



# PREDICT STUDENT TEST SCORE

AIO Conquer  
Warm Up Project Presentation

Presented by Team CONQ013

# MOTIVATION & OBJECTIVE

## Motivations

- To understand how a real ML project is built end-to-end
- To learn why each stage in the ML workflow is necessary
- Performance matters, but is not the primary goal

## Objectives

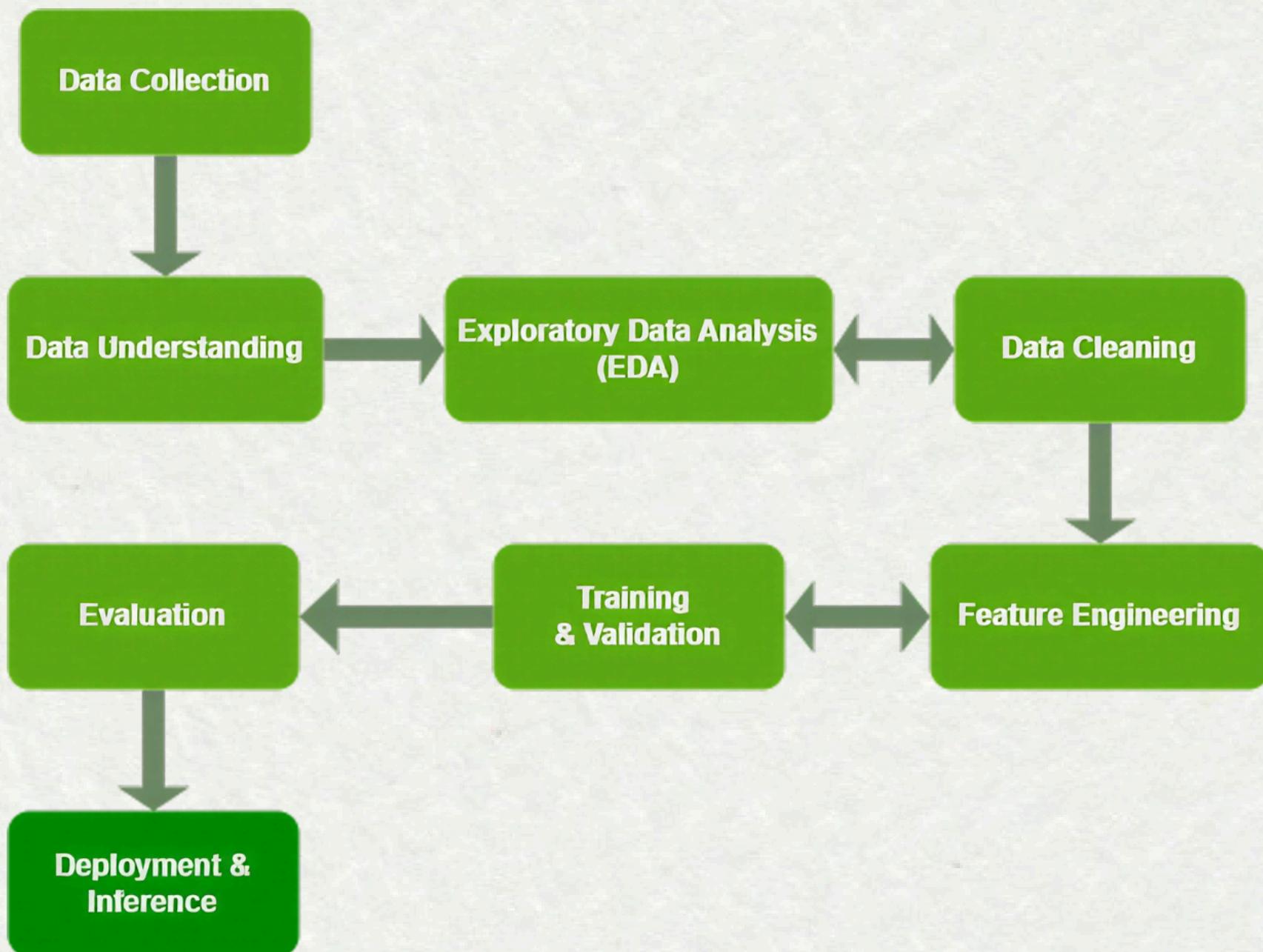
- Predict student exam score based on personal & academic factors using ML
- **Input:** study habits, sleep behavior
- **Output:** Predicted exam score (0-100)
- **Real-world value:** Help educators identify at-risk students early

