

Programação para Ciência de Dados

Pré-processamento de Dados - Limpeza e Tratamento

Arthur Casals

4 de Novembro de 2025

- ▶ Última entrega - prática de 14/11 (sexta-feira)

Projeto Integrador Final

Informações Gerais

O que é:

- ▶ Trabalho final que integra **todo o conteúdo do curso**
- ▶ Análise completa de dados reais
- ▶ Simulação de cenário profissional
- ▶ Trabalho **INDIVIDUAL**

Peso e Pontuação:

- ▶ **30% da nota final**
- ▶ Total: 100 pontos base

Projeto Integrador Final

Informações Gerais

Divulgação do enunciado:

- ▶ Enunciado completo: **Última aula prática (14/11/25)**
- ▶ Materiais de apoio no Moodle
- ▶ Slides detalhados
- ▶ Rubrica de avaliação

Resultado esperado:

- ▶ Notebook Jupyter (Google Colab)
- ▶ Tempo estimado: **5-6 horas**
- ▶ Entrega: **17/11/25, 18:00 (São Paulo)**
- ▶ Entrega feita via Moodle (arquivo .ipynb)

O que Será Avaliado

Competências Técnicas (70%):

- ▶ Python, NumPy, Pandas
- ▶ Análise exploratória de dados
- ▶ Pré-processamento e limpeza
- ▶ Feature engineering
- ▶ Análise multidimensional
- ▶ Visualização de dados
- ▶ Qualidade e eficiência do código

Soft Skills (30%):

▶ Documentação

- ▶ Clareza nas explicações
- ▶ Estrutura organizada
- ▶ Interpretação de resultados

▶ Visão de Negócio

- ▶ Comunicação para não-técnicos
- ▶ Priorização estratégica
- ▶ Recomendações acionáveis

⚠ Atenção

O projeto avalia não apenas **o que você sabe fazer**, mas também **como você comunica** seus achados e insights.

Como Funciona a Avaliação

Avaliação em Dois Níveis:

1. Verificação Automática

- ▶ Via Otter-Grader
- ▶ Verifica requisitos mínimos:
 - ▶ Estrutura correta
 - ▶ Elementos obrigatórios presentes
 - ▶ Notebook executa sem erros
- ▶ Garante que você atendeu o básico

2. Avaliação Manual

- ▶ Feita pelo professor
- ▶ Avalia qualidade e profundidade
- ▶ Diferencial entre notas

CrITÉrios de Excelência:

- ▶ Análises **profundas** (não superficiais)
- ▶ Insights **não óbvios**
- ▶ Interpretações **claras e conectadas**
- ▶ Código **limpo e eficiente**
- ▶ Comunicação **adaptada ao público**
- ▶ Recomendações **específicas e implementáveis**

💡 Nota Importante

Não basta fazer, tem que fazer **bem feito e comunicar bem.**

Como Se Preparar Agora

Revisão de Conteúdo:

- ▶ Revise notebooks das aulas anteriores
- ▶ Refaça exercícios das práticas independentes
- ▶ Pratique especialmente:
 - ▶ Pandas: groupby, pivot, merge
 - ▶ EDA completa
 - ▶ Visualizações com Seaborn
 - ▶ Feature engineering
 - ▶ Aulas 9-12

Habilidades Não-Técnicas:

- ▶ Pratique escrever textos explicativos
- ▶ Interprete visualizações em voz alta

Organize Seu Tempo:

- ▶ O projeto precisa de **5-6 horas**
- ▶ **Não deixe para última hora**

Recursos Disponíveis:

- ▶ Horário de atendimento (use!)
- ▶ Fórum do Moodle
- ▶ Material completo das aulas
- ▶ Documentação oficial

Agenda

- ▶ Introdução: Por que Limpar Dados?
- ▶ Bloco 1: Identificando Dados Faltantes
- ▶ Bloco 2: Estratégias de Imputação
- ▶ Bloco 3: Limpeza e Validação de Dados
- ▶ Bloco 4: Pipeline Completo e Documentação

Bloco 0

Introdução e Contexto

A Realidade dos Dados

Dados do mundo real são SEMPRE sujos!

