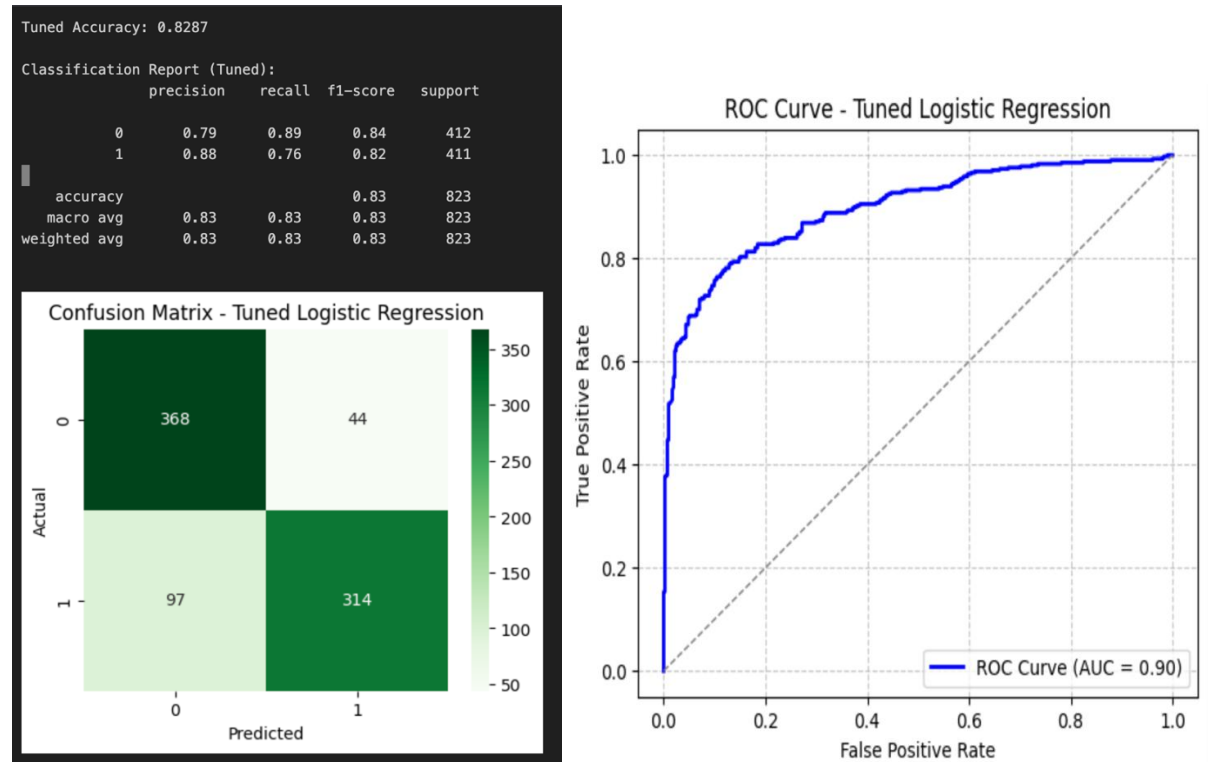


# Models & Their Scores

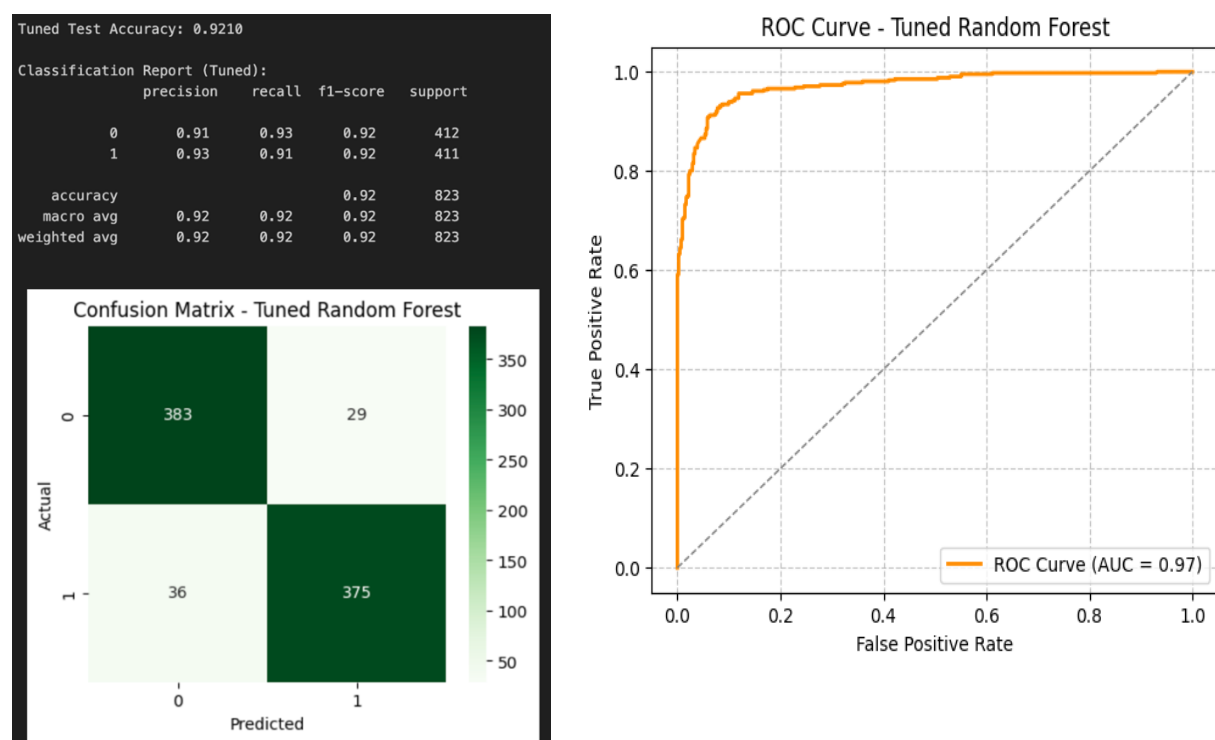
## Logistic Regression

We selected logistic regression because it's a simple and interpretable