## Machine Learning Framing

- Type: Supervised Learning
- Task: Regression



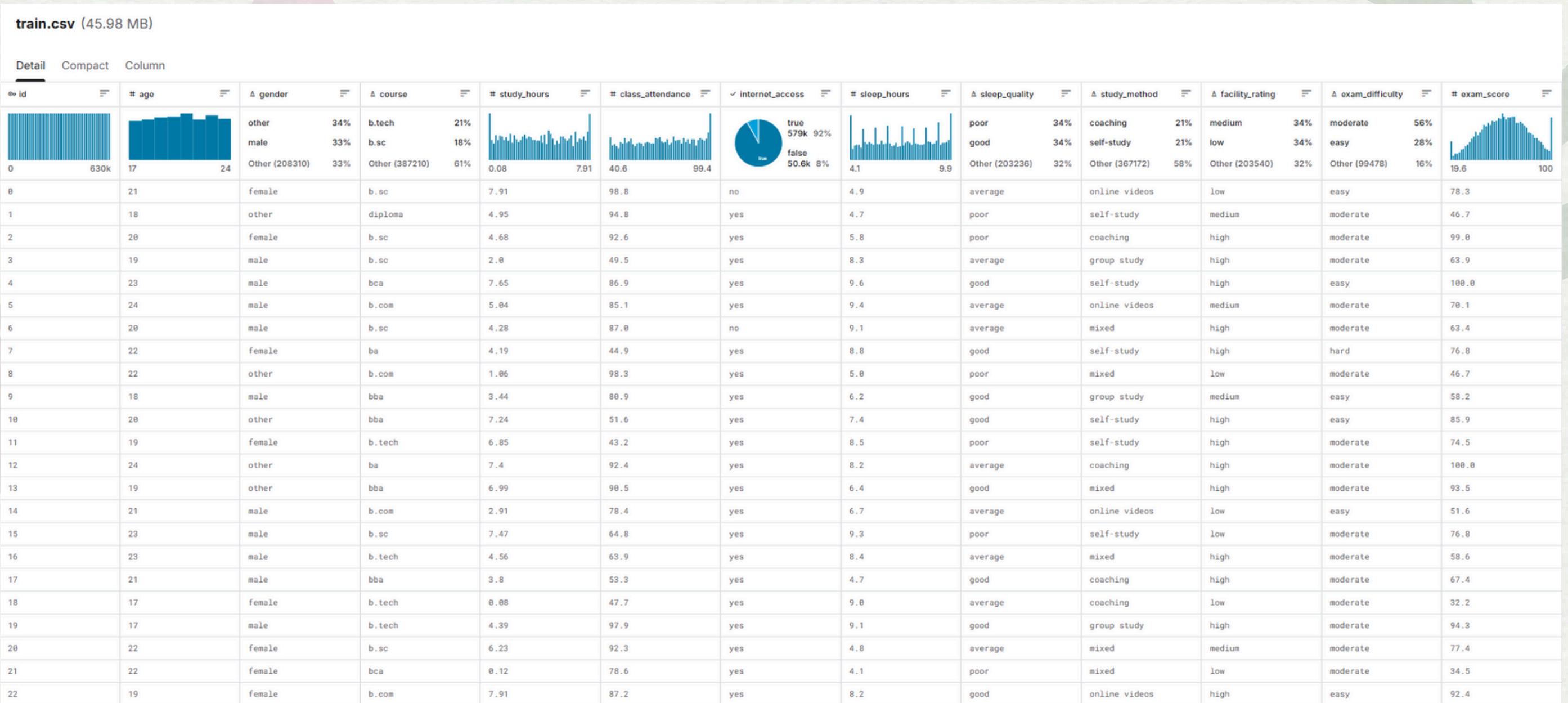
# OVERALL ML WORKFLOW



The workflow is inherently iterative, with continuous feedback between data, features, and models.

# DATA COLLECTION

Dataset: from Kaggle Tabular Playground Series Season 6 Episode 1: "Predicting Student Test Scores"



Synthetically-generated  
from deep learning model

# DATA UNDERSTANDING

Using `df.info()`, `df.describe()`, and `df.head()`, we initially discovered a remarkably clean dataset:

## Key observations:

- Dataset size: ~630,000 records
- Feature overview: 13 features, mix of numerical and categorical variables
- Target variable: `exam_score` (continuous)
- **Exam scores range widely from 19.6 to 100 points**
- **Average study time: ~4 hours daily**
- **Average sleep: ~7 hours nightly**

[4]:	df.head()												
[4]:	id	age	gender	course	study_hours	class_attendance	internet_access	sleep_hours	sleep_quality	study_method	facility_rating	exam_difficulty	exam_score
0	0	21	female	b.sc	7.91	98.8	no	4.9	average	online videos	low	easy	78.3
1	1	18	other	diploma	4.95	94.8	yes	4.7	poor	self-study	medium	moderate	46.7
2	2	20	female	b.sc	4.68	92.6	yes	5.8	poor	coaching	high	moderate	99.0
3	3	19	male	b.sc	2.00	49.5	yes	8.3	average	group study	high	moderate	63.9
4	4	23	male	bca	7.65	86.9	yes	9.6	good	self-study	high	easy	100.0