O que você vai encontrar:

- ▶ Valores faltantes (missing data)
- ▶ Erros de digitação
- ▶ Formatos inconsistentes
- ▶ Duplicatas
- ▶ Outliers problemáticos
- ▶ Dados impossíveis (idade = -5, área = 0)
- ▶ Unidades misturadas (metros e pés)

 **Atenção**

Estudos mostram: 80% do tempo em Data Science é gasto limpando dados!

Por que Limpeza de Dados Importa?

Garbage In, Garbage Out (GIGO)

Impacto de dados sujos:

▶ **Análises erradas**

- ▶ Conclusões incorretas
- ▶ Decisões ruins de negócio
- ▶ Perda de credibilidade

▶ **Modelos ruins**

- ▶ Baixa acurácia
- ▶ Previsões não confiáveis
- ▶ Overfitting em ruído

▶ **Perda de tempo**

- ▶ Retrabalho constante
- ▶ Debugging difícil
- ▶ Frustração da equipe

Impacto Financeiro: Exemplo Real

Caso: Empresa de E-commerce

Problema:

- ▶ Sistema de recomendação usando dados sujos
- ▶ 30% dos produtos com descrições incompletas
- ▶ Preços com erros de digitação (R\$ 1.0 em vez de R\$ 100)
- ▶ Categorias inconsistentes

Resultado:

- ▶ Recomendações irrelevantes
- ▶ Clientes frustrados
- ▶ Perda de vendas: R\$ 2 milhões/ano

Solução:

- ▶ Investimento em limpeza: R\$ 100 mil
- ▶ ROI: 20x em um ano!

Caso Real: Preços de Imóveis

Cenário: Startup de análise imobiliária

Problemas encontrados nos dados:

- ▶ **Missing:** 15% dos preços faltando
- ▶ **Formato:** Áreas como "N/A", "unknown", "???"
- ▶ **Impossíveis:** Quartos = 0, -1, 99
- ▶ **Inconsistência:** Datas DD/MM vs MM/DD
- ▶ **Duplicatas:** 5% registros repetidos
- ▶ **Outliers:** Preços R\$ 1 ou R\$ 999.999.999
- ▶ **Texto:** Bairros com variações de escrita

Resultado:

- ▶ Modelo inicial: erro de 80% (inutilizável)
- ▶ Após limpeza: erro de 15% (excelente)

Estatísticas Sobre Qualidade de Dados

Pesquisas revelam:

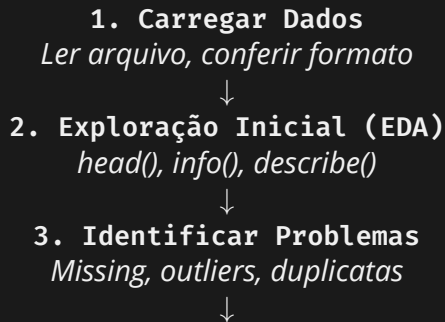
- ▶ **80-90%** do tempo é gasto em preparação de dados
- ▶ **1 em 3** analistas cita dados ruins como maior problema
- ▶ **40%** das iniciativas de BI falham por dados ruins
- ▶ Custo anual de dados ruins: **trilhões de dólares** globalmente

Tipos de problemas mais comuns:

1. Valores faltantes (presente em 90% dos datasets)
2. Inconsistências de formato (70%)
3. Duplicatas (50%)
4. Valores impossíveis (40%)
5. Erros de digitação (60%)

Pipeline de Pré-processamento

Ordem típica das operações:



Pipeline de Pré-processamento

3. Identificar Problemas

Missing, outliers, duplicatas



4. Tratar Valores Faltantes (HOJE)

Imputação, remoção



5. Limpar Inconsistências (HOJE)

Padronização, validação



6. Transformar (Próxima Aula)

Encoding, normalização

Tipos de Problemas de Qualidade

Framework de qualidade de dados:

1. Completude (Completeness)

- ▶ Dados faltantes, valores nulos
- ▶ **Exemplo:** 20% dos emails estão vazios

2. Validade (Validity)

- ▶ Valores fora do range esperado
- ▶ **Exemplo:** Idade = 200 anos

3. Consistência (Consistency)

- ▶ Formatos diferentes para mesma informação
- ▶ **Exemplo:** "SP", "São Paulo", "Sao Paulo"

4. Unicidade (Uniqueness)

- ▶ Registros duplicados
- ▶ **Exemplo:** Mesmo cliente cadastrado 3 vezes

Tipos de Problemas de Qualidade

5. Acurácia (Accuracy)

Definição:

- ▶ Dados refletem a realidade?
- ▶ Medições estão corretas?

