

RetentionClassifier:Model Training and Evaluation Report

1. Input Features

The dataset comprises the following features:

- **Age:** Continuous numerical feature representing employee age.
 - **Gender:** Categorical feature indicating male or female.
 - **Department:** Categorical feature representing the department in which the employee works.
 - **Job Title:** Categorical feature specifying the employee's role.
 - **Years at Company:** Numerical feature indicating the duration of employment.
 - **Satisfaction Level:** Numerical feature (likely normalized between 0 and 1) reflecting employee satisfaction.
 - **Average Monthly Hours:** Numerical feature showing average working hours per month.
 - **Promotion Last 5 Years:** Binary feature indicating if the employee was promoted in the last 5 years.
 - **Salary:** Ordinal categorical feature (e.g., low, medium, high).
 - **Attrition:** Target variable (binary: 0 = No Attrition, 1 = Attrition).
-

2. Models Trained

You trained several machine learning models:

1. Logistic Regression
2. Decision Tree
3. Random Forest
4. Support Vector Machine (SVM)
5. Naive Bayes
6. K-Nearest Neighbors (KNN)
7. XGBoost
8. Gradient Boosting

Best Model

- **Gradient Boosting** was identified as the best-performing model.
 - **Accuracy:** 52%
-

3. Best Model Performance

Gradient Boosting Results

Metric	Class 0 (No Attrition)	Class 1 (Attrition)	Macro Avg	Weighted Avg
Precision	53%	50%	52%	52%
Recall	51%	52%	52%	52%
F1-Score	52%	51%	51%	52%
Support (Count)	102	98	-	-
Overall Accuracy	-	-	-	52%

Insights:

- The model struggles to differentiate between the two classes.
 - The F1-Score for both classes is almost equal, indicating no bias toward a particular class.
 - The recall for both classes hovers around 51-52%, showing that the model is failing to identify many true positives for either class.
-

4. Observations

Low Accuracy and Metrics:

- **Accuracy of 52%** is marginally better than random guessing (50% for a binary classifier).
- Precision, recall, and F1-scores indicate that the model's performance is limited, potentially due to:
 - Insufficient features or poor feature quality.
 - Imbalanced target classes.
 - Data preprocessing issues (e.g., lack of proper encoding or scaling).
 - Lack of hyperparameter tuning for models.

Class Imbalance:

If the target classes (Attrition: 0 vs. 1) are imbalanced, the model might not generalize well. Check the class distribution in your dataset.

Feature Importance:

Gradient Boosting models can provide feature importance. Use this to determine which features are contributing most to the predictions. For instance, features like **Satisfaction Level** or **Years at Company** might be critical in predicting attrition.

5. Recommendations

Data Improvements:

1. **Feature Engineering:**
 - Combine or derive new features that better capture employee behavior.
 - For example, create a "Work-Life Balance Index" using **Satisfaction Level** and **Average Monthly Hours**.
2. **Data Balancing:**
 - If the classes are imbalanced, consider techniques like SMOTE (Synthetic Minority Oversampling Technique) to balance them.
3. **Categorical Encoding:**
 - Use one-hot encoding or ordinal encoding for categorical variables like **Department**, **Job Title**, and **Salary**.

Model Tuning:

1. Perform hyperparameter tuning for Gradient Boosting or other tree-based models using techniques like Grid Search or Random Search.
2. Experiment with ensemble models (e.g., combining Random Forest and XGBoost).

Alternative Approaches:

1. **Try Neural Networks:** If the dataset size is large enough, simple feed-forward neural networks may perform better.
2. **Consider Domain Knowledge:** Incorporate expert knowledge about the factors driving attrition.

Evaluation Metrics:

1. Use additional metrics like the Area Under the ROC Curve (AUC-ROC) to evaluate model discrimination power.
 2. Generate a confusion matrix to understand misclassifications.
-

6. Next Steps

1. **Conduct an EDA (Exploratory Data Analysis):**
 - Investigate relationships between features and the target variable.

- Identify outliers or unusual patterns.
- 2. **Analyze Feature Importance:**
 - Use `feature_importances_` from Gradient Boosting to rank the top predictors.
- 3. **Refine the Dataset:**
 - Address any missing data, outliers, or inconsistencies.
 - Normalize or scale numerical features for better model convergence.
- 4. **Iterate with Tuned Models:**
 - Train Gradient Boosting with optimized parameters.
 - Explore hybrid models or stacking approaches.