

ENTREGA 1 – PROJETO INTERDISCIPLINAR

Introdução

Este documento apresenta uma análise estruturada do código Python desenvolvido no ambiente Google Colab para a análise dos dados disponibilizados pela *PicMoney* em suas quatro tabelas (*base_players*, *transacoes_cupons*, *pedestres_paulista* e *lojas_valores*). O objetivo principal deste documento é documentar as etapas feitas para cumprir os requisitos propostos: entender os dados, coletar dados iniciais, descrever os dados, explorar os dados e verificar a qualidade dos dados.

Bloco 1: Importação de Bibliotecas e Configurações

Neste bloco inicial, foram importadas as bibliotecas fundamentais para a execução das tarefas de manipulação e visualização de dados. A biblioteca Pandas (pd) é utilizada para a criação e gestão de DataFrames, que são estruturas tabulares que organizam os dados. A biblioteca NumPy (np) oferece suporte a operações numéricas de alto desempenho. Para a visualização de dados, as bibliotecas Matplotlib (plt) e Seaborn (sns) foram importadas, permitindo a criação de gráficos e a customização de seus elementos visuais. Adicionalmente, foram aplicadas configurações para o estilo visual dos gráficos e para garantir que todas as colunas de um DataFrame sejam exibidas integralmente.

```
1 # Bloco 1: Imports e Configurações
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 # Configuração visual
8 plt.style.use("seaborn-v0_8")
9 sns.set_palette("pastel")
10 pd.set_option("display.max_columns", None)
```

Bloco 2: Carregamento dos Dados

A segunda etapa do processo consiste no carregamento dos conjuntos de dados a partir de arquivos CSV. O código carrega quatro arquivos distintos: *base_players.csv*, *lojas_valores.csv*, *transacoes_cupons.csv* e *pedestres_paulista.csv*. Cada arquivo é lido por meio da biblioteca Pandas, utilizando o ponto e vírgula (;) como delimitador de colunas. Após o carregamento, a forma (dimensões) de cada DataFrame é exibida, indicando o número de linhas e colunas de cada conjunto de dados.

players: 10.000 linhas e 11 colunas.

lojas: 10.000 linhas e 11 colunas.

transacoes: 100.000 linhas e 12 colunas.

pedestres: 100.000 linhas e 15 colunas.

```
1 # Bloco 2: Carregar os datasets
2 players = pd.read_csv("base_players.csv", sep=";")
3 lojas = pd.read_csv("lojas_valores.csv", sep=";")
4 transacoes = pd.read_csv("transacoes_cupons.csv", sep=";")
5 pedestres = pd.read_csv("pedestres_paulista.csv", sep=";")
6
7 # Mostrar dimensões
8 print("Players:", players.shape)
9 print("Lojas:", lojas.shape)
10 print("Transações:", transacoes.shape)
11 print("Pedestres:", pedestres.shape)
```

Bloco 3: Verificação da Qualidade dos Dados

Para assegurar a integridade dos dados, o terceiro bloco implementa uma função (qualidade_dados) que realiza uma análise detalhada de cada DataFrame. Esta função verifica as dimensões dos conjuntos, os tipos de dados de cada coluna, a presença de valores ausentes (NaN) e o número de linhas duplicadas.

A análise revelou:

players: As colunas cidade_trabalho e bairro_trabalho apresentam 3.969 valores ausentes, enquanto cidade_escola e bairro_escola possuem 6.932 valores ausentes. Não foram encontradas linhas duplicadas.

lojas: Não foram identificados valores ausentes ou linhas duplicadas.

pedestres: As colunas data_ultima_compra, ultimo_tipo_cupom, ultimo_valor_capturado e ultimo_tipo_loja apresentam 40.043 valores ausentes. Não foram encontradas linhas duplicadas.

```
1 # Bloco 3: Verificar qualidade dos dados
2
3 def qualidade_dados(df, nome):
4     print(f"\nDataset: {nome}")
5     print("Dimensões:", df.shape)
6     print("\nTipos de dados:")
7     print(df.dtypes)
8     print("\nValores ausentes:")
9     print(df.isnull().sum())
10    print("\nDuplicados:", df.duplicated().sum())
11    print("="*60)
12
13 qualidade_dados(players, "Players")
14 qualidade_dados(lojas, "Lojas")
15 qualidade_dados(transacoes, "Transações")
16 qualidade_dados(pedestres, "Pedestres")
```

