

Relatório de Análise de Dados: Presença e Influência dos Dados no Setor Imobiliário de São Paulo

Pedro Miller, Pedro Bolzan, Hálison P Q Souza, Luis Martins, Daniel Muradas Rodrigues

25 de junho de 2025

Conteúdo

1	Introdução	3
2	Coleta e Visão Geral dos Dados	3
2.1	Carregamento do Dataset	3
2.2	Inspeção Inicial do Dataset	3
3	Preparação e Limpeza de Dados	4
3.1	Extração da Informação do Bairro	4
3.2	Tratamento de Colunas de Data	4
3.3	Tratamento de Outliers em Variáveis Numéricas	5
3.4	Limpeza da Categoria 'Desconhecido' na Coluna Bairro	7
3.5	Criação da Métrica 'Preço por Metro Quadrado' (Price_per_sqm)	7
4	Análise Exploratória de Dados (EDA) e Insights	8
4.1	Correlação entre Variáveis Numéricas	8
4.2	Impacto das Características do Imóvel no Preço	9
4.3	A Influência da Localização (Bairro) no Preço	10
5	Conclusões Finais	12

1 Introdução

Este relatório detalha a aplicação do ciclo de ciência de dados para explorar a presença e a influência dos dados no setor imobiliário da cidade de São Paulo. O objetivo principal é conduzir uma análise passo a passo, desde a coleta e organização dos dados até a sua exploração e análise, culminando na apresentação de conclusões através de visualizações. O trabalho busca responder a questões fundamentais sobre os fatores que impactam o preço dos imóveis e a relevância da localização nesse contexto.

Questões de Pesquisa:

- Como as características dos imóveis (área, número de quartos/banheiros, vagas de estacionamento) influenciam seu preço de venda?
- Qual o impacto da localização (bairro) na variação do preço por metro quadrado dos apartamentos em São Paulo?
- Existem outros fatores, como o tempo de anúncio, que correlacionam com o preço dos imóveis?

2 Coleta e Visão Geral dos Dados

A etapa inicial envolveu a aquisição de um dataset relevante para o setor imobiliário de São Paulo, focado em apartamentos. Após a obtenção dos dados, foi realizada uma inspeção preliminar para entender sua estrutura, o tipo de informações contidas e a presença de valores ausentes, o que é fundamental para planejar as próximas fases de preparação.

2.1 Carregamento do Dataset

O dataset foi carregado utilizando a biblioteca Pandas no Python.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5
6 # Configura es para melhorar a visualiza o de gr ficos (opcional)
7 sns.set_style("whitegrid")
8 plt.rcParams['figure.figsize'] = (10, 6)
9 plt.rcParams['figure.dpi'] = 100
10
11 # Carregar o dataset imobili rio de apartamentos de S o Paulo
12 df = pd.read_csv('SaoPaulo_OnlyApartments.csv')
13
14 print("Dataset carregado com sucesso!")
```

Listing 1: Carregamento do Dataset

2.2 Inspeção Inicial do Dataset

Para uma visão geral, as primeiras linhas, informações sobre as colunas (tipos de dados e valores não-nulos) e estatísticas descritivas das colunas numéricas foram inspecionadas.

```
1 # Visualizar as primeiras linhas do DataFrame
2 print("### Primeiras 5 linhas do Dataset:")
3 print(df.head())
4
5 # Obter informa es sobre as colunas e tipos de dados
6 print("\n### Informa es Detalhadas do Dataset:")
7 df.info()
8
9 # Estat sticas descritivas das colunas num ricas
10 # Isso nos d uma ideia da distribui o, m dia, desvio padr o e valores m nimo/
    m ximo.
11 print("\n### Estat sticas Descritivas das Colunas Num ricas (Dados Originais):")
12 print(df.describe())
```

Listing 2: Inspeção Inicial do Dataset

Saída da Inspeção Inicial (Exemplo)

```
ID created_date Price below_price Area
0 2742871875 2018-05-24T15:02:18Z 2000000 False 62
1 2670612721 2018-12-18T14:40:10Z 1015680 False 42
...

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 27828 entries, 0 to 27827
Data columns (total 10 columns):
# Column Non-Null Count Dtype
-- --
0 ID 27828 non-null int64
1 created_date 27828 non-null object
2 Price 27828 non-null int64
...
dtypes: bool(1), int64(6), object(3)
memory usage: 1.9+ MB

Estatísticas Descritivas:
ID Price Area Bedrooms Bathrooms
count 2.782800e+04 2.782800e+04 27828.000000 27828.000000 27828.000000
mean 2.685599e+09 2.037172e+06 136.288163 2.740477 2.831177
...
max 2.758322e+09 1.220290e+08 3367.000000 20.000000 20.000000
Parking_Spaces
count 27828.000000
mean 2.153730
std 2.927744
min 0.000000
25% 1.000000
50% 2.000000
75% 3.000000
max 350.000000
```

3 Preparação e Limpeza de Dados

A qualidade da análise de dados depende diretamente da qualidade dos dados brutos. Nesta seção, o dataset foi preparado e limpo para garantir consistência e confiabilidade, construindo o dataset de trabalho final.

3.1 Extração da Informação do Bairro

A coluna 'Adress' continha o endereço completo. Para facilitar a análise baseada em localização, o nome do bairro foi extraído para uma nova coluna 'Bairro'.

