

AC2a

Nome: Thiago Faro Ribeiro

RA: 235553

Nome: Giovanna Andreoli

RA: 235831

Predição de Câncer

Informações do atributo:

a) ID Número de identificação

b) Diagnóstico (M = maligno, B = benigno) (a variável para predefinir)

Lista de exercícios:

1- Ler o dataset "Cancer_diagnostic.csv" e seguir os seguintes passos.

```
file_path = 'Cancer_diagnostic.csv'
data = pd.read_csv(file_path)
```

2- Qual é o tamanho do dataset?

```
print(f"Tamanho do dataset: {data.shape}")
```

✓ 0.0s

Tamanho do dataset: (569, 33)

3- Verifique colunas e tipos de dados

```
data.info()
```

✓ 0.0s

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 569 entries, 0 to 568
```

```
Data columns (total 33 columns):
```

#	Column	Non-Null Count	Dtype
0	id	569 non-null	int64
1	diagnosis	569 non-null	object
2	radius_mean	569 non-null	float64
3	texture_mean	569 non-null	float64
4	perimeter_mean	569 non-null	float64
5	area_mean	569 non-null	float64
6	smoothness_mean	569 non-null	float64
7	compactness_mean	569 non-null	float64
8	concavity_mean	569 non-null	float64
9	concave points_mean	569 non-null	float64
10	symmetry_mean	569 non-null	float64
11	fractal_dimension_mean	569 non-null	float64
12	radius_se	569 non-null	float64
13	texture_se	569 non-null	float64
14	perimeter_se	569 non-null	float64
15	area_se	569 non-null	float64
16	smoothness_se	569 non-null	float64
17	compactness_se	569 non-null	float64
18	concavity_se	569 non-null	float64
19	concave points_se	569 non-null	float64
...			
31	fractal_dimension_worst	569 non-null	float64
32	Unnamed: 32	0 non-null	float64

```
dtypes: float64(31), int64(1), object(1)
```

4- Mostra o resumo estatístico do dataset

