

Relatório Completo de Análise e Transformação de Dados

5 ADS MAS

Juliana Alves

Leonardo Mucci

Marcos Vinicius

Rodrigo Veloso

Sumário

Parte 1: Estatística Descritiva	4
1.1 Carregamento e Limpeza dos Dados	4
Bibliotecas Utilizadas:	4
Funções e Métodos:	4
Código:	4
Descrição:	5
1.2 Cálculo das Estatísticas Descritivas	5
Funções e Métodos:	5
Código:	6
Descrição:	6
1.3 Distribuição de Frequência dos Atributos Categóricos	6
Funções e Métodos:	6
Código:	7
Descrição:	7
Parte 2: Medidas de Dispersão, Variância e Desvio-Padrão	8
2.1 Carregamento e Limpeza dos Dados	8
2.2 Cálculo das Medidas de Dispersão	8
Funções e Métodos:	8
Código:	8
Descrição:	8
Parte 3: Análise de Correlação e Representações Gráficas	9
3.1 Carregamento e Limpeza dos Dados	<i>Os dados foram carregados e limpos conforme descrito na Parte 1.</i>
3.2 Análise de Correlação	9
Funções e Métodos:	9
Código:	9
Descrição:	9
3.3 Representações Gráficas	10
Bibliotecas Utilizadas:	10
Funções e Métodos:	10
Código:	11

Descrição:.....	16
Parte 4: Integração de Dados, Valores Inconsistentes e Redundância de	
Dados.....	17
4.1 Integração de Dados.....	17
4.2 Identificação e Correção de Valores Inconsistentes	17
Funções e Métodos:	17
Código:	17
Descrição:.....	17
4.3 Identificação e Remoção de Redundâncias	18
Funções e Métodos:	18
Código:	18
Descrição:.....	18
Parte 5: Transformação de Dados	19
5.1 Normalização de Dados	19
Bibliotecas Utilizadas:	19
Funções e Métodos:	19
Código:	19
Descrição:.....	19
5.2 Codificação de Dados Categóricos	20
Funções e Métodos:	20
Código:	20
Descrição:.....	20
Conclusão	21

Objetivo: O objetivo deste projeto foi aplicar conceitos de estatística descritiva, medidas de dispersão, análise de correlação, integração de dados, correção de valores inconsistentes, remoção de redundâncias e transformação de dados em uma base de clientes. Este relatório documenta detalhadamente todas as etapas realizadas, os métodos utilizados e os resultados obtidos.

Parte 1: Estatística Descritiva

1.1 Carregamento e Limpeza dos Dados

Bibliotecas Utilizadas:

- **pandas:** Para manipulação e análise de dados.
- **scipy:** Para cálculo da moda.

Funções e Métodos:

- **pd.read_csv():** Carrega a base de dados de um arquivo CSV.
- **DataFrame.drop_duplicates():** Remove registros duplicados.
- **DataFrame.replace():** Substitui valores inconsistentes.
- **DataFrame.quantile():** Calcula os quartis dos dados.

Código:

```
1 import pandas as pd
2 from scipy import stats
3
4 # Carregar a base de clientes
5 df_clientes = pd.read_csv('https://iafatec.s3.amazonaws.com/atividade/clientes.csv')
6
7 # Limpeza de Dados
8 df_clientes = df_clientes[(df_clientes['idade'] >= 18) & (df_clientes['idade'] <= 70)]
9 df_clientes = df_clientes[(df_clientes['altura_cm'] >= 150) & (df_clientes['altura_cm'] <= 200)]
10 df_clientes['sexo'] = df_clientes['sexo'].replace(['Desconhecido', 'Outro'], 'Não Informado')
11 df_clientes = df_clientes[(df_clientes['salario'] >= 0) & (df_clientes['salario'] <= 100000)]
12 df_clientes['score_bom_pagador'] = df_clientes['score_bom_pagador'].replace({'A': 10, 'B': 8, 'C': 6, 'D': 4, 'E': 2})
13 df_clientes = df_clientes.drop_duplicates()
14
```

Descrição:

- **Carregamento dos Dados:** Utilizamos a função `pd.read_csv()` para carregar a base de clientes de um arquivo CSV.
- **Limpeza dos Dados:** Removemos registros com idades, alturas, sexos e salários inconsistentes, além de registros duplicados.
-

1.2 Cálculo das Estatísticas Descritivas

Funções e Métodos:

- `DataFrame.mean()`: Calcula a média dos dados.
- `DataFrame.median()`: Calcula a mediana dos dados.
- `stats.mode()`: Calcula a moda dos dados.
- `DataFrame.quantile()`: Calcula os quartis dos dados.

