

## Predicting Student Exam Scores - Notes

### Dataset Overview

#### Columns:

- student\_id
- age
- gender
- course
- study\_hours
- class\_attendance
- internet\_access
- sleep\_hours
- sleep\_quality
- study\_method
- facility\_rating
- exam\_difficulty
- exam\_score

**Target:** exam\_score

**Dataset Size:** 20,000 entries

#### Data types:

- Float: study\_hours, class\_attendance, sleep\_hours, exam\_score
  - Integer: student\_id, age
  - ~~Object (numerical)~~: gender, course, internet\_access, sleep\_quality, study\_method, facility\_rating,
- 

### General Observations

- **Irrelevant features:** student\_id, age, course
- **Highly relevant features:**

- study\_hours: More study hours → higher marks
  - class\_attendance: Higher attendance → higher marks, but some variance exists
  - study\_method: Efficient methods → higher marks
  - facility\_rating: Better facilities → higher marks
- **Moderately relevant features:**
- sleep\_hours: Both too little and too much sleep can reduce performance
  - sleep\_quality: Good quality may slightly improve marks
  - exam\_difficulty: Harder exams → lower marks

- **Interactions & patterns:**

- High study\_hours + low sleep\_hours → possibly high marks
  - Low attendance + high study\_hours + easy exam → likely high marks
  - High sleep\_hours + high study\_hours → good marks
- 

## Correlation with Exam Score

Feature	Correlation	Strength
Study Hours	0.718	Strong
Class Attendance	0.309	Moderate
Sleep Hours	0.133	Weak

Interpretation: study\_hours is the most predictive; attendance has moderate influence; sleep hours alone has weak effect.

---

## Linear Regression

### Steps:

1. Drop student\_id

2. Encode categorical features
3. Split into training (80%) and testing (20%) sets
4. Train Linear Regression model
5. Evaluate using RMSE

### Results:

- Validation RMSE  $\approx 9.77$
- Interpretation:
- On average, predictions are off by  $\pm 10$  marks (reasonable for 0–100 range, but not strong)

### Limitations of Linear Regression:

- Cannot model non-linear relationships (e.g., optimal sleep hours, difficulty impact)
  - Cannot capture interactions ( $\text{StudyHours} \times \text{Attendance} \times \text{Method}$ )
  - Treats categorical features as additive
  - RMSE  $\sim 10$  is expected for this dataset
- 

### Random Forest Regression

- Captures non-linear patterns and feature interactions
- Train RMSE  $\approx 4.67$ , Validation RMSE  $\approx 10.55 \rightarrow$  overfitting observed
- Feature Importance (top 10):
  1. study\_hours: 58.5%
  2. class\_attendance: 15.9%
  3. sleep\_hours: 6.5%
  4. age: 2.5%
  5. sleep\_quality\_poor: 2.3%
  6. facility\_rating\_low: 2.2%
  7. sleep\_quality\_good: 1.96%

8. study\_method\_self-study: 1.09%
9. study\_method\_online videos: 0.93%
10. facility\_rating\_medium: 0.77%

### **Reason for overfitting:**

- Bounded exam scores (0–100)
  - Noise from human behavior (sleep, study habits)
  - Deep trees (max\_depth=15) overfit sparse patterns
- 

### **CatBoost Regressor**

- Handles categorical features natively
  - Train RMSE  $\approx 9.47$ , Validation RMSE  $\approx 9.78 \rightarrow$  minimal overfitting
  - Feature importance (CatBoost):
    - study\_hours: 50.7%
    - class\_attendance: 18.7%
    - sleep\_quality: 9.6%
    - study\_method: 8.3%
    - facility\_rating: 6.8%
    - sleep\_hours: 4.7%
    - Others <1%
  - Slight hyperparameter adjustments  $\rightarrow$  RMSE  $\approx 9.82$
- 

### **Ensemble Model**

- Combining predictions slightly improved RMSE
  - Ensemble RMSE  $\approx 9.77$
- 

### **Key Insights**

1. **Most important features:** study\_hours, class\_attendance, sleep\_quality, study\_method
2. **Linear models** are limited due to non-linear patterns and interactions
3. **Random forests** may overfit without careful tuning
4. **CatBoost** is effective for categorical-heavy datasets and generalizes well
5. Noise and bounded score ranges limit the achievable RMSE (~9–10)