

```
In [1]: 1 import os
        2 os.getcwd()#my current path
```

```
Out[1]: 'C:\\Users\\sumit\\Data Science\\Live Project\\Projects end to end\\ML Height Prediction'
```

```
In [2]: 1 #os.chdir('D:\\heightweight') #changing path
```

```
In [3]: 1 #importing eda libraries
        2 import numpy as np #math
        3 import pandas as pd #excellent for data manipulation
        4
        5 #visualization
        6 import matplotlib.pyplot as plt
        7 import seaborn as sns
        8
        9 #preprocessing
       10 from sklearn.preprocessing import StandardScaler
       11
       12 #splitting the data
       13 from sklearn.model_selection import train_test_split
       14
       15 # importing Algorithms
       16 from sklearn.linear_model import LinearRegression
       17 from sklearn.tree import DecisionTreeRegressor
       18 from sklearn.ensemble import RandomForestRegressor
       19
       20 #evaluation metrics
       21 from sklearn.metrics import mean_squared_error
       22
       23
```

```
In [4]: 1 df=pd.read_csv('SOCR-HeightWeight.csv') #reading csv file
```

```
In [5]: 1 df.head() #1 pound=453grams
```

```
Out[5]:
```

	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 [6]: 1 #converting weight pounds to kg
        2 df['Weight_kg']=df['Weight(Pounds)']*0.453592
        3
        4 # Convert inches to the desired format (feet.inches)
        5 df['Height(Feet.Inches)'] = df['Height(Inches)'] // 12 + (df['Height(Inches)'] % 12) / 12
```

In [7]: 1 df.describe()

Out[7]:

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

In [8]: 1 drop\_col=['Index', 'Height(Inches)', 'Weight(Pounds)'] # selecting column  
2  
3 #dropping columns  
4 df=df.drop(columns=drop\_col,axis=1)

In [9]: 1 df.sample(3)

Out[9]:

	Weight_kg	Height(Feet.Inches)
<b>15349</b>	59.688217	5.781337
<b>486</b>	57.838967	5.926560
<b>493</b>	57.390773	5.602496

In [10]: 1 df.shape #checking shape of the data

Out[10]: (25000, 2)

In [11]: 1 df.isna().any() #checking null values

Out[11]: Weight\_kg False  
Height(Feet.Inches) False  
dtype: bool

In [12]: 1 df.dtypes #checking dtypes for our dataframe  
2

Out[12]: Weight\_kg float64  
Height(Feet.Inches) float64  
dtype: object

In [13]: 1 df.corr() #correlation

Out[13]:

	Weight_kg	Height(Feet.Inches)
<b>Weight_kg</b>	1.000000	0.499192
<b>Height(Feet.Inches)</b>	0.499192	1.000000