model that works well for binary classification problems like predicting student graduation or dropout.



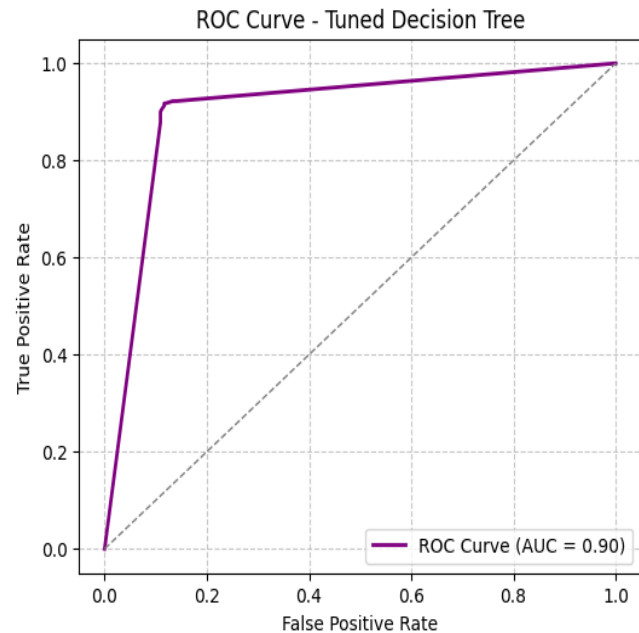
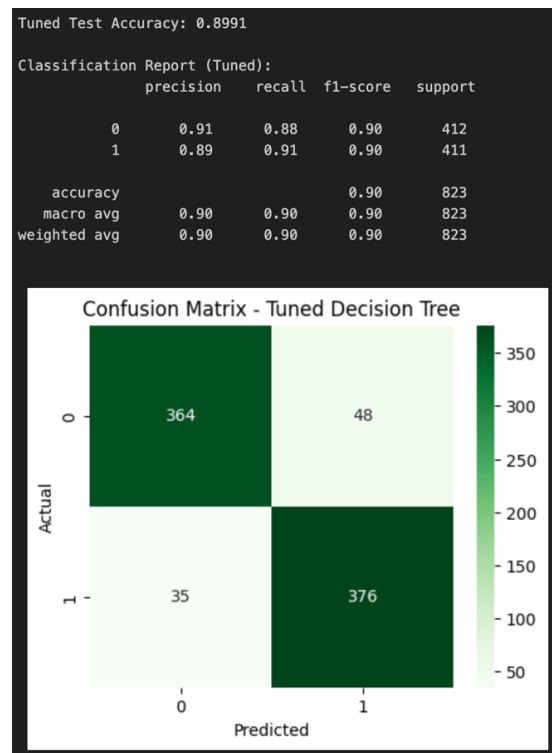
## Random Forest

