

Data Dictionary

- **Gender** : Sex of individual
- **Age** : How old is the individual
- **Height** : How tall is the individual
- **Weight** : How heavy is the individual
- **family_history_with_overweight** :Family history of obesity
- **FAVC** : Frequency of consumption of high-calorie foods
- **FCVC** : Frequency of vegetable consumption
- **NCP** : Number of main meals
- **CAEC** : Food consumption between meals
- **SMOKE** : Tobacco use
- **CH20** : Daily water consumption
- **SCC** : Monitoring calorie consumption.
- **FAF** : Frequency of physical activity
- **TUE** : Time for using technological devices
- **CALC** : Alcohol consumption
- **MTanANS** : Mode of transport used
- **NObesyedad** : Gives us the obesity status of the person

✓ Coletando Dados

Os dados foram disponibilizados em meu github público. Eu subi esses arquivos no github e estou usando esse script pronto para carregar esses arquivos no Google Colab.

```

1 import os
2 import zipfile
3 import requests
4 from io import BytesIO
5
6 def download_and_extract_github_folder(repo_url, branch, folder_path, target_path):
7     zip_url = f"{repo_url}/archive/{branch}.zip"
8
9     print("Downloading the zip file...")
10    response = requests.get(zip_url, stream=True)
11    zip_file = zipfile.ZipFile(BytesIO(response.content))
12
13    print("Extracting files...")
14    for file in zip_file.namelist():
15        if file.startswith(f"{repo_url.split('/')[-1]}-{branch}/{folder_path}"):
16            zip_file.extract(file, target_path)
17
18    extracted_folder = os.path.join(target_path, f"{repo_url.split('/')[-1]}-{branch}", folder_path)
19    os.system(f"mv {extracted_folder}/* {target_path}")
20    os.system(f"rm -r {target_path}/{repo_url.split('/')[-1]}-{branch}")
21
22    print("Files are ready to use at:", target_path)
23
24 repo_url = 'https://github.com/geeklicantropo/MESTRADO_PUBLICO'
25 branch = 'main'
26 folder_path = 'DATA_MINING/TAREFAS/Tarefa1/data/'
27 target_path = '/content/'
28
29 download_and_extract_github_folder(repo_url, branch, folder_path, target_path)
30
31
32 Downloading the zip file...
33 Extracting files...
34 Files are ready to use at: /content/

```

✓ Loading libraries

```
1 import pandas as pd
2 import numpy as np
3
4 #Visualização de Dados
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 #Machine Learning - Modelagem e Avaliação
9 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
10 from sklearn.preprocessing import StandardScaler, LabelEncoder
11 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
12
13 #Algoritmos preditivos
14 from sklearn.ensemble import RandomForestClassifier
15 from sklearn.naive_bayes import GaussianNB
16 from sklearn.svm import SVC
17 import xgboost as xgb
18 import lightgbm as lgb
19
20 #Biblioteca para mostrar barras de progresso em loops.
21 from tqdm import tqdm
22
23 #Estatísticas
24 from scipy.stats import chi2_contingency, f_oneway, spearmanr, skew

```

```
1 train = pd.read_csv("/content/train.csv")
2 test = pd.read_csv("/content/test.csv")
3
```

EDA

Aqui estou fazendo a análise exploratória dos dados, para investigar quaisquer informações implícitas nos dados que possa me ajudar a criar um modelo de classificação assertivo.

```
1 train.head()
```

	id	Gender	Age	Height	Weight	family_history_with_overweight	FAVC
0	0	Male	24.443011	1.699998	81.669950	yes	yes
1	1	Female	18.000000	1.560000	57.000000	yes	yes
2	2	Female	18.000000	1.711460	50.165754	yes	yes
3	3	Female	20.952737	1.710730	131.274851	yes	yes

Next steps: [View recommended plots](#)

```
1 test.head()
```

	id	Gender	Age	Height	Weight	family_history_with_overweight	FA
0	20758	Male	26.899886	1.848294	120.644178	yes	y
1	20759	Female	21.000000	1.600000	66.000000	yes	y
2	20760	Female	26.000000	1.643355	111.600553	yes	y
3	20761	Male	20.979254	1.553127	103.669116	yes	y

Next steps: [View recommended plots](#)

```
1 train.shape

(20758, 18)
```

```
1 test.shape
```

```
(13840, 17)
```