## df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 630000 entries, 0 to 629999
Data columns (total 13 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               630000 non-null   int64  
 1   age              630000 non-null   int64  
 2   gender            630000 non-null   object 
 3   course             630000 non-null   object 
 4   study_hours       630000 non-null   float64
 5   class_attendance 630000 non-null   float64
 6   internet_access   630000 non-null   object 
 7   sleep_hours        630000 non-null   float64
 8   sleep_quality      630000 non-null   object 
 9   study_method        630000 non-null   object 
 10  facility_rating    630000 non-null   object 
 11  exam_difficulty    630000 non-null   object 
 12  exam_score         630000 non-null   float64
dtypes: float64(4), int64(2), object(7)
memory usage: 62.5+ MB
```

## df.describe()

	id	age	study_hours	class_attendance	sleep_hours	exam_score
count	630000.000000	630000.000000	630000.000000	630000.000000	630000.000000	630000.000000
mean	314999.500000	20.545821	4.002337	71.987261	7.072758	62.506672
std	181865.479132	2.260238	2.359880	17.430098	1.744811	18.916884
min	0.000000	17.000000	0.080000	40.600000	4.100000	19.599000
25%	157499.750000	19.000000	1.970000	57.000000	5.600000	48.800000
50%	314999.500000	21.000000	4.000000	72.600000	7.100000	62.600000
75%	472499.250000	23.000000	6.050000	87.200000	8.600000	76.300000
max	629999.000000	24.000000	7.910000	99.400000	9.900000	100.000000

# EDA - DATA QUALITY

## Missing Value:

Use `df.isnull().sum()` to count missing entries.

- Less than 5%? Drop rows or impute.
- Between 5-30%? Use advanced imputation (mean, median, KNN).
- Over 30%? Consider dropping the column.

