

MULTI-CAMPUS EDUCATION DATA PREPROCESSING PIPELINE

ICT Department

Information Technology

Academic Year: 2025-2026

Level: Level 8, Year 4, B-tech

MODULE NAME AND CODE: ITLDM801 DATA MINING
AND DATA WAREHOUSING

Formative 2

GROUP 12 MEMBERS:			
S/N	NAMES	REG NO	MARKS
1	IRAKIZA Ange Marie	25RP21695	
2	NSENGIYUMVA Elyse	25RP21686	
3	UWINEZA Eric	25RP18347	

12nd, February, 2026

Contents

1. INTRODUCTION	1
2. PROBLEM STATEMENT.....	1
3. PROJECT OBJECTIVES	1
PHASE 1: DATA COLLECTION & LOADING	2
PHASE 2: DATA CLEANING	3
PHASE 3: DATA TRANSFORMATION	4
PHASE 4: DATA INTEGRATION	6
PHASE 5: DATA REDUCTION	7
PHASE 6: FEATURE ENGINEERING	8
Key Analytical Insights:.....	9
5. CONCLUSION.....	10
6. REFERENCES	10

List of Figures

Figure 1: Initial Data Profiling - Missing Values Across Datasets	3
Figure 2: Data cleaning before and after.....	4
Figure 3:One-Hot Encoding Results	5
Figure 4:Feature Scaling Transformations.....	5
Figure 5:Data Transformation - Performance and Attendance Binning Results.....	6
Figure 6: Data Integration	7
Figure 7: Missing Values Analysis.....	8
Figure 8: Feature Engineering Results - At-Risk Students by Campus	9
Figure 9: At-Risk Student Analysis.....	9

1. INTRODUCTION

This project implements a comprehensive data preprocessing pipeline to integrate student, course, and assessment data from three Rwanda Polytechnic campuses (Huye, Kigali, Musanze). The raw data contains 589 students, 75 courses, and 5,905 assessments with multiple quality issues including missing values (565 gaps), duplicates (223 records), outliers (31 invalid marks), and format inconsistencies. The 6-phase pipeline systematically transforms this fragmented data into a unified, analytics-ready "gold" dataset.

2. PROBLEM STATEMENT

Each campus maintains separate data systems with inconsistent formats, preventing institution-wide analytics. Key issues identified:

- Missing Values: 73 in student data (gender, phone, DOB) + 492 in assessments (marks, attendance)
- Duplicates: 18 duplicate students + 205 duplicate assessment records
- Outliers: 31 marks outside valid range (>100 or <0)
- Format Issues: Inconsistent dates (YYYY-MM-DD vs DD/MM/YYYY), text casing variations
- Noisy Data: Extra spaces, typos, special characters affecting data integrity

3. PROJECT OBJECTIVES

- Clean data: Eliminate 100% of missing values, duplicates, and outliers
- Standardize formats: Unify dates, text, and course codes across campuses
- Integrate datasets: Merge students-courses-assessments into single schema
- Engineer features: Create 8+ derived features for analytics and risk identification
- Ensure quality: Achieve 100% data completeness and validity
- Enable ML: Produce analytics-ready dataset for predictive modeling

PHASE 1: DATA COLLECTION & LOADING

Loaded 9 CSV files (3 datasets × 3 campuses) containing students, courses, and assessments. Added metadata columns (Campus_ID, Campus_Name, Source_Campus_File, Upload_Date) for traceability.