Bloco 4: Tratamento de Dados Duplicados e Ausentes

Neste bloco, é realizado um tratamento inicial para a remoção de possíveis linhas duplicadas e valores ausentes (NaN) nos quatro DataFrames. Embora a análise anterior tenha indicado que não há duplicados, a execução do método `drop_duplicates()` garante um processo de limpeza padronizado e robusto. Adicionalmente, o método `dropna()` é aplicado para remover as linhas que contêm valores ausentes. É importante notar que essa etapa de tratamento de valores ausentes pode reduzir significativamente o tamanho dos DataFrames, dependendo da quantidade de dados faltantes.

```
1 # Bloco 4: Tratamento básico dos dados
2
3 # Remover possíveis duplicados
4 players = players.drop_duplicates()
5 lojas = lojas.drop_duplicates()
6 transacoes = transacoes.drop_duplicates()
7 pedestres = pedestres.drop_duplicates()
8
9 # Remover possíveis valores ausentes
10 players = players.dropna()
11 lojas = lojas.dropna()
12 transacoes = transacoes.dropna()
13 pedestres = pedestres.dropna()
```

Bloco 5: Pré-visualização dos Dados

Este bloco apresenta uma pré-visualização das cinco primeiras linhas de cada DataFrame (`players`, `lojas`, `transacoes` e `pedestres`) para verificar o formato e o conteúdo dos dados após as etapas de carregamento e tratamento. Esta etapa é crucial para uma inspeção visual rápida e para confirmar que o processamento inicial foi bem-sucedido. A função `display()` é utilizada para exibir os DataFrames de forma formatada.

```
1 # Bloco 5: Pré-visualização
2 print("Players:")
3 display(players.head())
4
5 print("Lojas:")
6 display(lojas.head())
7
8 print("Transações:")
9 display(transacoes.head())
10
11 print("Pedestres:")
12 display(pedestres.head())
```

Bloco 6: Análise de Dados Agrupados

Este bloco de código realiza análises estatísticas dos dados, agrupando-os por categoria e calculando a soma e a média das transações. Os resultados são exibidos em um formato tabulado, revelando insights sobre a frequência e o valor das transações em cada categoria.

```
1 # Bloco 6: Estatísticas descritivas
2 print("Descrição - Players")
3 display(players.describe(include="all").T)
4
5 print("Descrição - Lojas")
6 display(lojas.describe(include="all").T)
7
8 print("Descrição - Transações")
9 display(transacoes.describe(include="all").T)
10
11 print("Descrição - Pedestres")
12 display(pedestres.describe(include="all").T)
```

Bloco 7: KPIs e Análises

```
1 # Bloco 7: KPIs e análises
2
3 # KPI 1: Número de clientes únicos
4 clientes_totais = players.shape[0]
5
6 # KPI 2: Número de lojas parceiras
7 lojas_unicas = transacoes['nome_estabelecimento'].unique()
8
9 # KPI 3: Volume total transacionado
10 volume_total = lojas["valor_compra"].sum()
11
12 # KPI 4: Valor médio dos cupons
13 valor_medio_cupom = lojas["valor_cupom"].mean()
14
15 # KPI 5: Ticket médio por compra
16 ticket_medio = lojas["valor_compra"].mean()
17
18 # KPI 6: Número de pedestres únicos
19 pedestres_unicos = pedestres.shape[0]
20
21
22 print("🌟 KPIs para CEO & CFO")
23 print(f"Clientes únicos: {clientes_totais}")
24 print(f"Lojas parceiras: {len(lojas_unicas)}")
25 print(f"Volume total transacionado: R$ {volume_total:,.2f}")
26 print(f"Valor médio dos cupons: R$ {valor_medio_cupom:,.2f}")
27 print(f"Ticket médio: R$ {ticket_medio:,.2f}")
28 print(f"Pedestres únicos registrados: {pedestres_unicos}")
29
```

```
30
31 # Visualizações
32 fig, axs = plt.subplots(1, 2, figsize=(12,5))
33
34 # Distribuição de valores de compra
35 sns.histplot(lojas["valor_compra"], bins=30, ax=axs[0], kde=True)
36 axs[0].set_title("Distribuição dos Valores de Compra")
37
38 # Relação cupom vs compra
39 sns.scatterplot(data=lojas.sample(5000), x="valor_cupom", y="valor_compra", alpha=0.5, ax=axs[1])
40 axs[1].set_title("Cupom vs Valor da Compra")
41
42 plt.tight_layout()
43 plt.show()
```