```
data.describe()
```

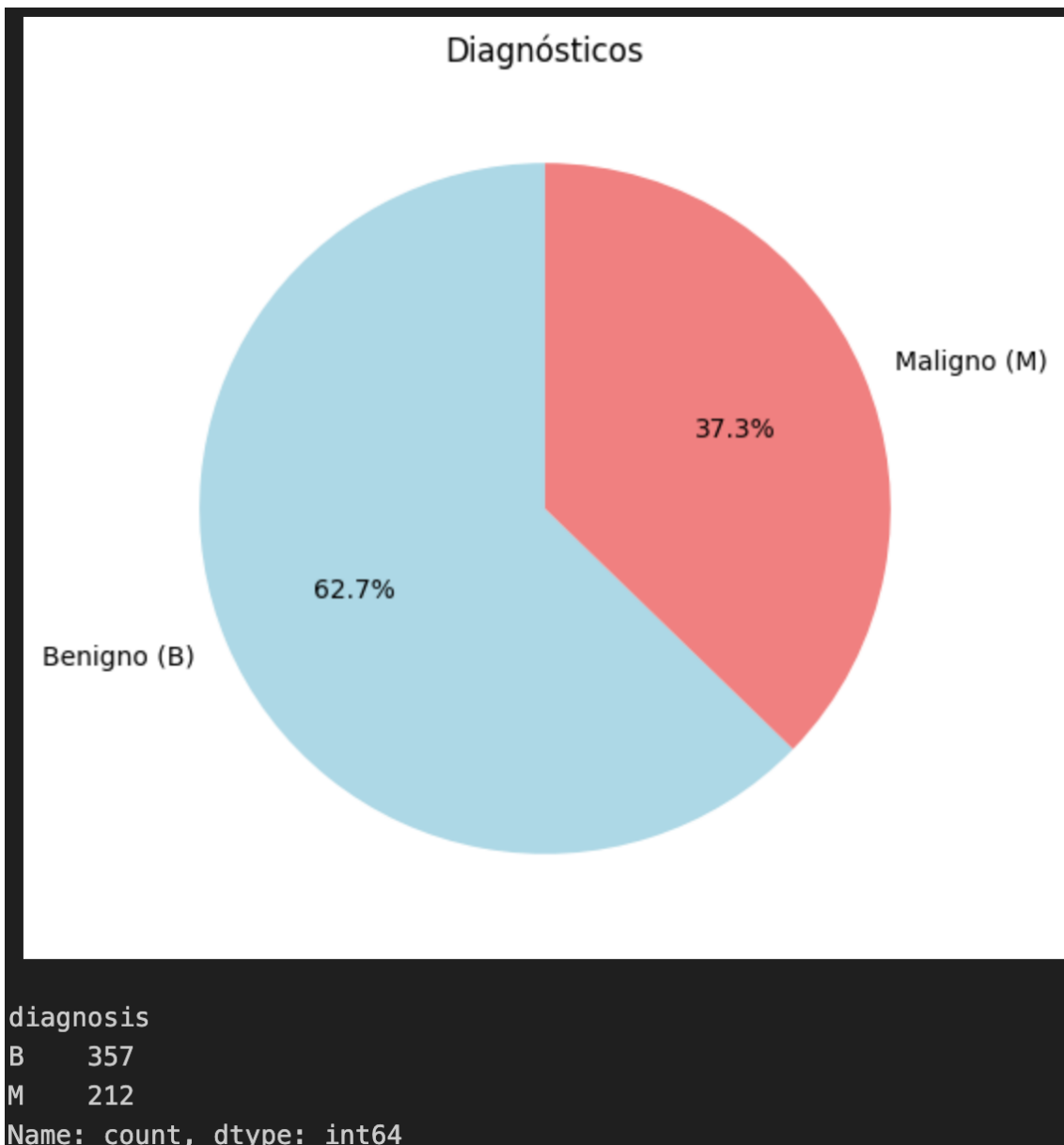
	id	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave_points_mean	symmetry_mean	...	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	concave_points_worst	symmetry_worst	fractal_dimension_worst	Unnamed: 32
count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	...	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	5.0
mean	3.027184e+02	14.127220	19.248454	91.869313	664.488414	0.049380	0.104341	0.048789	0.248879	0.181102	...	25.477223	107.281219	660.548328	0.132288	0.264385	0.227783	0.148500	0.390576	0.203848	NaN
std	1.250206e+08	0.524049	4.301036	24.288881	351.874129	0.014054	0.052813	0.078720	0.038863	0.027414	...	6.146258	33.802542	568.358883	0.022832	0.167136	0.208614	0.048732	0.041887	0.014081	NaN
min	8.470000e+01	6.863000	9.710000	43.790000	143.800000	0.002000	0.018000	0.000000	0.000000	0.100000	...	12.000000	85.410000	184.200000	0.097000	0.027596	0.000000	0.000000	0.000000	0.000000	NaN
