

SemiAuto Regression Report

Dataset: GYM

Generated on: 2025-05-16 12:02:31

Project Flow



Table of Contents

1. Data Ingestion
2. Data Preprocessing
3. Feature Engineering
4. Model Building
5. Model Evaluation
6. Model Optimization (if performed)
7. Final Evaluation Results

1. Data Ingestion

This step involves loading and analyzing the original dataset to understand its structure and characteristics.

Dataset Overview

Dataset: gym

Train samples: 778, **Test samples:** 195

Target column: BMI

Column Types

Original Columns:

Age, Gender, Weight (kg), Height (m), Max_BPM, Avg_BPM, Resting_BPM, Session_Duration (hours), Calories_Burned, Workout_Type, Fat_Percentage, Water_Intake (liters), Workout_Frequency (days/week), Experience_Level, BMI

Numerical Columns:

Age, Weight (kg), Height (m), Max_BPM, Avg_BPM, Resting_BPM, Session_Duration (hours), Calories_Burned, Fat_Percentage, Water_Intake (liters), Workout_Frequency (days/week), Experience_Level, BMI

Categorical Columns:

Gender, Workout_Type

Skewed Columns:

Weight (kg), Fat_Percentage, BMI

Normal Columns:

Age, Height (m), Max_BPM, Avg_BPM, Resting_BPM, Session_Duration (hours), Calories_Burned, Water_Intake (liters), Workout_Frequency (days/week), Experience_Level

Columns with Nulls:

None

Columns with Outliers:

Weight (kg), Calories_Burned, BMI

Highly Correlated Features

Weight (kg):

- BMI: 0.8532

Session_Duration (hours):

- Calories_Burned: 0.9081

- Experience_Level: 0.7648

Calories_Burned:

- Session_Duration (hours): 0.9081
- Experience_Level: 0.6941

Fat_Percentage:

- Experience_Level: -0.6544

Workout_Frequency (days/week):

- Experience_Level: 0.8371

Experience_Level:

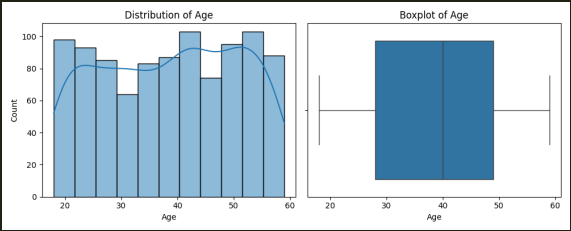
- Session_Duration (hours): 0.7648
- Calories_Burned: 0.6941
- Fat_Percentage: -0.6544
- Workout_Frequency (days/week): 0.8371

BMI:

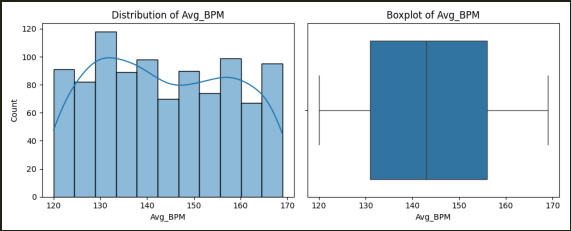
- Weight (kg): 0.8532

Feature Distributions

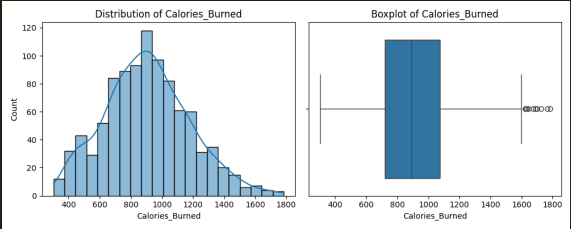
Age



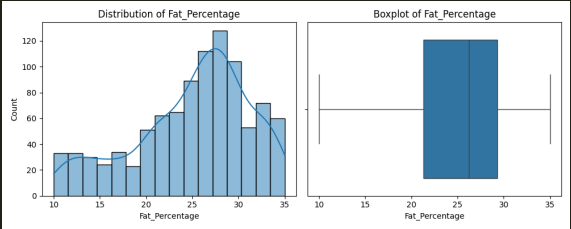
Avg_BPM



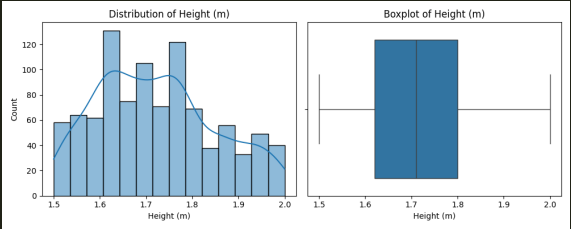
Calories_Burned



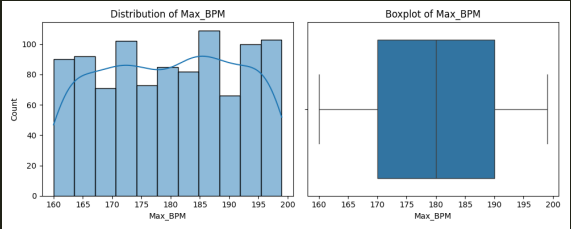
Fat_Percentage



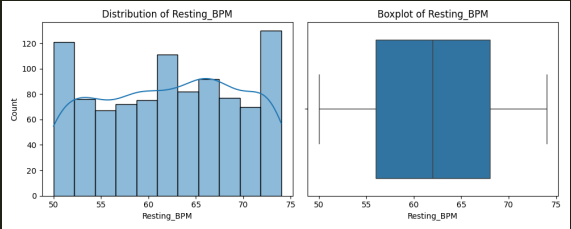
Height (m)