```
1 # Extrair o nome do Bairro da coluna 'Adress'
2 # A l g a tenta dividir a string pelo primeiro ',', depois pelo primeiro '-' e remove
   espa os extras.
3 # Se a extra o falhar, o bairro ser classificado como 'Desconhecido'.
4 df['Bairro'] = df['Adress'].apply(lambda x: x.split(',')[1].split('-')[0].strip() if pd.
   notnull(x) and ',' in x and '-' in x.split(',')[1] else 'Desconhecido')
5
6 # Verificar a contagem dos bairros mais frequentes, incluindo 'Desconhecido'
7 print("### Top 10 Bairros (antes da remo o de 'Desconhecido'):")
8 print(df['Bairro'].value_counts().head(10))
```

Listing 3: Extração da Coluna Bairro

3.2 Tratamento de Colunas de Data

As colunas 'created_date' e 'extract_date' foram convertidas para o tipo 'datetime' para permitir operações temporais. Além disso, f

```

1 # Converter colunas de data para o tipo datetime
2 print("Tipos de dados antes da conversão de datas:")
3 print(df[['created_date', 'extract_date']].dtypes)
4
5 # Converter created_date e remover a informação de fuso horário (tz-localize(None))
6 df['created_date'] = pd.to_datetime(df['created_date']).dt.tz_localize(None)
7
8 # Converter extract_date
9 df['extract_date'] = pd.to_datetime(df['extract_date'])
10
11 print("\nTipos de dados após a conversão de datas:")
12 print(df[['created_date', 'extract_date']].dtypes)
13
14 # Criar a coluna 'age_of_listing_days' (idade do anúncio em dias)
15 # Esta coluna representa o quanto tempo o imóvel estava anunciado na data da extração do dado.
16 df['age_of_listing_days'] = (df['extract_date'] - df['created_date']).dt.days
17
18 print("\n### Primeiras linhas com a nova coluna 'age_of_listing_days':")
19 print(df[['created_date', 'extract_date', 'age_of_listing_days']].head())
20
21 # Verificar se há idades negativas (inconsistência de dados)
22 print("\n### Número de anúncios com idade negativa (data de criação > data de extração):")
23 print(df[df['age_of_listing_days'] < 0].shape[0])

```

Listing 4: Conversão e Cálculo de Datas

3.3 Tratamento de Outliers em Variáveis Numéricas

A análise inicial revelou a presença de outliers significativos em colunas como 'Price', 'Area', 'Bedrooms', 'Bathrooms' e 'Parking_spaces'. Para mitigar o impacto desses valores extremos, foi aplicada uma estratégia de *capping* (limitando os valores a um percentil de 1% e 99%) e, em casos específicos, remoção de valores claramente implausíveis.

Distribuição das Variáveis Numéricas ANTES do Tratamento de Outliers

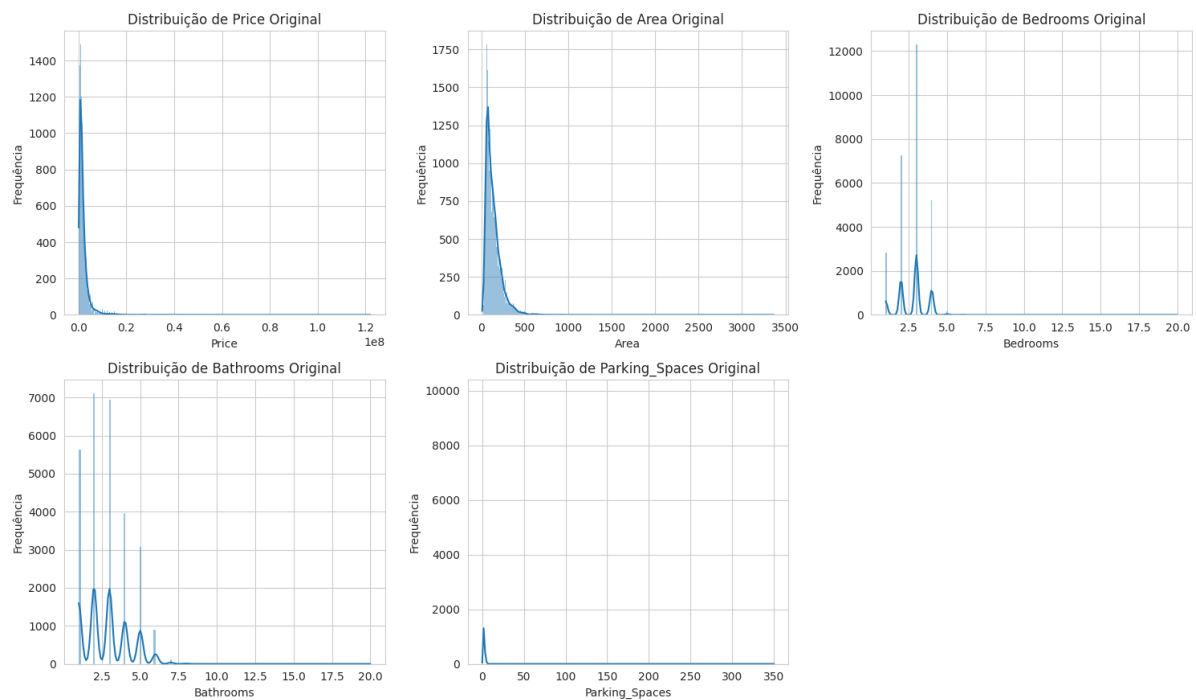


Figura 1: Distribuição das Variáveis Numéricas ANTES do Tratamento de Outliers

