# Model Development and Training Report.

**INTRODUCTION**

This report details the development and assessment of various machine learning models utilized for classification tasks. It emphasizes the comparison of model performance, analysis of learning curves, and diagnostics. The models examined include Logistic Regression, Decision Tree Classifier, Random Forest Classifier, Support Vector Machine (SVM), and XGBoost. The objective is to determine the best-performing model based on accuracy and other pertinent metrics, while also taking into account interpretability and generalization capabilities.

**DESCRIPTION OF MODEL ARCHITECTURE**

**1. Logistic Regression**

- **Use case**: As a simple baseline model for binary classification (dropout vs graduate).
- **Architecture**: A single-layer model with a logistic (sigmoid) activation function. It predicts the probability of a class based on a linear combination of input features.
- **Advantages**:
  - Simple and interpretable.
  - Works well when the relationship between features and outcome is linear.
- **Disadvantages**:
  - May underperform if the relationships are non-linear or complex.

**2. Decision Trees**

- **Use case**: Can be used when feature interactions and non-linearity are important.
- **Architecture**: A tree structure where each internal node splits the data based on a feature, and each leaf represents a predicted class.
- **Advantages**:
  - Captures non-linear relationships.
  - Easy to interpret and visualize.
- **Disadvantages**:
  - Prone to overfitting on complex datasets unless properly tuned.

**3. Random Forest**

- **Use case**: A more robust model than a single decision tree, useful when higher accuracy is desired and the dataset has many features.
- **Architecture**: An ensemble of decision trees. Each tree is built from a random sample of the data, and the final output is based on the majority vote (classification).

- **Advantages**:
  - Reduces overfitting compared to a single decision tree.
  - Handles large numbers of features well.
- **Disadvantages**:
  - Less interpretable than individual decision trees.

## 4. XGBoost (Extreme Gradient Boosting)

- **Use case**: XGBoost is a specific implementation of Gradient Boosting Machines (GBM) optimized for speed and performance. It's commonly used for structured/tabular data and is known for achieving high accuracy in classification tasks like predicting student dropout.
- **Architecture**: Similar to standard gradient boosting, XGBoost trains decision trees in a sequential manner, where each new tree attempts to correct the errors of the previous ones. It uses advanced regularization techniques to prevent overfitting and employs efficient handling of missing values.
- **Advantages**:
  - High predictive performance with tabular data.
  - Robust to overfitting due to its regularization.
  - Handles large datasets and features efficiently.
  - Flexible with custom loss functions.
- **Disadvantages**:
  - Less interpretable compared to simpler models like logistic regression or single decision trees.
  - Requires careful tuning of hyperparameters for optimal performance.

## 5. Support Vector Machines (SVM)

- **Use case**: SVM is a strong model for classification tasks, especially when the classes (dropout vs. graduate) are separable with a clear boundary. It works well with both linear and non-linear relationships, depending on the kernel used.
- **Architecture**: SVM tries to find a hyperplane that maximally separates the two classes (dropout vs. graduate) in the feature space. If the data is not linearly separable, kernel functions (such as RBF, polynomial, or sigmoid) can be used to transform the data into higher dimensions where a separating hyperplane exists.
- **Advantages**:
  - Effective in high-dimensional spaces.
  - Works well when the data has clear separation between classes.
  - Can model non-linear relationships using kernel trick.
- **Disadvantages**:
  - Not as easy to interpret as models like logistic regression or decision trees.
  - Computationally expensive, especially for large datasets.

        o    Requires careful tuning of hyperparameters (like the kernel type and regularization parameter) for optimal performance.

## COMPARISON OF MODEL PERFORMANCES

The following machine learning models: Logistic Regression, Decision Tree, Random Forest, Support Vector Machine (SVM), and XGBoost were evaluated based the following metrics: accuracy, precision, recall, F1 score and ROC-AUC. Cross Validation was used to compare the model performance and that aided to select the most appropriate model.

| Model | Accuracy | Precision | Recall | F1 Score | ROC-AUC |
|---|---|---|---|---|---|
| Logistic Regression | 0.8804 | 0.9406 | 0.6690 | 0.7819 | 0.9274 |
| Decision Tree | 0.8036 | 0.6797 | 0.7324 | 0.7051 | 0.7848 |
| Random Forest | 0.8849 | 0.8760 | 0.7465 | 0.8061 | 0.9274 |
| Support Vector Machine | 0.8691 | 0.9200 | 0.6479 | 0.7603 | 0.9214 |
| XGBoost | 0.8894 | 0.8720 | 0.7676 | 0.8165 | 0.9193 |

### Key Insights

- XGBoost has the highest accuracy at 0.8894, indicating it performs best in correctly classifying instances as Random Forest and Logistic Regression follow closely.
- Logistic Regression leads with the highest precision (0.9406), effectively identifying true positives.
- Decision Tree has the highest recall (0.7324), showing its ability to capture a significant number of actual positive cases.
- XGBoost achieves the best F1 Score (0.8165), balancing precision and recall well.
- Logistic Regression and Random Forest both excel with high ROC-AUC values (0.9274), indicating strong performance in distinguishing between classes.

### In Summary,

- XGBoost is the best-performing model overall, particularly in accuracy and F1 Score.
- Random Forest and Logistic Regression are also strong contenders, especially in precision and ROC-AUC.
- Decision Tree may not be suitable for this dataset due to lower accuracy and precision.

## ANALYSIS OF LEARNING CURVES AND MODEL DIAGNOSTICS

### Key Insights from Learning Curves

- Logistic Regression: Strong training performance (accuracy ~0.88) but a validation score of ~0.86 indicates some overfitting. Future improvements should focus on regularization and cross-validation.

- Decision Tree: High training accuracy (~1.00) with a validation score of ~0.80 suggests significant overfitting. Techniques like pruning and cross-validation could enhance generalization.
- Random Forest: Consistently high training score (~1.00) and a validation score of ~0.86 indicate good performance but potential overfitting. Further tuning may improve validation results.
- Support Vector Machine (SVM): Training score (~0.88) and validation score (~0.85) suggest reasonable generalization with minimal overfitting. Hyperparameter tuning could enhance performance.
- XGBoost: Perfect training performance (1.00) with a validation score around 0.87 indicates potential overfitting, although it shows good generalization. Further tuning could improve validation performance.

**Bias-Variance Tradeoff.**

The Random Forest and Support Vector Machine (SVM) models demonstrate the best bias-variance tradeoff. Random Forest maintains moderate bias and lower variance, effectively balancing complexity and generalization. SVM also offers a good balance with moderate bias and low to moderate variance, making it suitable for various datasets when properly tuned. Both models are robust against overfitting while capturing complex patterns.

## CONCLUSION

In this report, I evaluated several machine learning models for classification tasks, including Logistic Regression, Decision Tree, Random Forest, Support Vector Machine (SVM), and XGBoost. The analysis revealed that XGBoost emerged as the best-performing model, achieving the highest accuracy and F1 score. Random Forest and Logistic Regression also demonstrated strong performance, particularly in precision and ROC-AUC metrics. Decision Trees, while easy to interpret, showed lower accuracy and precision, making them less suitable for this dataset.

## RECOMMENDATION

Based on the findings, I recommend using XGBoost for optimal performance in classification tasks due to its high accuracy and ability to handle complex relationships in data. For scenarios requiring interpretability, Logistic Regression is advisable, especially for linear relationships. Random Forest is also recommended for its robustness and moderate interpretability, making it a solid choice for datasets with many features. Further tuning of hyperparameters and techniques like cross-validation should be employed to enhance model performance and generalization.