25%	6.692300e+01	11.700000	16.170000	75.170000	400.300000	0.003000	0.044800	0.020000	0.020000	0.100000	...	17.000000	84.100000	195.300000	0.100000	0.147200	0.150000	0.004900	0.020000	0.027400	NaN
50%	0.000000e+00	13.710000	18.840000	86.240000	551.100000	0.008000	0.050000	0.000000	0.030000	0.100000	...	25.410000	97.680000	666.500000	0.131000	0.211000	0.228700	0.099300	0.282200	0.080400	NaN
75%	8.810200e+01	15.780000	21.800000	104.100000	782.700000	0.016000	0.104000	0.100000	0.070000	0.100000	...	33.700000	125.800000	1084.000000	0.148000	0.338000	0.388000	0.181600	0.377000	0.280000	NaN
max	9.113200e+02	26.710000	39.280000	148.800000	2507.000000	0.163400	0.345400	0.428800	0.201000	0.304000	...	46.640000	311.200000	4264.000000	0.222800	1.048000	1.252000	0.291000	0.618000	0.207000	NaN

5 rows x 22 columns

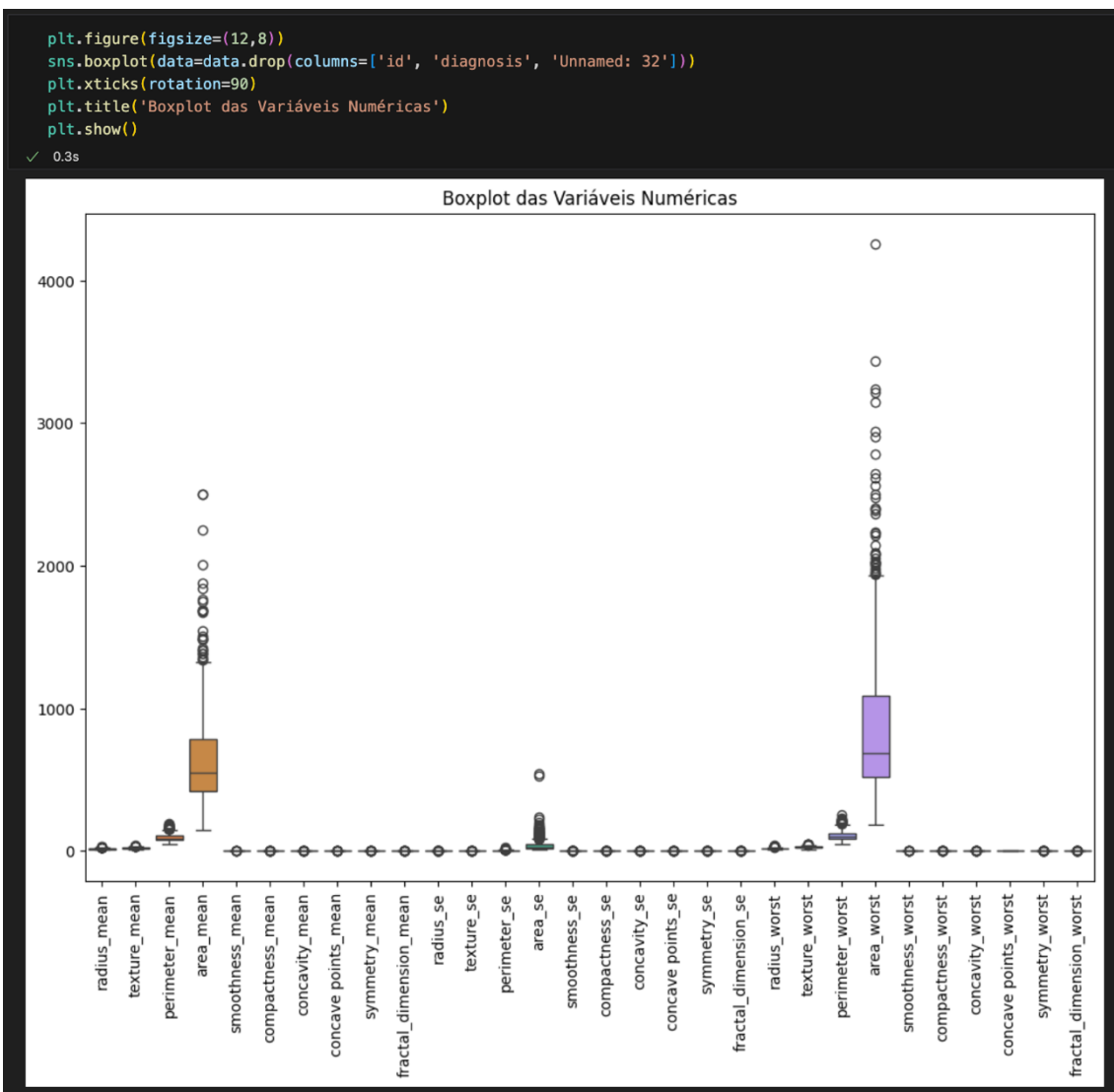
5- O dataset está balanceado?

O dataset não está balanceado.

6- Mostrar um gráfico de pizza do dataset



7- Mostra o gráfico de boxplot do dataset



8- Tem valores ausentes? Se sim quais?

