

# A Data-Driven Multimodal Analysis of How Attendance and Behavioral Factors Influence Academic Success in Higher Education

**Course:** MSDS696\_S70 – Data Science Practicum 2

**Author:** Balarama Raju Saripalli

**Email:** [bsaripalli@regis.edu](mailto:bsaripalli@regis.edu)

**GitHub:** <https://github.com/varma1234/MSDS696-S70-Practicum-Project-2>

## Table of Contents

- Project Overview
- Problem Definition and Motivation
- Data Description and Collection Effort
- Data Collection & Cleaning (pipeline summary)
- Analysis Approach and Methodology
- EDA, Feature Engineering, Modeling & Evaluation, Analysis and Improvement
- Key Tools and Libraries Used
- Practicum Deliverables and Evaluation Rubric Alignment
- Week-by-Week Timeline and Progress
- Roadblocks and Learnings
- Online Presence and Documentation
- Visualizations and Dashboard Summary
- Project Structure (drive / repo layout)
- References

## Project Overview

This project examines how attendance patterns, homework completion, classroom engagement, and resilience-related behavioral measures influence academic success. By merging multiple data sources — attendance logs, homework records, class participation scores, assessment performance, and mental-health/resilience measures — into a single analytical dataset (~12,156 students), the work aims to identify early risk markers and provide actionable guidance for educators to design targeted interventions.

## Problem Definition and Motivation

The project frames student success prediction as a supervised classification problem: predict target\_pass (binary: pass vs. not pass) using multimodal features that capture attendance behavior, homework completion, engagement, and resilience. Motivation includes: earlier identification of at-risk students, data-driven resource allocation, and evidence-based intervention design. Unlike single-source studies, this project fuses academic, behavioral, and resilience signals to improve prediction and interpretability.

## Data Description and Collection Effort

### Datasets used (stored in Google Drive under Practicum\_Project\_2/Data):

- Data/CleanData/ — cleaned component files (attendance\_cleaned.csv, homework\_cleaned.csv, performance\_cleaned.csv, mental\_health\_cleaned.csv, students\_cleaned.csv, mendeley\_raw\_cleaned.csv)
- Data/Final/ — merged master files (student\_master.csv, student\_master\_updated.csv, student\_master\_updated\_1.csv)
- Data/RawData/ — original sources (SEL dataset, Mendeley student dataset, mental health/resilience dataset, student performance & attendance dataset)

### Master dataset summary (final cleaned file used for modeling):

- File loaded: /content/drive/My Drive/Practicum\_Project\_2/Data/Final/student\_master\_updated\_1.csv
- Raw dataframe (post-merge) shape: **(12156, 31)** → after dropping all-NaN/time-id fields and IDs: **(12156, 16)** used for modeling
- Typical features: grade\_level, age, attendance\_rate, hw\_done\_rate, avg\_score, attendance\_pct, homework\_completion, class\_participation, engagement\_index, stress\_level, coping\_score, sleep\_hours, resilience\_score, attendance\_cat, engage\_x\_resilience
- Target: target\_pass (binary).

**Collection effort notes:** multiple datasets were joined by student identifier; several columns were engineered (attendance categories, engagement × resilience), and noisy or empty sensor-like fields were removed.

## Data Collection & Cleaning (pipeline summary)

1. **Identify & load raw files** from Data/RawData and Data/CleanData.
2. **Drop ID/text columns & all-NaN channels** (e.g., communication-text features that had no usable data).
3. **Generate target:** If raw avg\_score was scaled, target was created via median split; where raw pass threshold existed, that was used instead. (Final dataset uses target\_pass.)
4. **Imputation & normalization:** numeric fields imputed (median) and scaled (StandardScaler) as needed during modeling pipelines; categorical fields one-hot encoded.
5. **Save cleaned master** to /content/drive/My Drive/Practicum\_Project\_2/outputs/data/clean\_master.csv and record processing steps in README.md.

## Analysis Approach and Methodology

### High-level stages:

- Preprocessing: drop useless columns, impute, scale, and one-hot encode necessary categoricals (e.g., grade\_level).