We chose random forest as it's an ensemble of decision trees that provides higher accuracy, reduces overfitting, handles numerical and categorical data effectively, and offers feature importance scores for student retention factors.



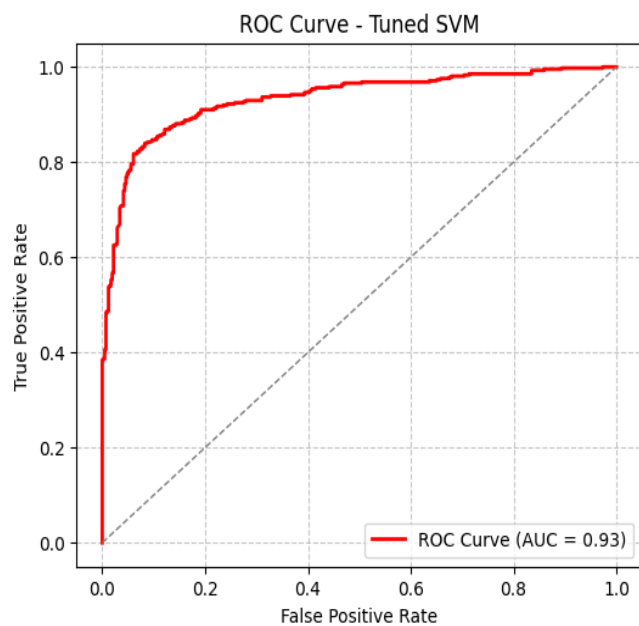
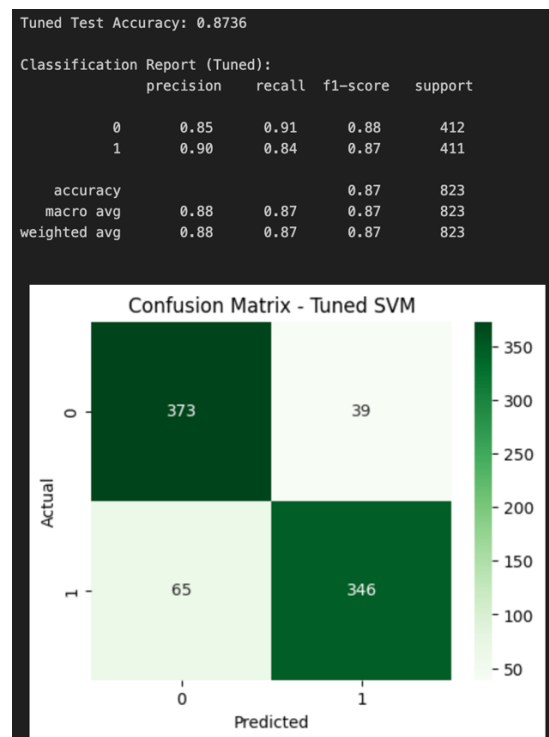
## Decision Tree

We picked decision tree because it's easy to understand and visualize, helping identify key factors influencing student outcomes in this binary classification task.