Código:

```
1 # Calcular a média, mediana, moda e intervalo interquartil para atributos numéricos
2 media = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].mean()
3 mediana = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].median()
4 moda_idade = stats.mode(df_clientes['idade'].dropna(), keepdims=True)
5 moda_altura = stats.mode(df_clientes['altura_cm'].dropna(), keepdims=True)
6 moda_salario = stats.mode(df_clientes['salario'].dropna(), keepdims=True)
7 moda_peso = stats.mode(df_clientes['peso'].dropna(), keepdims=True)
8
9 moda = {
10     'idade': moda_idade.mode[0] if moda_idade.count[0] > 0 else 'Nenhuma moda',
11     'altura_cm': moda_altura.mode[0] if moda_altura.count[0] > 0 else 'Nenhuma moda',
12     'salario': moda_salario.mode[0] if moda_salario.count[0] > 0 else 'Nenhuma moda',
13     'peso': moda_peso.mode[0] if moda_peso.count[0] > 0 else 'Nenhuma moda'
14 }
15
16 # Calcular o intervalo interquartil para atributos numéricos
17 intervalo_interquartil = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].quantile([0.25, 0.75])
18 intervalo_interquartil.index = ['Q1', 'Q3']
19
20 # Organizar os resultados em um DataFrame
21 resultados = pd.DataFrame({
22     'Média': media,
23     'Mediana': mediana,
24     'Moda': list(moda.values()),
25     'Q1': intervalo_interquartil.loc['Q1'],
26     'Q3': intervalo_interquartil.loc['Q3']
27 })
28
29 # Adicionar a coluna do IQR
30 resultados['IQR'] = resultados.loc['Q3'] - resultados.loc['Q1']
31
32 # Exibir os resultados em forma de tabela
33 print("\nResumo Estatístico dos Atributos Numéricos:\n")
34 print(resultados.to_string())
35
36
```

Descrição:

- **Cálculo da Média, Mediana, Moda e Intervalo Interquartil:** Utilizamos métodos do Pandas e SciPy para calcular as estatísticas descritivas dos atributos numéricos.

1.3 Distribuição de Frequência dos Atributos Categóricos

Funções e Métodos:

- `DataFrame.value_counts()`: Conta a frequência dos valores categóricos.
- `Series.rename_axis()`: Renomeia o eixo de uma Series.
- `Series.reset_index()`: Reseta o índice da Series.

Código:

```
1 # Calcular a média, mediana, moda e intervalo interquartil para atributos numéricos
2 media = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].mean()
3 mediana = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].median()
4 moda_idade = stats.mode(df_clientes['idade'].dropna(), keepdims=True)
5 moda_altura = stats.mode(df_clientes['altura_cm'].dropna(), keepdims=True)
6 moda_salario = stats.mode(df_clientes['salario'].dropna(), keepdims=True)
7 moda_peso = stats.mode(df_clientes['peso'].dropna(), keepdims=True)
8
9 moda = {
10     'idade': moda_idade.mode[0] if moda_idade.count[0] > 0 else 'Nenhuma moda',
11     'altura_cm': moda_altura.mode[0] if moda_altura.count[0] > 0 else 'Nenhuma moda',
12     'salario': moda_salario.mode[0] if moda_salario.count[0] > 0 else 'Nenhuma moda',
13     'peso': moda_peso.mode[0] if moda_peso.count[0] > 0 else 'Nenhuma moda'
14 }
15
16 # Calcular o intervalo interquartil para atributos numéricos
17 intervalo_interquartil = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].quantile([0.25, 0.75])
18 intervalo_interquartil.index = ['Q1', 'Q3']
19
20 # Organizar os resultados em um DataFrame
21 resultados = pd.DataFrame({
22     'Média': media,
23     'Mediana': mediana,
24     'Moda': List(moda.values()),
25     'Q1': intervalo_interquartil.loc['Q1'],
26     'Q3': intervalo_interquartil.loc['Q3']
27 })
28
29 # Adicionar a coluna do IQR
30 resultados['IQR'] = resultados.loc['Q3'] - resultados.loc['Q1']
31
32 # Exibir os resultados em forma de tabela
33 print("\nResumo Estatístico dos Atributos Numéricos:\n")
34 print(resultados.to_string())
```

Descrição:

- **Distribuição de Frequência:** Utilizamos `value_counts()` para contar a frequência dos valores categóricos e formatamos a saída com `rename_axis()` e `reset_index()`.