```
missing_values = data.isnull().sum()
print("Valores ausentes:")
print(missing_values[missing_values > 0])
```

✓ 0.0s

Valores ausentes:
Unnamed: 32 569
dtype: int64

9- Realiza a padronização (standardization)

Clique para adicionar um ponto de interrupção

'Unnamed: 32', 'id']]

```
# Separando os dados em X (features) e y (alvo)
X = data_cleaned.drop(columns=['diagnosis'])
y = data_cleaned['diagnosis']

# Convertendo o rótulo de diagnóstico para valores binários (M = 1, B = 0)
y = y.map({'M': 1, 'B': 0})

# Padronizando os dados
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Exibindo as primeiras linhas dos dados padronizados
pd.DataFrame(X_scaled, columns=X.columns).head()
```

✓ 0.0s

radius_mean

texture_mean

perimeter_mean

area_mean

smoothness_mean

compactness_mean

concavity_mean

concave points_mean

symmetry_mean

fractal_dimension_mean

radius_worst

texture_worst

perimeter_worst

area_worst

smoothness_worst

compactness_worst

concavity_worst

concave points_worst

symmetry_worst

fractal_dimension_worst

id

diagnosis

0	17.09544	-0.271591	2.064914	0.844470	1.544648	0.266910	2.652474	0.121740	2.215110	2.201547	1.860890	1.016293	2.201001	2.201237	1.567166	2.404604	2.105158	2.299076	2.709122	1.945110	1	M
1	1.623821	-0.301802	1.688961	1.008709	-0.828962	-0.487072	-0.023848	0.548154	0.001380	-0.868852	1.809307	-0.368203	1.531128	1.880489	-0.375613	-0.430444	-0.166749	1.097084	-0.243890	0.281980	1	M
2	1.573688	0.486387	1.666203	1.516884	0.842270	1.002828	1.301479	2.037231	0.839685	-0.388208	1.517870	-0.022874	1.547470	1.468285	0.227607	1.082932	0.854874	1.950500	1.102255	0.201381	1	M
3	-0.766809	0.320332	-0.582887	-0.361444	2.281013	0.402829	1.891687	1.481707	3.867383	-0.971019	-0.281464	0.328384	-0.248926	-0.550051	3.246170	3.881380	1.895448	1.705186	6.546461	0.826010	1	M
4	1.760287	-1.181816	1.778573	1.826229	0.280372	0.838340	1.377011	1.428483	-0.008560	-0.861450	1.288576	-1.468770	1.388139	1.220724	0.220668	-0.310369	0.813179	0.792259	-0.868363	-0.387700	1	M

5 rows x 20 columns

10- Se tiver outlier fazer a remoção deles e mostrar o boxplot do dataset.

```

# Calculando os quartis e o IQR
Q1 = data.drop(columns=['id', 'diagnosis', 'Unnamed: 32']).quantile(0.25)
Q3 = data.drop(columns=['id', 'diagnosis', 'Unnamed: 32']).quantile(0.75)
IQR = Q3 - Q1

