	Obesity Level Predictor 1. Descriptive Statistics, Simple Exploration and Data Cleaning
In []	import pandas as pd import matplotlib.pyplot as plt import seaborn as sns import numpy as np from sklearn.cluster import KMeans from sklearn.model_selection import train_test_split from sklearn.ensemble import RandomForestClassifier from sklearn.metrics import classification_report
In []	<pre># Reading the file df = pd.read_csv('/Users/azamrahman/Desktop/Projects/Capstone/ObesityDataSet.csv') # Display the first 5 rows df.head()</pre>
Out[]	Figure 1. Sender
In [] Out[]	df.describe() Age Height Weight FCVC NCP CH2O FAF TUE
	count 2111.00000 </td
In []	75% 26.00000 1.768464 107.430682 3.00000 3.00000 2.477420 1.666678 1.000000 max 61.00000 1.980000 173.000000 3.000000 3.000000 3.000000 2.000000 # Display the total null values in each attribute df.isnull().sum()
Out[]	Age 0 Height 0 Weight 0 family_history_with_overweight 0 FAVC 0
	FCVC NCP CAEC CAEC SMOKE CH2O SCC FAF TUE CALC CALC CALC CALC CALC CALC CALC CAL
In []	MTRANS 0 NObeyesdad 0 dtype: int64 # Number of total duplicated rows duprows = df.duplicated(subset = None, keep = 'first').sum() print("There are", duprows, "duplicated rows") #Drop the duplicated rows
	<pre>df = df.drop_duplicates() # Re-count the total number of rows and display the type of data attributes df.info() # Total number of rows dropped from 2110 to 2087 rows There are 24 duplicated rows <class 'pandas.core.frame.dataframe'=""></class></pre>
	Int64Index: 2087 entries, 0 to 2110 Data columns (total 17 columns): # Column
	5 FAVC 2087 non-null object 6 FCVC 2087 non-null float64 7 NCP 2087 non-null object 8 CAEC 2087 non-null object 9 SMOKE 2087 non-null object 10 CH2O 2087 non-null float64 11 SCC 2087 non-null object 12 FAF 2087 non-null float64 13 TUE 2087 non-null float64
In []	14 CALC 2087 non-null object 15 MTRANS 2087 non-null object 16 NObeyesdad 2087 non-null object dtypes: float64(8), object(9) memory usage: 293.5+ KB
Out[]	/Library/Frameworks/Python.framework/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 a rerror or misinterpretation. FutureWarning <axessubplot:xlabel='height'></axessubplot:xlabel='height'>
In []	15 16 17 18 19 20 Height # Display the boxplot for Weight to identify any outliers sns.boxplot(df['Weight'])
Out[]	/Library/Frameworks/Python.framework/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 a nerror or misinterpretation. FutureWarning
In []	# Display the boxplot for Age to identify any outliers sns.boxplot(df['Age']) /Library/Frameworks/Python.framework/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 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 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 a keyword arg: x. From version 0.12, the only valid positional argument will be 'data'.
Out[]	/Library/Frameworks/Python.framework/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 a rerror or misinterpretation. FutureWarning <axessubplot:xlabel='age'></axessubplot:xlabel='age'>
In []	# Display the count of outliers for each attribute Q1 = df.quantile(0.25) Q3 = df.quantile(0.75) IQR = Q3 - Q1 outliers = ((df < (Q1 - 1.5 * IQR)) (df > (Q3 + 1.5 * IQR))).sum()
	Age 167 CAEC 0 CALC 0 CH2O 0 FAF 0 FAF 0 FAVC 0
	FCVC 0 Gender 0 Height 1 MTRANS 0 NCP 577 NObeyesdad 0 SCC 0 SMOKE 0 TUE 0 0
	Weight 1 family_history_with_overweight 0 dtype: int64 /Users/azamrahman/Library/Python/3.7/lib/python/site-packages/ipykernel_launcher.py:5: FutureWarning: 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` """ We can see that NCP has a very high count 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.