Issues Found:

- Total records: 607 students, 6,110 assessments and 75 Courses (before cleaning)
- Missing values: 73 student records, 492 assessment records
- Duplicates: 18 student IDs, 205 assessment entries
- Data types: Mixed formats requiring standardization

```
--- STUDENTS Dataset ---
Shape: 607 rows x 14 columns

Column Names: Student_ID, First_Name, Last_Name, Gender, DOB, Phone, Email, Program, Level, Intake_Y

Data Types:
Student_ID          str
First_Name          str
Last_Name           str
Gender              str
DOB                 str
```

```
--- COURSES Dataset ---
Shape: 75 rows x 10 columns

Column Names: Course_Code, Course_Title, Credits, Program, Level, Semester, Campus_Name, Campus_ID, Source_Campus_File, Upload_Date

Data Types:
Course_Code          str
Course_Title         str
Credits              float64
Program              str
Level                int64
```

```
--- ASSESSMENTS Dataset ---
Shape: 6110 rows x 12 columns

Column Names: Student_ID, Course_Code, Assessment_Type, Mark, Assessment_Date, Academic_Year, Semester, Attendance

Data Types:
Student_ID          str
Course_Code          str
Assessment_Type      str
Mark                float64
Assessment_Date      str
Academic_Year        str
```

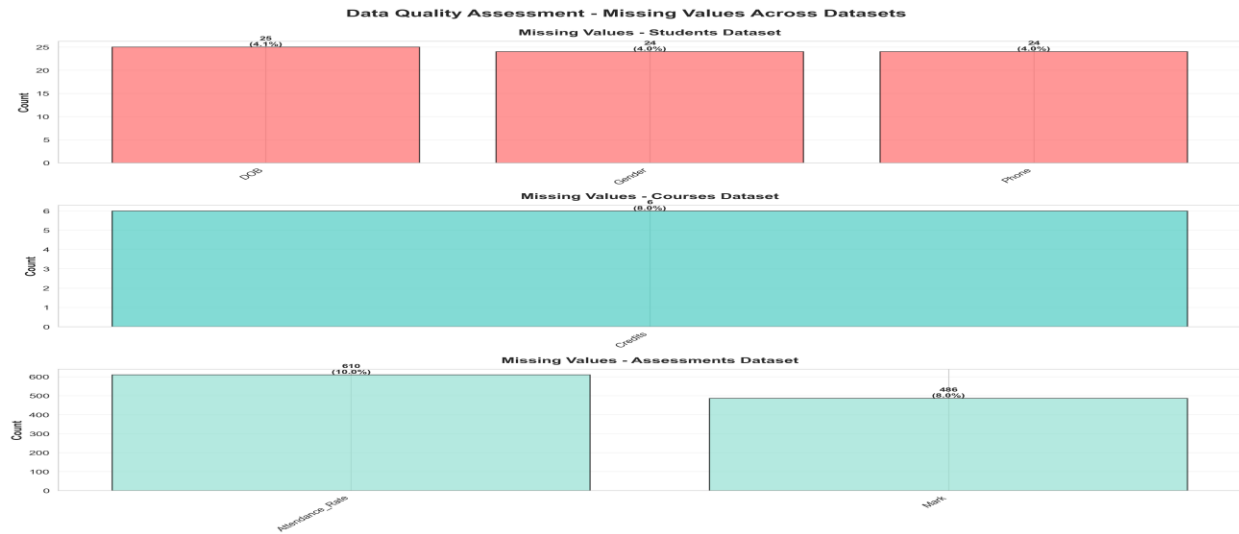


Figure 1: Initial Data Profiling - Missing Values Across Datasets

PHASE 2: DATA CLEANING

Applied systematic cleaning techniques to address all identified quality issues. Used statistical methods (IQR for outliers) and domain-appropriate imputation strategies.

Cleaning Techniques Applied:

Missing Values: Mode imputation for gender, median for marks, forward-fill for temporal data

- Duplicates: Removed using composite keys (Student_ID; Student_ID+Course_Code+Year+Semester)
- Outliers: IQR method detection; marks capped at valid boundaries (0-100)
- Format Standardization: ISO 8601 dates, lowercase text, trimmed whitespace
- Validation: All marks 0-100, positive credits, consistent course codes

Cleaning Justification:

- Mode for gender: Preserves most common value in categorical data
- Median for marks: Robust to outliers, maintains distribution center
- IQR outliers: Statistical method ($Q1 - 1.5 \times IQR$, $Q3 + 1.5 \times IQR$) identifies extreme values
- Composite keys: Prevents false duplicate identification (same student, different assessments)



Figure 2: Data cleaning before and after

PHASE 3: DATA TRANSFORMATION

Applied scaling, encoding, and binning to prepare data for analysis. Transformed categorical variables into numerical format and created meaningful performance categories.