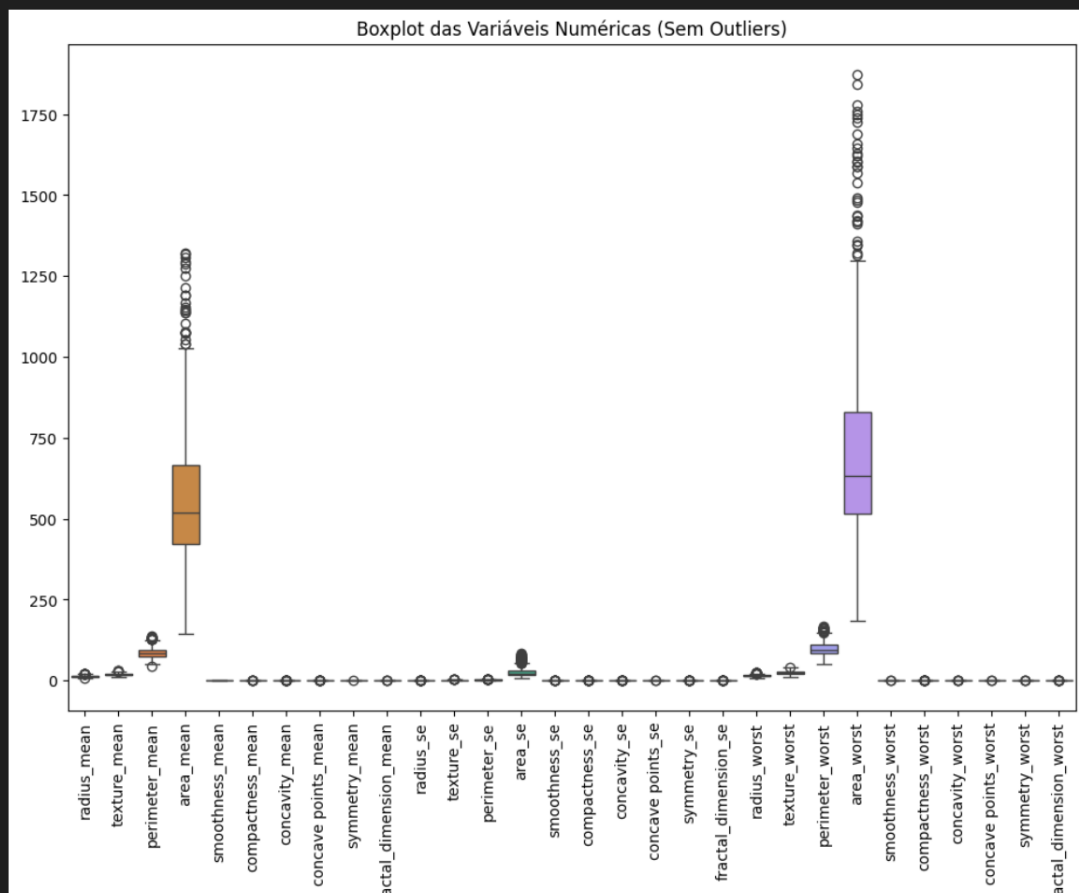
# Definindo os limites
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Removendo os outliers
data_no_outliers = data.drop(columns=['id', 'diagnosis', 'Unnamed: 32']).mask(
    (data.drop(columns=['id', 'diagnosis', 'Unnamed: 32']) < lower_bound) |
    (data.drop(columns=['id', 'diagnosis', 'Unnamed: 32']) > upper_bound)
).dropna()

# Exibindo o novo boxplot
plt.figure(figsize=(12,8))
sns.boxplot(data=data_no_outliers)
plt.xticks(rotation=90)
plt.title('Boxplot das Variáveis Numéricas (Sem Outliers)')
plt.show()

```

✓ 0.3s



11- Dividir o dataset em training (80%) e testing (20%)

```

X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)

```

✓ 0.0s

12- Faz a classificação usando o algoritmo KNN com K = 3, K = 5, K = 7 e mostrar as métricas de accuracy, recall e precision para cada K.

```
# Função para treinar e avaliar o modelo KNN com diferentes valores de K
def evaluate_knn(k):
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    y_pred = knn.predict(X_test)

    accuracy = accuracy_score(y_test, y_pred)
    recall = recall_score(y_test, y_pred)
    precision = precision_score(y_test, y_pred)

    print(f"K={k}: Accuracy={accuracy:.4f}, Recall={recall:.4f}, Precision={precision:.4f}")

# Avaliando o modelo KNN com K=3, K=5, e K=7
for k in [3, 5, 7]:
    evaluate_knn(k)
```

✓ 0.0s

K=3: Accuracy=0.9474, Recall=0.9302, Precision=0.9302
K=5: Accuracy=0.9474, Recall=0.9302, Precision=0.9302
K=7: Accuracy=0.9474, Recall=0.9302, Precision=0.9302