Exemplos de problemas:

- ▶ Endereço cadastrado está desatualizado
- ▶ Sensor descalibrado gerando medições erradas
- ▶ Dados transcritos incorretamente
- ▶ Conversão de unidades errada (km \rightarrow m)

Como detectar:

- ▶ Validação com fonte externa
- ▶ Verificação amostral manual
- ▶ Comparação com benchmarks conhecidos

Dataset de Hoje: Housing Prices

Dados de preços de imóveis com problemas reais

Variáveis principais:

- ▶ **price:** Preço do imóvel (R\$)
- ▶ **area:** Área total em m²
- ▶ **bedrooms:** Número de quartos
- ▶ **bathrooms:** Número de banheiros
- ▶ **age:** Idade do imóvel (anos desde construção)
- ▶ **neighborhood:** Bairro/região
- ▶ **condition:** Estado de conservação
- ▶ **garage_spaces:** Vagas de garagem
- ▶ **has_garden:** Possui jardim (sim/não)

Dataset de Hoje: Housing Prices

Tamanho:

- ▶ 3000 registros
- ▶ 9 variáveis

Origem:

- ▶ The Ames Housing Dataset
- ▶ Subconjunto de features (82 \rightarrow 9)
- ▶ Poluído propositalmente

Problemas Inseridos no Dataset

Problemas propositalmente incluídos: Valores Faltantes:

- ▶ price: 12% missing
- ▶ area: 8% missing
- ▶ age: 15% missing
- ▶ neighborhood: 5% missing

Inconsistências:

- ▶ Bairros com nomes variados: "Centro", "centro", "CENTRO"
- ▶ condition: "Excelente", "Muito Bom", "Bom", etc.
- ▶ Espaços extras em strings

Problemas Inseridos no Dataset

Problemas propositalmente incluídos:

Valores Impossíveis:

- ▶ bedrooms = -1, 0, 99
- ▶ area = 0, 5, 10000 m²
- ▶ age = -5, 250 anos
- ▶ price = R\$ 1, R\$ 9999999999

Objetivos de Aprendizado

Conhecimento:

- ▶ Identificar tipos de problemas de qualidade
- ▶ Classificar dados faltantes (MCAR, MAR, MNAR)
- ▶ Conhecer múltiplas estratégias de imputação

Habilidades:

- ▶ Aplicar técnicas de imputação apropriadas
- ▶ Limpar inconsistências de formato
- ▶ Validar ranges e tipos de dados
- ▶ Construir pipeline de limpeza reproduzível

Atitudes:

- ▶ Sempre documentar decisões
- ▶ Ser conservador com remoção de dados
- ▶ Validar resultados sistematicamente

Pandas: Nossa ferramenta principal

Detecção de problemas:

- ▶ `isnull()`, `notnull()`, `isna()`
- ▶ `duplicated()`
- ▶ `describe()`, `value_counts()`

Tratamento de missing:

- ▶ `fillna()`, `dropna()`
- ▶ `interpolate()`
- ▶ `groupby().transform()`

Ferramentas de Hoje

Limpeza de strings:

- ▶ `str.strip()`, `str.lower()`, `str.upper()`
- ▶ `str.replace()`, `str.extract()`

Validação e conversão:

- ▶ `astype()`, `to_numeric()`, `to_datetime()`
- ▶ `between()`, `isin()`

Filosofia da Limpeza de Dados

1. Sempre documente suas decisões

- ▶ Por que removeu ou alterou algo?
- ▶ Qual critério utilizou?
- ▶ Quantos registros foram afetados?