```
df.isnull().sum()
```

```
age           0  
gender        0  
course         0  
study_hours    0  
class_attendance 0  
internet_access 0  
sleep_hours     0  
sleep_quality    0  
study_method     0  
facility_rating   0  
exam_difficulty 0  
exam_score       0  
dtype: int64
```

## Duplicate Rows

Check with `df.duplicated().sum()`.

Duplicates can cause overfitting and data leakage. Remove them carefully, keeping first occurrence with `keep='first'`.

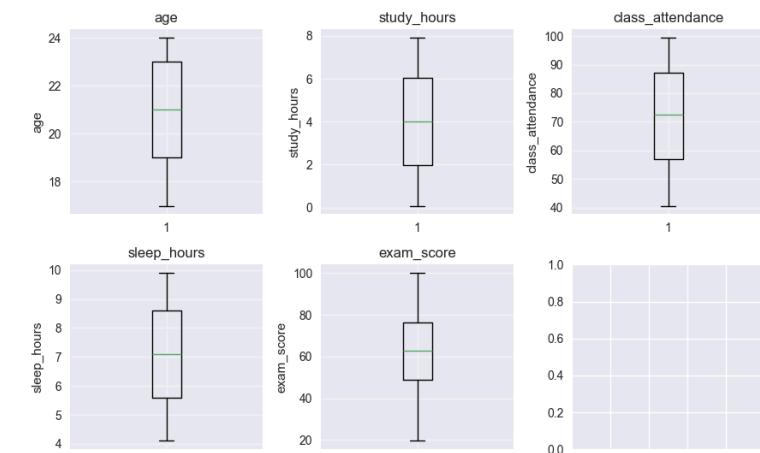
```
df.duplicated().sum()
```

0

## Outliers

Apply IQR method: Calculate Q1, Q3, and IQR. Outliers fall below  $Q1 - 1.5 \times IQR$  or above  $Q3 + 1.5 \times IQR$ .

Don't automatically delete—some outliers are valuable insights!

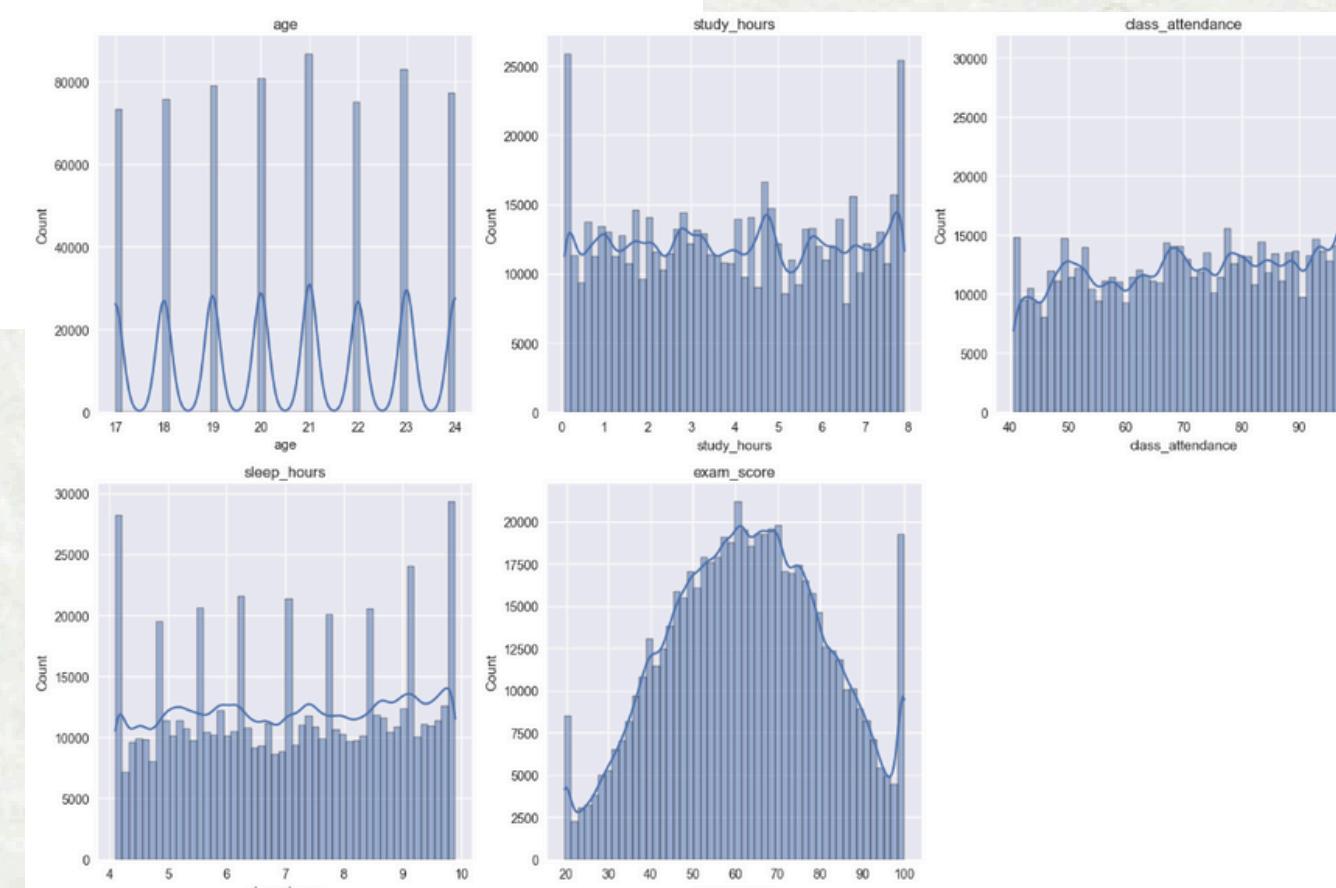
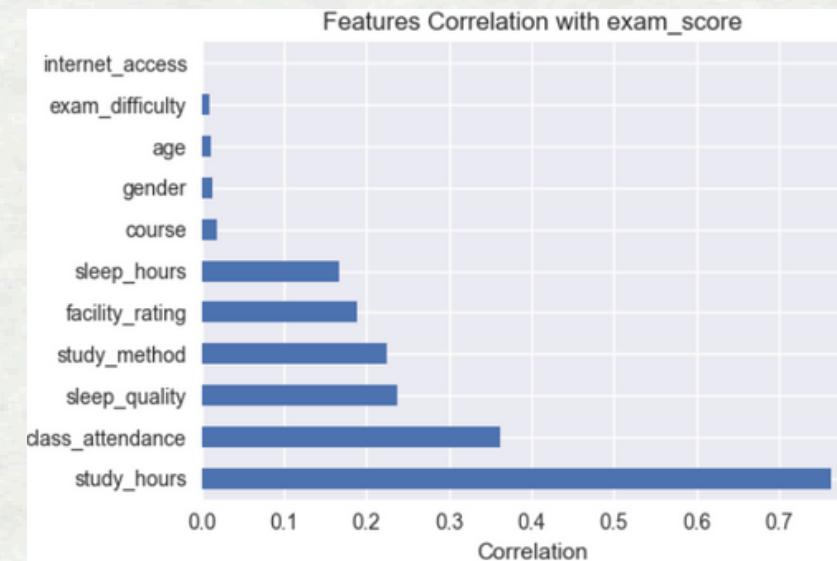
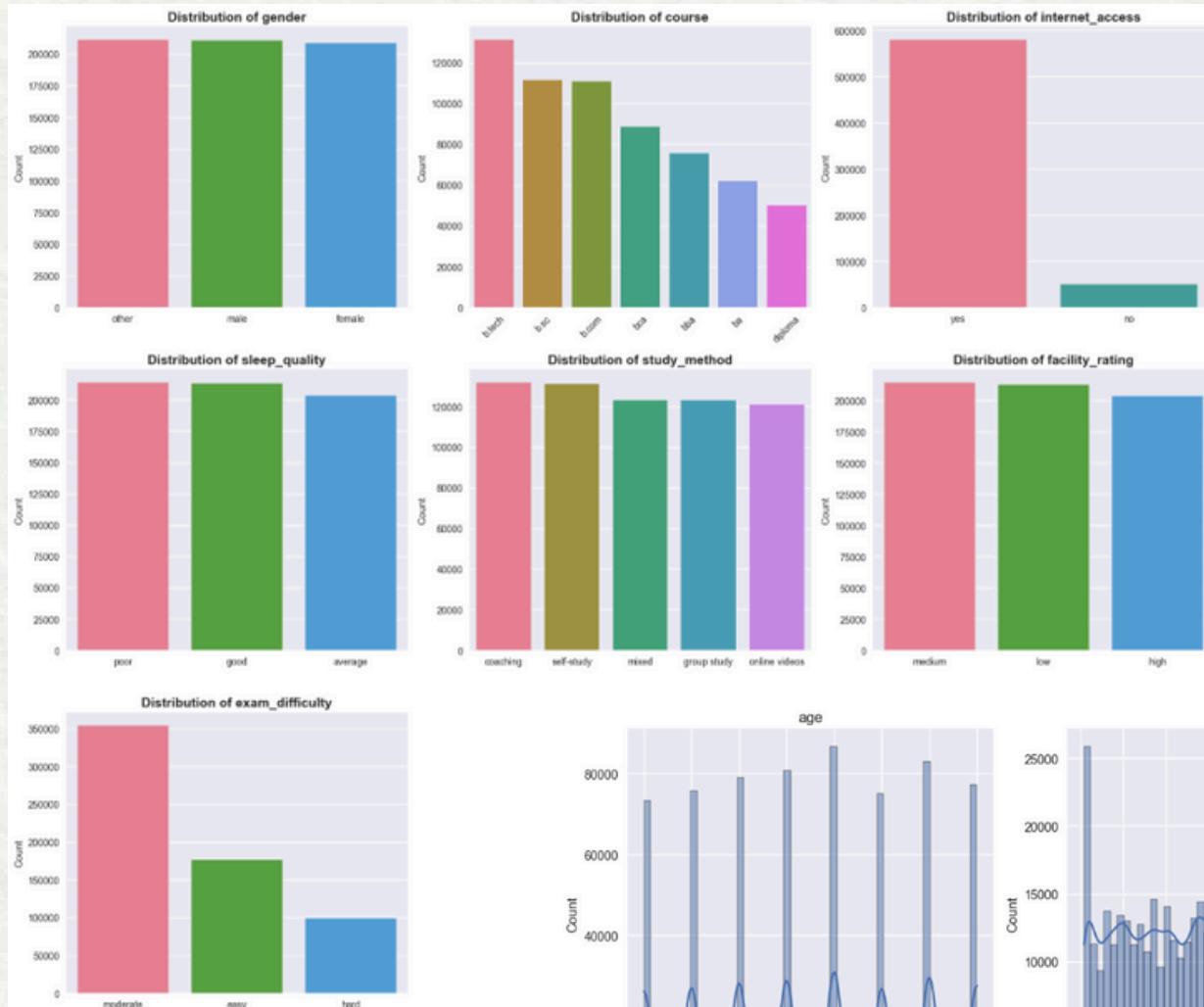


## Data Logic

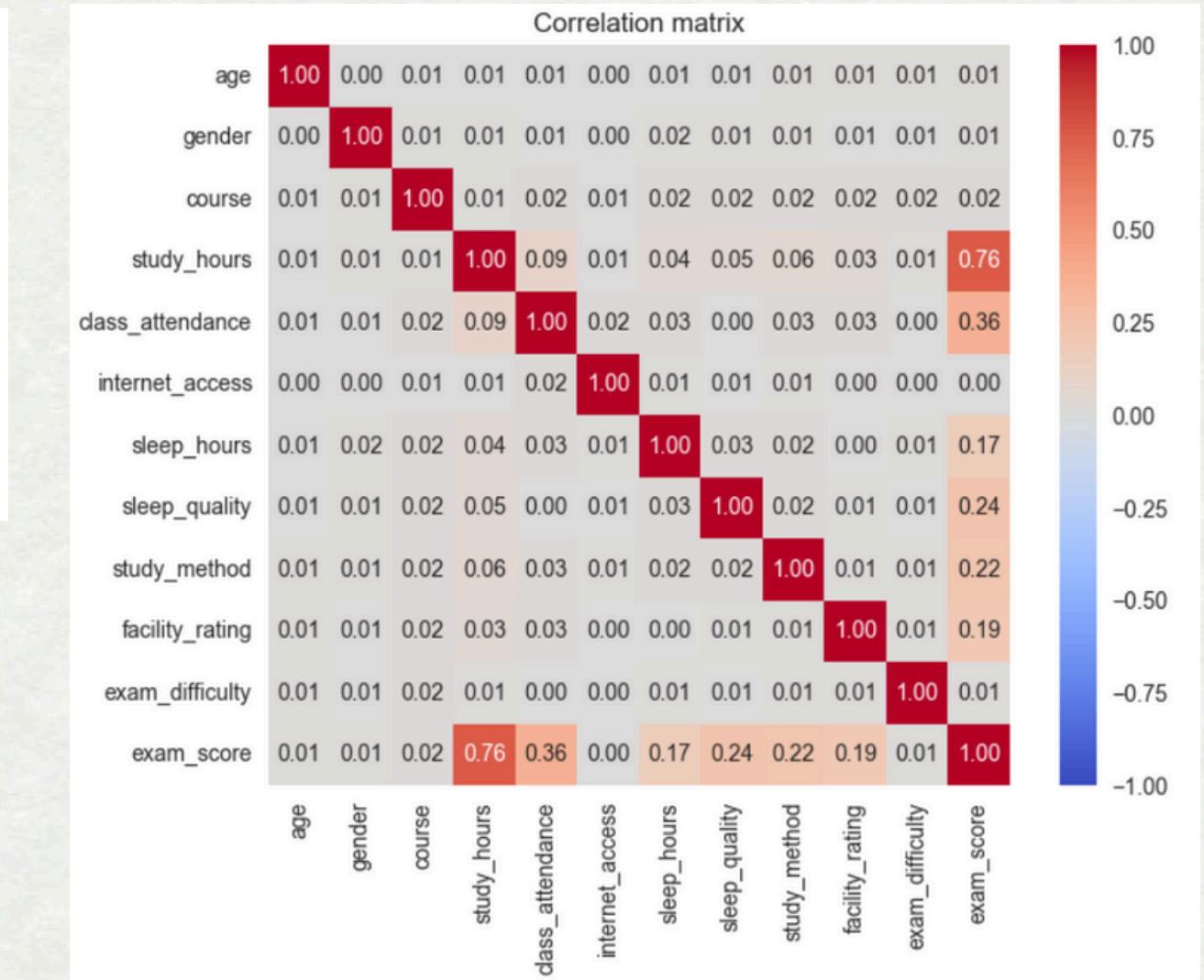
Verify data types match expected formats (ex: datetime, numeric vs ). Check score range (0-100), sleep quality (0-24). Fix encoding issues and standardize formats.

# EDA - DATA DISTRIBUTION

## Distribution Analysis



## Relationship Analysis

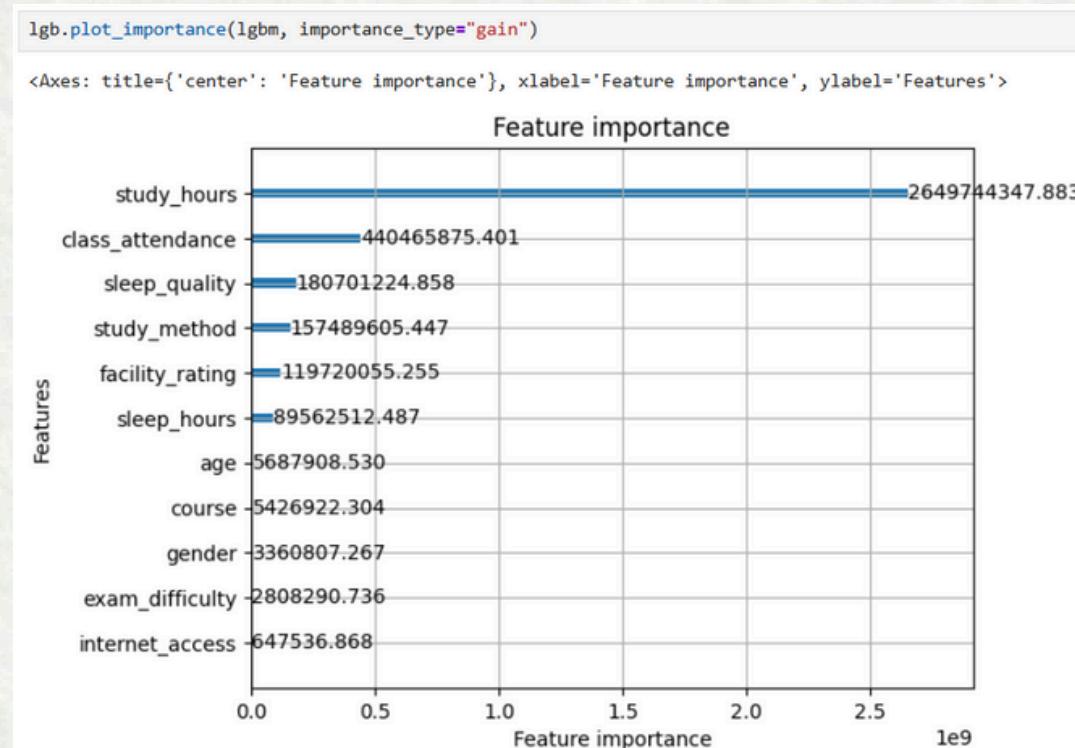


# FEATURE ENGINEERING

## Remove Irrelevant Features

Drop columns like IDs with no predictive value.