In []	2. Deep Data Exploration and Preparation for Cluster Analysis # plotting correlation heatmap sns.heatmap(df.corr(),annot=True, cmap='coolwarm') # displaying heatmap plt.show()
	Age - 1 0 0.032 0.2 0.014 0.056 0.044 0.15 0.3 - 0.04 0.15 0.15 0.15 0.15 0.15 0.15 0.15 0.15
	CH20 -0.044 0.22 0.2 0.081 0.075 1 0.17 0.021 FAF -0.15 0.29 0.056 0.022 0.13 0.17 1 0.059 TUE - 0.3 0.042 0.079 0.1 0.016 0.021 0.059 10.2 The correlation beatment indicates that some of the variables have a signifficant correlation to other that needs to be dealt with Unight and Weight have the description of 46%. Cince we know both variables are as wish in determining the Obesital and it will appear in the Obesital and it wi
In []	The correlation heatmap indicates that none of the variables have a significant correlation to each other that needs to be dealt with. Height and Weight have the strongest positive correlation of 46%. Since we know both variables are crucial in determining the Obesity Level, it will remain in the dataset. # Create a copy of the original dataframe for data exploration df_explore = df.copy() # Compute BMI using the original formula and Weight and Height columns df_explore["Meight"] = df_explore["Weight"] / df_explore["Height"]**2 # Display first 5 rows df_explore.head()
Out[]	Gender Age Height Weight family_history_with_overweight FAVE VEV NO CALC MTRANS NO beyesdad BMI Permale 21.0 1.62 64.0 1.02 4.03 1.02 5.03 1.02 5.03 1.02 5.03 1.02 5.03 1.02 5.03 1.02 5.03 1.02 5.03 1.02 5.03
In []	4 Male 22.0 1.78 89.8 no no 2.0 1.0 Sometimes no 2.0 no 0.0 0.0 Sometimes Public_Transportation Overweight_Level_II 28.342381 **# Apply logic to create another column named "Truth" to display if the BMI level falls into the correct category. df_explore["Truth"] = np.where(((df_explore["BMI"] < 18.5) & (df_explore["NObeyesdad"] == "Insufficient_Weight")) ((df_explore["BMI"] >= 18.5) & (df_explore["BMI"] < 25) & (df_explore["Nobeyesdad"] == "Normal_Weight")) ((df_explore["BMI"] >= 25) & (df_explore["Nobeyesdad"] == "Overweight_Level_I") (df_explore["Nobeyesdad"] == "Overweight_Level_I") ((df_explore["BMI"] >= 30) & (df_explore["BMI"] < 30) & (df_explore["Nobeyesdad"] == "Obesity_Type_I"))
	((df_explore["BMI"] >= 35) & (df_explore["NObeyesdad"] == "Obesity_Type_II")) ((df_explore["BMI"] >= 40) & (df_explore["NObeyesdad"] == "Obesity_Type_III")) , True, False) # Calculate the accuracy of BMI when compared to the class label. Acc_of_BMI = (df_explore["Truth"].sum() / len(df_explore))*100 print("The Accuracy of BMI and Nobeyesdad is", Acc_of_BMI) The Accuracy of BMI and Nobeyesdad is 95.11260182079539 This shows that BMI is not directly correlated to the Obesity level since not all levels are True. This means that some of the other columsn come intp the equation into determining the proper diagnosis for Obesity Level. The class label is not completely determined by the BMI (Height vs Weight) but have a very accurate representation.
	# Break the Data into Gender and Class label to make sure there is accurate representation df.groupby(["NObeyesdad", "Gender"]).size().unstack(level=1).plot(kind='bar') <pre></pre>
	250 - 200 - 150 - 100 - 70
	Normal_Weight - Normal_Weight Obesity_Type_III Obesity_Type_III Overweight_LeveI Overweight_DeveI
In []	plt.scatter(df_explore["Height"], df_explore["Weight"]) plt.xlabel("Height (m)")
Out[]	plt.ylabel("Weight vs. Weight") # Add best line of fit plt.plot(np.unique(df_explore["Height"]), np.polyld(np.polyfit(df_explore["Weight"], 1))(np.unique(df_explore["Height"])), color = 'black') [<matplotlib.lines.line2d 0x7fde192caad0="" at="">] Height vs. Weight</matplotlib.lines.line2d>
	160 - 140 - (g) 120 - (g) 100 - (g)
	Completing the Cluster Analysis
In []	<pre># How to check the right number of clusters for annual income vs spending score: x = df.iloc[:,[2,3]].values wcss = [] for i in range(1,11): km = KMeans(n_clusters = i, init = 'k-means++', max_iter=300, n_init = 10, random_state=0) km.fit(x) wcss.append(km.inertia_)</pre>
	<pre>plt.plot(range(1,11), wcss) plt.xlabel("k value") plt.ylabel("wcss") plt.title("Elbow Analysis") plt.show() # Best number of clusters is 4 as the curve smoothens after 4. # Keep in mind that this is only for height vs weight</pre>
	14 - 12 - 10 - 10 - 10 - 10 - 10 - 10 - 10
	0.4 0.2 0.0 2 4 6 8 10 k value
In []	According to the elbow method, the best values 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 clustering. # Number of Clusters: K = 7
	<pre>y_means = km.fit_predict(x) # Select random observations as centroids fig=plt.figure(figsize=(12,6)) plt.scatter(x[y_means == 5, 0], x[y_means == 5, 1], s = 100, c = 'red', label = 'Insufficient_Weight') plt.scatter(x[y_means == 0, 0], x[y_means == 0, 1], s = 100, c = 'pink', label = 'Normal_Weight') plt.scatter(x[y_means == 2, 0], x[y_means == 2, 1], s = 100, c = 'cyan', label = 'Overweight_Level_I') plt.scatter(x[y_means == 6, 0], x[y_means == 6, 1], s = 100, c = 'grey', label = 'Overweight_Level_II') plt.scatter(x[y_means == 1, 0], x[y_means == 1, 1], s = 100, c = 'yellow', label = 'Obesity_Type_II') plt.scatter(x[y_means == 4, 0], x[y_means == 4, 1], s = 100, c = 'orange', label = 'Obesity_Type_II')</pre>
	<pre>plt.scatter(x[y_means == 3, 0], x[y_means == 3, 1], s = 100, c = 'magenta', label = 'Obesity_Type_III') plt.scatter(km.cluster_centers_[:, 0], km.cluster_centers_[:, 1], s = 50, c = 'blue', label = 'Centeroid') plt.ylabel("Height (kg)") plt.ylabel("Weight (kg)") plt.grid() plt.show()</pre>
	Insufficient_Weight Normal_Weight Overweight_Level II Obesity_Type_II Obesity_Type_III Centeroid
In []	Comparing original Class Label and Clustering to measure Performance # Export the clustered results to a columns
Out[]	df["Cluster"] = y_means df.head() Gender Age Height Weight family_history_with_overweight FAVC FCVC NCP CAEC SMOKE CH20 SCC FAF TUE CALC MTRANS NObeyesdad Cluster Female 21.0 1.62 64.0 yes no 2.0 3.0 Sometimes no 2.0 no 0.0 1.0 no Public_Transportation Normal_Weight 1 Female 21.0 1.52 56.0 yes no 3.0 3.0 Sometimes yes 3.0 yes 3.0 yes 3.0 ves 3.0 yes 3.0 Normal_Weight 3
In []	# Apply logic to create another column hamed fluth to display if the cluster falls into the correct category.