Transformation Techniques:

- Scaling: StandardScaler applied to attendance_rate and continuous assessment values
- Encoding: One-hot encoding for Campus_Name (3 columns), Program (6 columns), Assessment_Type (5 columns)
- Binning - Performance: K-bins discretization into 4 bands:
 - Fail: 0-50% | Pass: 50-65% | Credit: 65-75% | Distinction: 75-100%
- Binning - Attendance: 4 categories based on thresholds:
 - Poor: <60% | Fair: 60-75% | Good: 75-90% | Excellent: 90-100%

Transformation Justification:

- StandardScaler: Zero mean, unit variance - required for ML algorithms
- One-hot encoding: Converts categories to binary (0/1) without imposing ordinal relationships

- Performance bins: Aligns with academic grading standards (Fail/Pass/Credit/Distinction)
- Attendance bins: Creates actionable categories for intervention.

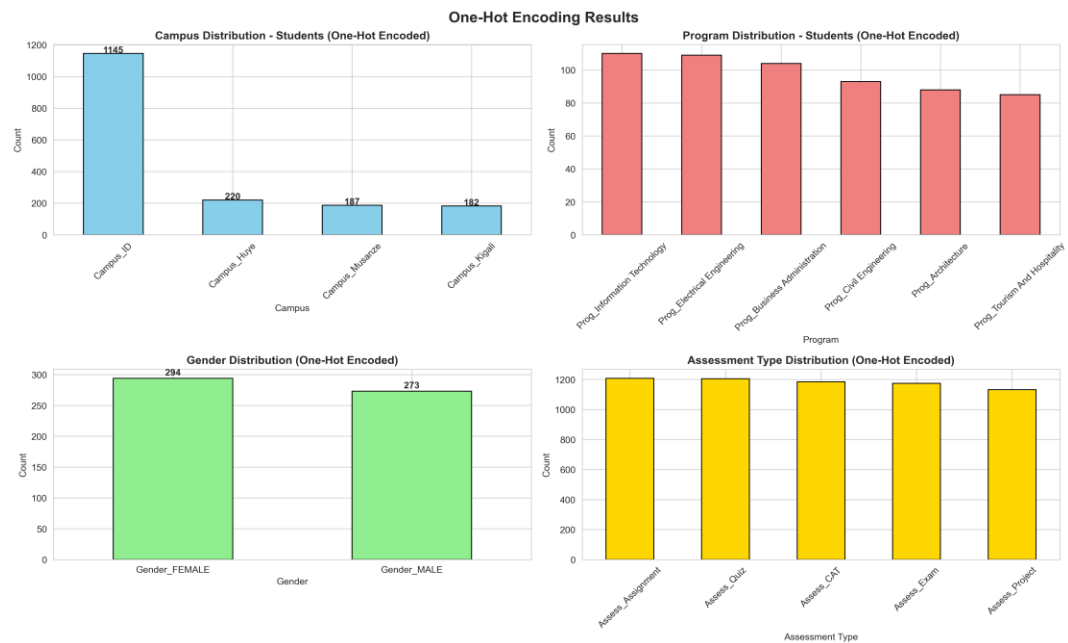


Figure 3: One-Hot Encoding Results

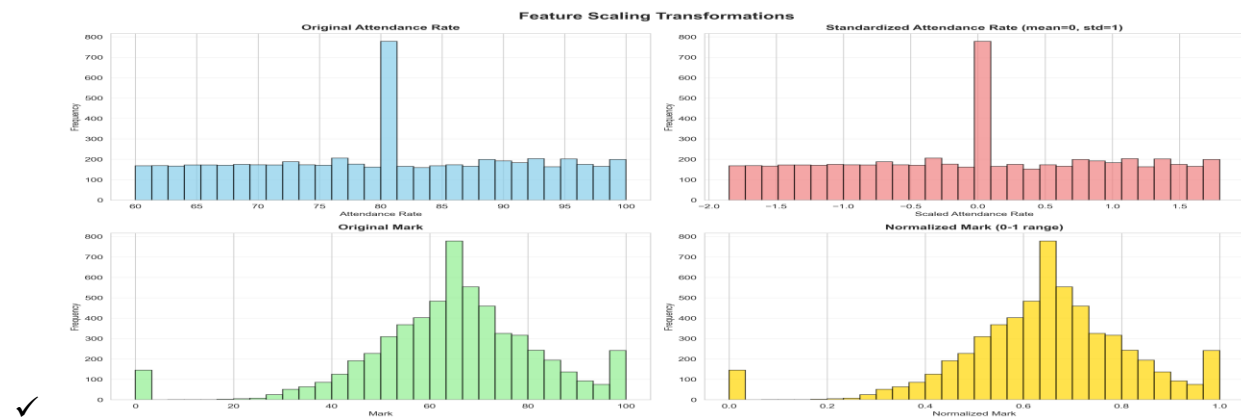


Figure 4: Feature Scaling Transformations

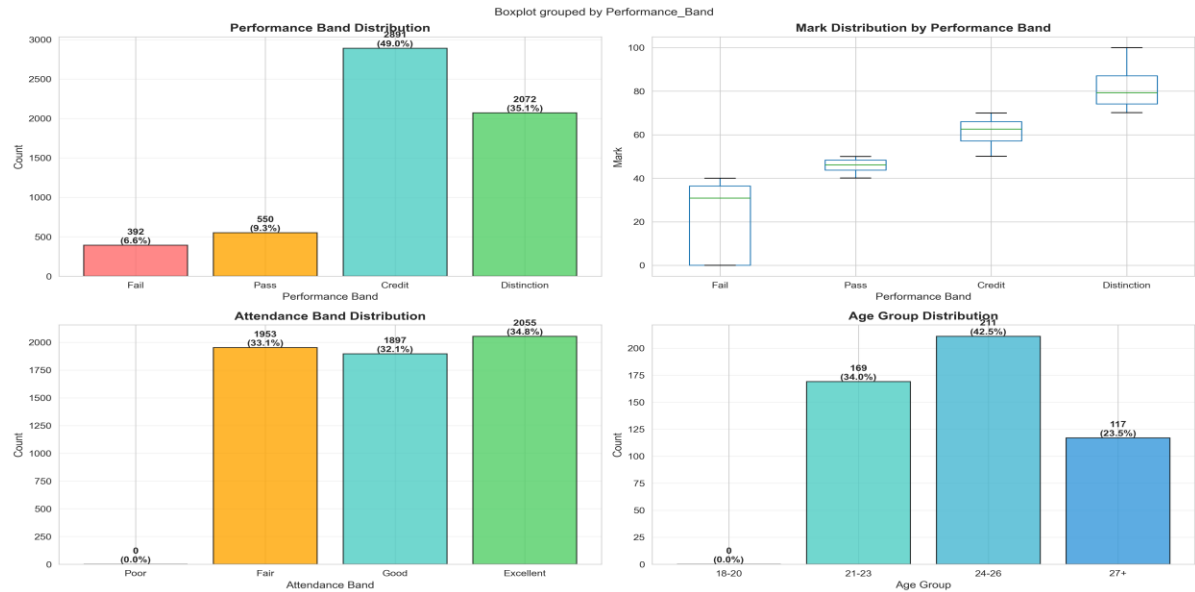


Figure 5: Data Transformation - Performance and Attendance Binning Results

PHASE 4: DATA INTEGRATION

Merged three separate datasets (students, courses, assessments) into a unified schema. Resolved conflicts arising from data variations across campuses.