```
In [14]: 1 df.describe()
```

```
Out[14]:
```

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

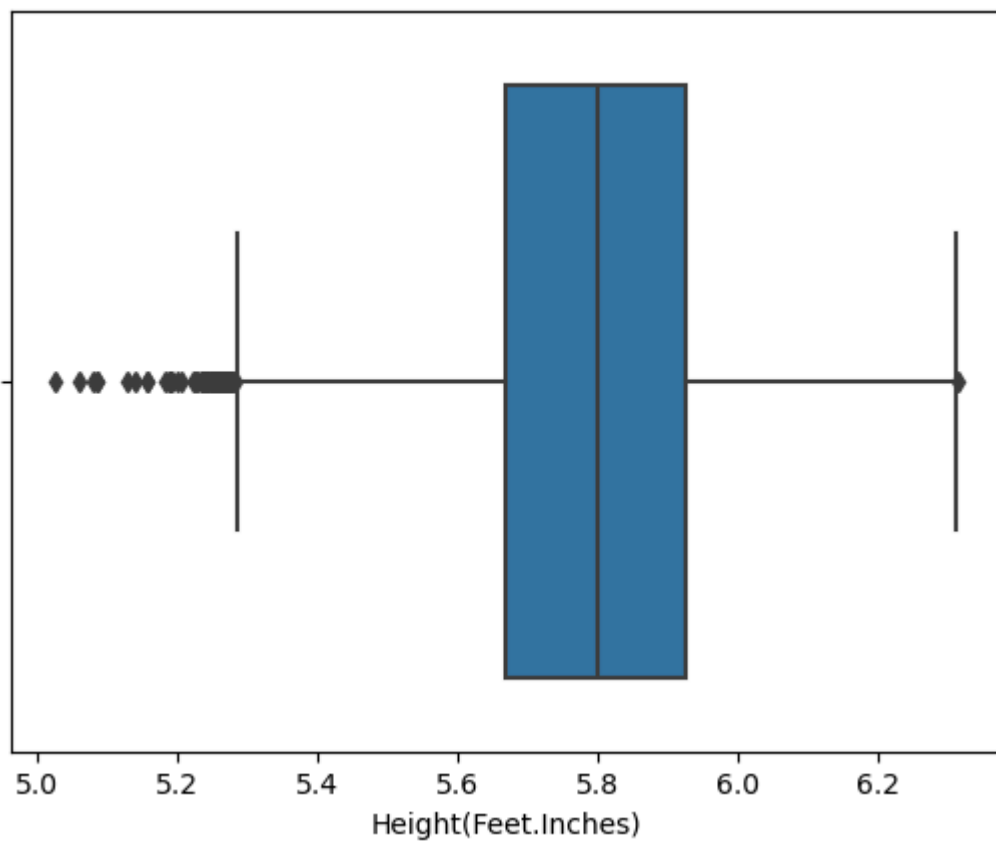
The minimum height recorded is approximately 60.28 inches, and the maximum height is approximately 75.15 inches. The minimum weight recorded is approximately 78.01 pounds, and the maximum weight is approximately 170.92 pounds. Percentiles (25th, 50th, and 75th):

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

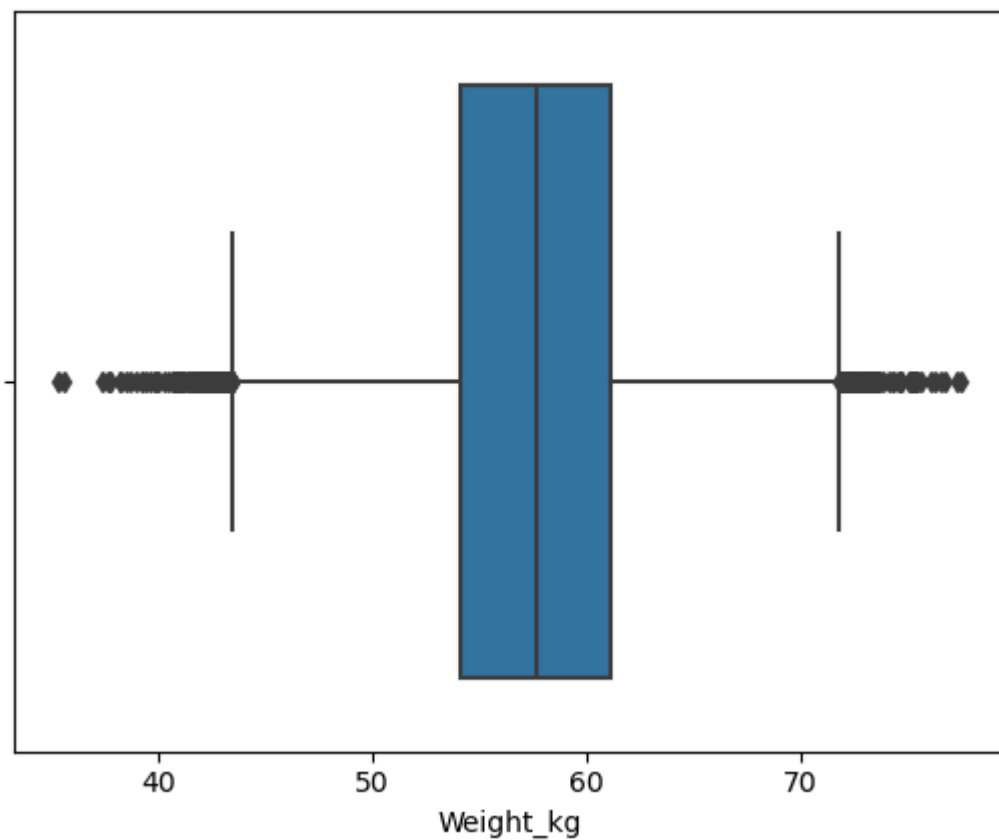
```
In [15]: 1 sns.boxplot(x=df['Height(Feet.Inches)'])
```