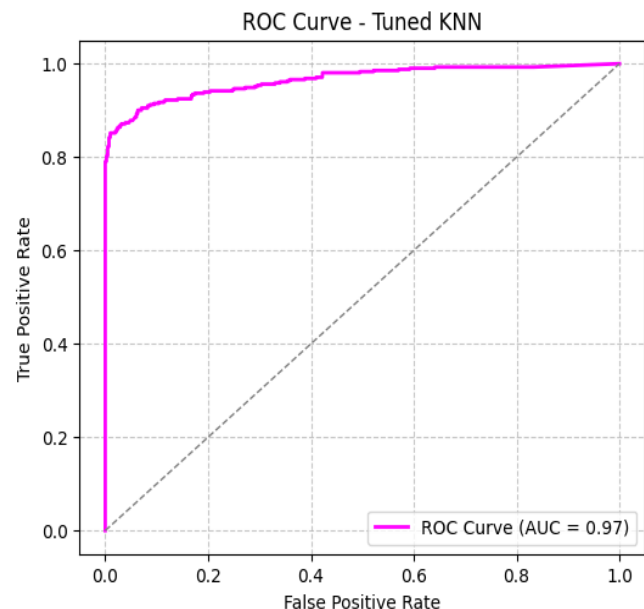
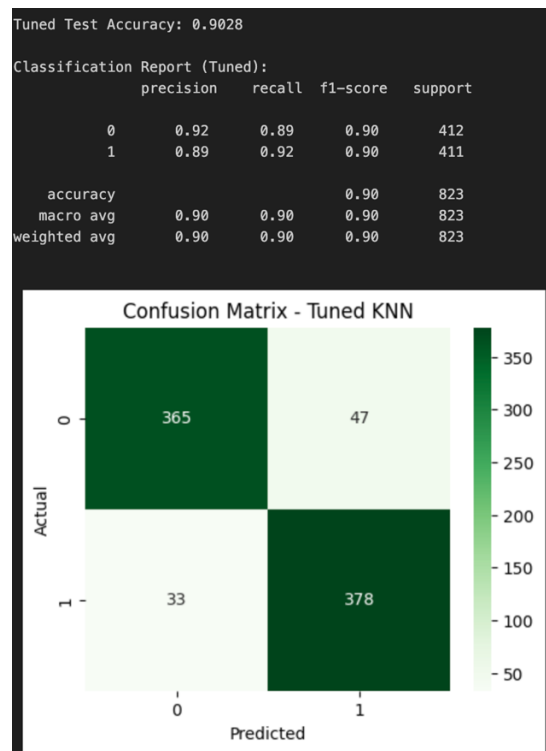
## SVM (Support Vector Machine)

We selected SVM for its usefulness in binary classification, as it effectively separates classes with a hyperplane, working well with both linear and non-linear data when margins are clear.