Criando uma classe para múltiplas tarefas de EDA

Aqui estou encapsulando múltiplas tarefas da análise exploratória dos dados em uma única classe.

```

1 class EDA:
2     def __init__(self, data, target_col):
3         self.data = data
4         self.target_col = target_col
5         self.target_data = data[target_col]
6
7     #Principais análises estatísticas
8     def data_summary(self):
9         display(self.data.describe(include='all').transpose())
10
11     #Verificar valores nulos
12     def missing_value_analysis(self):
13         display(
14             pd.DataFrame(
15                 self.data.isnull().sum(axis = 0) / len(self.data) * 100, columns = ['Null Value Ratio']
16             ).applymap(lambda x: f"{np.round(x, decimals = 5)}%")
17         )
18
19     #Checar os tipos de dados
20     def checking_datatypes(self):
21         display(self.data.dtypes)
22
23     #Checar linhas duplicadas
24     def verify_duplicates(self):
25         n_duplicates = self.data.duplicated().sum()
26         print(f"Number of duplicates: {n_duplicates}.")
27
28     #Checar a distribuição dos dados de maneira visual
29     def data_distribution(self):
30         numerical_cols = self.data.select_dtypes(include=['float64', 'int64']).columns
31         categorical_cols = self.data.select_dtypes(include=['object', 'category']).columns
32
33         #Determina o número de linhas necessárias para subplots
34         total_cols = len(numerical_cols) + len(categorical_cols)
35         num_rows = (total_cols + 3) // 4 #Garanta linhas suficientes adicionando 3 antes da divisão inteira por 4
36
37         #Cria um esboço de imagem maior que irá conter os subplots
38         fig, axes = plt.subplots(nrows=num_rows, ncols=4, figsize=(20, num_rows * 4))
39         axes = axes.flatten() #Estratégia usada para evitar erros
40
41         #Plotando historigramas para features numéricas
42         for idx, col in enumerate(numerical_cols):
43             if idx < len(axes): #Estratégia usada para evitar erros
44                 sns.histplot(self.data[col], kde=True, bins=30, color='blue', ax=axes[idx])
45                 axes[idx].set_title(f'Distribution of {col}')
46
47         #Plotando gráficos de barras para features categóricas
48         offset = len(numerical_cols) #Começa a plotar as variáveis categóricas depois das numéricas
49         for idx, col in enumerate(categorical_cols, start=offset):
50             if idx < len(axes): #Estratégia usada para evitar erros
51                 sns.countplot(x=self.data[col], ax=axes[idx])
52                 axes[idx].set_title(f'Distribution of {col}')
53                 axes[idx].tick_params(axis='x', rotation=45)
54
55         #Esconde qualquer eixo não usado
56         for ax in axes[total_cols:]:
57             ax.axis('off')
58
59         plt.tight_layout()
60         plt.show()
61
62     #Identifica outliers usando o intervalo interquartil
63     def identify_outliers(self):
64         numerical_features = self.data.select_dtypes(include=['float64', 'int64']).columns
65         outliers_info = {}
66
67         # Dicionário com as descrições das variáveis
68         descriptions = {
69             'FAF': 'Frequency of physical activity',
70             'Gender': 'Sex of individual',
71             'Age': 'How old is the individual',
72             'Height': 'How tall is the individual',
73             'Weight': 'How heavy is the individual',
74             'family_history_with_overweight': 'Family history of obesity',
75             'FAVC': 'Frequency of consumption of high-calorie foods',
76             'FCVC': 'Frequency of vegetable consumption',
77             'NCP': 'Number of main meals',

```



```
78     'CAEC' : 'Food consumption between meals',
79     'SMOKE' : 'Tobacco use',
80     'CH20' : 'Daily water consumption',
81     'SCC' : 'Monitoring calorie consumption',
82     'FAF' : 'Frequency of physical activity',
83     'TUE' : 'Time for using technological devices',
84     'CALC' : 'Alcohol consumption',
85     'MTanANS' : 'Mode of transport used',
86     'NObesyedad' : 'Gives us the obesity status of the person'
87 }
88
89 #Determine o número de linhas necessárias para subplots
90 total_cols = len(numerical_features)
91 num_rows = (total_cols + 3) // 4
92
93 #Cria um conjunto de subplots onde em cada linha, teremos no máximo 4 subplots
94 fig, axes = plt.subplots(nrows=num_rows, ncols=4, figsize=(20, num_rows * 4))
95 axes = axes.flatten()
96
97 #Loop sobre todas as features numéricas para calcular e plotar os outliers
98 for idx, feature in enumerate(numerical_features):
99     Q1 = self.data[feature].quantile(0.25)
100     Q3 = self.data[feature].quantile(0.75)
101     IQR = Q3 - Q1
102
103     lower_bound = Q1 - 1.5 * IQR
104     upper_bound = Q3 + 1.5 * IQR
105
106     outliers = self.data[(self.data[feature] < lower_bound) | (self.data[feature] > upper_bound)]
107     outliers_info[feature] = {'count': len(outliers), 'outliers': outliers}
108
109     #Plota cada boxplot em seu respectivo subplot
110     sns.boxplot(y=self.data[feature], color='orange', ax=axes[idx])
111     full_title = descriptions.get(feature, feature) #Pega a descrição completa do dicionário de descriptions
112     axes[idx].set_title(f'Outliers in {feature} ({full_title})')
113
114 #Oculta quaisquer eixos não utilizados se o número de features for menor que o número de subplots
115 for i in range(idx + 1, len(axes)):
116     axes[i].axis('off')
117
118 plt.tight_layout()
119 plt.show()
120
121 return outliers_info
122
123 #Analisa a assimetria das features
124 def analyze_skewness(self, plot=False):
125     numerical_features = self.data.select_dtypes(include=['float64', 'int64']).columns
126     skew_info = {}
127     if plot:
128         total_features = len(numerical_features)
129         num_rows = (total_features + 3) // 4
130         fig, axes = plt.subplots(nrows=num_rows, ncols=4, figsize=(20, num_rows * 4))
131         axes = axes.flatten()
132
133         for idx, feature in enumerate(numerical_features):
134             skewness = skew(self.data[feature].dropna())
135             skew_info[feature] = skewness
136             if plot:
137                 sns.histplot(self.data[feature], kde=True, color='purple', ax=axes[idx])
138                 axes[idx].set_title(f'Skewness of {feature}: {skewness:.2f}')
139
140         if plot:
141             for i in range(idx + 1, len(axes)):
142                 axes[i].axis('off')
143             plt.tight_layout()
144             plt.show()
145
146     #Esse retorno é necessário para ser usado na plotagem da correlação logo após.
147     return skew_info
148
149
150 #Cria o gráfico de correlação (Heatmap)
151 def correlation_analysis_heatmap(self, skew_threshold=1):
152     '''
153     Verifica a assimetria usando a função de analisar assimetria. Isso é importante porque
154     se a feature for assimétrica (skewness acima de 1), usar a correlação de Pearson não irá funcionar.
```

```
155     Como alternativa, a correlação de Spearman será usada.
156     ...
157     skew_info = self.analyze_skewness()
158     skewed_features = any(abs(skew) > skew_threshold for skew in skew_info.values())
159
160     corr_data = self.data.copy()
161     for column in corr_data.select_dtypes(include=['object', 'category']).columns:
162         corr_data[column] = corr_data[column].astype('category').cat.codes
163
164     #Define o método de correlação baseado na assimetria.
165     correlation_method = 'spearman' if skewed_features else 'pearson'
166     correlation_matrix = corr_data.corr(method=correlation_method)
167
168     plt.figure(figsize=(12, 10))
169     sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', cbar=True)
170     plt.title(f'Correlation Matrix using {correlation_method.capitalize()}')
171     plt.show()
172
173     return correlation_matrix
174
```

```
1 #Aqui estou criando uma cópia do dataframe do trem apenas para deletar a coluna id e fazer algum EDA
2 df_train = train.copy()
3 df_train = df_train.drop(labels=["id"], axis=1)
```

```
1 #Cria uma instância da classe de EDA para a análise exploratória de dados
2 analyzer = EDA(df_train, 'NObeyesdad')
```

```
1 #Análises estatísticas principais
2 print("Resumo dos dados")
3 analyzer.data_summary()
```