```
Out[15]: <Axes: xlabel='Height(Feet.Inches)'>
```



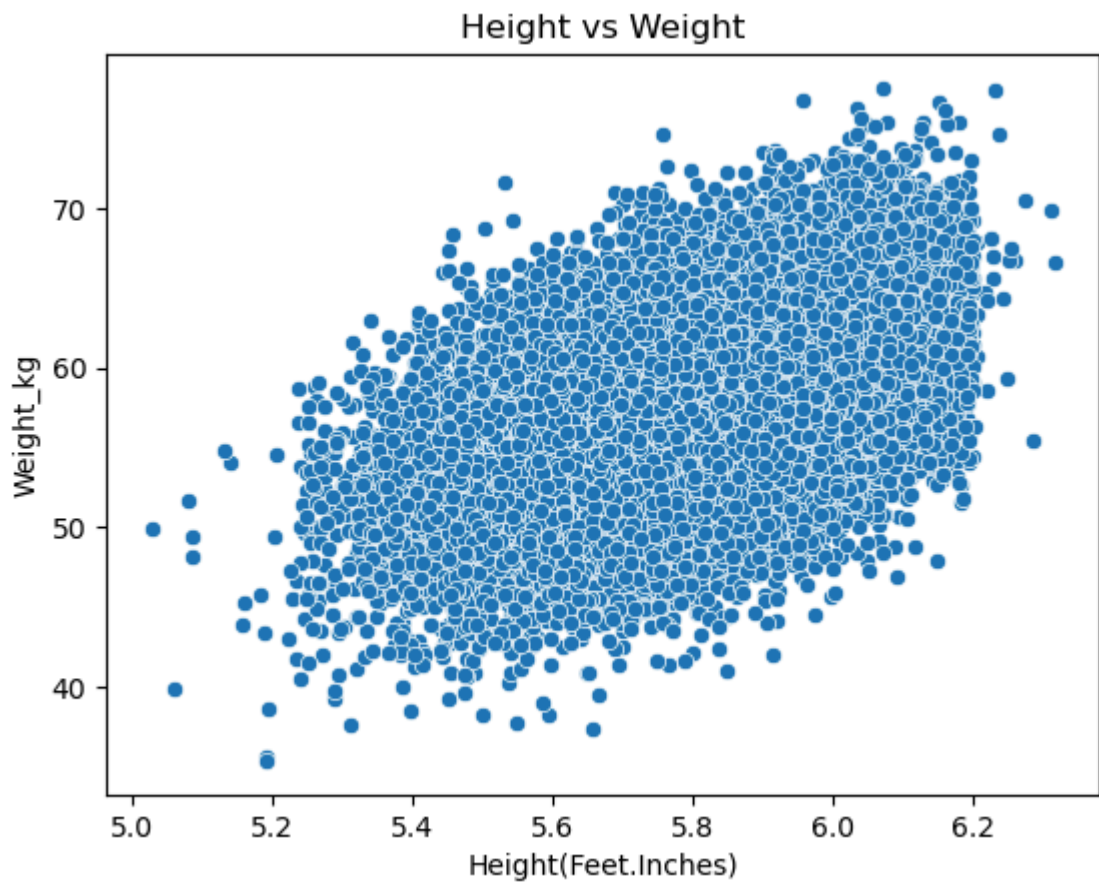
```
In [16]: 1 sns.boxplot(x=df['Weight_kg']) #checking outliers for
```

```
Out[16]: <Axes: xlabel='Weight_kg'>
```



A correlation coefficient of 0.502859 suggests a moderate positive correlation between height and weight

```
In [17]: 1 x=df['Height(Feet.Inches)']
2 y=df['Weight_kg']
3
4 sns.scatterplot(x=x,y=y)
5 plt.title('Height vs Weight')
6 plt.xlabel('Height(Feet.Inches)')
7 plt.ylabel('Weight_kg')
8 plt.show()
```



```
In [18]: 1 df.sample(3)
```

Out[18]:

	Weight_kg	Height(Feet.Inches)
7148	50.461067	5.765882
12746	62.735448	6.000383
10106	55.045566	5.750667

