

```
In [3]: import os
os.getcwd()
```

```
Out[3]: 'c:\\Users\\Prachi\\Documents\\VS Code Files\\ML CAPSTONE PROJECT\\Height Weight'
```

```
In [4]: os.chdir(r"C:\Users\Prachi\Documents\VS Code Files\ML CAPSTONE PROJECT\Height Weight")
```

```
In [25]: # importing eda Libraries
import numpy as np
import pandas as pd

# Visualization
import matplotlib.pyplot as plt
import seaborn as sns

# preprocessing
from sklearn.preprocessing import StandardScaler

# splitting the data
from sklearn.model_selection import train_test_split

# Splitting the data
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor

# evaluation metrics
from sklearn.metrics import mean_squared_error
```

```
In [26]: df = pd.read_csv('SOCR-HeightWeight.csv')
```

```
In [6]: df.head() # 1 pound = 453 grams
```

```
Out[6]:
```

	Index	Height(Inches)	Weight(Pounds)
0	1	65.78331	112.9925
1	2	71.51521	136.4873
2	3	69.39874	153.0269
3	4	68.21660	142.3354
4	5	67.78781	144.2971

```
In [7]: #converting weight pounds to kg
df['Weight_kg']=df['Weight(Pounds)']*0.453592

# Convert inches to the desired format (feet.inches)
df['Height(Feet.Inches)'] = df['Height(Inches)'] // 12 + (df['Height(Inches)'] % 12
```

```
In [8]: df.describe()
```

```
Out[8]:
```

	Index	Height(Inches)	Weight(Pounds)	Weight_kg	Height(Feet.Inches)
count	25000.000000	25000.000000	25000.000000	25000.000000	25000.000000
mean	12500.500000	67.993114	127.079421	57.642209	5.795967
std	7217.022701	1.901679	11.660898	5.289290	0.183513
min	1.000000	60.278360	78.014760	35.386871	5.027836
25%	6250.750000	66.704397	119.308675	54.117461	5.670440
50%	12500.500000	67.995700	127.157750	57.677738	5.799570
75%	18750.250000	69.272958	134.892850	61.186318	5.927296
max	25000.000000	75.152800	170.924000	77.529759	6.315280

```
In [9]: drop_col=['Index','Height(Inches)','Weight(Pounds)'] # selecting columns to del it
#dropping columns
df=df.drop(columns=drop_col,axis=1)
```

```
In [10]: df.sample(3) #it will give random row information
```

```
Out[10]:
```

	Weight_kg	Height(Feet.Inches)
312	52.304056	5.821327
8164	67.614555	5.803615
24822	60.672148	5.836479

```
In [11]: df.shape #checking shape of the data
```

```
Out[11]: (25000, 2)
```

```
In [12]: df.isna().any() #checking null values
```

```
Out[12]: Weight_kg      False
Height(Feet.Inches)  False
dtype: bool
```

```
In [13]: df.dtypes #checking dtypes for our dataframe
```

```
Out[13]: Weight_kg      float64
Height(Feet.Inches)  float64
dtype: object
```

```
In [14]: df.corr() #correlation
```

```
Out[14]:
```

	Weight_kg	Height(Feet.Inches)
Weight_kg	1.000000	0.499192
Height(Feet.Inches)	0.499192	1.000000

```
In [15]: df.describe()
```

Out[15]:

	Weight_kg	Height(Feet.Inches)
count	25000.000000	25000.000000
mean	57.642209	5.795967
std	5.289290	0.183513
min	35.386871	5.027836
25%	54.117461	5.670440
50%	57.677738	5.799570
75%	61.186318	5.927296
max	77.529759	6.315280

Mean:

The mean height is approximately 67.99 inches. The mean weight is approximately 127.08 pounds. Standard Deviation (Std):

The standard deviation for height is approximately 1.90 inches, indicating the spread or dispersion of heights around the mean. The standard deviation for weight is approximately 11.66 pounds, indicating the spread or dispersion of weights around the mean. Minimum and Maximum Values:

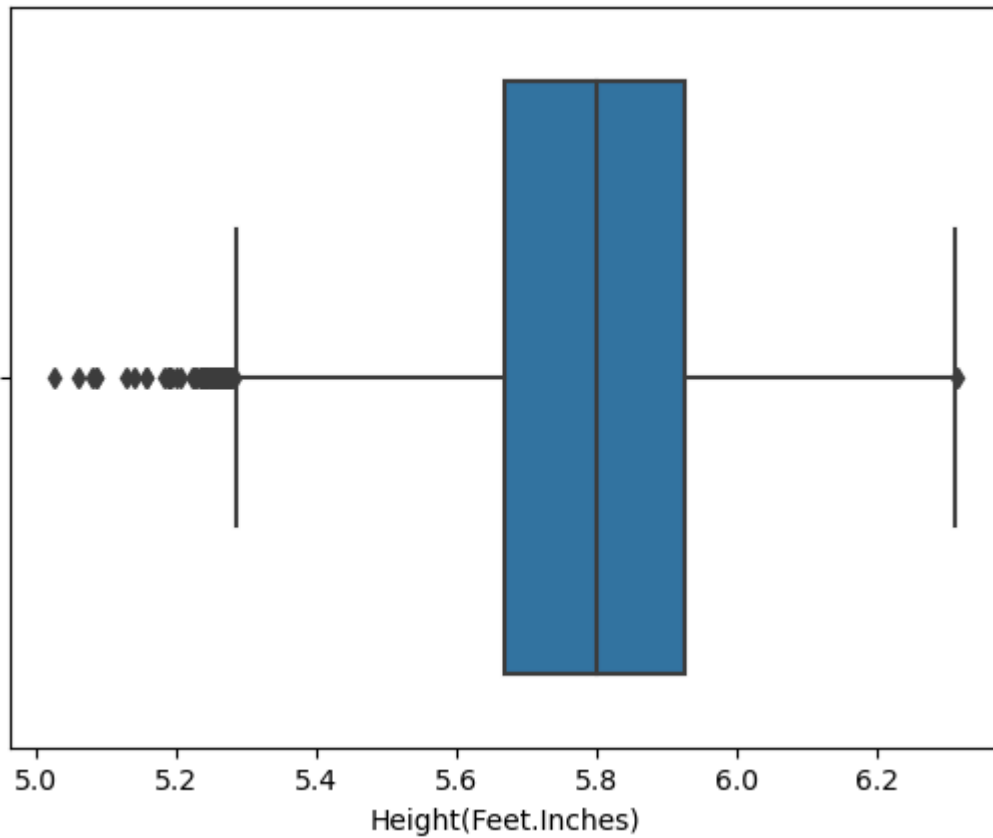
The minimum height recorded is approximately 60.28 inches, and the maximum height is approximately 75.15 inches. The minimum weight recorded is approximately 78.01 pounds, and the maximum weight is approximately 170.92 pounds. Percentiles (25th, 50th, and 75th):

The 25th percentile (Q1) indicates that 25% of the data falls below a height of approximately 66.70 inches and a weight of approximately 119.31 pounds. The 50th percentile (median) indicates that 50% of the data falls below a height of approximately 67.99 inches and a weight of approximately 127.16 pounds. The 75th percentile (Q3) indicates that 75% of the data falls below a height of approximately 69.27 inches and a weight of approximately 134.89 pounds.

