Obesity Level Predicto 1. Descriptive Statistics, Sir []: import pandas as pd	or mple Exploration and Data Cleaning
<pre>import pandas as pd import matplotlib.pyplot as plt import seaborn as sns import numpy as np from sklearn.cluster import KMeans</pre>	n/Desktop/Projects/Capstone/ObesityDataSet.csv')
0 Female 21.0 1.62 64.0 1 Female 21.0 1.52 56.0	istory_with_overweight FAVC FCVC NCP CAEC SMOKE CH20 SCC FAF TUE CALC MTRANS NObeyesdad yes no 2.0 3.0 Sometimes no 2.0 no 0.0 1.0 no Public_Transportation Normal_Weight yes no 3.0 3.0 Sometimes yes 3.0 yes 3.0 yes 3.0 verifies Public_Transportation Normal_Weight
2 Male 23.0 1.80 77.0 3 Male 27.0 1.80 87.0 4 Male 22.0 1.78 89.8 []: # Display descriptive statistics of df.describe()	yes no 2.0 3.0 Sometimes no 2.0 no 2.0 no 2.0 1.0 Frequently Public_Transportation Normal_Weight no no 3.0 3.0 Sometimes no 2.0 no 2.0 no 2.0 0.0 Frequently Walking Overweight_Level_I no no 2.0 1.0 Sometimes no 2.0 no 0.0 Sometimes Public_Transportation Overweight_Level_II f each numerical attribute
mean 24.312600 1.701677 86.58605 std 6.345968 0.093305 26.19117 min 14.000000 1.450000 39.00000	00 2111.000000 21
25% 19.947192 1.630000 65.47334 50% 22.777890 1.700499 83.00000 75% 26.000000 1.768464 107.43068 max 61.000000 1.980000 173.000000	00 2.385502 3.000000 2.000000 1.000000 0.625350 82 3.000000 3.000000 2.477420 1.666678 1.000000 00 3.000000 3.000000 3.000000 2.000000
# There are no missing values preserved: Gender 0 Age 0 Height 0 Weight 0 family_history_with_overweight 0	
FAVC 0 FCVC 0 NCP 0 CAEC 0 SMOKE 0 CH2O 0 SCC 0 FAF 0 TUE 0 CALC 0	
MTRANS 0 NObeyesdad 0 dtype: int64 # Number of total duplicated rows duprows = df.duplicated(subset = No print("There are", duprows, "duplic #Drop the duplicated rows	one, keep = first).sum()
<pre>df.info() # Total number of rows dropped from There are 24 duplicated rows <class 'pandas.core.frame.dataframe'<="" pre=""></class></pre>	'>
<pre>Int64Index: 2087 entries, 0 to 2110 Data columns (total 17 columns): # Column 0 Gender 1 Age 2 Height 3 Weight 4 family_history_with_overweight 5 FAVC</pre>	Non-Null Count Dtype
6 FCVC 7 NCP 8 CAEC 9 SMOKE 10 CH2O 11 SCC 12 FAF 13 TUE 14 CALC	2087 non-null float64 2087 non-null object 2087 non-null object 2087 non-null float64 2087 non-null float64 2087 non-null float64 2087 non-null object 2087 non-null float64 2087 non-null float64 2087 non-null object 2087 non-null object 2087 non-null object 2087 non-null object 2088 non-null object 2089 non-null object 2089 non-null object
<pre>15 MTRANS 16 NObeyesdad dtypes: float64(8), object(9) memory usage: 293.5+ KB []: # Display the boxplot for Height to sns.boxplot(df['Height']) /Library/Frameworks/Python.framework</pre>	2087 non-null object 2087 non-null object obj
n error or misinterpretation. FutureWarning <axessubplot:xlabel='height'></axessubplot:xlabel='height'>	
Height []: # Display the boxplot for Weight to sns.boxplot(df['Weight']) /Library/Frameworks/Python.framework n error or misinterpretation.	o identify any outliers k/Versions/3.7/lib/python3.7/site-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in the control of the co
FutureWarning <axessubplot:xlabel='weight'></axessubplot:xlabel='weight'>	
40 60 80 100 120 140 Weight	
<pre># Display the boxplot for Age to id sns.boxplot(df['Age']) /Library/Frameworks/Python.framework n error or misinterpretation. FutureWarning <axessubplot:xlabel='age'></axessubplot:xlabel='age'></pre>	dentify any outliers k/Versions/3.7/lib/python3.7/site-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in the control of the contr
-	* *
20 30 40 50 Age # Display the count of outliers for	
Q1 = df.quantile(0.25) Q3 = df.quantile(0.75) IQR = Q3 - Q1 outliers = ((df < (Q1 - 1.5 * IQR)) print(outliers)	
CALC CH2O FAF FAVC FCVC Gender Height MTRANS NCP 57	0 0 0 0 0 0 0 1 0 0
	O O O O Automatic reindexing on DataFrame vs Series comparisons is deprecated and will raise ValueError in a future version. Do `left, right = left.align(right, axis=1, copy=False)` before e.g. `left == right` of outliers. Since the attributes acts as an ordinal variable, we can keep it for now and remove the attribute if needed when developing the model. The most important attributes are Height and Weight for clustering, and we can see that it only has one outlier which shows that it does not need manipulation at the moment.