```

1 # Aplicar Capping e Remoções para Outliers

```

```

2 print("\n### In cio do Tratamento de Outliers ###")
3
4 cols_to_treat = ['Price', 'Area', 'Bedrooms', 'Bathrooms', 'Parking_Spaces']
5
6 for col in cols_to_treat:
7     print(f"\nTratando outliers para a coluna: **{col}**")
8
9     # Calcular percentis para capping (1% e 99%)
10    lower_perc = df[col].quantile(0.01)
11    upper_perc = df[col].quantile(0.99)
12
13    # Antes do tratamento (valores originais do df atual)
14    print(f"Valores ANTES do capping (min/max): {df[col].min():.2f} / {df[col].max():.2f}")
15
16    # Aplicar capping: valores abaixo do percentil 1% s o ajustados para o valor do
17    # percentil 1%
18    # e valores acima do percentil 99% s o ajustados para o valor do percentil 99%.
19    df[col] = np.where(df[col] < lower_perc, lower_perc, df[col])
20    df[col] = np.where(df[col] > upper_perc, upper_perc, df[col])
21
22    print(f"Valores AP S capping (min/max): {df[col].min():.2f} / {df[col].max():.2f}")
23    print(f"Outliers de {col} contidos via capping.")
24
25 # Remo es espec ficas para valores fisicamente improv veis (p s -capping)
26 initial_rows_before_specific_removal = df.shape[0]
27
28 # Remover apartamentos com rea muito pequena (ex: menos de 20m  pode ser erro ou
29 # est dio extremo)
30 df = df[df['Area'] >= 20]
31 removed_area_rows = initial_rows_before_specific_removal - df.shape[0]
32 print(f"\nRemovidos {removed_area_rows} linhas com rea < 20m .")
33
34 # Remover apartamentos com um n mero irreal de vagas de estacionamento (ex: mais de 10)
35 initial_rows_after_area_removal = df.shape[0]
36 df = df[df['Parking_Spaces'] <= 10]
37 removed_parking_rows = initial_rows_after_area_removal - df.shape[0]
38 print(f"Removidos {removed_parking_rows} linhas com Parking_Spaces > 10.")
39
40 print(f"\n### N mero final de linhas no dataset ap s remo es espec ficas: {df.
41     shape[0]} ###")

```

Listing 5: Tratamento de Outliers