Remove features with near-zero correlation (<0.05), also confirmed by tree-based model importance analysis.



## Feature Encoding

### Numerical features:

Used in their original scale without transformation.

### Categorical features:

- Ordinal encoding for ordered variables (sleep\_quality, facility\_rating, exam\_difficulty).
- One-hot encoding for nominal variables (study\_method, course).

## Create New Features

### Interaction features:

- $\text{study\_efficiency} = \text{study\_hours} \times \text{class\_attendance}$ ,
- $\text{sleep\_attend\_product} = \text{sleep\_hours} \times \text{class\_attendance}$ .

### Deviation feature:

- $\text{sleep\_deviation} = |\text{sleep\_hours} - 7|$  to capture deviation from optimal sleep duration.

# MODEL TRAINING

## Model Selection Strategy

- Baseline: Linear Regression
- Regularized Models: Ridge (L2), Lasso (L1)
- Final Model: LightGBM (tree-based)

## Training Setup

- Train/Validation split: 80/20
- Early stopping

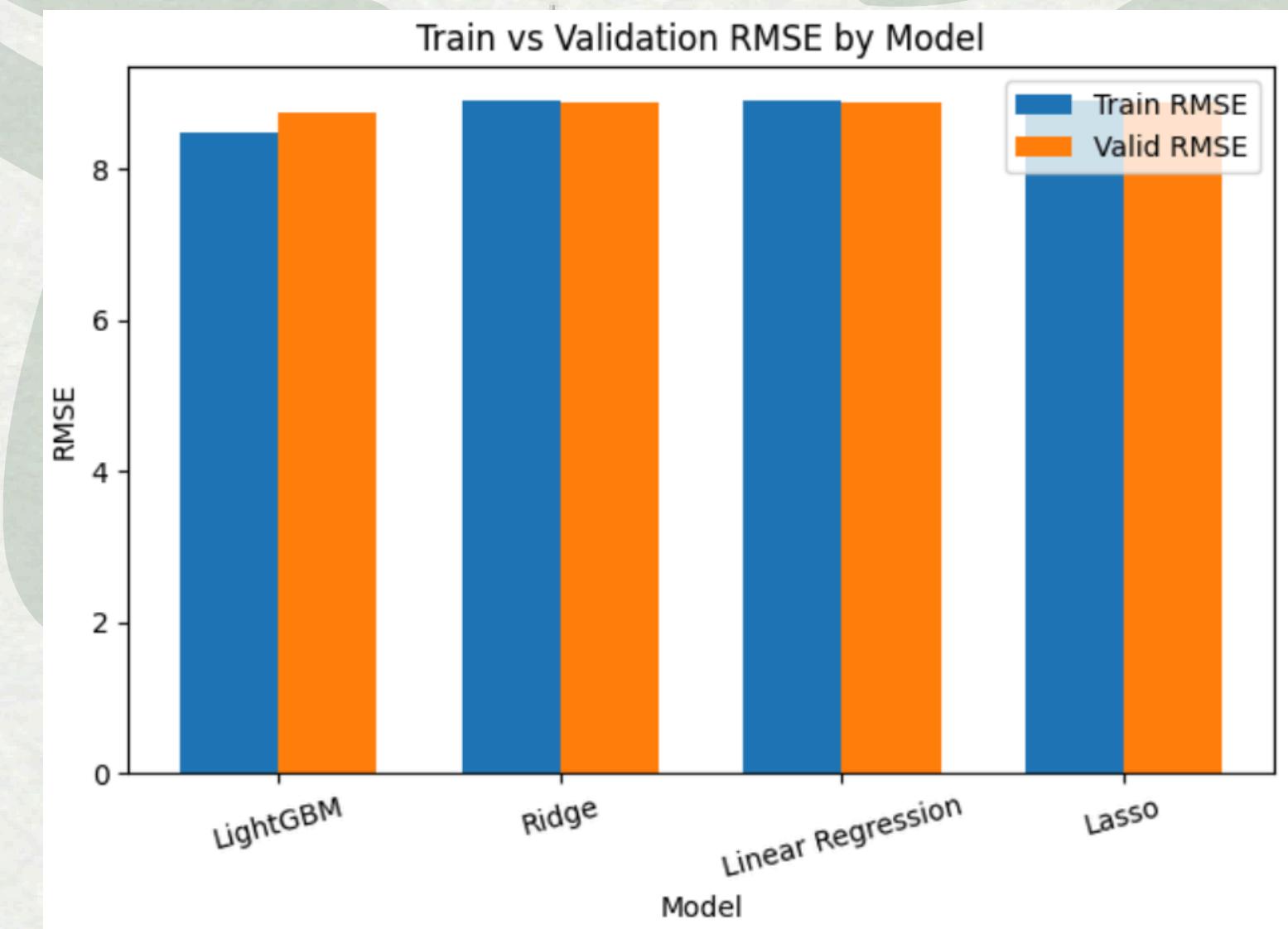
## Why LightGBM?

- Captures non-linear relationships and feature interactions
- Fast training on large dataset (630K records)
- Best performance: RMSE  $\approx$  8.75 on validation set