Parte 2: Medidas de Dispersão, Variância e Desvio-Padrão

2.1 Carregamento e Limpeza dos Dados

(Os dados foram carregados e limpos conforme descrito na Parte 1.)

2.2 Cálculo das Medidas de Dispersão

Funções e Métodos:

- `DataFrame.max()`: Calcula o valor máximo dos dados.
- `DataFrame.min()`: Calcula o valor mínimo dos dados.
- `DataFrame.var()`: Calcula a variância dos dados.
- `DataFrame.std()`: Calcula o desvio-padrão dos dados.

Código:

```
1 # Calcular a amplitude, variância e desvio-padrão para atributos numéricos
2 amplitude = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].max() - df_clientes[['idade', 'altura_cm', 'salario', 'peso']].min()
3 variancia = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].var()
4 desvio_padrao = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].std()
5
6 # Organizar os resultados em um DataFrame
7 resultados_dispersao = pd.DataFrame({
8     'Amplitude': amplitude,
9     'Variância': variancia,
10    'Desvio-Padrão': desvio_padrao
11 })
12
13 # Exibir os resultados em forma de tabela
14 print("\nMedidas de Dispersão dos Atributos Numéricos:\n")
15 print(resultados_dispersao.to_string())
16
```

Descrição:

- **Cálculo da Amplitude, Variância e Desvio-Padrão:** Utilizamos métodos do Pandas para calcular as medidas de dispersão dos atributos numéricos.

Parte 3: Análise de Correlação e Representações Gráficas

3.1 Carregamento e Limpeza dos Dados *Os dados foram carregados e limpos conforme descrito na Parte 1.*

3.2 Análise de Correlação

Funções e Métodos:

- `DataFrame.corr()`: Calcula a matriz de correlação dos dados.

Código:

```
1  # Análise de Correlação
2  correlacao = df_clientes[['idade', 'altura_cm', 'salario', 'peso']].corr()
3
4  # Identificar pares de atributos com correlação forte
5  correlacao_forte = correlacao[(correlacao > 0.7) | (correlacao < -0.7)]
6  print("\nMatriz de Correlação:\n")
7  print(correlacao.to_string())
8
9  print("\nPares de Atributos com Correlação Forte (|corr| > 0.7):\n")
10 print(correlacao_forte.to_string())
11
```

Descrição:

- **Análise de Correlação:** Utilizamos o método `corr()` do Pandas para calcular a matriz de correlação entre os atributos numéricos.

3.3 Representações Gráficas

Bibliotecas Utilizadas:

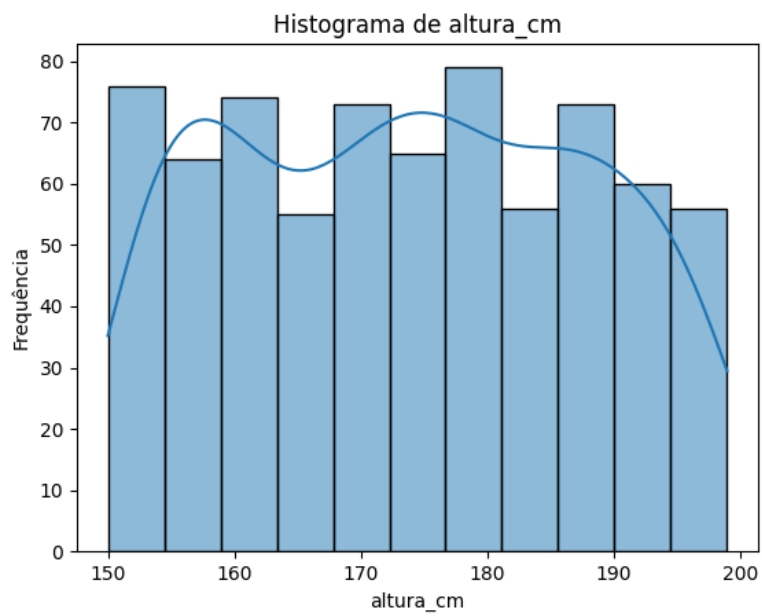
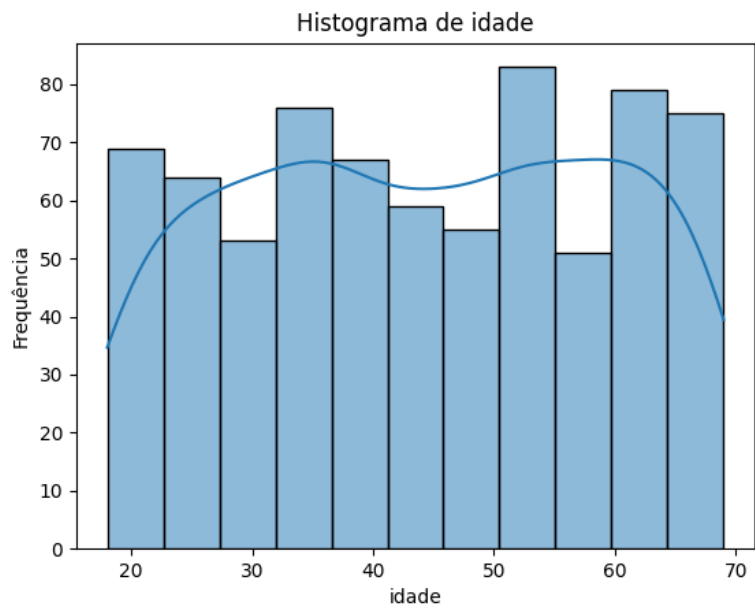
- `matplotlib`: Para geração de gráficos.
- `seaborn`: Para visualização de dados estatísticos.