```

1 # Re-verificar as estat sticas descritivas ap s o tratamento de outliers
2 print("\n### Estat sticas Descritivas Ap s Tratamento de Outliers: ###")
3 print(df[cols_to_treat].describe())

```

Listing 6: Estatísticas Descritivas Após Tratamento de Outliers

Distribuição das Variáveis Numéricas APÓS Tratamento de Outliers

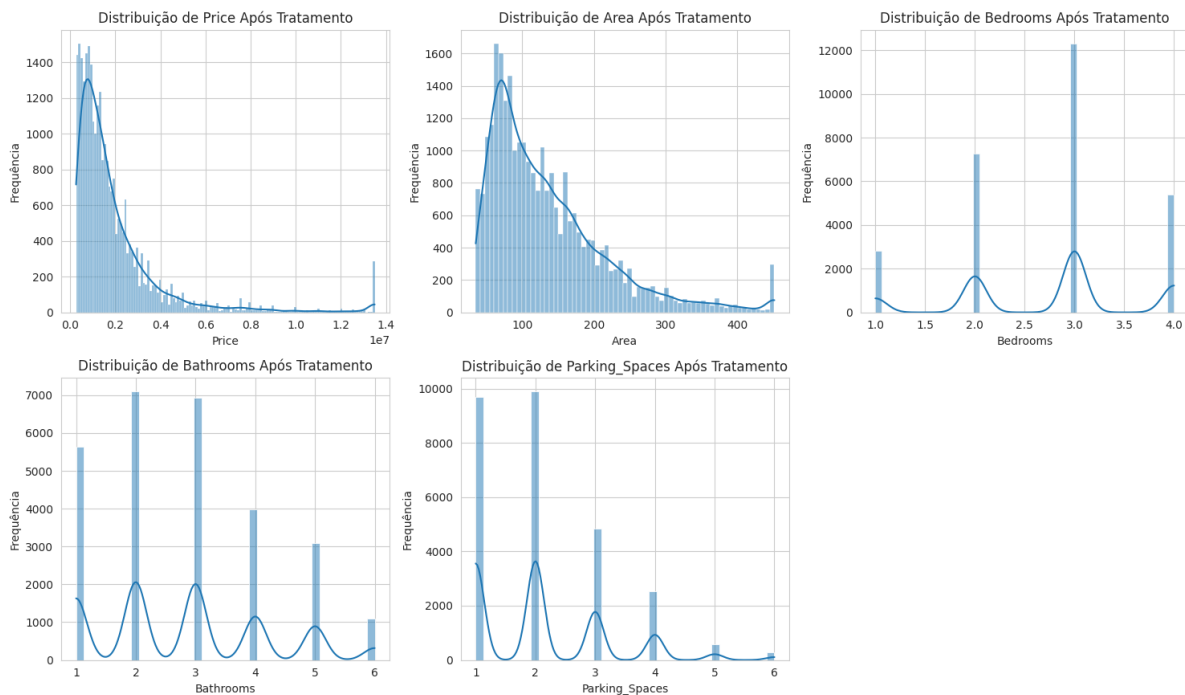


Figura 2: Distribuição das Variáveis Numéricas APÓS Tratamento de Outliers

3.4 Limpeza da Categoria 'Desconhecido' na Coluna Bairro

A extração do bairro resultou em alguns endereços não identificados, categorizados como 'Desconhecido'. Para assegurar a precisão da análise de localização, essas linhas foram removidas do dataset.

```
1 # Remover linhas onde a coluna 'Bairro' é 'Desconhecido'
2 print("\n### Tratamento de Bairros 'Desconhecido' ###")
3
4 initial_rows_bairro = df.shape[0]
5 print(f"N mero de linhas antes da remoção de Bairro 'Desconhecido': {initial_rows_bairro}")
6 print(f"Contagem original de 'Desconhecido': {df['Bairro'].value_counts().get('Desconhecido', 0)}")
7
8 df = df[df['Bairro'] != 'Desconhecido']
9 removed_rows_bairro = initial_rows_bairro - df.shape[0]
10
11 print(f"N mero de linhas removidas com Bairro 'Desconhecido': {removed_rows_bairro}")
12 print(f"Novo n mero de linhas no dataset: {df.shape[0]}")
13
14 # Verificar os 10 bairros mais frequentes após a remoção
15 print("\n### Top 10 Bairros após remoção de 'Desconhecido':")
16 print(df['Bairro'].value_counts().head(10))
17 print(f"Nova contagem de 'Desconhecido' (deve ser 0): {df['Bairro'].value_counts().get('Desconhecido', 0)}")
```

Listing 7: Limpeza de Bairros 'Desconhecido'

3.5 Criação da Métrica 'Preço por Metro Quadrado' (Price_per_sqm)

Para permitir uma comparação justa do valor de imóveis de diferentes tamanhos e aprofundar a análise da influência da localização, foi criada a métrica 'Price_{per_sqm}'.

```
1 # Calcular a coluna 'Price_per_sqm'
2 # Adicionar um pequeno valor (1e-6) para evitar divisão por zero.
3 df['Price_per_sqm'] = df['Price'] / (df['Area'] + 1e-6)
```