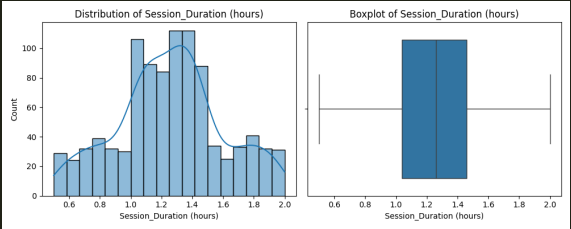
Max_BPM



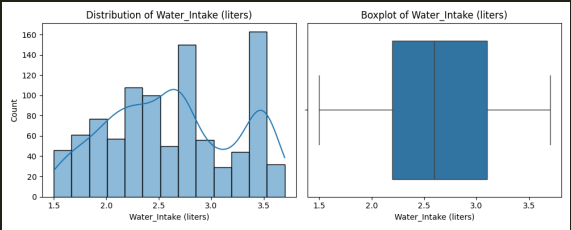
Resting_BPM



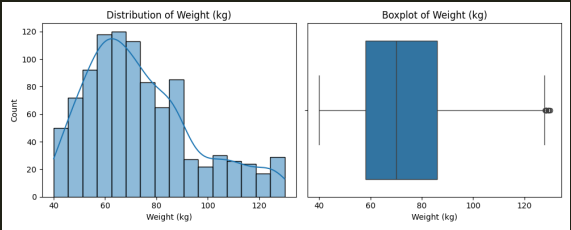
Session_Duration (hours)



Water_Intake (liters)



Weight (kg)



Correlation Heatmap

Age

Weight (kg)

Height (m)

Max_BPM

Avg_BPM

Resting_BPM

Session_Duration (hours)

Calories_Burned

Fat_Percentage

Water_Intake (liters)

Workout_Frequency (days/week)

Experience_Level

BMI

0.8

0.6

0.4

0.2

0.0

-0.2

-0.4

-0.6

2. Data Preprocessing

This step involves cleaning the dataset and preparing it for model training.

Preprocessing Configuration

Duplicate handling:

Remove duplicates: True

Outlier Treatment:

Method: IQR

Applied to columns:

- Weight (kg)
- Calories_Burned

Skewed Data Transformation:

Method: box-cox

Applied to columns:

- Weight (kg)
- Fat_Percentage

Numerical Scaling:

Method: standard

Applied to columns:

- Age
- Weight (kg)
- Height (m)
- Max_BPM
- Avg_BPM
- Resting_BPM
- Session_Duration (hours)
- Calories_Burned
- Fat_Percentage
- Water_Intake (liters)
- Workout_Frequency (days/week)
- Experience_Level

Categorical Encoding:

Method: onehot

Drop first: True

Applied to columns:

- Gender
- Workout_Type

Preprocessed Data Preview

Training Data Sample (First 5 rows):

Age	Weight (kg)	Height (m)	Max_BPM	Avg_BPM	Resting_BPM	Session_Du..	Calories_B..	Fat_Percen..	Water_Inta..
-1.5330451..	-3.5324010..	-0.9700883..	-0.6810859..	1.70030808..	0.64940678..	-0.4037490..	0.48828893..	-4.0568311..	0.93198659..
0.03035418..	-3.5144073..	-0.5824515..	-1.5471938..	1.63084945..	0.10774130..	-0.9977108..	-0.5071114..	-3.9546773..	-1.5637186..
1.26461680..	-3.4756539..	-0.0397601..	1.22435159..	0.72788720..	-0.2985078..	-0.8492203..	-0.8650082..	-3.9902220..	-1.2309579..
-1.2039084..	-3.5196223..	-1.6678344..	0.61807601..	0.93626310..	0.78482315..	-0.6413337..	-0.3281630..	-3.8220239..	-0.2326758..
-1.5330451..	-3.5006984..	-1.5903071..	0.44485442..	-1.1474959..	-1.6526715..	-0.6710318..	-0.9321139..	-3.9424981..	-0.7318169..

Test Data Sample (First 5 rows):

Age	Weight (kg)	Height (m)	Max_BPM	Avg_BPM	Resting_BPM	Session_Du..	Calories_B..	Fat_Percen..	Water_Inta..
-0.3810666..	-3.4841766..	-0.6599789..	-0.8543075..	-0.6612854..	0.64940678..	0.01202429..	0.09311124..	-4.1245912..	-0.8981972..
0.77091175..	-3.4584482..	0.73551330..	-0.3346427..	0.38059403..	1.05565589..	1.91270222..	1.85277042..	-4.2832452..	1.43112764..
0.85319592..	-3.4678900..	-0.0397601..	1.22435159..	-1.0780372..	0.92023952..	0.60598614..	0.07819887..	-3.8823926..	1.09836694..
-0.7102033..	-3.4940894..	0.19282187..	0.01180044..	1.70030808..	0.64940678..	-0.0176738..	0.93565991..	-3.9081537..	0.76560624..
-0.5456350..	-3.4546463..	-0.9700883..	-0.3346427..	0.10275949..	-0.2985078..	2.17998505..	2.54619544..	-4.2612287..	1.43112764..