Resumo dos dados

	count	unique	top	freq	mean	
Gender	20758	2	Female	10422	NaN	
Age	20758.0	NaN	NaN	NaN	23.841804	5.0
Height	20758.0	NaN	NaN	NaN	1.700245	0.0
Weight	20758.0	NaN	NaN	NaN	87.887768	26.0
family_history_with_overweight	20758	2	yes	17014	NaN	
FAVC	20758	2	yes	18982	NaN	
FCVC	20758.0	NaN	NaN	NaN	2.445908	0.0
NCP	20758.0	NaN	NaN	NaN	2.761332	0.0
CAEC	20758	4	Sometimes	17529	NaN	
SMOKE	20758	2	no	20513	NaN	
CH2O	20758.0	NaN	NaN	NaN	2.029418	0.0
SCC	20758	2	no	20071	NaN	
FAF	20758.0	NaN	NaN	NaN	0.981747	0.0
TUE	20758.0	NaN	NaN	NaN	0.616756	0.0
CALC	20758	3	Sometimes	15066	NaN	
MTRANS	20758	5	Public_Transportation	16687	NaN	



```
1 #Verifica valores nulos
2 analyzer.missing_value_analysis()
```

	Null Value Ratio	
Gender	0.0%	
Age	0.0%	
Height	0.0%	
Weight	0.0%	
family_history_with_overweight	0.0%	
FAVC	0.0%	
FCVC	0.0%	
NCP	0.0%	
CAEC	0.0%	
SMOKE	0.0%	
CH2O	0.0%	
SCC	0.0%	
FAF	0.0%	
TUE	0.0%	
CALC	0.0%	
MTRANS	0.0%	
NObeyesdad	0.0%	

```
1 #Verifica os tipos de dados
2 analyzer.checking_datatypes()

Gender      object
Age          float64
Height       float64
Weight       float64
family_history_with_overweight  object
FAVC         object
FCVC         float64
NCP          float64
CAEC         object
SMOKE        object
CH2O         float64
SCC          object
FAF          float64
TUE          float64
CALC         object
MTRANS       object
NObeyesdad   object
dtype: object
```