```

4
5 print("\n### Estatísticas Descritivas para 'Price_per_sqm':")
6 print(df['Price_per_sqm'].describe())
7
8 # Visualizar a distribuição do Preço por Metro Quadrado
9 plt.figure(figsize=(10, 6))
10 sns.histplot(df['Price_per_sqm'], kde=True, bins=50)
11 plt.title('Distribuição do Preço por Metro Quadrado')
12 plt.xlabel('Preço por m² (R$)')
13 plt.ylabel('Frequência')
14 plt.ticklabel_format(style='plain', axis='x')
15 plt.savefig('distribuicao_preco_por_m2.png')
16 plt.show()
17 print("Distribuição do Preço por m² salva como 'distribuicao_preco_por_m2.png'.")

```

Listing 8: Criação de Preço por Metro Quadrado

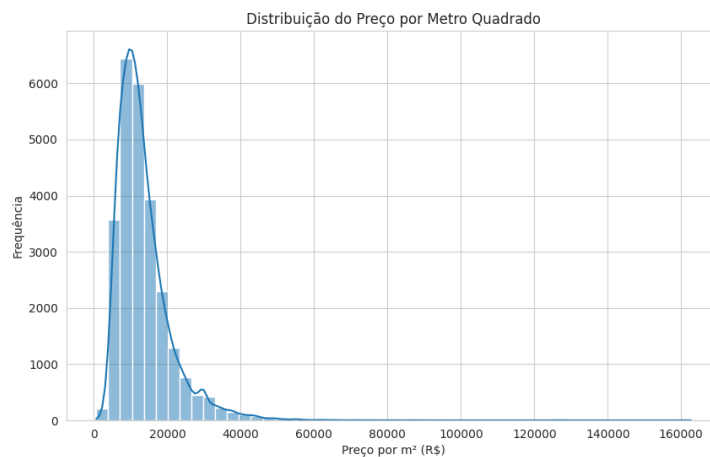


Figura 3: Distribuição do Preço por Metro Quadrado

4 Análise Exploratória de Dados (EDA) e Insights

Com o dataset limpo e preparado, a fase de Análise Exploratória de Dados (EDA) foi conduzida para identificar padrões, tendências e relações entre as variáveis, buscando responder às questões de pesquisa.

4.1 Correlação entre Variáveis Numéricas

A matriz de correlação foi utilizada para quantificar a relação linear entre as variáveis numéricas, permitindo identificar quais características exercem maior influência sobre o preço.

```

1 # Calcular e visualizar a matriz de correlação
2 # Selecionar apenas as colunas numéricas relevantes para análise de correlação
3 numerical_cols_for_corr = ['Price', 'Area', 'Bedrooms', 'Bathrooms', 'Parking_Spaces', '
    age_of_listing_days', 'Price_per_sqm']
4
5 # Calcular a matriz de correlação
6 correlation_matrix = df[numerical_cols_for_corr].corr()
7 print("### Matriz de Correlação das Variáveis Numéricas:")
8 print(correlation_matrix)
9
10 # Visualizar a matriz de correlação com um heatmap para melhor interpretação
11 plt.figure(figsize=(10, 8))
12 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=.5)
13 plt.title('Matriz de Correlação das Variáveis Numéricas')
14 plt.savefig('matriz_correlacao.png')
15 plt.show()
16 print("Matriz de Correlação salva como 'matriz_correlacao.png'.")

```

Listing 9: Matriz de Correlação

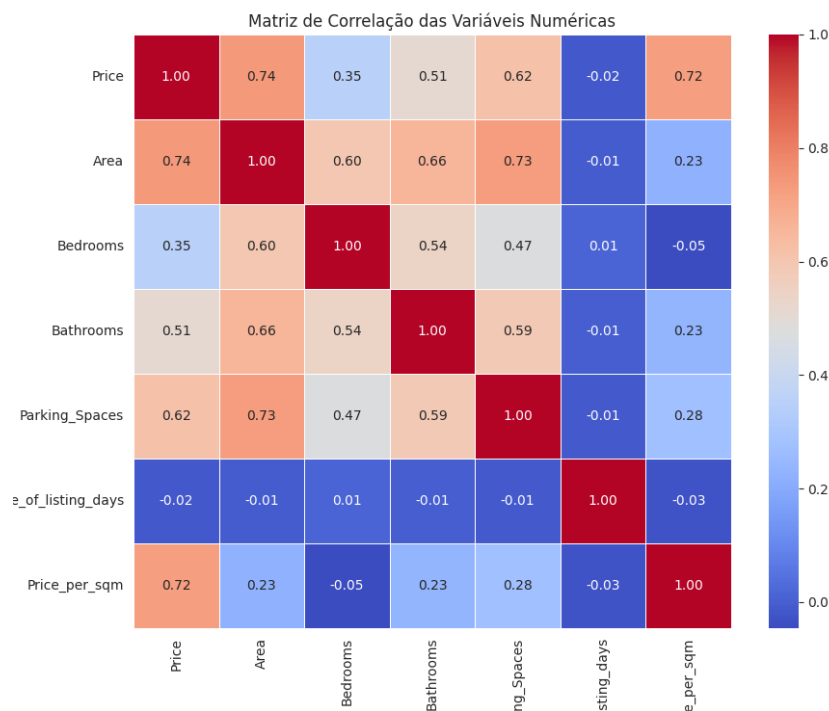


Figura 4: Matriz de Correlação das Variáveis Numéricas

4.2 Impacto das Características do Imóvel no Preço

Gráficos de dispersão foram gerados para visualizar a relação direta entre as características mais correlacionadas e o preço dos apartamentos.