- Feature selection: preserve meaningful numeric features and derived aggregates; remove leakage (avg\_score removal experiments performed).
- Balance handling: SMOTE used on training set to handle class imbalance where appropriate.
- Modeling: baseline and tuned tree models (Random Forest typically primary), logistic regression baseline, optional XGBoost experiments.
- Explainability: feature importances (RF), permutation importance, SHAP TreeExplainer for best tree model.
- Validation: stratified train/test split, cross-validation for hyperparameter tuning, ablation studies (e.g., remove age) and permutation importance to ensure robustness.

**Evaluation metrics:** Accuracy, F1-score, ROC-AUC, confusion matrix, precision/recall at class level.

## EDA, Feature Engineering, Modeling & Evaluation, Analysis and Improvement

### Exploratory Data Analysis (key findings)

- **Target distribution:** target\_pass ≈ 60% pass vs 40% fail.
- **Important patterns:** attendance rate, engagement index, and homework completion show meaningful separation by target in boxplots/summary statistics.
- **Correlations:** low–moderate correlations found between attendance/engagement/resilience indicators; no single feature dominated after leakage removal.
- **Missingness:** several feature columns were fully NaN and dropped; remaining features have full non-null in cleaned master.

### Feature Engineering

- Engineered features include attendance\_cat (Low/High), engage\_x\_resilience (engagement × resilience), and normalized versions of raw academic scores.
- Categorical encoding: grade\_level → one-hot (Grade 1–5).
- Leakage handling: experiments removed avg\_score from predictors and re-evaluated model performance.

### Modeling & Evaluation (selected results)

- **Pipelines used:** ColumnTransformer (numeric imputer + scaler, categorical imputer + one-hot) → SMOTE → Classifier (RandomForest).
- **Model performance (examples):**
  - RandomForest after leak-removal: Accuracy ≈ 0.504, F1 ≈ 0.556, ROC-AUC ≈ 0.496 (test set).
  - Logistic Regression baseline: Accuracy ≈ 0.507, F1 ≈ 0.551.
  - Earlier (leaky) runs produced inflated perfect metrics (1.0) which were diagnosed and corrected by removing leakage features.

- **Feature importance:** After rebuilding the preprocessor, RF importances consistently highlight age, grade\_level dummies, engagement\_index, and engage\_x\_resilience. Permutation importance confirmed relative ordering (negative permuted decreases indicate usefulness).
- **Explainability:** SHAP TreeExplainer applied on the top-performing tree model; key SHAP plots and force plots were produced for interpretation.

### **Analysis & Improvement Steps Taken**

- Detected target leakage, removed offending features (notably scaled score duplicates), rebuilt pipeline and re-ran.
- Performed ablation studies (e.g., remove age) to test feature influence; age removal decreased performance slightly, indicating limited but non-dominant effect.
- Reconstructed a robust feature-name alignment pipeline to ensure feature importances and permutation importances map correctly to transformed columns.
- Saved artifacts: models, feature importance CSVs, permutation\_importances\_fixed.csv, confusion matrices, and plots in /outputs.

### **Key Tools and Libraries Used**

- Python (pandas, numpy)
- scikit-learn (preprocessing, models, ColumnTransformer, RandomizedSearchCV)
- imbalanced-learn (SMOTE, imblearn.pipeline)
- xgboost (optional)
- shap (for explainability)
- matplotlib, seaborn (visualizations)
- joblib (model persistence)
- Google Colab environment & Google Drive for storage

### **Practicum Deliverables and Evaluation Rubric Alignment**

This project targets each rubric item: problem sourcing, data cleaning complexity, rigorous EDA, model building & optimization, handling class imbalance, interpretability, and polished presentation. Specific deliverables include:

- Cleaned master dataset saved under /outputs/data/clean\_master.csv
- Modeling notebook(s) and hyperparameter search artifacts (/outputs/models/)
- Visualizations: distribution plots, correlation heatmap, feature importances, ROC/PR curves (/outputs/charts/)
- README and notebook documentation for reproducibility.