```
In [19]: 1 X=df.iloc[:,1]
2 y=df.iloc[:,0]
```

In [20]:

1 X

Out[20]:

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 [21]:

1 df.columns[1] *#X variable column name*

Out[21]: 'Height(Feet.Inches)'

In [22]:

1 df.columns[0] *# y variable*

Out[22]: 'Weight\_kg'

In [23]:

1 *#Data scaling(preprocessing data)*

In [24]:

```
1 scaler_X = StandardScaler()
2 X_scaled = scaler_X.fit_transform(X.values.reshape(-1,1))
3
4
5 scaler_y = StandardScaler()
6 y_scaled = scaler_y.fit_transform(y.values.reshape(-1, 1))
```

### ***splitting data into 70% 30% radio***

In [25]:

1 X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2)

In [26]:

```
1 print('Shape of trining data')
2 print(X_train.shape)
3 print(y_train.shape)
4
5 print('Shpae of testing data')
6 print(X_test.shape)
7 print(y_test.shape)
```

Shape of trining data

(20000,)

(20000,)

Shpae of testing data

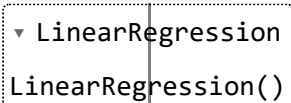
(5000,)

(5000,)

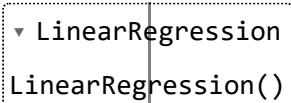
```
In [27]: 1 #Linear regression model X should be 2d array so we are reshaping it to 2d array
2
3 # Reshape training data
4 X_train_2d = X_train.values.reshape(-1, 1)
5 y_train_2d = y_train.values.reshape(-1, 1)
6
7 # Reshape testing data
8 X_test_2d = X_test.values.reshape(-1, 1)
9 y_test_2d = y_test.values.reshape(-1, 1)
10
11 print("Shape of training data (X):", X_train_2d.shape)
12 print("Shape of training data (y):", y_train_2d.shape)
13 print("Shape of testing data (X):", X_test_2d.shape)
14 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 [28]: 1 lr=LinearRegression() #Linear Regression
2 lr
```

Out[28]: A dashed box containing a dropdown arrow pointing to 'LinearRegression' and the text 'LinearRegression()' below it.

```
In [29]: 1 lr.fit(X_train_2d,y_train_2d)
```

Out[29]: A dashed box containing a dropdown arrow pointing to 'LinearRegression' and the text 'LinearRegression()' below it.

```
In [30]: 1 y_pred=lr.predict(X_test_2d)
2 y_pred[:10]
```

Out[30]: array([[58.30995889],  
[57.55379225],  
[57.21368631],  
[59.15305375],  
[59.14955828],  
[57.21726773],  
[54.13829253],  
[57.67365521],  
[60.77133975],  
[58.75644733]])



```
In [31]: 1 y_test_2d[:10]
```

```
Out[31]: array([[60.91082852],
                [50.78715403],
                [61.68266066],
                [57.32504768],
                [45.44130015],
                [48.93509253],
                [56.34655902],
                [60.25330155],
                [59.08262596],
                [63.5877017 ]])
```

```
In [32]: 1 mean_squared_error(y_pred,y_test_2d)
```

```
Out[32]: 21.232939817362524
```

```
In [33]: 1 model_dtr=DecisionTreeRegressor()
2 model_dtr
```

```
Out[33]: ▼ DecisionTreeRegressor
DecisionTreeRegressor()
```

```
In [34]: 1 model_dtr.fit(X_train_2d, y_train_2d)
```

```
Out[34]: ▼ DecisionTreeRegressor
DecisionTreeRegressor()
```

```
In [35]: 1 y_pred_dtr=model_dtr.predict(X_test_2d)
2 y_pred_dtr[:10]
```

```
Out[35]: array([53.40614888, 47.98372867, 58.31968422, 52.68566334, 57.0185102 ,
                55.81358842, 54.73145398, 50.68808953, 58.78946945, 50.35941677])
```

```
In [36]: 1 mean_squared_error(y_pred_dtr,y_test_2d)
```

```
Out[36]: 40.98078515849597
```

### ***RandomForestRegressor***

```
In [37]: 1 model_rfr=RandomForestRegressor()
2 model_rfr.fit(X_train_2d,y_train_2d)
```