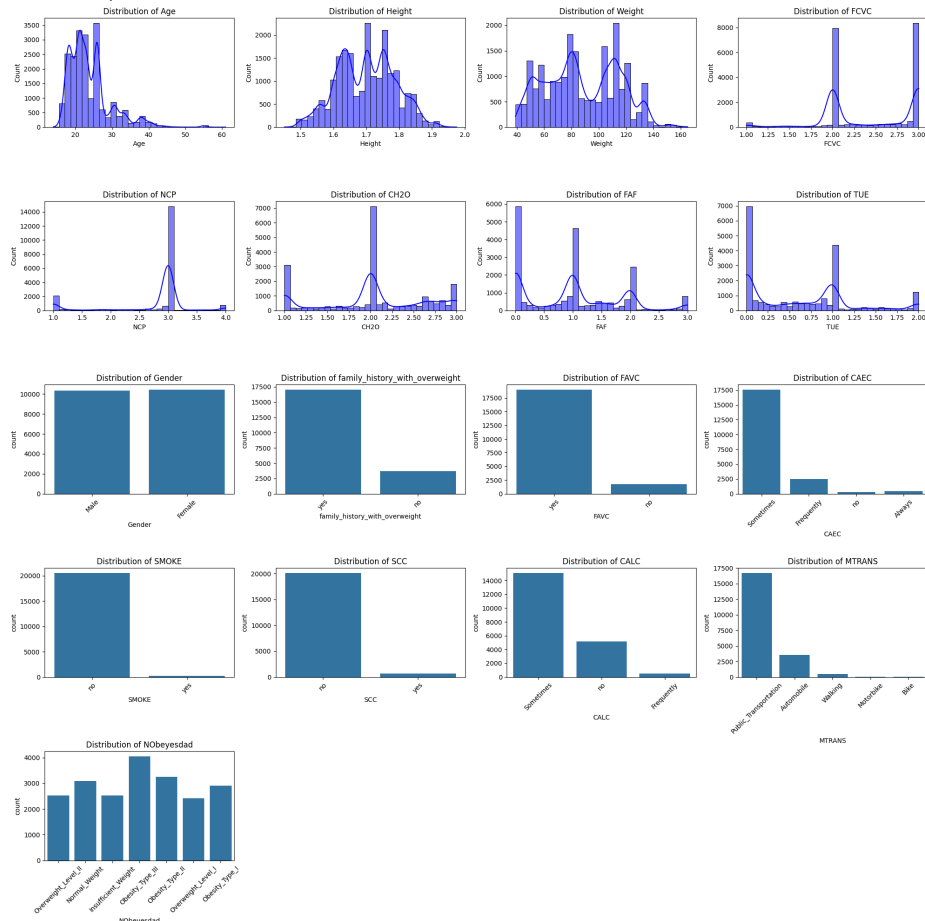
```
1 #Verifica linhas duplicadas
2 analyzer.verify_duplicates()

Number of duplicates: 0.
```

- A obesidade tipo 3 (o pior tipo) é a mais comum, infelizmente
- Existem alguns valores discrepantes para algumas colunas, especialmente Age
- A coluna Age é um valor do tipo Float. Ela precisa ser alterada para um inteiro
- Existem algumas features muito desbalanceadas
- Não temos problemas com valores nulos
- Não temos problemas com duplicatas também.

```
1 #Verifica a distribuição dos dados
2 print("Distribuição de Dados:")
3 analyzer.data_distribution()
```


Distribuição de Dados:



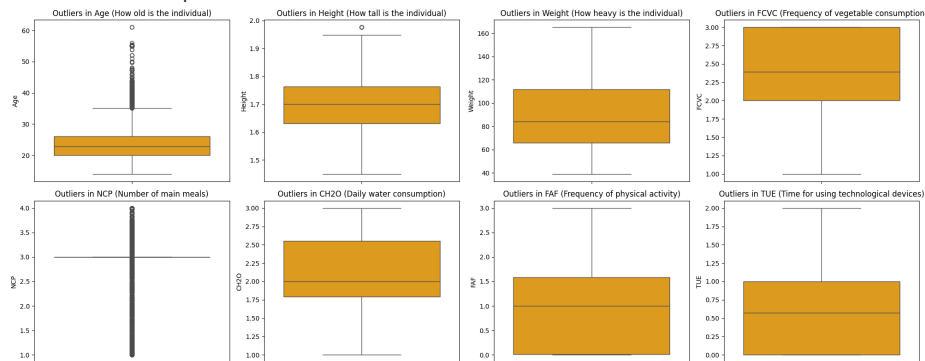
- O histograma que mostra gráficos numéricos está skewed na maioria das features.
- Apenas Height e Weight e Age têm uma distribuição que parece uma distribuição normal.
- Precisamos aplicar a normalização a esses dados.
- Os valores categóricos são em sua maioria desequilibrados. Apenas a coluna Gender está bem equilibrada

```

1 # Identificar e imprimir informações sobre outliers
2 print("\nIdentificando e plotando outliers:")
3 outliers_info = analyzer.identify_outliers()
4 for feature, details in outliers_info.items():
5     print(f"{feature} tem {details['count']} outliers")

```

Identificando e plotando outliers:

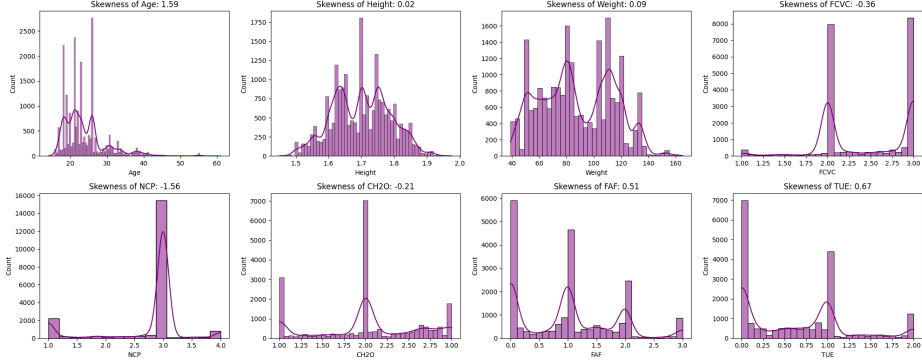


Age tem 1074 outliers
 Height tem 4 outliers
 Weight tem 0 outliers
 FCVC tem 0 outliers
 NCP tem 6052 outliers
 CH2O tem 0 outliers
 FAF tem 0 outliers
 TUE tem 0 outliers