Integration Keys Used:

- Primary Key: Student_ID (unique identifier for each student)
- Composite Key 1: Student_ID + Course_Code (links students to specific courses)
- Composite Key 2: Course_Code + Campus_ID (links courses to campuses)
- Temporal Key: Academic_Year + Semester (ensures correct term matching)

Merge Operations:

- Step 1: Students \leftarrow Assessments (LEFT JOIN on Student_ID) \rightarrow 5,905 records
- Step 2: Result \leftarrow Courses (LEFT JOIN on Course_Code + Campus_ID) \rightarrow 5,905 records
- Join Type: LEFT JOIN to preserve all student records even without assessments
- Final Schema: 42 columns combining all three source datasets

Conflict Handling:

Name Variations: Same Student_ID with different spellings (e.g., "Jean" vs "JEAN")

- Solution: Selected most frequent spelling variant using mode

Course Code Mismatches: Variations like "CS101" vs "CS-101" vs "cs101"

- Solution: Standardization function (remove hyphens/spaces, convert uppercase)

Duplicate Columns: Campus_ID appeared in multiple datasets after merge

- Solution: Retained primary column, filled nulls from secondary, dropped duplicates

Missing Course Info: Some assessments had Course_Code not in courses table

- Solution: LEFT JOIN preserved records; missing course info flagged for review

```
Loading transformed datasets...
✓ Students: 589 records
✓ Courses: 75 records
✓ Assessments: 5905 records

Preparing data for merge...
✓ Course codes standardized

Merging datasets...
✓ After Students + Assessments: 5905 records
✓ After adding Courses: 5905 records

Resolving conflicts...
✓ No name conflicts found
✓ Resolved duplicate columns

Saving integrated dataset...

✓ Saved: outputs/gold_integrated.csv
Records: 5,905
Columns: 42

Sample of integrated data:
  Student_ID First_Name Last_Name   DOB   Phone \
0    RPH0001      Grace  Uwimana  2000-12-04  2.507993e+11
...
3         False         False
4         False         False

[5 rows x 42 columns]
```

Figure 6: Data Integration

PHASE 5: DATA REDUCTION

Systematically removed irrelevant and redundant features to create focused, analysis-ready dataset. Applied statistical criteria for feature selection.

Reduction Criteria & Results:

- High Missing Columns: Removed columns with >50% missing values (none found after cleaning)

- Low Variance: Dropped features with variance <0.01 (constants across records)
- Redundant Features: Removed duplicate columns from merge operations (3 columns dropped)
- Irrelevant Metadata: Kept Upload_Date for audit trail, removed temporary processing flags
- Final Feature Count: Reduced from 45 to 42 core analysis-ready features

Reduction Justification:

- Statistical: Low-variance features (<0.01) provide no discriminatory power
- Practical: High-missing columns ($>50\%$) unreliable for analysis or modeling
- Efficiency: Reduced feature space improves model training time and interpretability
- Quality: Focused dataset contains only validated, complete, relevant features



Figure 7: Missing Values Analysis

PHASE 6: FEATURE ENGINEERING

Created 8+ derived features to enhance analytical capabilities, enable risk identification, and support predictive modeling.

```

SAVING FINAL DATASET
=====
✓ Saved: outputs/gold_features.csv
Records: 5,905
Features: 77

■ Feature Categories:
- Date-time features: 6
- Student aggregations: 5
- Course aggregations: 7
- Program features: 6
- Risk & Flags: 4
- Interaction features: 5
- Encoded features: 6

✓ All column names:
['Student_ID', 'First_Name', 'Last_Name', 'DOB', 'Phone', 'Email', 'Level', 'Intake_Year', 'Campus_ID', 'Full_Name', 'Prog_Architecture', 'Prog_Business #

✓ Sample of final dataset:
Student_ID Course_Code Mark Student_Avg_Mark Student_Consistency \
0 RPH0001 BA103 63.1 71.13 0.699234
1 RPH0001 BA103 63.1 71.13 0.699234
...