Checking outliers using boxplot

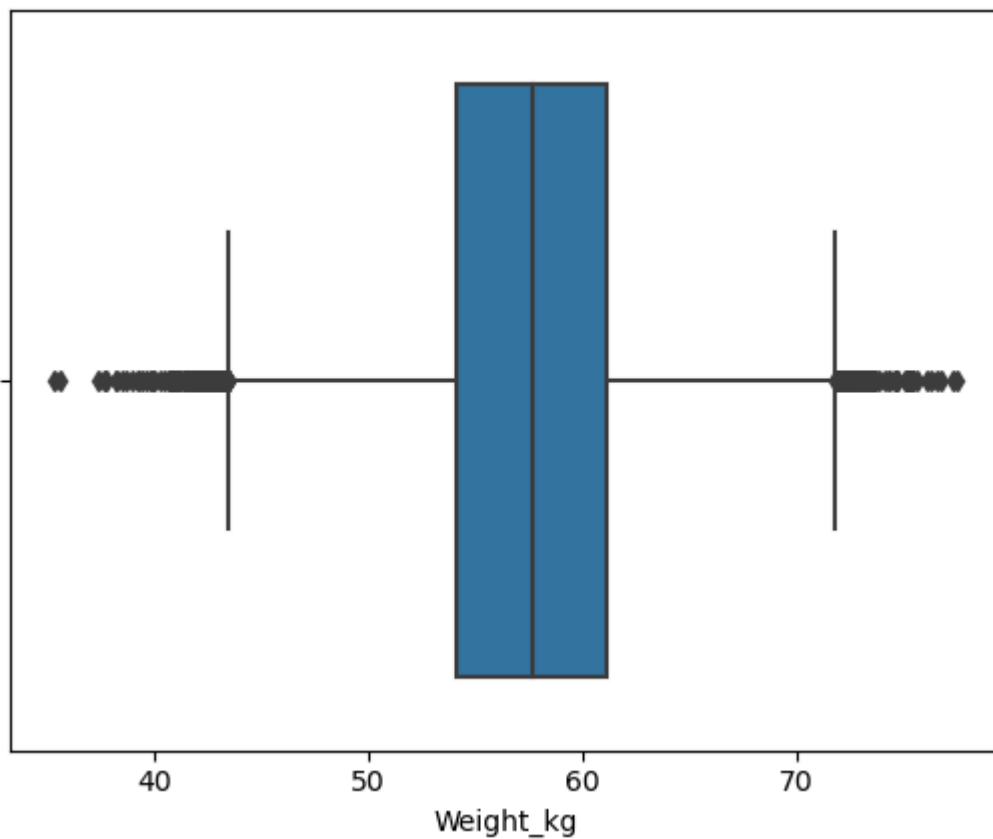
```
In [16]: sns.boxplot(x=df['Height(Feet.Inches)'])
```

```
Out[16]: <Axes: xlabel='Height(Feet.Inches)'>
```



```
In [17]: sns.boxplot(x=df['Weight_kg']) #checking outliers for
```

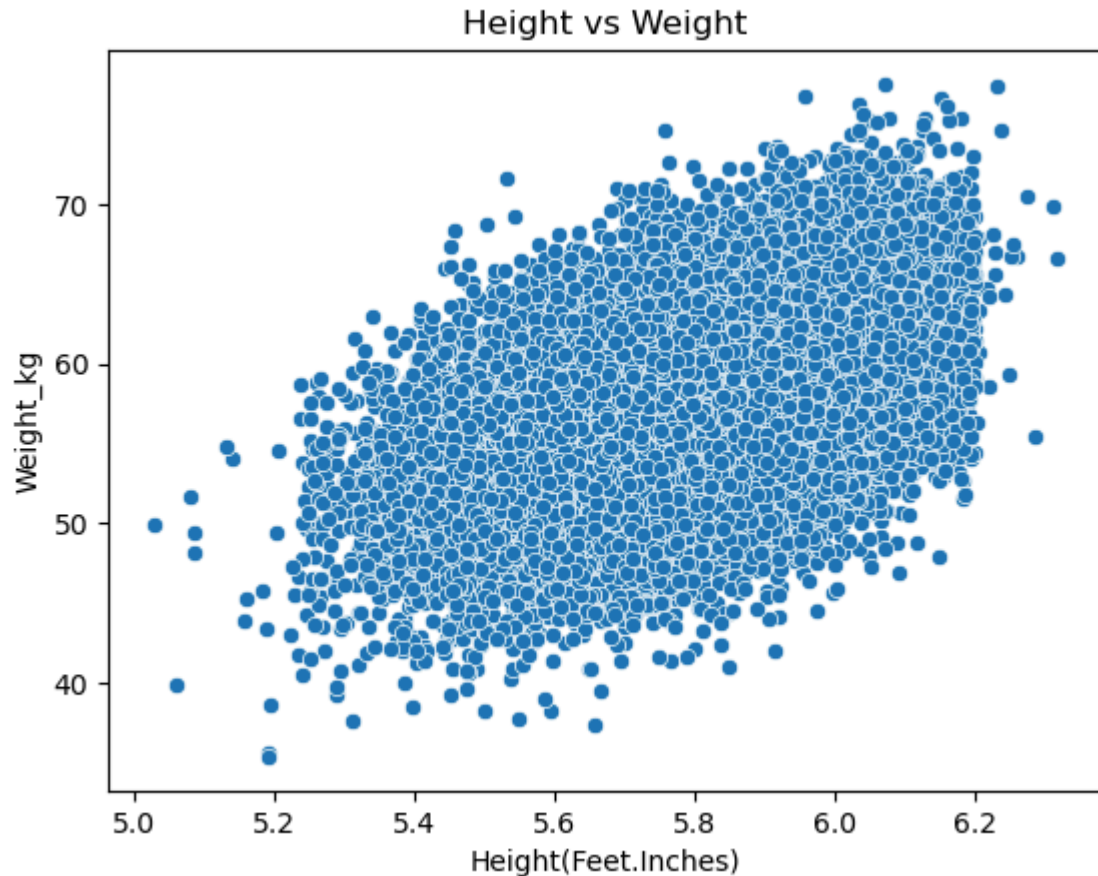
```
Out[17]: <Axes: xlabel='Weight_kg'>
```



```
In [18]: x=df['Height(Feet.Inches)']
          y=df['Weight_kg']

          sns.scatterplot(x=x,y=y)
          plt.title('Height vs Weight')
```

```
plt.xlabel('Height(Feet.Inches)')
plt.ylabel('Weight_kg')
plt.show()
```



```
In [19]: df.sample(3)
```

```
Out[19]:
```

