# Predicting Employee Attrition with ML

Using machine learning to predict employee attrition.

Leveraging HR data and advanced models to identify at-risk employees.

Objective: reduce attrition and improve employee retention.

# Data Cleaning and Preparation

## Data Cleaning

- Remove irrelevant columns.
- Handle missing values.
- Ensure data quality.

## Label Encoding

- Convert categories to numbers.
- Prepare for model training.
- Improve model compatibility.

# Exploratory Data Analysis (EDA)

### Correlation Matrix
Visualize feature correlations.

### Multicollinearity Detection
Reduce redundant variables.

### Key Insights
Age vs. Monthly Income.

Keep Attrition as Dependent variable and Run the models

# Machine Learning Models

## Logistic Regression

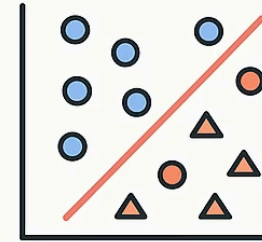Binary classification baseline.

## SVM

Separating hyperplanes.

## Decision Tree

Interpretable rules.

## Random Forest
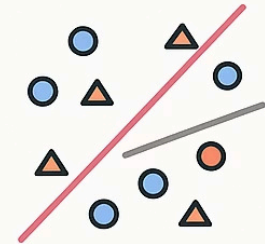
Ensemble accuracy.

# MACHINE LEARNING

### LOGISTIC REGRESSION

Binary classifier that models the probability of class membership

### SUPPORT VECTOR MACHINE

Classifies data by finding the hyperplane that best separates the classes

### DECISION TREE

Tree-like model of decisions based on feature values

### RANDOM FOREST

Ensemble of decision trees that improves classification accuracy

# Model Evaluation Metrics

✓ **Accuracy**

**Confusion Matrix**

**ROC-AUC**

Comprehensive metrics to compare models.

# Key Results and Insights

**Performance**

SVM and Random Forest excel.

**Interpretation**

Logistic Regression is clear.

**Identification**

Pinpoint high-risk attrition.

```
[ ]  # Fit logistic regression model
     log_model = LogisticRegression(max_iter=1000)
     log_model.fit(X_train, y_train)

     # Predictions
     y_pred = log_model.predict(X_test)
     y_prob = log_model.predict_proba(X_test)[:, 1]
```

## Logistic regression

Accuracy: 89%

```
[ ]  # Train SVM model
     svm_model = SVC(kernel="rbf", random_state=42)
     svm_model.fit(X_train, y_train)

     # Predictions
     y_pred = svm_model.predict(X_test)
```

## Support Vector Machine

Accuracy: 89%

```python
# Train Decision Tree model
dt_model = DecisionTreeClassifier(random_state=42)
dt_model.fit(X_train, y_train)

# Predictions
y_pred_dt = dt_model.predict(X_test)

# Evaluate model
accuracy_dt = accuracy_score(y_test, y_pred_dt)
report_dt = classification_report(y_test, y_pred_dt)

accuracy_dt, report_dt

print(accuracy_dt)
print(report_dt)
```