# EVALUATION

## Performance Metric: Validation RMSE (*lower is better*)

- LightGBM achieves the lowest validation RMSE, indicating better performance.
- Our model achieved RMSE  $\approx 8.75$ , quite close to Kaggle leaderboard's 8.53—validating our EDA and feature engineering pipeline.



Predicting Student Test Scores

Overview Data Code Models Discussion Leaderboard Rules Team Submissions Late Submission ...

This leaderboard is calculated with approximately 20% of the test data. The final results will be based on the other 80%, so the final standings may be different.

#	Team	Members	Score	Entries	Last	Solution
1	Mahog		8.53096	51	3d	
2	Yew Jin Lim		8.53303	61	3d	
3	Traiko Dinev		8.53475	155	2d	

=====

MODEL COMPARISON SUMMARY

=====

Model	Train RMSE	Valid RMSE	Valid R <sup>2</sup>	Valid MAE	Overfit Gap
LightGBM	8.476879	8.749616	0.784742	6.972499	-0.272737
Ridge	8.896482	8.886395	0.777960	7.093235	0.010088
Linear Regression	8.896457	8.886428	0.777958	7.093134	0.010029
Lasso	8.898044	8.887389	0.777910	7.094680	0.010655

BEST MODEL: LightGBM

Valid RMSE: 8.7496

Valid R<sup>2</sup>: 0.7847

Valid MAE: 6.9725

# DEPLOYMENT



 **PredictScore.AI**  
Student Score Predictor