- Como identificado antes, Age tem a maior quantidade de outliers.
- NCP possui valores entre 1 e 4. Se todos ou a maioria dos dados estão concentrados em um único valor ou em poucos valores muito próximos, será difícil enxergar o boxplot. Neste caso, o primeiro quartil (Q1), a mediana (Q2) e o terceiro quartil (Q3) podem ser todos iguais ou muito próximos.
- Eu entendo que o NCP deveria ser uma variável inteira e não um float.

```
1 # Analisar e imprimir informações sobre a assimetria (skewness)
2 print("\nAnalisando visualmente a assimetria dos dados:")
3 skew_info = analyzer.analyze_skewness(plot=True)
4 print("Informações de Skewness:")
5 for feature, skewness in skew_info.items():
6     print(f"{feature}: Skewness = {skewness:.2f}")
```

Analizando visualmente a assimetria dos dados:



Informações de Skewness:
Age: Skewness = 1.59
Height: Skewness = 0.02
Weight: Skewness = 0.09
FCVC: Skewness = -0.36
NCP: Skewness = -1.56
CH2O: Skewness = -0.21
FAF: Skewness = 0.51
TUE: Skewness = 0.67

O `skew_threshold` é um parâmetro que define o limite para considerar uma distribuição como significativamente assimétrica ou normal.

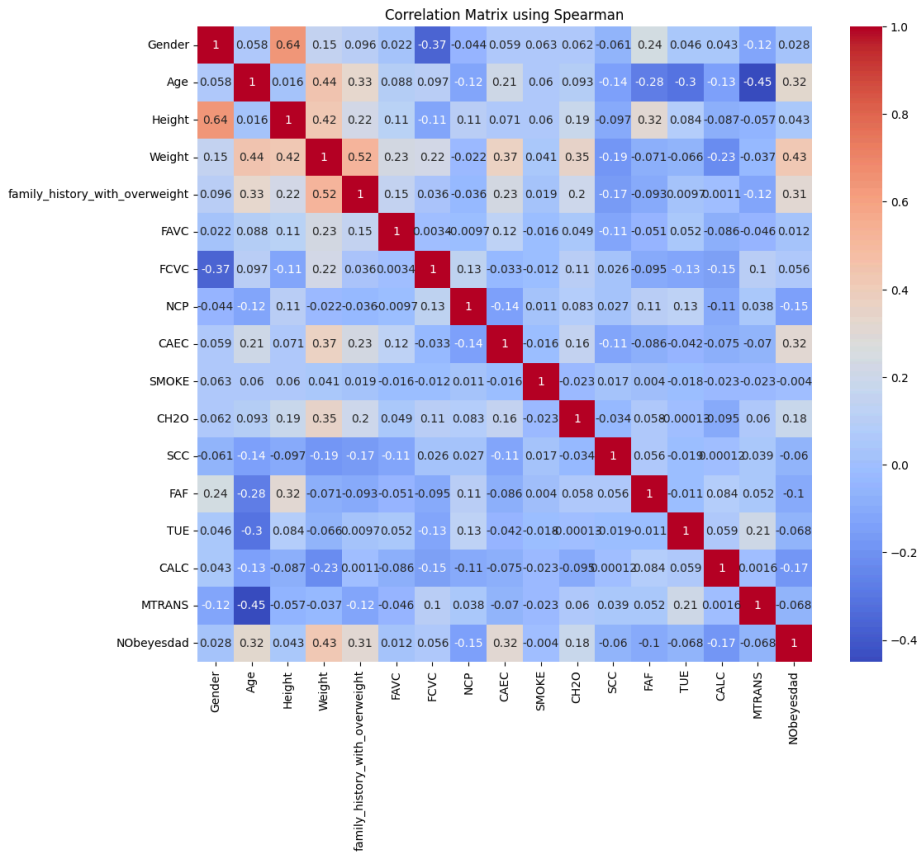
Se o skewness é um valor:

- Próximo a 0 indica uma distribuição simétrica (normal).
- Maior que 0 indica uma assimetria positiva (cauda longa à direita).
- Menor que 0 indica uma assimetria negativa (cauda longa à esquerda).
- Valor maior que 1 sugere uma distribuição altamente assimétrica.

Somente Height e Weight possuem um valor de skill mais ou menos simétrico.

```
1 #Realizar análise de correlação e imprimir a matriz de correlação
2 print("\nHeatmap que mostra a correlação dentre todas as features:")
3 correlation_matrix = analyzer.correlation_analysis_heatmap()
```

Heatmap que mostra a correlação dentre todas as features:



- Como era de se esperar, a variável weight é a variável que possui maior correlação com a variável alvo.
- CAEC : Food consumption between meals, também possui uma significativa correlação com a variável alvo
- Age and family_history_with_overweight é altamente correlacionada com a variável alvo.
- Gender não tem uma correlação muito alta com a variável alvo, ao contrário do que o senso comum possa dizer.

Feature Engineering

Baseado nas informações coletadas pela análise exploratória de dados, algumas mudanças nas features se fazem necessárias.