## Decision tree Classifier

Accuracy: 80%

```python
# Train Random Forest model
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)

# Predictions
y_pred = rf_model.predict(X_test)

# Evaluate model
accuracy = accuracy_score(y_test, y_pred)
report = classification_report(y_test, y_pred)

accuracy, report

print(accuracy)
print(report)
```
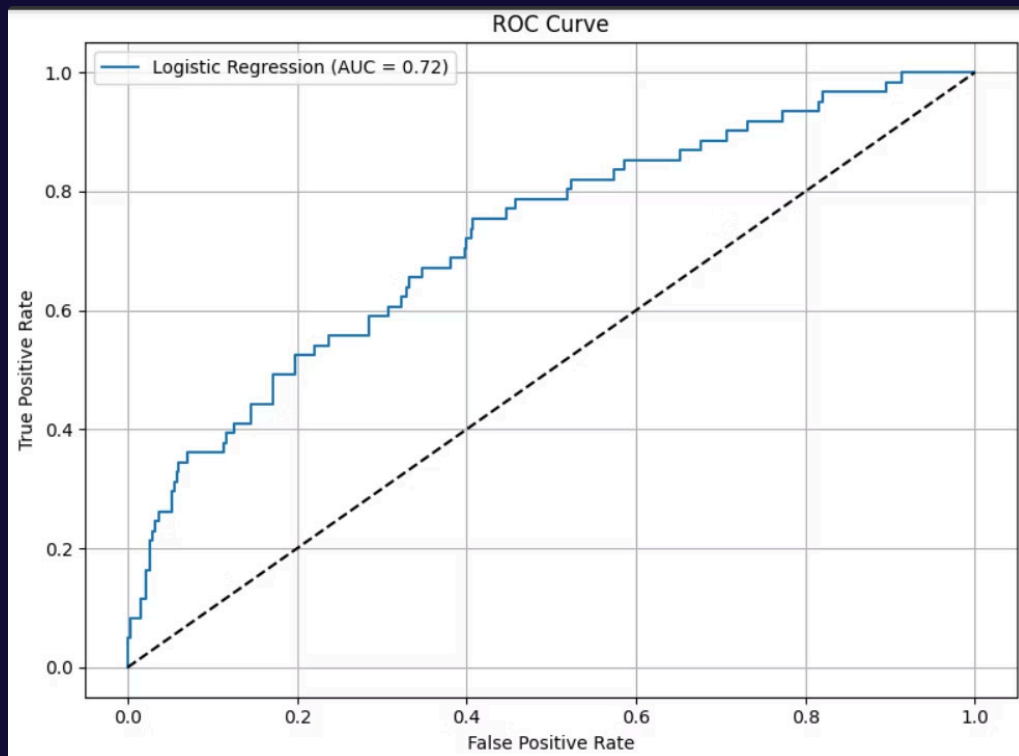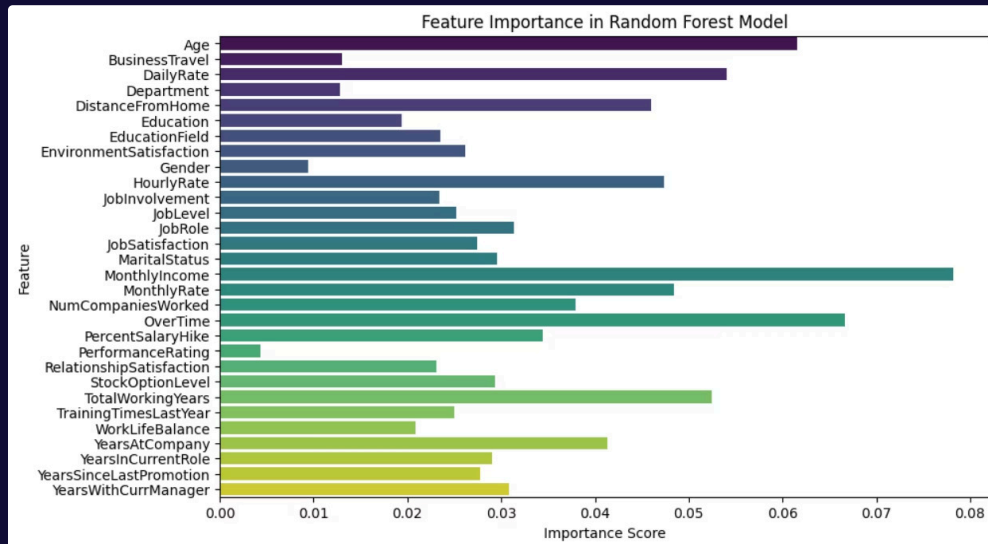
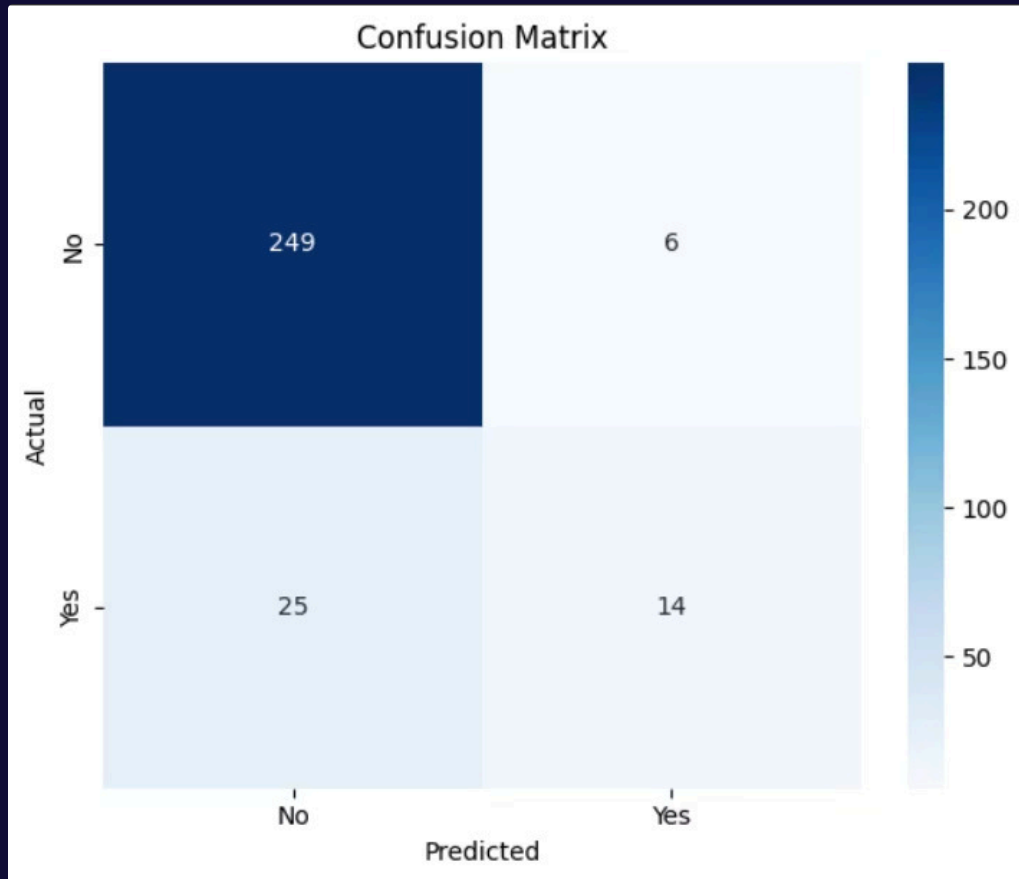## Random Forest Classifier

