

Data Collection and Preprocessing Phase

Date	15 June 2024
Team ID	739820
Project Title	Predicting the Unpredictable: A Look into the World of Powerlifting
Maximum Marks	6 Marks

Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

Section	Description
Data Overview	<u>Dimension:</u> 614 rows × 13 columns
	<u>Descriptive statistics:</u>
Univariate Analysis	-

Outliers and Anomalies	-																																																						
Data Preprocessing Code Screenshots																																																							
Loading Data	<div>#importing the data which is in csv file</div> <div>data1=pd.read_csv("/content/X_train.csv",header='infer')</div> <div>data2=pd.read_csv("/content/y_train.csv",header='infer')</div> <div>data1.head()</div> <table><thead><tr><th></th><th>playerId</th><th>Name</th><th>Sex</th><th>Equipment</th><th>Age</th><th>BodyweightKg</th><th>BestSquatKg</th><th>BestDeadliftKg</th></tr></thead><tbody><tr><td>0</td><td>19391.0</td><td>Carlos Ceron</td><td>M</td><td>Raw</td><td>23.0</td><td>87.30</td><td>205.0</td><td>235.0</td></tr><tr><td>1</td><td>15978.0</td><td>Tito Herrera</td><td>M</td><td>Wraps</td><td>23.0</td><td>73.48</td><td>220.0</td><td>260.0</td></tr><tr><td>2</td><td>27209.0</td><td>Levi Lehman</td><td>M</td><td>Raw</td><td>26.0</td><td>112.40</td><td>142.5</td><td>220.0</td></tr><tr><td>3</td><td>27496.0</td><td>Stacy Hayford</td><td>F</td><td>Raw</td><td>35.0</td><td>59.42</td><td>95.0</td><td>102.5</td></tr><tr><td>4</td><td>20293.0</td><td>Brittany Hirt</td><td>F</td><td>Raw</td><td>26.5</td><td>61.40</td><td>105.0</td><td>127.5</td></tr></tbody></table>		playerId	Name	Sex	Equipment	Age	BodyweightKg	BestSquatKg	BestDeadliftKg	0	19391.0	Carlos Ceron	M	Raw	23.0	87.30	205.0	235.0	1	15978.0	Tito Herrera	M	Wraps	23.0	73.48	220.0	260.0	2	27209.0	Levi Lehman	M	Raw	26.0	112.40	142.5	220.0	3	27496.0	Stacy Hayford	F	Raw	35.0	59.42	95.0	102.5	4	20293.0	Brittany Hirt	F	Raw	26.5	61.40	105.0	127.5
	playerId	Name	Sex	Equipment	Age	BodyweightKg	BestSquatKg	BestDeadliftKg																																															
0	19391.0	Carlos Ceron	M	Raw	23.0	87.30	205.0	235.0																																															
1	15978.0	Tito Herrera	M	Wraps	23.0	73.48	220.0	260.0																																															
2	27209.0	Levi Lehman	M	Raw	26.0	112.40	142.5	220.0																																															
3	27496.0	Stacy Hayford	F	Raw	35.0	59.42	95.0	102.5																																															
4	20293.0	Brittany Hirt	F	Raw	26.5	61.40	105.0	127.5																																															
Handling Missing Data	<div>[] data['Age'].fillna(data['Age'].mean(),inplace=True)</div>																																																						

Data Transformation	<pre>data1['Name'] = LabelEncoder().fit_transform(data1['Name'])</pre> <pre>data['Sex'] = data['Sex'].map({'M':1, 'F':0}) from sklearn.preprocessing import LabelEncoder data['Equipment'] = LabelEncoder().fit_transform(data['Equipment'])</pre> <pre>data['BestSquatKg'] = LabelEncoder().fit_transform(data['BestSquatKg'])</pre>
Feature Engineering	Attached the codes in final submission.
Save Processed Data	-