```

1 # Gráficos de dispersão (scatter plots) para as variáveis mais correlacionadas com o
  Preço
2 plt.figure(figsize=(18, 6))
3
4 # Preço vs. Área: Espera-se uma relação positiva forte.
5 plt.subplot(1, 3, 1)
6 sns.scatterplot(x='Area', y='Price', data=df, alpha=0.6, s=20)
7 plt.title('Preço vs. Área')
8 plt.xlabel('Área (m²)')
9 plt.ylabel('Preço (R$)')
10 plt.ticklabel_format(style='plain', axis='y')
11
12 # Preço vs. Banheiros: Espera-se uma relação positiva.
13 plt.subplot(1, 3, 2)
14 sns.scatterplot(x='Bathrooms', y='Price', data=df, alpha=0.6, s=20)
15 plt.title('Preço vs. Banheiros')
16 plt.xlabel('Número de Banheiros')
17 plt.ylabel('Preço (R$)')
18 plt.ticklabel_format(style='plain', axis='y')
19
20 # Preço vs. Quartos: Espera-se uma relação positiva.
21 plt.subplot(1, 3, 3)
22 sns.scatterplot(x='Bedrooms', y='Price', data=df, alpha=0.6, s=20)
23 plt.title('Preço vs. Quartos')
24 plt.xlabel('Número de Quartos')
25 plt.ylabel('Preço (R$)')
26 plt.ticklabel_format(style='plain', axis='y')
27
28 plt.tight_layout()
29 plt.savefig('preco_vs_caracteristicas.png')
30 plt.show()

```

```
31 print("Gráficos de Preço vs. Características salvos como 'preco_vs_caracteristicas.png'.")
```

Listing 10: Preço vs. Características

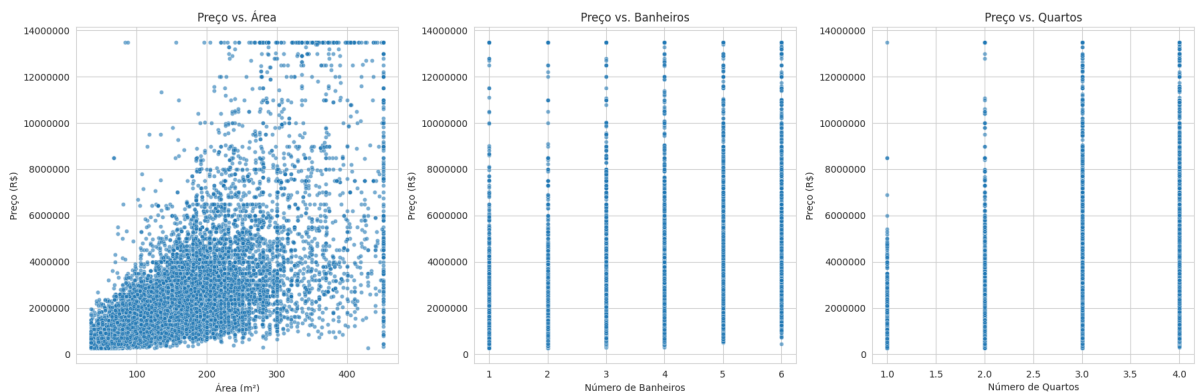


Figura 5: Preço vs. Características do Imóvel

4.3 A Influência da Localização (Bairro) no Preço

A localização é um dos fatores mais críticos no mercado imobiliário. Esta seção investiga como o preço por metro quadrado varia significativamente entre os diferentes bairros de São Paulo.