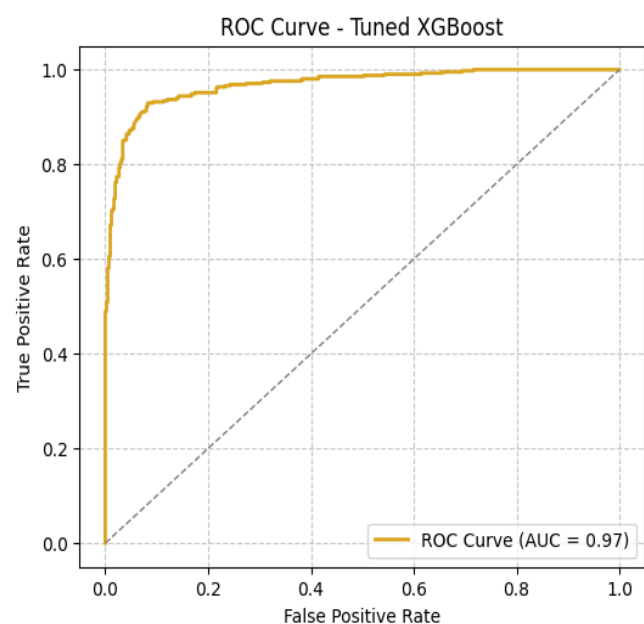
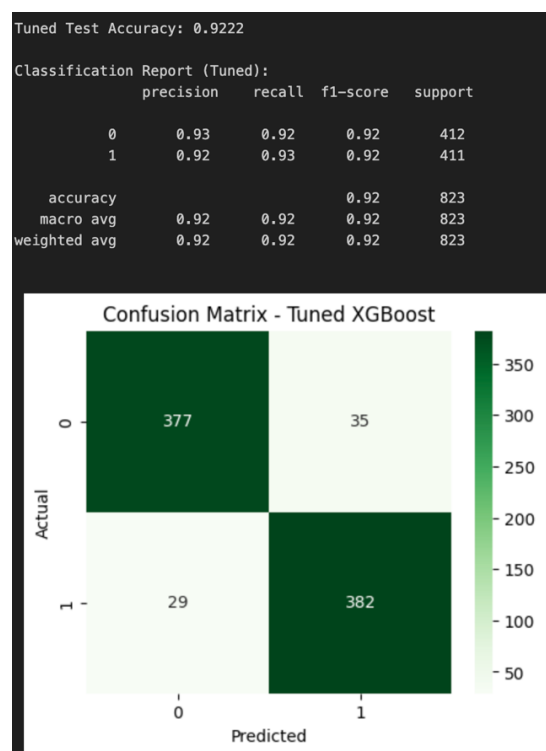
## KNN (K-Nearest Neighbors)

We chose KNN as a simple instance-based method that classifies based on similar student profiles, useful for data with clear clusters or patterns in retention prediction.