- Em alguns casos, os seus tipos de dados precisam ser alterados.
- Em outros casos, algumas features precisa ter os seus valores normalizados.
- Algumas novas colunas precisam ser criadas.

```
1 train.head(2)
```

	id	Gender	Age	Height	Weight	family_history_with_overweight	FAVC	FC
0	0	Male	24.443011	1.699998	81.66995		yes	yes

Next steps: [View recommended plots](#)

```

1 #Função para aplicar transformações nos dados
2 def feature_engineering(data):
3     #Transformar todas as colunas numéricas para float, exceto 'Age'
4     numerical_cols = data.select_dtypes(include=['int64', 'float64']).columns
5
6     #Converte as colunas para os seus tipos adequados
7     data[numerical_cols] = data[numerical_cols].astype(float)
8     data['id'] = data['id'].astype(int)
9     data['Age'] = data['Age'].astype(int)
10    data['CH20'] = data['CH20'].astype(int)
11
12    #Criar nova coluna para Age em escalas de 10 em 10
13    data['Age_Scale'] = data['Age'] // 10 + 1
14
15    #Criar nova coluna para Weight em escalas de 20 em 20
16    data['Weight_Scale'] = data['Weight'] // 20 + 1
17
18    #Criar nova coluna para Height em escalas de 10 em 10 centímetros
19    #Encontrar o mínimo e o máximo da coluna 'Height'
20    min_height = data['Height'].min()
21    max_height = data['Height'].max()
22    #Arredondar para baixo o mínimo e para cima o máximo para o múltiplo de 0.1 mais próximo
23    min_height = np.floor(min_height * 10) / 10
24    max_height = np.ceil(max_height * 10) / 10
25    #Criar os bins a partir do mínimo ao máximo com incrementos de 0.1 metros
26    bins = np.arange(min_height, max_height + 0.1, 0.1) #Usar np.arange para garantir inclusão do último intervalo
27    #Usar pd.cut para categorizar os valores de 'Height' e obter os códigos dos bins
28    data['Height_bin'] = pd.cut(data['Height'], bins=bins, right=False, labels=False)
29    #Calcular o ponto médio de cada bin para atribuir à nova coluna 'Height_scale'
30    data['Height_scale'] = [(bins[i] + bins[i+1])/2 for i in data['Height_bin']]
31    #Dropar coluna Height_bin
32    data = data.drop(labels=["Height_bin"], axis=1)
33
34    #Calcular o Índice de Massa Corporal (BMI)
35    data['BMI'] = data['Weight'] / (data['Height'] ** 2)
36
37    #Arredondar todas as colunas float para 2 casas decimais
38    for column in data.select_dtypes(include=['float64']).columns:
39        data[column] = data[column].round(2)
40
41    #Transformar colunas binárias
42    binary_columns = ['Gender', 'family_history_with_overweight', 'FAVC', 'SMOKE', 'SCC']
43    binary_mapping = {'yes': 1, 'no': 0, 'true': 1, 'false': 0, 'male': 1, 'female': 0}
44    for column in binary_columns:
45        #Certificar-se de que a coluna é de tipo string e lida com valores ausentes
46        data[column] = data[column].astype(str).str.lower().replace(binary_mapping)
47
48    #Transformar colunas categóricas em números, considerando sua ordem hierárquica (na maioria das features)
49    data['CAEC'] = data['CAEC'].map({'no': 0, 'Sometimes': 1, 'Frequently': 2, 'Always': 3})
50    data['CALC'] = data['CALC'].map({'no': 0, 'Sometimes': 1, 'Frequently': 2})
51    data['MTRANS'] = data['MTRANS'].astype('category').cat.codes #Sem hierarquia específica, encoding ordinal
52
53    if "NObeyesdad" in data.columns:
54        data['NObeyesdad'] = data['NObeyesdad'].map({
55            'Insufficient_Weight': 0,
56            'Normal_Weight': 1,
57            'Overweight_Level_I': 2,
58            'Overweight_Level_II': 3,
59            'Obesity_Type_I': 4,
60            'Obesity_Type_II': 5,
61            'Obesity_Type_III': 6
62        })
63
64    return data
65
66

```

```
1 # Aplicar a função de feature engineering
2 df_transformed = train.copy()
3 df_transformed = feature_engineering(df_transformed)
4 df_transformed = df_transformed.drop(labels=["id"], axis=1)
```

```
1 df_transformed.head(2)
```

	Gender	Age	Height	Weight	family_history_with_overweight	FAVC	FCVC	NCP	CAE
0	1	24	1.70	81.67	1	1	2.0	2.98	
1	0	18	1.56	57.00	1	1	2.0	3.00	

Encontrando as features mais importantes

Um processo será rodado abaixo com 2 objetivos:

- Encontrar a quantidade ideal de features para criar um modelo de classificação
- Encontrar quais são as features mais importantes para criar o modelo de classificação

```

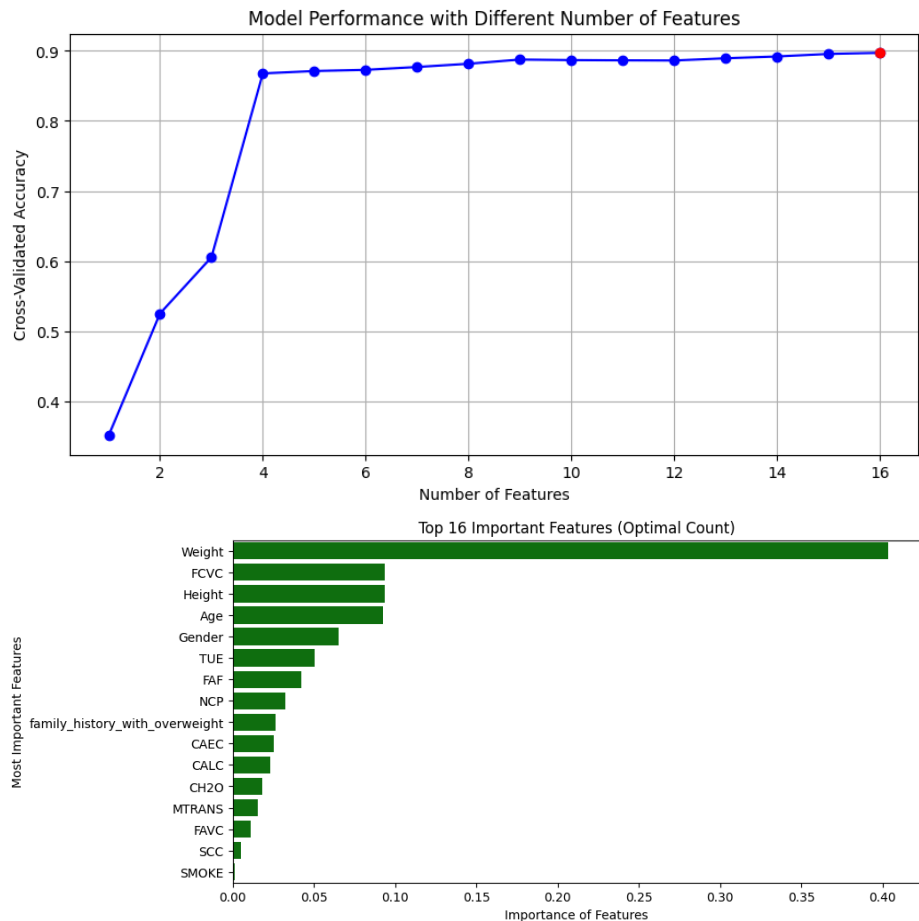
1 '''
2 Determina as características ótimas removendo características altamente correlacionadas e
3 avaliando o desempenho com um número crescente de características.
4 '''
5 #Plota a performance do modelo dado números diferentes de features
6 def plot_feature_performance(max_features, scores, best_feature_count, best_score):
7     plt.figure(figsize=(10, 5))
8     plt.plot(range(1, max_features + 1), scores, marker='o', linestyle='-', color='blue')
9     plt.plot(best_feature_count, best_score, 'ro')
10    plt.title('Model Performance with Different Number of Features')
11    plt.xlabel('Number of Features')
12    plt.ylabel('Cross-Validated Accuracy')
13    plt.grid(True)
14    plt.show()
15
16 #Plota a importância das features selecionados após identificar o número ideal de recursos
17 def plot_feature_importance(X, best_features, y, num_features):
18     model = RandomForestClassifier(n_estimators=100, random_state=42)
19     model.fit(X[best_features], y)
20     feature_importances = pd.Series(model.feature_importances_, index=best_features).sort_values(ascending=False)
21
22     plt.figure(figsize=(10, 5))
23     sns.barplot(x=feature_importances.values, y=feature_importances.index, color='green')
24     plt.title(f'Top {num_features} Important Features (Optimal Count)')
25     plt.xlabel('Importance of Features')
26     plt.ylabel('Most Important Features')
27     plt.show()
28
29 def find_optimal_features(data, target_col, max_features=None, corr_threshold=0.85):
30     #Separando as variáveis independentes da variável alvo
31     X = data.drop(target_col, axis=1)
32     y = data[target_col]
33
34     #Normalizando os dados
35     scaler = StandardScaler()
36     X_scaled = scaler.fit_transform(X)
37     X_scaled_df = pd.DataFrame(X_scaled, columns=X.columns)
38
39     #Calculando a matriz de correlação e removendo features altamente correlacionadas
40     corr_matrix = X_scaled_df.corr().abs()
41     upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k=1).astype(bool))
42     to_drop = [column for column in upper.columns if any(upper[column] > corr_threshold)]
43     X_reduced = X_scaled_df.drop(columns=to_drop)
44
45     if max_features is None or max_features > len(X_reduced.columns):
46         max_features = len(X_reduced.columns)
47
48     model = RandomForestClassifier(n_estimators=100, random_state=42)
49     scores = []
50
51     # Modelos de treinamento com um número crescente de features importantes
52     for i in range(1, max_features + 1):
53         selected_features = X_reduced.columns[:i]
54         score = np.mean(cross_val_score(model, X_reduced[selected_features], y, cv=5))
55         scores.append(score)
56
57     #Identifica o número otimizado de features
58     best_feature_count = np.argmax(scores) + 1
59     best_score = max(scores)
60     best_features = X_reduced.columns[:best_feature_count]
61
62     #Plotagem a performance do modelo e a importância das features
63     plot_feature_performance(max_features, scores, best_feature_count, best_score)
64     plot_feature_importance(X_reduced, best_features, y, best_feature_count)
65
66     return list(best_features)
67
68
69
70
71 def evaluate_models_with_gridsearch(data, target_col):
72     best_features = find_optimal_features(data, target_col)
73     X_best = data[best_features]
74     y = data[target_col]
75
76     #Normaliza os valores das features.
77     scaler = StandardScaler()
78     X_scaled = scaler.fit_transform(X_best)
79
80     #Gridsearch para encontrar o melhor modelo
81     param_grid = {
82         'n_estimators': [50, 100, 200],
83         'max_depth': [None, 10, 20, 30],
84         'min_samples_split': [2, 5, 10],
85         'min_samples_leaf': [1, 2, 5]
86     }
87     model = RandomForestClassifier()
88     grid_search = GridSearchCV(model, param_grid, cv=5)
89     grid_search.fit(X_scaled, y)
90     best_model = grid_search.best_estimator_
91     best_score = grid_search.best_score_
92
93     #Printando o melhor modelo encontrado
94     print(f"Melhor modelo encontrado com {best_model.n_estimators} estimadores e profundidade máxima de {best_model.max_depth}."
95           f"O melhor score de validação cruzada é {best_score} com as seguintes features: {best_features}")
96
97     return best_model, best_score, best_features