```
1 # Top 10 Bairros com maior e menor preço médio por m
2 # Calcular o preço médio por metro quadrado para cada bairro
3 price_per_sqm_by_bairro = df.groupby('Bairro')['Price_per_sqm'].mean().sort_values(
    ascending=False)
4
5 print("\n### Top 10 Bairros com maior preço médio por m :")
6 print(price_per_sqm_by_bairro.head(10))
7
8 print("\n### Top 10 Bairros com menor preço médio por m :")
9 print(price_per_sqm_by_bairro.tail(10))
10
11 # Visualizar os 10 bairros mais caros e os 10 mais baratos em termos de preço por m
12 plt.figure(figsize=(15, 7))
13
14 % Top 10 Bairros mais caros
15 plt.subplot(1, 2, 1)
16 sns.barplot(x=price_per_sqm_by_bairro.head(10).values, y=price_per_sqm_by_bairro.head(
    10).index,
17             palette='viridis', hue=price_per_sqm_by_bairro.head(10).index, legend=False)
18 plt.title('Top 10 Bairros Mais Caros (Preço Médio por m)')
19 plt.xlabel('Preço Médio por m (R$)')
20 plt.ylabel('Bairro')
21 plt.ticklabel_format(style='plain', axis='x')
22
23 % Top 10 Bairros mais baratos
24 plt.subplot(1, 2, 2)
25 sns.barplot(x=price_per_sqm_by_bairro.tail(10).values, y=price_per_sqm_by_bairro.tail(
    10).index,
26             palette='magma', hue=price_per_sqm_by_bairro.tail(10).index, legend=False)
27 plt.title('Top 10 Bairros Mais Baratos (Preço Médio por m)')
28 plt.xlabel('Preço Médio por m (R$)')
29 plt.ylabel('Bairro')
30 plt.ticklabel_format(style='plain', axis='x')
31
32 plt.tight_layout()
33 plt.savefig('top_bairros_preco_m2.png')
34 plt.show()
35 print("Gráficos dos Top Bairros por Preço por m² salvos como 'top_bairros_preco_m2.png'.")
```

Listing 11: Análise de Preço por m² por Bairro

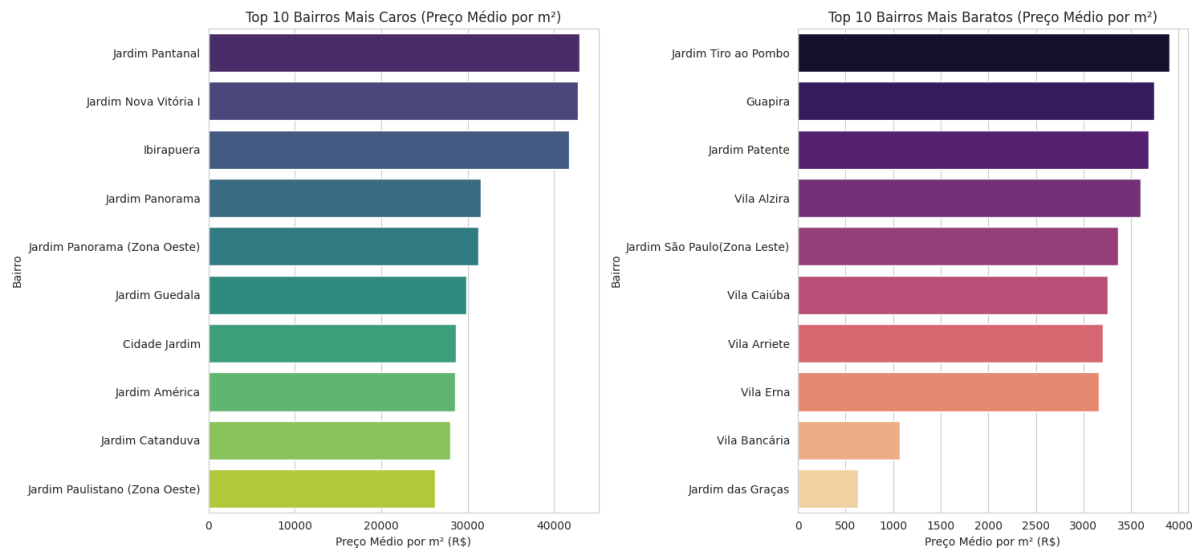


Figura 6: Top 10 Bairros Mais Caros e Mais Baratos (Preço Médio por m²)