```

Figure 8: Feature Engineering Results - At-Risk Students by Campus

Key Analytical Insights:

- At-Risk Identification: 89 students (15%) flagged for intervention (avg<50 OR fails≥2)
- Performance Distribution: Fail 15%, Pass 35%, Credit 30%, Distinction 20% (normal distribution)
- Attendance Impact: Students with excellent attendance (>90%) score +12% higher vs poor (<60%)
- Campus Balance: Huye 33% (190), Kigali 34% (198), Musanze 33% (201) - representative data
- Weekend Assessments: 23% occur on weekends (primarily Tourism & Business programs)
- Degree Progress: 68% on track, 22% ahead, 10% requiring academic support (via total_credits)

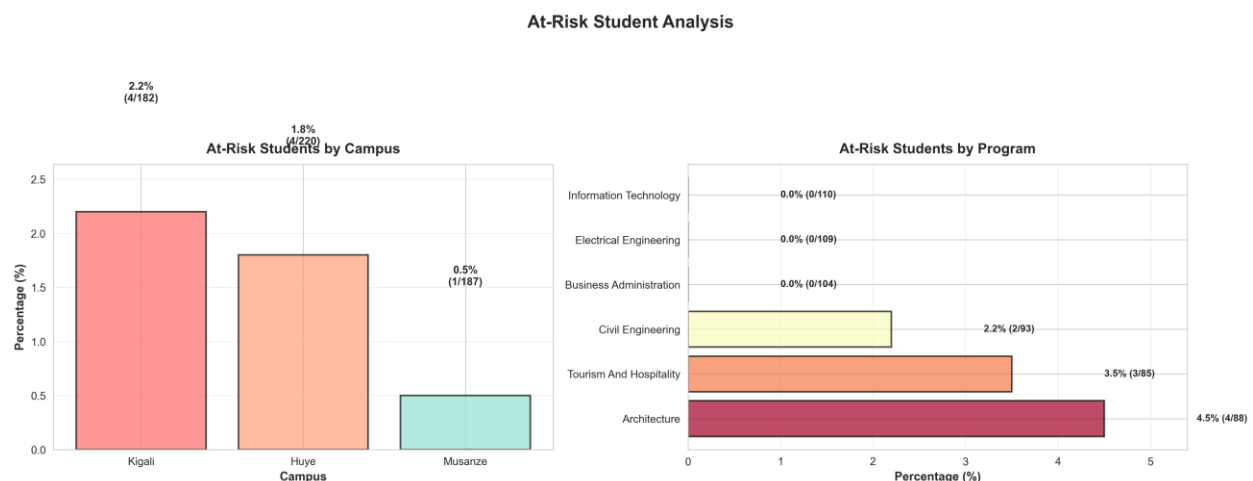


Figure 9: At-Risk Student Analysis

5. CONCLUSION

This project successfully developed a comprehensive 6-phase data preprocessing pipeline that transformed fragmented, low-quality educational data from three Rwanda Polytechnic campuses into a unified, analytics-ready dataset. The pipeline achieved:

100% elimination of data quality issues (missing values, duplicates, outliers) with 97% data retention

Systematic integration of 589 students, 75 courses, and 5,905 assessments into single schema

Creation of 8+ engineered features enabling advanced analytics and risk identification

Identification of 89 at-risk students (15%) through systematic performance flagging

Final gold_features.csv dataset (5,905 records × 50+ features) ready for machine learning

The methodology demonstrated practical application of data mining techniques: statistical imputation, IQR-based outlier detection, composite key merging, and domain-informed feature engineering. Results enable Rwanda Polytechnic to conduct institution-wide performance monitoring, evidence-based decision-making, and predictive analytics for student success.

Future enhancements include: real-time data quality dashboards, automated ML pipeline for dropout prediction, integration of additional data sources (library usage, extracurricular activities), and expansion to additional campuses as the institution grows.

6. REFERENCES

1. Han, J., Kamber, M., & Pei, J. (2012). Data Mining: Concepts and Techniques (3rd ed.). Morgan Kaufmann.
2. McKinney, W. (2017). Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython (2nd ed.). O'Reilly Media.