2. Deep Data Exploration and provided and pr	and Preparation for Cluster Analysis
Plt.show() Age - 1	0.042 - 0.8 5 - 0.079 - 0.6
NCP -0.056 0.23 0.092 0.035 1 0.075 0.13 CH2O -0.044 0.22 0.2 0.081 0.075 1 0.17 FAF - 0.15 0.29 0.056 0.022 0.13 0.17 1 TUE - 0.3 0.042 0.079 -0.1 0.016 0.021 0.059	0.021 0.059 1 -0.2
<pre>[]: # Create a copy of the original dat df_explore = df.copy()</pre>	ormula and Weight and Height columns
Gender Age Height Weight family_his O Female 21.0 1.62 64.0 1 Female 21.0 1.52 56.0 2 Male 23.0 1.80 77.0 3 Male 27.0 1.80 87.0	istory_with_overweight FAVC RVC RVC RVC RVC RVC RVC RVC RVC RVC R
<pre>df_explore["Truth"] = np.where(((df ((df_explore["BMI"] >= 18.5) & (df_ ((df_explore["BMI"] >= 25) & (df_explore["BMI"] >= 25)</pre>	no no 2.0 1.0 Sometimes no 2.0 no 0.0 0.0 Sometimes Public_Transportation Overweight_Level_II 28.342381 lumn named "Truth" to display if the BMI level falls into the correct category. f_explore["BMI"] < 18.5) & (df_explore["Nobeyesdad"] == "Insufficient_Weight")) explore["BMI"] < 25) & (df_explore["Nobeyesdad"] == "Normal_Weight")) vexplore["BMI"] < 30) & ((df_explore["Nobeyesdad"] == "Overweight_Level_I") (df_explore["Nobeyesdad"] == "Overweight_Level_II"))) xplore["BMI"] < 35) & (df_explore["Nobeyesdad"] == "Obesity_Type_I"))
<pre>((df_explore["BMI"] >= 35) & (df_ex ((df_explore["BMI"] >= 40) & (df_ex # Calculate the accuracy of BMI whe Acc_of_BMI = (df_explore["Truth"].s print("The Accuracy of BMI and Nobe</pre> The Accuracy of BMI and Nobeyesdad in	<pre>xplore["BMI"] < 40) & (df_explore["NObeyesdad"] == "Obesity_Type_II")) xplore["NObeyesdad"] == "Obesity_Type_III")) , True, False) en compared to the class label. sum() / len(df_explore))*100 eyesdad is", Acc_of_BMI)</pre>
	lass label to make sure there is accurate representation [].size().unstack(level=1).plot(kind='bar') Gender Female Female Male
250 - 200 - 150 - 100 - 50 -	Male Male Male Male Male Male Male Male
Normal Weight Normal Weight Obesity Type II Obesity Type III	Overweight Level II -
contribute into lowering the errors of the oth	e for clustering. Having a visual can show us if there is any cluster clearly depicted.