```

```

8     x_best_scaled = scaler.fit_transform(x_best)
9
10    models = {
11        'Naive Bayes': (GaussianNB(), {}),
12        'SVM': (SVC(probability=True), {'C': [0.1, 1, 10], 'gamma': ['scale', 'auto'], 'kernel': ['linear', 'rbf']}),
13        'Random Forest': (RandomForestClassifier(), {'n_estimators': [100, 200, 300], 'max_depth': [10, 20, 30, None]}),
14        'XGBoost': (xgb.XGBClassifier(use_label_encoder=False, eval_metric='mlogloss'), {'n_estimators': [100, 200], 'max_de
15        'LightGBM': (lgb.LGBMClassifier(), {'n_estimators': [100, 200], 'num_leaves': [31, 64, 128], 'learning_rate': [0.1,
16    }
17
18    results = []
19    for name, (model, params) in tqdm(models.items(), desc="Training Models"):
20        grid_search = GridSearchCV(model, params, cv=5, scoring='accuracy', verbose=2, n_jobs=-1)
21        grid_search.fit(X_best_scaled, y)
22        best_model = grid_search.best_estimator_
23        best_params = grid_search.best_params_
24        scores = cross_val_score(best_model, X_best_scaled, y, cv=5, scoring='accuracy')
25        accuracy = np.mean(scores)
26        results.append({'Model': name, 'Best Accuracy': accuracy, 'Best Parameters': best_params})
27
28    results_df = pd.DataFrame(results)
29    print(results_df)
30
31    # Plot the model performances for visual comparison
32    plt.figure(figsize=(10, 6))
33    sns.barplot(x='Model', y='Best Accuracy', data=results_df, palette='viridis')
34    plt.title('Performance of Various Models')
35    plt.ylabel('Accuracy')
36    plt.show()
37
38    return results_df
39
40
41 results_df = evaluate_models_with_gridsearch(df_transformed, 'NObeyesdad')

```

```
Training Models: 0%|          | 0/5 [00:00<?, ?it/s]Fitting 5 folds for each of 1
Training Models: 20%|██        | 1/5 [00:02<00:08, 2.11s/it]Fitting 5 folds for ea
Training Models: 40%|██████    | 2/5 [37:10<1:05:34, 1311.55s/it]Fitting 5 folds fo
Training Models: 60%|██████████| 3/5 [41:03<27:18, 819.19s/it] Fitting 5 folds fo
Training Models: 80%|██████████| 4/5 [44:07<09:28, 568.45s/it]Fitting 5 folds for e
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing wa
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 1064
[LightGBM] [Info] Number of data points in the train set: 20758, number of used feat
[LightGBM] [Info] Start training from score -2.107483
[LightGBM] [Info] Start training from score -1.907353
[LightGBM] [Info] Start training from score -2.146276
[LightGBM] [Info] Start training from score -2.107879
[LightGBM] [Info] Start training from score -1.964779
[LightGBM] [Info] Start training from score -1.854892
[LightGBM] [Info] Start training from score -1.635203
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing wa
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 1055
[LightGBM] [Info] Number of data points in the train set: 16606, number of used feat
[LightGBM] [Info] Start training from score -2.107657
[LightGBM] [Info] Start training from score -1.907167
[LightGBM] [Info] Start training from score -2.146046
[LightGBM] [Info] Start training from score -2.108153
[LightGBM] [Info] Start training from score -1.964755
[LightGBM] [Info] Start training from score -1.855022
[LightGBM] [Info] Start training from score -1.635117
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing wa
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 1061
[LightGBM] [Info] Number of data points in the train set: 16606, number of used feat
[LightGBM] [Info] Start training from score -2.107657
[LightGBM] [Info] Start training from score -1.907167
[LightGBM] [Info] Start training from score -2.146561
[LightGBM] [Info] Start training from score -2.108153
[LightGBM] [Info] Start training from score -1.964755
```