**AI-Powered Predictor**

## Predict Your Exam Score

Enter your study habits, attendance, and other factors to get an AI-powered prediction of your exam performance. Understand what impacts your grades most.

**Data-Driven**  
Trained on large student data

**ML Based**  
LightGBM Regressor

**Instant Results**  
Get predictions fast

**Your Predicted Exam Score**

**86.0**  
out of 100  
PREDICTED SCORE

Based on the factors you provided, this is your estimated performance. Excellent! Keep up the great work!

**Predict Score**

**Reset Form**

The interface is a web-based form for predicting exam scores. It features a header with the app name and a sub-header "AI-Powered Predictor". Below this is a main title "Predict Your Exam Score" with a descriptive subtitle about entering study habits and attendance. Three cards below the title provide details: "Data-Driven" (trained on large student data), "ML Based" (LightGBM Regressor), and "Instant Results" (get predictions fast). To the right is a detailed input form with sliders for Class Attendance (75%), Study Hours (per day) (8), Sleep Hours (per night) (6), and dropdown menus for Study Method (mixed), Course (b.tech), Facility Rating (medium), Sleep Quality (average), and Exam Difficulty (moderate). A prominent green circular progress bar on the right displays a predicted score of 86.0 out of 100, labeled as the "PREDICTED SCORE". Below the progress bar, a message states: "Based on the factors you provided, this is your estimated performance. Excellent! Keep up the great work!". At the bottom are two buttons: a blue "Predict Score" button and a grey "Reset Form" button.

# FURTHER IMPROVEMENT

## 01 Iterative EDA & Feature Engineering

- Revisit EDA after feature creation to validate new feature distributions

## 02 Model & Training Enhancements

- Try advanced models (e.g., Gradient Boosting, Ensemble methods)
- Perform systematic hyperparameter tuning (Grid / Random / Bayesian Search)

## 03 Deployment & Monitoring (Future Work)

- Monitor model performance and data drift in production
- Automate the pipeline using MLOps practices



A stylized illustration of a person with short brown hair, wearing a teal hoodie and blue pants, sitting cross-legged on the floor. They are holding a stack of yellow-lined paper in their hands. In front of them is an open teal laptop, and to the right are several stacked books. The background features a light gray surface with faint, concentric circular patterns.

# THANK YOU FOR YOUR ATTENTION!

AIO Conquer  
Warm Up Project Presentation

Presented by Team CONQ013