```
C:\Users\sumit\AppData\Local\Temp\ipykernel_9532\3179291635.py:2: 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().
model_rfr.fit(X_train_2d,y_train_2d)
```

```
Out[37]: ▼ RandomForestRegressor
RandomForestRegressor()
```

```
In [38]: 1 y_pred_rfr=(X_test_2d)
          2 y_pred_rfr[:10]
```

```
Out[38]: array([[5.842447],
                 [5.789663],
                 [5.765922],
                 [5.901299],
                 [5.901055],
                 [5.766172],
                 [5.551245],
                 [5.79803 ],
                 [6.014263],
                 [5.873614]])
```

```
In [39]: 1 mean_squared_error(y_pred_rfr,y_test_2d)
```

```
Out[39]: 2715.7568554601735
```

### Hyperparameter tuning

```
In [40]: 1 from sklearn.model_selection import GridSearchCV
          2 from sklearn.linear_model import LinearRegression
          3
          4 # Define hyperparameters to tune
          5 param_grid = {
          6     'fit_intercept': [True, False],
          7     'copy_X': [True, False]
          8 }
          9
         10 # Create a Linear Regression model
         11 model_lr = LinearRegression()
         12
         13 # Initialize GridSearchCV
         14 grid_search = GridSearchCV(model_lr, param_grid, cv=5, scoring='neg_mean_squared_error')
         15
         16 # Fit the model
         17 grid_search.fit(X_train_2d, y_train_2d)
         18
         19 # Print the best parameters and best MSE score
         20 print("Best Parameters:", grid_search.best_params_)
         21 print("Best Negative MSE Score:", grid_search.best_score_)
         22
```

```
Best Parameters: {'copy_X': True, 'fit_intercept': True}
Best Negative MSE Score: -20.951832946201606
```

```
In [41]: 1 from sklearn.model_selection import cross_val_score
2 from sklearn.linear_model import LinearRegression
3
4 # Create a Linear Regression model
5 model_lr = LinearRegression()
6
7 # Perform 10-fold cross-validation
8 accuracy_scores = cross_val_score(model_lr, X_train_2d, y_train_2d, cv=
9
10 # Convert negative mean squared error to positive
11 mse_scores = -accuracy_scores
12
13 # Print the MSE scores
14 print("MSE Scores:", mse_scores)
15
```

MSE Scores: [20.89965582 21.67357185 22.39972494 20.26920314 20.56598084 20.7416008 20.67197687 20.67941739 21.2313776 20.38123774]

## final model

```
In [42]: 1 from sklearn.linear_model import LinearRegression
2
3 # Initialize the Linear Regression model with the best parameters
4 final_model = LinearRegression(fit_intercept=False, copy_X=True)
5
6 # Fit the model to the entire training data
7 final_model.fit(X_train_2d, y_train_2d)
8
9 # Now you can use final_model to make predictions on new data
10
```

```
Out[42]: LinearRegression
LinearRegression(fit_intercept=False)
```

## converting to pickle file

```
In [43]: 1 import pickle
2
3 # Define the filename for the pickle file
4 filename = 'final_model.pkl'
5
6 # Save the final_model to a pickle file
7 with open(filename, 'wb') as file:
8     pickle.dump(final_model, file)
9
```

```
In [44]: 1 os.path.abspath('final_model.pkl')
```

```
Out[44]: 'C:\\Users\\sumit\\Data Science\\Live Project\\Projects end to end\\ML Height Prediction\\final_model.pkl'
```

```
In [45]: 1 import pickle
2 import numpy as np
3
4 # Load the saved model from the file
5 filename = 'final_model.pkl'
6 with open(filename, 'rb') as file:
7     loaded_model = pickle.load(file)
8
9 # Input height for prediction
10 height_input = 6.0
11
12 # Reshape the input height to match the shape expected by the model (2D)
13 height_input_2d = np.array(height_input).reshape(1, -1)
14
15 # Use the loaded model to make predictions
16 predicted_weight = loaded_model.predict(height_input_2d)
17
18 # Print the predicted weight
19 print("Predicted weight:", predicted_weight[0, 0])
20
```

Predicted weight: 59.69879152974141