3. Feature Engineering

This step involves creating new features or selecting the most important ones.

Feature Engineering Configuration

Applied Techniques:

Automated Feature Engineering: No
SHAP-based Feature Selection: No

Transformed Data Preview

Transformed Training Data Sample (First 5 rows):

Age	Weight (kg)	Height (m)	Max_BPM	Avg_BPM	Resting_BPM	Session_Du..	Calories_B..	Fat_Percen..	Water_Inta..
-1.5330451..	-3.5324010..	-0.9700883..	-0.6810859..	1.70030808..	0.64940678..	-0.4037490..	0.48828893..	-4.0568311..	0.93198659..
0.03035418..	-3.5144073..	-0.5824515..	-1.5471938..	1.63084945..	0.10774130..	-0.9977108..	-0.5071114..	-3.9546773..	-1.5637186..
1.26461680..	-3.4756539..	-0.0397601..	1.22435159..	0.72788720..	-0.2985078..	-0.8492203..	-0.8650082..	-3.9902220..	-1.2309579..
-1.2039084..	-3.5196223..	-1.6678344..	0.61807601..	0.93626310..	0.78482315..	-0.6413337..	-0.3281630..	-3.8220239..	-0.2326758..
-1.5330451..	-3.5006984..	-1.5903071..	0.44485442..	-1.1474959..	-1.6526715..	-0.6710318..	-0.9321139..	-3.9424981..	-0.7318169..

Transformed Test Data Sample (First 5 rows):

Age	Weight (kg)	Height (m)	Max_BPM	Avg_BPM	Resting_BPM	Session_Du..	Calories_B..	Fat_Percen..	Water_Inta..
-0.3810666..	-3.4841766..	-0.6599789..	-0.8543075..	-0.6612854..	0.64940678..	0.01202429..	0.09311124..	-4.1245912..	-0.8981972..
0.77091175..	-3.4584482..	0.73551330..	-0.3346427..	0.38059403..	1.05565589..	1.91270222..	1.85277042..	-4.2832452..	1.43112764..
0.85319592..	-3.4678900..	-0.0397601..	1.22435159..	-1.0780372..	0.92023952..	0.60598614..	0.07819887..	-3.8823926..	1.09836694..
-0.7102033..	-3.4940894..	0.19282187..	0.01180044..	1.70030808..	0.64940678..	-0.0176738..	0.93565991..	-3.9081537..	0.76560624..
-0.5456350..	-3.4546463..	-0.9700883..	-0.3346427..	0.10275949..	-0.2985078..	2.17998505..	2.54619544..	-4.2612287..	1.43112764..

4. Model Building

This step involves training the regression model on the transformed data.

Model Selection

Selected Model:

CatBoost

Training timestamp: 2025-05-16 11:48:56

5. Model Evaluation

This step involves evaluating the performance of the trained model.

Performance Metrics

Original Model Performance:

Evaluation timestamp: 2025-05-16 11:59:07

Metric	Value
R ² Score	0.99258
Explained Variance Score	0.99258
Mean Squared Error	0.36174
Root Mean Squared Error	0.60145
Mean Absolute Error	0.37558
Mean Absolute Percentage Error	0.01376
Max Error	3.84596

6. Model Optimization

This step involves tuning the hyperparameters of the model to improve performance.

Optimized Hyperparameters

Parameter	Value
border_count	104
depth	4
iterations	193
l2_leaf_reg	1
learning_rate	0.17459696479273032
verbose	False

Optimization timestamp: 2025-05-16 11:57:50

7. Final Evaluation Results

This section presents the final performance of the optimized model.

Optimized Model Performance

Metric	Value
R ² Score	0.99839
Explained Variance Score	0.99840
Mean Squared Error	0.07824
Root Mean Squared Error	0.27971
Mean Absolute Error	0.21122
Mean Absolute Percentage Error	0.00844
Max Error	1.36023

Evaluation timestamp: 2025-05-16 11:59:07

Performance Comparison

Metric	Original Model	Optimized Model	Improvement
R ² Score	0.99258	0.99839	+0.59%
RMSE	0.60145	0.27971	+53.49%
MAE	0.37558	0.21122	+43.76%

Conclusion

Summary of the regression model development and performance.

This report summarizes the development of a regression model to predict BMI using the gym dataset. A CatBoost regression model was trained and optimized using hyperparameter tuning. The optimization process improved the model's R^2 score from 0.99258 to 0.99839, representing a 0.59% improvement.

This automatic report was generated to provide insights into the model development process and performance metrics. It includes details about data preprocessing, feature engineering, model selection, and evaluation results.