

# Assignment: Feature Importance Analysis using SHAP

Domain: Education (Student Performance)

Student Name: \_\_\_\_\_

Roll Number: \_\_\_\_\_

Date: \_\_\_\_\_

Dataset: student\_performance\_cleaned.csv

Rows (original): 649

Rows (after de-duplication): 649

Task Type: Classification

Target Variable: pass

## Preprocessing Steps

- Removed duplicates
- Missing values imputed: numeric=median, categorical=most\_frequent
- Categorical encoding: One-Hot Encoding
- Train/test split: 80/20 (stratified by target)
- Feature count after encoding: 39
- Categorical columns: 0
- Numeric columns: 39

## Model & Performance

Model: RandomForestClassifier (n\_estimators=300, random\_state=42)

Metrics:

Task: Classification

Target: pass

Accuracy: 0.8077

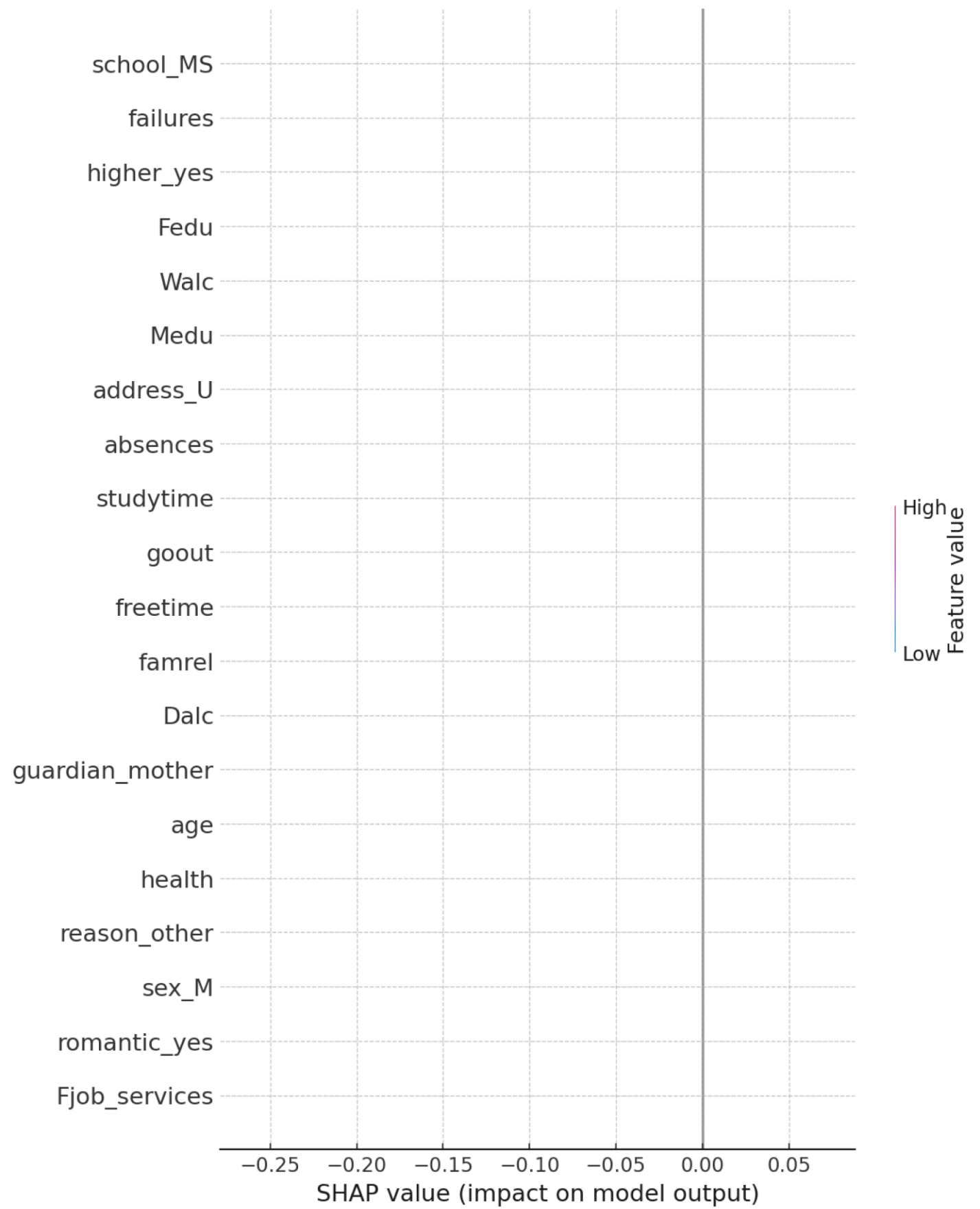
Precision: 0.8571

Recall: 0.9273

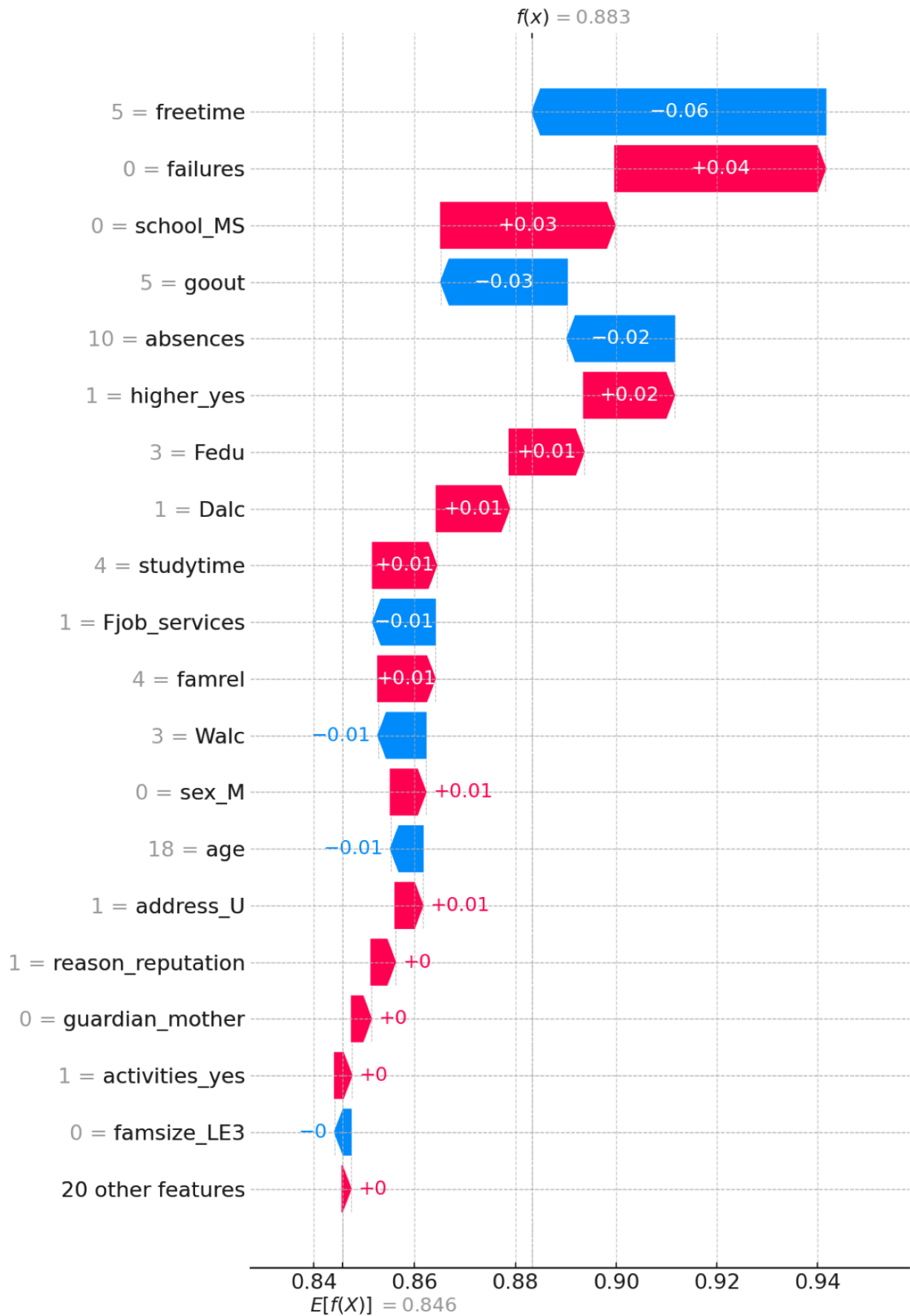
F1-score: 0.8908

ROC AUC: 0.6705

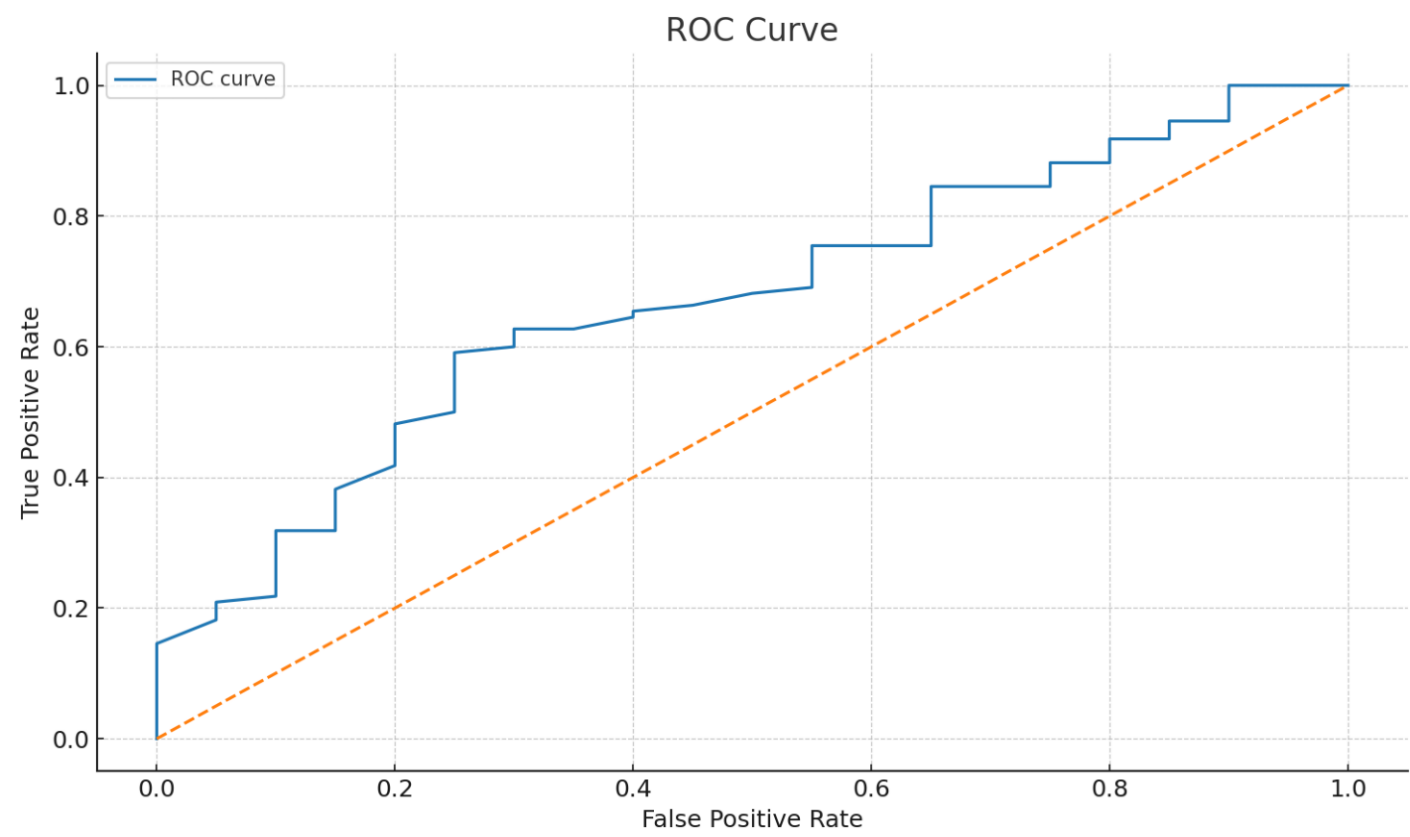
SHAP Summary Plot (overall feature importance)



SHAP Waterfall Plot (single prediction explanation)



ROC Curve



## SHAP Analysis & Interpretation

Top 5 features by SHAP importance:

1. school\_MS | SHAP=0.05524 | RF=0.06437
2. failures | SHAP=0.05392 | RF=0.11896
3. higher\_yes | SHAP=0.03450 | RF=0.06523
4. Fedu | SHAP=0.01734 | RF=0.04028
5. Walc | SHAP=0.01393 | RF=0.04465

Interpretation (template):

- These features show the strongest influence on predictions (by mean |SHAP|)
- Agreement with RF importance suggests consistent signal; disagreement may indicate non-linear effects
- Discuss domain sense and any surprises; note data/model limitations and next steps