	Weight_kg	Height(Feet.Inches)
18013	57.265854	5.815961
19622	54.457893	5.661994
6136	59.588562	5.907214

```
In [20]: # split the data into dependent & independent variable
X=df.iloc[:,1]
y=df.iloc[:,0]
```

```
In [21]: X
```

```
Out[21]:
```

0	5.578331
1	6.151521
2	5.939874
3	5.821660
4	5.778781
	...
24995	5.950215
24996	5.454826
24997	5.469855
24998	5.752918
24999	5.887761

Name: Height(Feet.Inches), Length: 25000, dtype: float64

```
In [22]: df.columns[1] #X variable column name
```

```
Out[22]: 'Height(Feet.Inches)'
```

```
In [23]: df.columns[0]    # y variable
```

```
Out[23]: 'Weight_kg'
```

Data Scaling (Preprocessing Data)

```
In [ ]: scaler_X = StandardScaler()
X_scaled = scaler_X.fit_transform(X.values.reshape(-1,1))

scaler_y = StandardScaler()
y_scaled = scaler_y.fit_transform(y.values.reshape(-1, 1))
```

Splitting data into 80% 20% ratio

```
In [27]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
In [28]: print('Shape of training data')
print(X_train.shape)
print(y_train.shape)

print('Shape of testing data')
print(X_test.shape)
print(y_test.shape)
```

```
Shape of training data
(20000,)
(20000,)
Shape of testing data
(5000,)
(5000,)
```

```
In [29]: #Linear regression model X should be 2d array so we are reshaping it to 2d array

# Reshape training data
X_train_2d = X_train.values.reshape(-1, 1)
y_train_2d = y_train.values.reshape(-1, 1)

# Reshape testing data
X_test_2d = X_test.values.reshape(-1, 1)
y_test_2d = y_test.values.reshape(-1, 1)

print("Shape of training data (X):", X_train_2d.shape)
print("Shape of training data (y):", y_train_2d.shape)
print("Shape of testing data (X):", X_test_2d.shape)
print("Shape of testing data (y):", y_test_2d.shape)
```

```
Shape of training data (X): (20000, 1)
Shape of training data (y): (20000, 1)
Shape of testing data (X): (5000, 1)
Shape of testing data (y): (5000, 1)
```

```
In [30]: lr=LinearRegression() #Linear Regression
lr
```

Out[30]:

▼ LinearRegression ⓘ ?
► Parameters

In [31]: lr.fit(X_train_2d,y_train_2d)

Out[31]:

▼ LinearRegression ⓘ ?
► Parameters

In [32]: y_pred=lr.predict(X_test_2d)
y_pred[:10]

Out[32]: array([[55.94425481],
[60.91226889],
[56.56867714],
[56.42643564],
[51.52547113],
[52.93798976],
[60.30463034],
[60.27256006],
[62.74472434],
[63.0616341]])

In [33]: y_test_2d[:10]