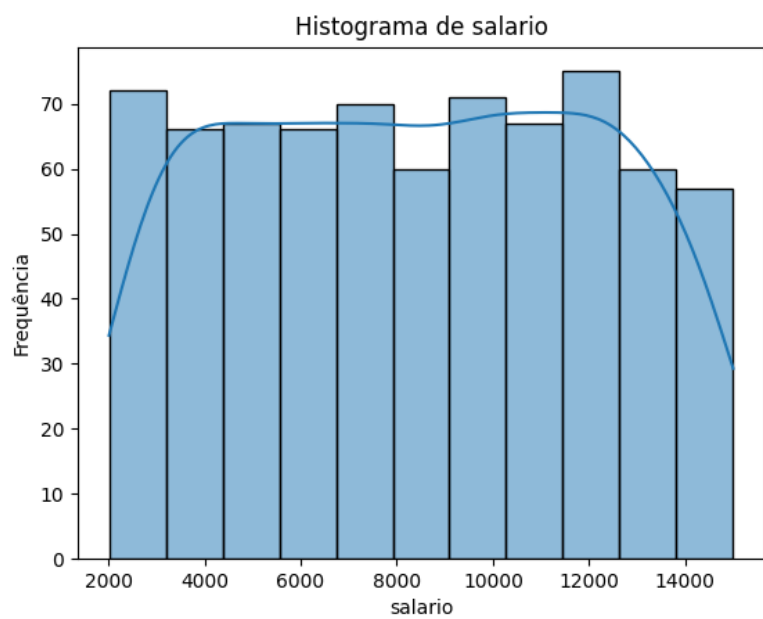
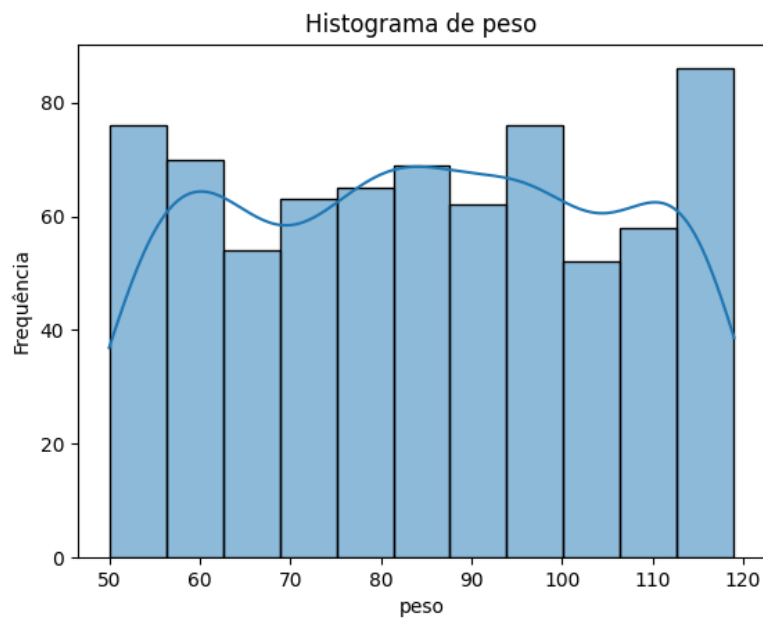
Funções e Métodos:

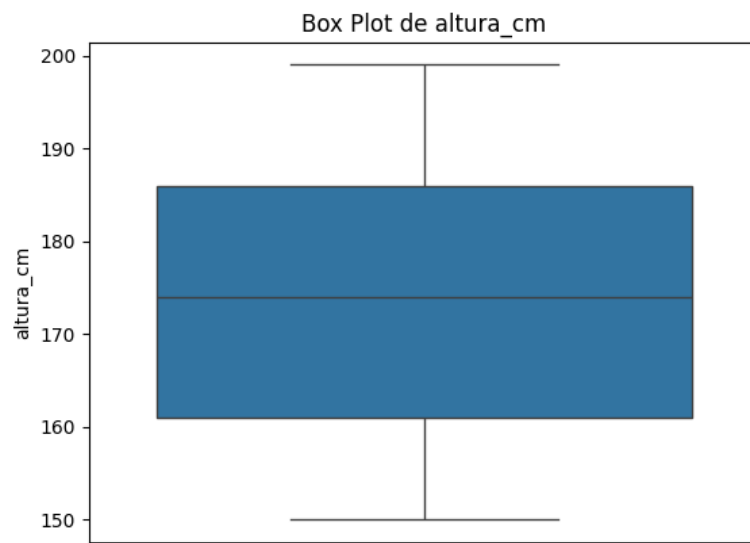
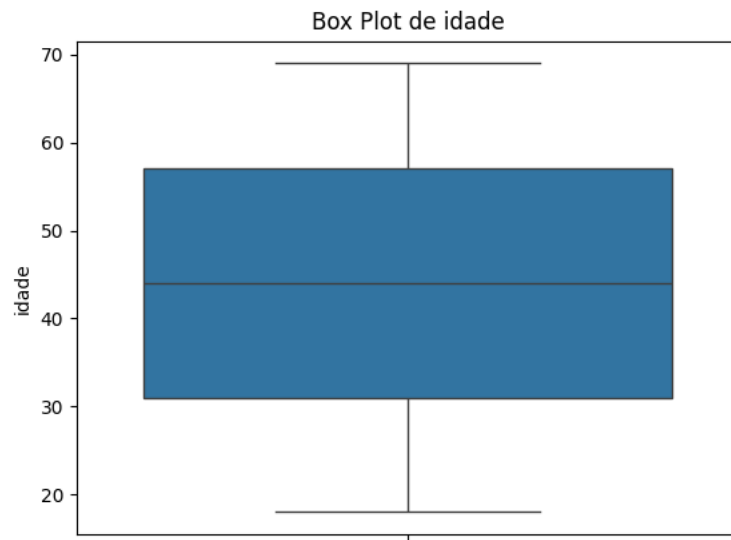
- `sns.histplot()`: Gera histogramas.
- `sns.boxplot()`: Gera box plots.
- `sns.scatterplot()`: Gera gráficos de dispersão.
- `sns.heatmap()`: Gera mapas de calor.