Accuracy: 89%

# ROC Curve

Visualizes the trade-off between true positive rate and false positive rate across thresholds.
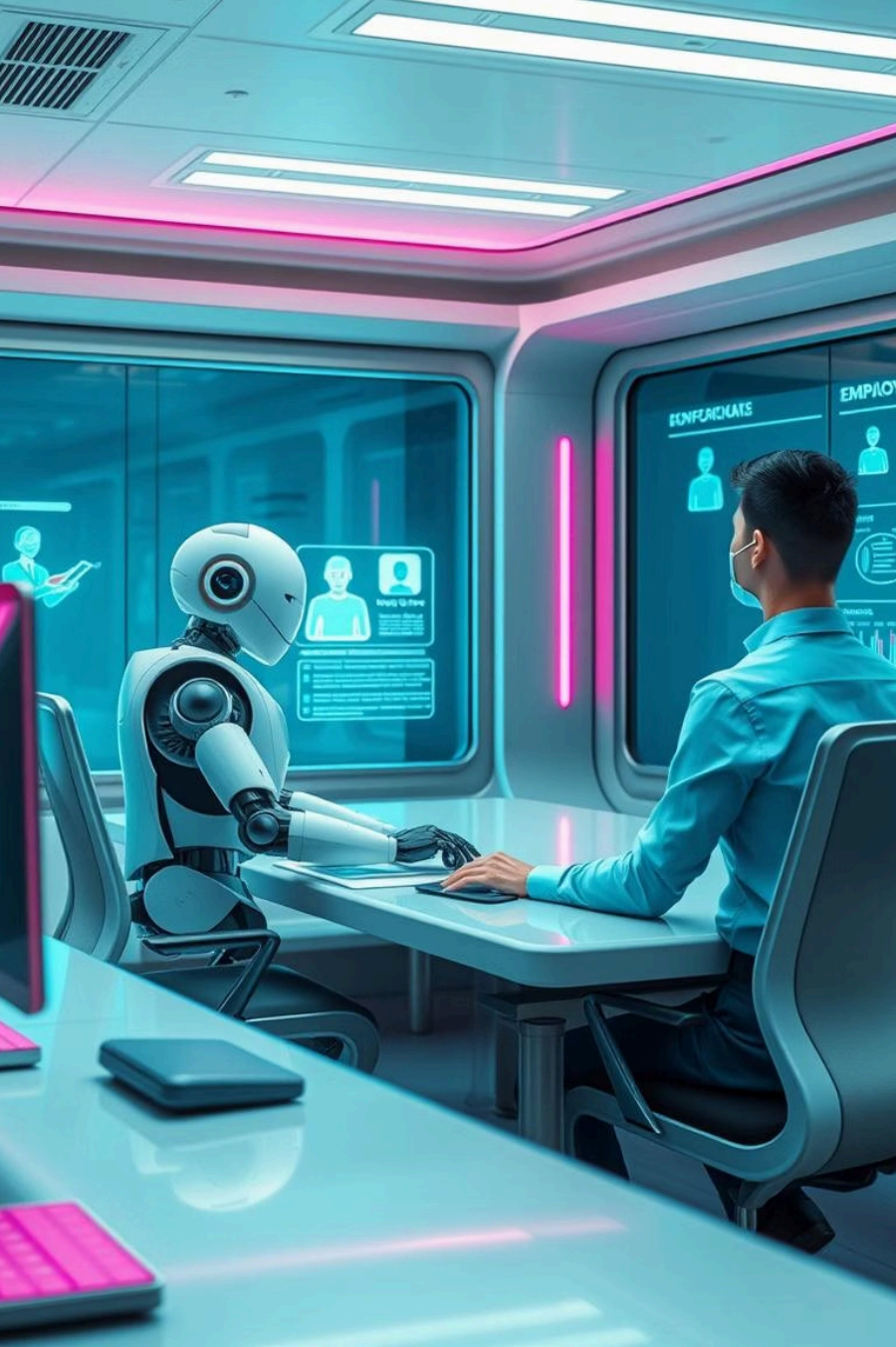
Feature Importance in Random Forest Model

# Feature Importance graph

Shows which features have the greatest impact on the model's predictions.

# Confusion Matrix

Summarizes the model's classification performance by showing correct and incorrect predictions.

# Conclusion and Future Steps

### Impact

Precise attrition prediction.

### Integration

HR dashboards ready.

### Next Steps

- Hyperparameter tuning.
- XGBoost for boost.
- Real-time deployment.