Out[33]: array([[60.87349789],
[64.25661383],
[50.63170805],
[53.62895327],
[46.5397639],
[48.20970821],
[55.81821505],
[55.03481631],
[76.60307055],
[55.98708736]])

In [34]: mean_squared_error(y_pred,y_test_2d)

Out[34]: 21.69730652290755

In [35]: model_dtr=DecisionTreeRegressor()
model_dtr

Out[35]:

▼ DecisionTreeRegressor ⓘ ?
► Parameters

In [36]: model_dtr.fit(X_train_2d, y_train_2d)

Out[36]:

▼ DecisionTreeRegressor ⓘ ?
► Parameters

In [37]: y_pred_dtr=model_dtr.predict(X_test_2d)
y_pred_dtr[:5]

```
Out[37]: array([63.37542065, 56.46639802, 56.7162365 , 64.71347169, 57.84078178])
```

```
In [39]: print(y_test_2d[:5])
```

```
[[60.87349789]
 [64.25661383]
 [50.63170805]
 [53.62895327]
 [46.5397639 ]]
```

```
In [40]: mean_squared_error(y_pred_dtr, y_test_2d)
```

```
Out[40]: 41.50751860513505
```

RandomForestRegressor

```
In [41]: model_rfr=RandomForestRegressor()
model_rfr.fit(X_train_2d, y_train_2d)
```

c:\Users\Prachi\anaconda3\Lib\site-packages\sklearn\base.py:1365: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
return fit_method(estimator, *args, **kwargs)

```
Out[41]: ▼ RandomForestRegressor ⓘ ?
```

► Parameters

```
In [42]: y_pred_rfr=(X_test_2d)
y_pred_rfr[:10]
```

```
Out[42]: array([[5.675233],
 [6.023626],
 [5.719022],
 [5.709047],
 [5.365356],
 [5.464412],
 [5.981014],
 [5.978765],
 [6.152131],
 [6.174355]])
```

```
In [43]: mean_squared_error(y_pred_rfr, y_test_2d)
```

```
Out[43]: 2708.30376658499
```

Hyperparameter Tuning

```
In [44]: from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LinearRegression
```

```
# Define hyperparameters to tune
param_grid = {
    'fit_intercept': [True, False],
    'copy_X': [True, False]
}
```



```

# Create a Linear Regression model
model_lr = LinearRegression()

# Initialize GridSearchCV
grid_search = GridSearchCV(model_lr, param_grid, cv=5, scoring='neg_mean_squared_er

# Fit the model
grid_search.fit(X_train_2d, y_train_2d)

# Print the best parameters and best MSE score
print("Best Parameters:", grid_search.best_params_)
print("Best Negative MSE Score:", grid_search.best_score_)

```

Best Parameters: {'copy_X': True, 'fit_intercept': True}
Best Negative MSE Score: -20.836260216566203

Final Model

```

In [45]: from sklearn.linear_model import LinearRegression

# Initialize the Linear Regression model with the best parameters
final_model = LinearRegression(fit_intercept=False, copy_X=True)

# Fit the model to the entire training data
final_model.fit(X_train_2d, y_train_2d)

# Now you can use final_model to make predictions on new data

```

Out[45]:

▼ LinearRegression ⓘ ?

► Parameters

Converting to pickle file

```

In [46]: import pickle

# Define the filename for the pickle file
filename = 'final_model.pkl'

# Save the final_model to a pickle file
with open(filename, 'wb') as file:
    pickle.dump(final_model, file)

```

```

In [47]: os.path.abspath('final_model.pkl')

```

Out[47]: 'C:\\Users\\Prachi\\Documents\\VS Code Files\\ML CAPSTONE PROJECT\\Height Weight\\final_model.pkl'

```

In [48]: import pickle
import numpy as np

# Load the saved model from the file
filename = 'final_model.pkl'
with open(filename, 'rb') as file:
    loaded_model = pickle.load(file)

# Input height for prediction
height_input = 6.0

# Reshape the input height to match the shape expected by the model (2D array)

```

```
height_input_2d = np.array(height_input).reshape(1, -1)

# Use the loaded model to make predictions
predicted_weight = loaded_model.predict(height_input_2d)

# Print the predicted weight
print("Predicted weight:", predicted_weight[0, 0])
```

Predicted weight: 59.7202378537033