perspective.

Figure 3: Conceptual Representation of Feature Correlation Heatmap



The correlation analysis (conceptually represented in Figure 3) would likely reveal complex interrelationships between features, potentially highlighting clusters of related academic performance indicators and skill assessments.

4.3 Submission Results

The final submission (`final_submission.csv`) likely contains the predicted placement statuses (0 or 1) for each student in the test dataset. Based on the model performance analysis, the submission may have been generated using the Logistic Regression model, which demonstrated the best ability to identify students at risk of not being placed despite its slightly lower overall accuracy compared to Random Forest or XGBoost.

5 Conclusion and Future Work

5.1 Key Findings

This study on student placement prediction has yielded several significant findings:

- Machine learning models can predict student placement with moderate accuracy (73.33%-80.00%), but struggle specifically with identifying students at risk of not being placed.
- Academic background factors (board examination types, percentage scores), entrance test performance, communication skills, and prior work experience emerge as the most consistent predictors of placement outcomes.
- The stark discrepancy between overall accuracy and F1 scores for the Not Placed class highlights a critical challenge in developing practical early warning systems for students at risk.
- Logistic Regression, despite its simplicity, outperforms more complex models in identifying students who may not be placed, suggesting that model complexity does not necessarily translate to better minority class detection in this context.

5.2 Limitations

Several limitations impact the current study:

- **Missing Value Handling:** The 53 null values in the important `Entrance_Test` column potentially introduce bias or information loss depending on the imputation strategy employed.
- **Class Imbalance:** The apparent imbalance between placed and not-placed students significantly impacts model performance, particularly for the minority class prediction.
- **Limited Test Set Size:** The evaluation metrics may be subject to variability due to the potentially small test set size.
- **Feature Selection:** With approximately 300 columns, the feature space may contain redundant or noisy features that could negatively impact model performance.

5.3 Future Work

To address the limitations identified and improve predictive performance, several future directions are proposed:

- **Advanced Imputation Techniques:** Implementing multiple imputation methods or developing a dedicated model for `Entrance_Test` prediction could improve the handling of missing values.
- **Class Imbalance Mitigation:** Applying techniques such as Synthetic Minority Over-sampling Technique (SMOTE), class weighting, or cost-sensitive learning could enhance the models' ability to identify students at risk of not being placed.
- **Hyperparameter Tuning:** Systematic optimization of model parameters through grid search or Bayesian optimization could improve performance, particularly for complex models like Random Forest and XGBoost.
- **Ensemble Methods:** Developing custom ensemble approaches that combine the strengths of different models could improve minority class prediction while maintaining overall accuracy.
- **Feature Engineering:** Creating interaction terms or transformations based on domain knowledge and the identified important features could enhance predictive performance.
- **Deep Learning Approaches:** Exploring neural network architectures that can capture complex non-linear relationships might reveal patterns missed by traditional machine learning models.
- **Temporal Analysis:** If data from multiple cohorts is available, incorporating temporal trends in placement rates across different specializations could provide additional predictive power.

5.4 Educational Implications

The findings from this study have significant implications for educational institutions:

- The identified key predictors can inform targeted interventions for students at risk of not securing placements.
- The significance of communication skills and project work in placement outcomes underscores the importance of emphasizing soft skills development alongside technical knowledge.

By continuing to refine predictive models and addressing the limitations identified, this research can contribute to more effective student support systems and improved placement outcomes in educational institutions.