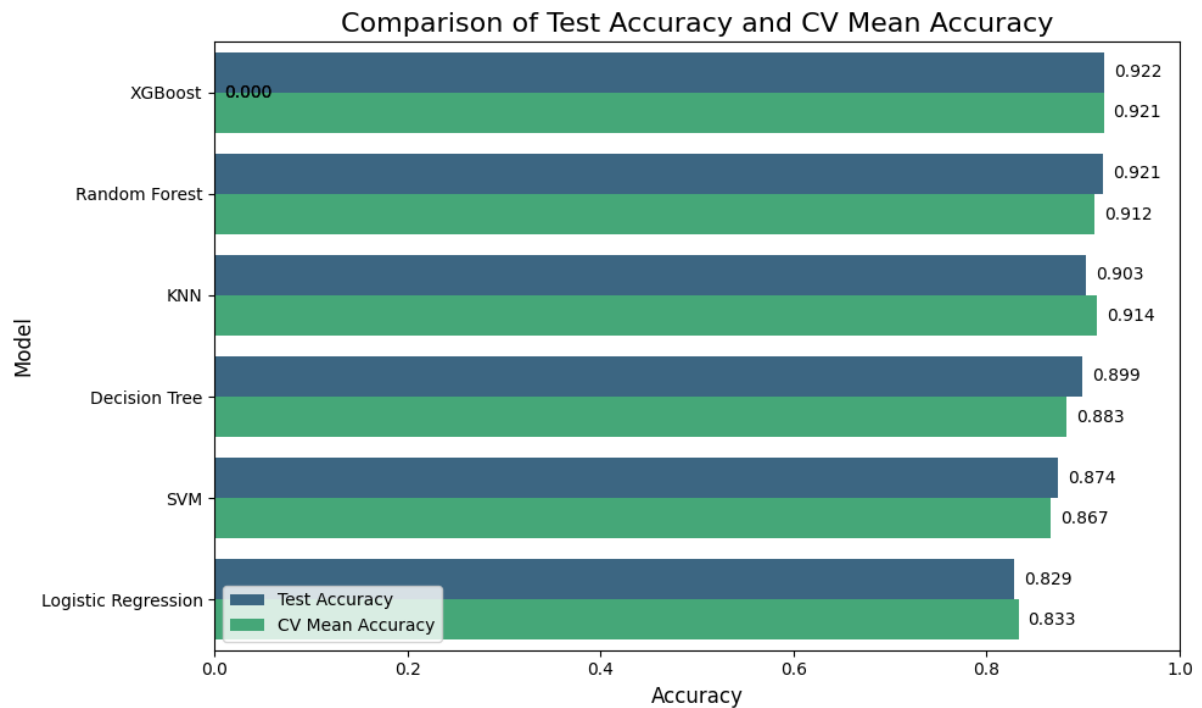
## XGBoost

We selected XGBoost for its efficiency as a gradient boosting algorithm that delivers high accuracy on structured tabular datasets like ours for student outcome prediction.



## Model Comparison

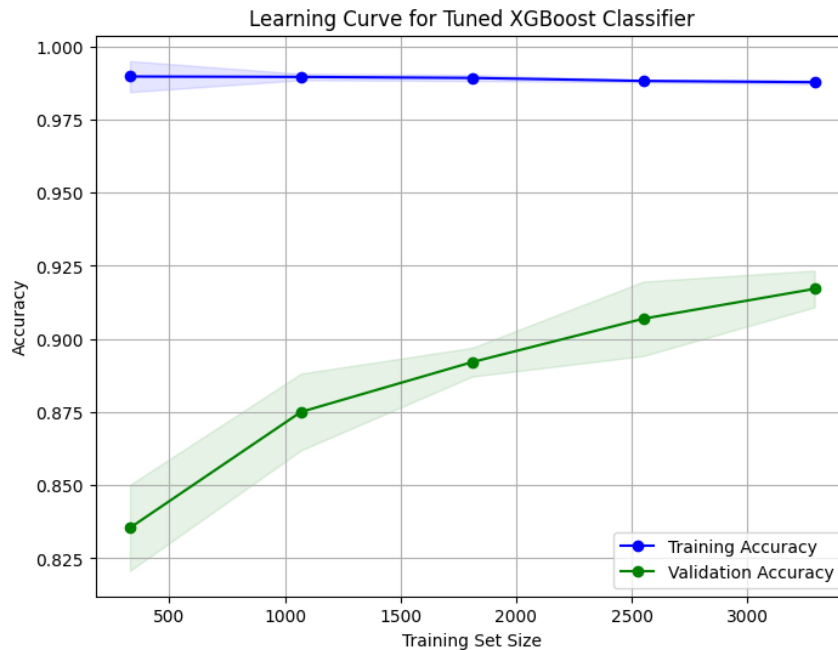
Model	Test Accuracy	CV Mean Accuracy	CV std
XGBoost	0.9222	0.9213	0.0089
Random Forest	0.9210	0.9116	0.0073
KNN	0.9028	0.9143	0.0054
Decision Tree	0.8991	0.8827	0.0126
SVM	0.8736	0.8669	0.0064
Logistic Regression	0.8287	0.8329	0.0129



## Final Model Selected – XGBoost

After evaluating all developed models based on accuracy, cross-validation results, confusion matrices, and ROC curves, the **tuned XGBoost (Extreme Gradient Boosting)** model was selected as the final predictive model for student graduation and dropout classification.

XGBoost demonstrated the **highest overall accuracy** among all models, along with excellent **cross-validation stability** and a **high AUC value** in the ROC analysis. The confusion matrix showed that the model correctly identified most dropout and graduate cases, with minimal misclassifications.



The learning curve of the tuned XGBoost model further confirmed that it achieved a **good balance between bias and variance**, indicating strong generalization to unseen data without overfitting.

The main advantages of XGBoost that contributed to its superior performance include:

- Ability to handle complex nonlinear relationships in the data
- Built-in regularization techniques to prevent overfitting
- Efficient computation and scalability for large datasets

Therefore, the tuned XGBoost model was finalized as the most suitable and reliable algorithm for predicting student retention outcomes. Its strong predictive capability and stable performance make it an effective choice for practical use in educational analytics and early intervention systems.