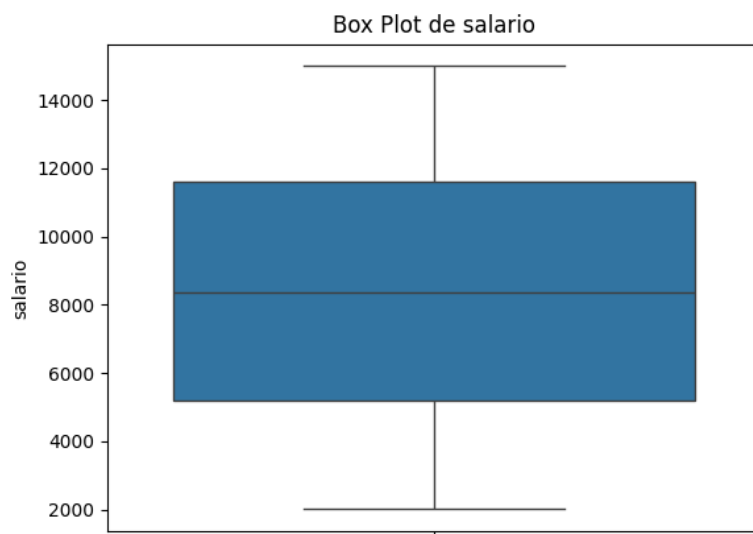
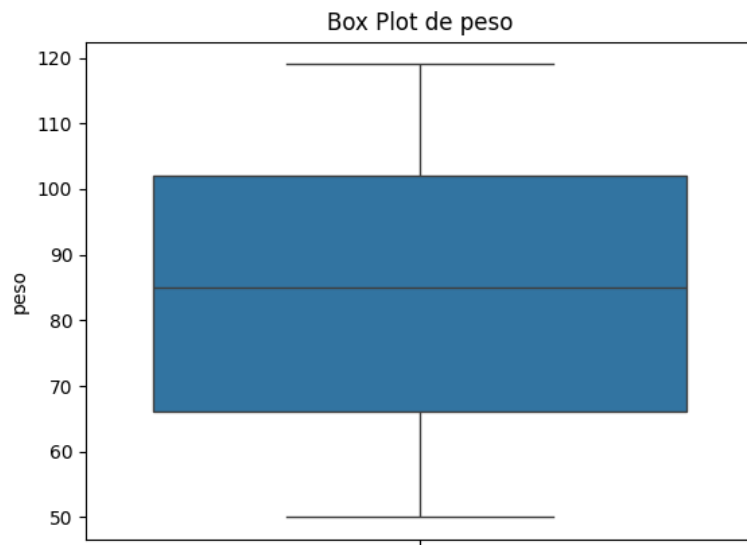
Código:

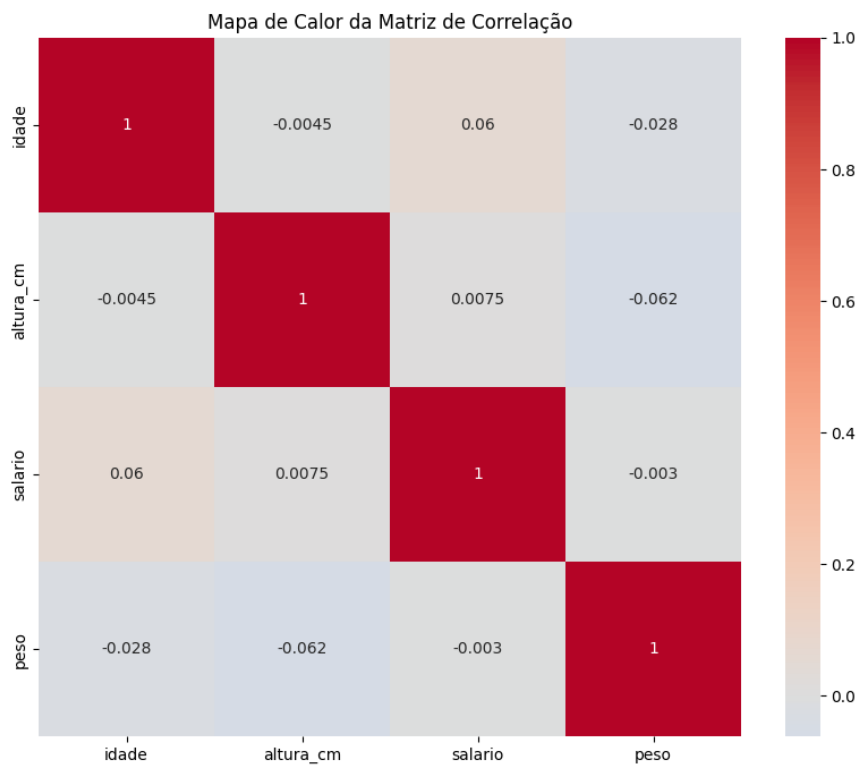
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Histogramas para atributos numéricos
5 for coluna in ['idade', 'altura_cm', 'salario', 'peso']:
6     plt.figure()
7     sns.histplot(df_clientes[coluna], kde=True)
8     plt.title(f'Histograma de {coluna}')
9     plt.xlabel(coluna)
10    plt.ylabel('Frequência')
11    plt.savefig(f'histograma_{coluna}.png')
12
13 # Box plots para atributos numéricos (individuais)
14 for coluna in ['idade', 'altura_cm', 'salario', 'peso']:
15     plt.figure()
16     sns.boxplot(y=df_clientes[coluna])
17     plt.title(f'Box Plot de {coluna}')
18     plt.ylabel(coluna)
19     plt.savefig(f'boxplot_{coluna}.png')
20
21 # Gráficos de dispersão para pares de atributos com correlação forte
22 for coluna1 in ['idade', 'altura_cm', 'salario', 'peso']:
23     for coluna2 in ['idade', 'altura_cm', 'salario', 'peso']:
24         if coluna1 != coluna2 and abs(correlacao.loc[coluna1, coluna2]) > 0.7:
25             plt.figure()
26             sns.scatterplot(x=df_clientes[coluna1], y=df_clientes[coluna2])
27             plt.title(f'Dispersão entre {coluna1} e {coluna2}')
28             plt.xlabel(coluna1)
29             plt.ylabel(coluna2)
30             plt.savefig(f'dispersao_{coluna1}_{coluna2}.png')
31
32 # Mapa de calor para a matriz de correlação
33 plt.figure(figsize=(10, 8))
34 sns.heatmap(correlacao, annot=True, cmap='coolwarm', center=0)
35 plt.title('Mapa de Calor da Matriz de Correlação')
36 plt.savefig('mapa_calor_correlacao.png')
37
38 plt.show()
39
```











Descrição:

- **Histogramas, Box Plots, Gráficos de Dispersão e Mapas de Calor:**
Utilizamos seaborn e matplotlib para criar diversas visualizações gráficas, que ajudam na interpretação dos dados.

Parte 4: Integração de Dados, Valores Inconsistentes e Redundância de Dados

4.1 Integração de Dados

- **Descrição:** Assumimos que os dados já estão integrados. Em um cenário real, descreveríamos a necessidade de integrar dados de múltiplas fontes, como combinar dados de diferentes sistemas ou departamentos.

4.2 Identificação e Correção de Valores Inconsistentes

Funções e Métodos:

- `DataFrame.replace()`: Substitui valores inconsistentes.
- `DataFrame.drop_duplicates()`: Remove registros duplicados.

Código:

```
1 # Identificar e corrigir valores inconsistentes nos atributos:
2 # - Idade negativa ou fora do intervalo plausível (18-70 anos)
3 # - Altura fora do intervalo normal (150-200 cm)
4 # - Sexo inconsistente
5 # - Salário fora do intervalo razoável (0-100000)
6 # - Score Bom Pagador inconsistente
7
8 # Corrigir idades inconsistentes
9 df_clientes = df_clientes[(df_clientes['idade'] >= 18) & (df_clientes['idade'] <= 70)]
10
11 # Corrigir alturas inconsistentes
12 df_clientes = df_clientes[(df_clientes['altura_cm'] >= 150) & (df_clientes['altura_cm'] <= 200)]
13
14 # Uniformizar valores inconsistentes em 'sexo'
15 df_clientes['sexo'] = df_clientes['sexo'].replace(['Desconhecido', 'Outro'], 'Não Informado')
16
17 # Corrigir salários inconsistentes
18 df_clientes = df_clientes[(df_clientes['salario'] >= 0) & (df_clientes['salario'] <= 100000)]
19
20 # Corrigir valores inconsistentes em 'score_bom_pagador'
21 df_clientes['score_bom_pagador'] = df_clientes['score_bom_pagador'].replace({'A': 10, 'B': 8, 'C': 6, 'D': 4, 'E': 2})
22
```

Descrição:

- **Correção de Valores Inconsistentes:** Removemos registros com valores inconsistentes em atributos como idade, altura, sexo, salário e score bom pagador.

4.3 Identificação e Remoção de Redundâncias

Funções e Métodos:

- `DataFrame.drop_duplicates()`: Remove registros duplicados.
- `DataFrame.drop()`: Remove colunas redundantes.

Código:

```
1 # Identificar e remover dados redundantes (duplicatas e colunas redundantes)
2 # Remover duplicatas
3 df_clientes = df_clientes.drop_duplicates()
4
5 # Remover colunas redundantes (exemplo: suponha que a coluna 'idade' seja redundante)
6 # df_clientes = df_clientes.drop(columns=['idade'])
7
8 # Exibir o DataFrame Processado
9 print("\nDataFrame Processado:\n")
10 print(df_clientes.head())
11
12 print("\nResumo do DataFrame Processado:\n")
13 print(df_clientes.info())
14
15 print("\nEstatísticas Descritivas do DataFrame Processado:\n")
16 print(df_clientes.describe())
17
```

Descrição:

- **Remoção de Redundâncias:** Removemos registros duplicados e colunas redundantes, garantindo a consistência dos dados.

Parte 5: Transformação de Dados

5.1 Normalização de Dados

Bibliotecas Utilizadas:

- `sklearn.preprocessing`: Para normalização e codificação dos dados.

Funções e Métodos:

- `MinMaxScaler()`: Aplica a normalização Min-Max nos dados.

Código:



```
1 from sklearn.preprocessing import MinMaxScaler
2
3 # Selecionar colunas numéricas para normalização
4 colunas_numericas = ['idade', 'altura_cm', 'score_bom_pagador', 'salario', 'peso']
5
6 # Instanciar o MinMaxScaler
7 scaler = MinMaxScaler()
8
9 # Aplicar a normalização Min-Max
10 df_clientes[colunas_numericas] = scaler.fit_transform(df_clientes[colunas_numericas])
11
```

Descrição:

- **Normalização de Dados:** Utilizamos `MinMaxScaler` para normalizar os atributos numéricos no intervalo `[0, 1]`.

5.2 Codificação de Dados Categóricos

Funções e Métodos:

- `OneHotEncoder()`: Aplica a codificação One-Hot nos dados categóricos.

Código:

```
1 from sklearn.preprocessing import OneHotEncoder
2
3 # Instanciar o OneHotEncoder
4 encoder = OneHotEncoder(sparse_output=False)
5
6 # Codificar o atributo 'sexo'
7 sexo_encoded = encoder.fit_transform(df_clientes[['sexo']])
8 sexo_encoded_df = pd.DataFrame(sexo_encoded, columns=encoder.get_feature_names_out(['sexo']))
9
10 # Codificar o atributo 'genero_musical_favorito'
11 genero_encoded = encoder.fit_transform(df_clientes[['genero_musical_favorito']])
12 genero_encoded_df = pd.DataFrame(genero_encoded, columns=encoder.get_feature_names_out(['genero_musical_favorito']))
13
14 # Concatenar as colunas codificadas ao DataFrame original
15 df_clientes = pd.concat([df_clientes, sexo_encoded_df, genero_encoded_df], axis=1)
16
17 # Remover colunas originais categóricas
18 df_clientes = df_clientes.drop(columns=['sexo', 'genero_musical_favorito'])
19
20 # Exibir o DataFrame Processado
21 print("\nDataFrame Processado:\n")
22 print(df_clientes.head())
23
24 print("\nResumo do DataFrame Processado:\n")
25 print(df_clientes.info())
26
27 print("\nEstatísticas Descritivas do DataFrame Processado:\n")
28 print(df_clientes.describe())
29
```

Descrição:

- **Codificação de Dados Categóricos:** Utilizamos `OneHotEncoder` para transformar atributos categóricos em variáveis dummy, permitindo seu uso em modelos de aprendizado de máquina.

Conclusão

Neste projeto, aplicamos uma abordagem abrangente para a análise e transformação de uma base de dados de clientes, utilizando uma variedade de técnicas estatísticas e computacionais. As etapas seguidas incluíram a limpeza e preparação dos dados, o cálculo de estatísticas descritivas, a análise de correlação, a visualização gráfica dos dados, a correção de valores inconsistentes, a remoção de redundâncias e a transformação de dados para fins de modelagem.

Os principais resultados obtidos incluem:

- **Limpeza e Preparação de Dados:** Conseguimos identificar e corrigir valores inconsistentes e duplicados, garantindo a qualidade dos dados para análises subsequentes.
- **Estatísticas Descritivas:** Foram calculadas medidas centrais (média, mediana, moda) e medidas de dispersão (variância, desvio-padrão, amplitude), fornecendo uma visão clara das características dos dados.
- **Análise de Correlação:** Identificamos relações significativas entre diferentes atributos numéricos, o que pode guiar futuras análises e decisões baseadas em dados.
- **Visualizações Gráficas:** Criamos diversas representações visuais, como histogramas, box plots, gráficos de dispersão e mapas de calor, que ajudaram a interpretar os dados de maneira intuitiva e informativa.
- **Transformação de Dados:** Normalizamos os dados numéricos e codificamos os atributos categóricos, tornando os dados prontos para uso em modelos de aprendizado de máquina e outras análises avançadas.

Este projeto demonstrou a importância de um processo bem-estruturado de análise e transformação de dados para extrair insights valiosos e garantir a qualidade dos dados. As técnicas aplicadas aqui são fundamentais para qualquer análise de dados robusta e servem como base para futuras análises mais complexas e modelagens preditivas.