	<pre>df["Truth"] = np.where(((df["Cluster"] == 0) & (df["Nobeyesdad"] == "Insufficient_Weight")) ((df["Cluster"] == 1) & (df["Nobeyesdad"] == "Normal_Weight")) ((df["Cluster"] == 2) & (df["Nobeyesdad"] == "Overweight_Level_I")) ((df["Cluster"] == 3) & (df["Nobeyesdad"] == "Overweight_Level_II")) ((dff["Cluster"] == 4) & (dff["Nobeyesdad"] == "Obesity_Type_I")) ((dff["Cluster"] == 5) & (dff["Nobeyesdad"] == "Obesity_Type_II")) ((dff["Cluster"] == 6) & (dff["Nobeyesdad"] == "Obesity_Type_III")) , True, False) # Calculate the accuracy of Cluster when compared to the class label. Acc_of_Cluster = (dff["Truth"].sum() / len(df))*100 print("The Accuracy of Cluster when compared to original class label is", Acc_of_Cluster, "%")</pre>
	The Accuracy of Cluster when compared to original class label is 10.11020603737422 % When comparing the cluster to the original class label, we can see that the K-means cluster analysis did not do a good job predicting the level of obesity correctly using 7 clusters. This is perhaps the amount of clusters are too much, or the height and weight were not the best indicators to diagnose obesity. In the Final results submission. I will include a code that runs the K-means Algorithm with just 4 clusters to see if it can develop a more accurate class label. For example, the 4 clusters can be Underweight, Normal, Overweight, and Obese. Since the elbow method, it will be interesting to see how well it will do. Data Preparation for Supervised Learning
In []	<pre># Columns of all numerical variables: num_cols = dfget_numeric_data().columns.tolist() # Columns of all categorical variables: cat_cols = df.select_dtypes(include=['object']).columns.tolist() cat_cols.remove('NObeyesdad')</pre>
In [] Out[]	# Creating a copy of the dataframe df_onehot=df.copy() # Converting categorical variables to dummy variables to complete encoding df_onehot = pd.get_dummies(df, columns=cat_cols, prefix = cat_cols) df_onehot df_onehot
- L J	Age Feight Weight FeV CV
In []	2087 rows × 34 columns Model Building for Random Forest to predict Obesity Level # Make sure that the class label is not included
In []	class_col_name = 'Nobeyesdad' one_hot_feature_names=df_onehot.columns[df_onehot.columns != class_col_name] # Split dataset into training set and test set: 70% training and 30% test X_train, X_test, y_train, y_test = train_test_split(df_onehot.loc[:, one_hot_feature_names], df_onehot[class_col_name], test_size=0.3,random_state=109)
	<pre># Create Random Forest Model clf = RandomForestClassifier(n_estimators=100) # Train the model clf.fit(X_train, y_train) # Create the Predictions y_pred = clf.predict(X_test)</pre>
In []	Initial Performance Measures for Random Forest # Show the classification Report to view the results print(classification_report(y_test, y_pred)) precision recall f1-score support Insufficient Weight 0.96 0.95 0.96 82
	Normal_Weight 0.82 0.87 0.85 70 Obesity_Type_I 0.98 0.96 0.97 122 Obesity_Type_III 0.99 0.98 0.98 84 Obesity_Type_III 0.99 1.00 0.99 98 Overweight_Level_I 0.94 0.88 0.91 95 Overweight_Level_II 0.89 0.96 0.92 76 accuracy 0.95 627
	macro avg 0.94 0.94 0.94 6.27 weighted avg 0.95 0.95 0.95 6.27 In our initial Results, we can see that the Random Forest model can predict the obesity level at an accuracy of 95%. However, we can see through the precision of Normal_Weight shows that it does have some trouble predicting the true positives of the obesity level. In the final results, we will investigate feature engineering, overtraining, and the best possible parameters for the model with in detailed classification results.