# Add best line of fit	<pre>ght"]), np.polyId(np.polyfit(df_explore["Height"], df_explore["Weight"], 1))(np.unique(df_explore["Height"])), color = 'black') 84c3c110>]</pre>
140 - (6) 120 - 80 - 60 -	
Completing the Cluster An # How to check the right number of	nalysis clusters for annual income vs spending score:
<pre>x = df.1loc[:,[2,3]].values wcss = [] for i in range(1,11):</pre>	<pre>it = 'k-means++', max_iter=300, n_init = 10, random_state=0)</pre>
<pre>plt.ylabel("wcss") plt.title("Elbow Analysis") plt.show() # Best number of clusters is 4 as t # Keep in mind that this is only fo le6 Elbow Analysis</pre>	
1.4 - 1.2 - 1.0 - 85 0.8 - 0.6 - 0.4 -	
0.2 - 0.0 - 2 4 6 k value According to the elbow method, the best value clustering.	slues of k is 4 clusters as the curve smoothens. However, we will still use 7 since we want to cluster based on each Obesity Level. On a more realistic note, having only 4 clusters that represent Underweight, Normal Weight, Overweight, and Obese would be the most effective way to develop the most accurate model based on
<pre>y_means = km.fit_predict(x) # Select random observations as cen fig=plt.figure(figsize=(12,6))</pre>	
<pre>plt.scatter(x[y_means == 0, 0], x[y plt.scatter(x[y_means == 2, 0], x[y plt.scatter(x[y_means == 6, 0], x[y plt.scatter(x[y_means == 1, 0], x[y plt.scatter(x[y_means == 4, 0], x[y plt.scatter(x[y_means == 3, 0], x[y plt.scatter(x[y_means == 3, 0], x[y plt.scatter(km.cluster_centers_[:,0 plt.xlabel("Height (m)")</pre>	<pre>y_means == 5, 1], s = 100, c = 'red', label = 'Insufficient_Weight') y_means == 0, 1], s = 100, c = 'pink', label = 'Normal_Weight') y_means == 2, 1], s = 100, c = 'cyan', label = 'Overweight_Level_I') y_means == 6, 1], s = 100, c = 'grey', label = 'Overweight_Level_II') y_means == 1, 1], s = 100, c = 'yellow', label = 'Obesity_Type_I') y_means == 4, 1], s = 100, c = 'orange', label = 'Obesity_Type_I') y_means == 3, 1], s = 100, c = 'magenta', label = 'Obesity_Type_III') y_means == 3, 1], s = 100, c = 'magenta', label = 'Centeroid')</pre>
<pre>plt.ylabel("Weight (kg)") plt.legend() plt.grid() plt.show()</pre> <pre>Insufficient_Weight</pre>	
Obesity_Type_II Obesity_Type_III Obesity_Type_III Centeroid 80	
15 16	17 18 19 20 Height (m)
<pre># Export the clustered results to a df["Cluster"] = y_means df.head()</pre>	_abel and Clustering to measure Performance a columns istory_with_overweight FAVC FCVC NCP CAEC SMOKE CH20 SCC FAF TUE CALC MTRANS NObeyesdad Cluster
Gender Age Height Weight family_his 0 Female 21.0 1.62 64.0 1 Female 21.0 1.52 56.0 2 Male 23.0 1.80 77.0 3 Male 27.0 1.80 87.0 4 Male 22.0 1.78 89.8	istory_with_overweight FAVC PCV NCP CACC SMOKE CH20 SCC FAF TUE CALC MTRANS NObeyesdad Cluster yes no 2.0 3.0 Sometimes no 2.0
<pre>df["Truth"] = np.where(((df["Cluste ((df["Cluster"] == 1) & (df["NObeye ((df["Cluster"] == 2) & (df["NObeye ((df["Cluster"] == 3) & (df["NObeye ((df["Cluster"] == 4) & (df["NObeye</pre>	esdad"] == "Overweight_Level_I")) esdad"] == "Overweight_Level_II")) esdad"] == "Obesity_Type_I"))
# Calculate the accuracy of Cluster Acc_of_Cluster = (df["Truth"].sum() print("The Accuracy of Cluster when The Accuracy of Cluster when compare When comparing the cluster to the original of	esdad"] == "Obesity_Type_III")) , True, False) r when compared to the class label.
	rvised Learning (To be continued)