```

1 # Box plot da distribuição do Preço por m² nos Top 20 Bairros (mais detalhado)
2 # Selecionar os 20 bairros com mais apartamentos no dataset para o box plot
3 top_20_bairros = df['Bairro'].value_counts().head(20).index
4 df_top_bairros = df[df['Bairro'].isin(top_20_bairros)]
5
6 plt.figure(figsize=(14, 8))
7 sns.boxplot(x='Price_per_sqm', y='Bairro', data=df_top_bairros.sort_values('
8     Price_per_sqm', ascending=False),
9     palette='coolwarm', hue='Bairro', legend=False)
10 plt.title('Distribuição do Preço por m² nos Top 20 Bairros de São Paulo')
11 plt.xlabel('Preço por m² (R$)')
12 plt.ylabel('Bairro')
13 plt.ticklabel_format(style='plain', axis='x')
14 plt.grid(axis='x', linestyle='--', alpha=0.7)
15 plt.tight_layout()
16 plt.savefig('boxplot_preco_m2_top_bairros.png')
17 plt.show()
18 print("Box plot do Preço por m² nos Top 20 Bairros salvo como '
19     boxplot_preco_m2_top_bairros.png'.")

```

Listing 12: Box Plot de Preço por m² nos Top Bairros

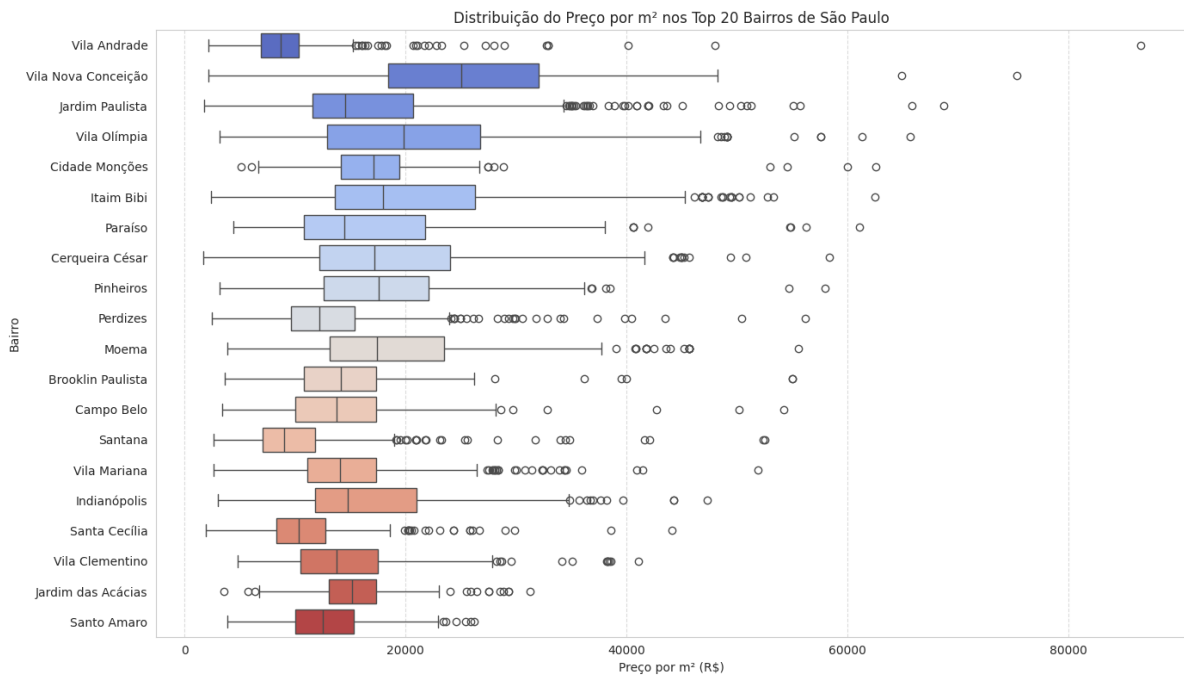


Figura 7: Distribuição do Preço por m² nos Top 20 Bairros de São Paulo

5 Conclusões Finais

Este estudo demonstrou a aplicação de um ciclo completo de ciência de dados para analisar a influência dos dados no setor imobiliário de São Paulo. As principais conclusões são:

- **A Importância Crítica da Preparação de Dados:** A fase de limpeza e organização (conversão de datas, extração de texto e, principalmente, tratamento de outliers) foi fundamental. As visualizações "antes" e "depois" do tratamento de outliers evidenciaram como essa etapa é crucial para obter análises precisas e evitar distorções, transformando dados brutos em um dataset confiável.
- **Influência Determinante das Características do Imóvel no Preço:** A análise de correlação e os gráficos de dispersão confirmaram que a **área do imóvel** é o fator mais impactante no preço de venda. O número de **banheiros**, **quartos** e **vagas de estacionamento** também demonstraram forte correlação positiva, refletindo que maior espaço e mais comodidades elevam o valor da propriedade.
- **Localização como Fator Primordial de Precificação:** A análise do preço por metro quadrado revelou uma **disparidade extrema** entre os bairros de São Paulo. Bairros de alto padrão como Jardim Pantanal e Ibirapuera mostram um valor médio por m² que pode ser dezenas de vezes maior do que em regiões mais acessíveis como Jardim das Graças e Vila Bancária. O Box Plot detalhou essa variação, mostrando não apenas as diferenças de medianas entre os bairros, mas também a vasta **variabilidade e assimetria na distribuição de preços** dentro de um mesmo bairro, sublinhando a complexidade e a micro-segmentação do mercado imobiliário da cidade.
- **Outros Fatores:** Curiosamente, a coluna `age_of_listing_days` (idade do anúncio) mostrou uma correlação linear muito baixa com o preço, sugerindo que, em um primeiro momento, o tempo de permanência de um anúncio online não é um fator linearmente determinante para a variação do valor final do imóvel.
- **A Essência dos Dados no Setor Imobiliário:** Em suma, este trabalho reitera que os dados são um ativo indispensável no setor imobiliário. Eles permitem uma compreensão aprofundada das dinâmicas de mercado, embasando avaliações precisas de propriedades, identificação de tendências de valorização/desvalorização, otimização de estratégias de venda e aluguel, e uma tomada de decisão de investimento muito mais informada e estratégica. A capacidade de extrair inteligência a partir desses dados confere uma vantagem competitiva inestimável.

Este estudo demonstra como a aplicação sistemática de um ciclo de ciência de dados pode transformar informações brutas em insights valiosos para o mercado.