## Week-by-Week Timeline and Progress

### Week Activities & Status

- 1 Proposal, problem definition — **Done**
- 2 Dataset review and project plan — **Done**
- 3 Data sources finalized and folder structure — **Done**
- 4 Data cleaning, merging, initial EDA — **Done**
- 5 Final master dataset & feature engineering — **Done**
- 6 Baseline modeling, deeper EDA, clustering experiments — **Done**
- 7 Final evaluation, visualization polishing, dashboard prototyping — **Done**
- 8 Final report, presentation recording, repo packaging — **Done**

### Roadblocks and Learnings

- **Target leakage** was detected (e.g., scaled avg\_score duplicates) and required removing features and re-running models — a useful lesson in careful feature audits.
- **Feature-name mismatch** between transformed arrays and feature importance outputs required building robust feature-name mapping logic.
- **Clustering sensitivity** to scaling demanded multiple normalization strategies and qualitative checks.
- **Compute time:** permutation importance and SHAP were computationally heavy; sampling strategies were used for explainability stages.  
Lessons learned: build reproducible preprocessing pipelines, validate for leakage early, and always align transformed features to importance outputs.

### Online Presence and Documentation

- A professional GitHub repo will include: polished README, notebooks organized by stage (data cleaning, EDA, modeling, explainability), saved outputs under /outputs, and an inference helper script under /src.
- Executive summary and visuals will be placed on the repo landing page for quick consumption by non-technical stakeholders.
- A small text file with the repo link will be uploaded to WorldClass by the Week 7 deadline.

### Visualizations and Dashboard Summary

Outputs saved to /content/drive/My Drive/Practicum\_Project\_2/outputs/charts/ include:

- Distributions: age, attendance\_rate, hw\_done\_rate, avg\_score
  - Boxplots by target\_pass for key predictors
  - Correlation heatmap and cluster visualizations
  - RF importances and permutation importance bar charts
  - Confusion matrix, ROC and PR curves
- Planned dashboard: lightweight interactive summarization (Tableau Public or Power BI) showing top risk signals and cluster profiles to help educators triage outreach.

### **Project Structure (repo / drive layout)**

Practicum\_Project\_2/

```

├── Data/
│   ├── CleanData/
│   ├── Final/
│   └── RawData/
└── notebooks/
    ├── 01_data_cleaning.ipynb
    ├── 02_eda.ipynb
    ├── 03_modeling.ipynb
    └── 04_explainability.ipynb
└── outputs/
    ├── charts/
    ├── data/
    └── models/
└── README.md
└── presentation/

```

### **Quick Start Guide**

**Prerequisites:** Python 3.8+, Google Colab or local Jupyter, required packages listed in requirements.txt.

#### **Mount Drive in Colab:**

```

from google.colab import drive
drive.mount('/content/drive')

```

**Load and run main notebook:**

Open notebooks/03\_modeling.ipynb and run sequentially. Saved artifacts will appear under /content/drive/My Drive/Practicum\_Project\_2/outputs/.

**Run inference example (after saving final\_preprocessor.joblib and final\_classifier.joblib):**

```
from src.inference import load_models, predict_df  
clf, preproc = load_models()  
preds, proba = predict_df(df_new, preproc, clf)
```

**Authors & Contact**

**Balarama Raju Saripalli**

Email: [bsaripalli@regis.edu](mailto:bsaripalli@regis.edu)

**References**

Kaggle Student Performance and Attendance Dataset

<https://www.kaggle.com/datasets/marvyaymanhalim/student-performance-and-attendance-dataset>

Mendeley Student Dataset

<https://data.mendeley.com/datasets/4mvgvtcx8s/1>

CORE Social-Emotional Learning (SEL) Dataset <https://dataqualitycampaign.org/wp-content/uploads/2018/03/DQC-CORE-CaseStudy-2018Mar22.pdf>

Kaggle Student Mental Health & Resilience Dataset

[https://www.kaggle.com/datasets/zoya07/student-mental-health-and-resilience-dataset](https://www.kaggle.com/datasets/ziya07/student-mental-health-and-resilience-dataset)

Why scientists need to be better at data visualization ([Betsy Mason](#) 11.12.2019)

<https://knowablemagazine.org/content/article/mind/2019/science-data-visualization>