2. Mantenha os dados originais intactos

- ▶ Sempre trabalhe em uma cópia
- ▶ Possibilidade de reverter decisões
- ▶ Comparar antes/depois

3. Seja conservador

- ▶ Dados são valiosos
- ▶ Não remova sem razão forte
- ▶ Prefira transformar a deletar

Filosofia da Limpeza de Dados (cont.)

4. Entenda o contexto do negócio

- ▶ Consulte especialistas do domínio
- ▶ O que faz sentido no negócio?
- ▶ Há regras específicas do setor?

5. Seja transparente

- ▶ Reporte tudo que foi feito
- ▶ Mostre impacto das decisões
- ▶ Permita reprodução do processo

6. Valide, valide, valide

- ▶ Sempre verifique resultados
- ▶ Compare distribuições antes/depois
- ▶ Testes de sanidade (sanity checks)

Fluxo Sistemático de Limpeza

Processo iterativo em 4 fases:

FASE 1: Auditoria Inicial (Assess)

- ▶ Carregar dados
- ▶ Dimensões: quantos registros? Features?
- ▶ Tipos de dados
- ▶ Primeira inspeção visual

FASE 2: Diagnóstico (Diagnose)

- ▶ Identificar todos os problemas
- ▶ Quantificar missing, duplicatas, outliers
- ▶ Classificar severidade
- ▶ Priorizar o que atacar

Fluxo Sistemático de Limpeza (cont.)

FASE 3: Planejamento (Plan)

- ▶ Definir estratégias para cada problema
- ▶ Estabelecer critérios de validação
- ▶ Criar regras de negócio
- ▶ Decidir ordem de execução

FASE 4: Execução e Validação (Execute & Validate)

- ▶ Aplicar transformações
- ▶ Verificar resultados após cada passo
- ▶ Comparar antes/depois
- ▶ Documentar mudanças
- ▶ Gerar relatório final

Fluxo Sistemático de Limpeza (cont.)

Nota Importante

Este é um processo iterativo - você voltará às fases anteriores!

Setup Inicial: Imports

</> Python

```
1 # Imports necessarios
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 # Configuracoes de display
7 pd.set_option('display.max_columns', None)
8 pd.set_option('display.max_rows', 100)
9 pd.set_option('display.float_format', '{:.2f}'.format)
10 # Configuracoes de plots
11 plt.style.use('seaborn-v0_8-darkgrid')
12 sns.set_palette("husl")
13 print("Bibliotecas carregadas com sucesso!")
```

Carregar Dados: Primeira Etapa

</> Python

```
1 # Carregar dataset
2 df = pd.read_csv('housing_prices_dirty.csv')
3 print("="*50)
4 print("DADOS CARREGADOS")
5 print("="*50)
6 print(f"\nShape: {df.shape}")
7 print(f"Linhas: {df.shape[0]}")
8 print(f"Colunas: {df.shape[1]}")
9 print(f"\nMemoria utilizada: {df.memory_usage(deep=True).sum() /
    1024**2:.2f} MB")
10
11 # SEMPRE salve uma copia dos dados originais
12 df_original = df.copy()
13 print("\nCopia de seguranca criada!")
```


Carregar Dados: Primeira Etapa (cont.)

</> Python

```
1 # =====
2 # DADOS CARREGADOS
3 # =====
4 #
5 # Shape: (3076, 9)
6 # Linhas: 3076
7 # Colunas: 9
8 #
9 # Memoria utilizada: 0.74 MB
10 #
11 # Copia de seguranca criada!
```

Inspeção Visual: Primeiras Linhas

</> Python

```
1 # Ver primeiras linhas
2 print("=== PRIMEIRAS 5 LINHAS ===")
3 print(df.head())
4
5 # Ver ultimas linhas (detectar problemas no fim)
6 print("\n=== ULTIMAS 5 LINHAS ===")
7 print(df.tail())
8
9 # Amostra aleatoria (visao geral melhor)
10 print("\n=== AMOSTRA ALEATORIA (10 linhas) ===")
11 print(df.sample(10, random_state=42))
12
13 # Sempre olhe head(), tail() E sample() - cada um revela coisas diferentes
```

Info: Visão Geral do Dataset

</> Python

```
1 # Informacoes gerais
2 print("=== INFORMACOES DO DATASET ===")
3 df.info()
4
5 # Interpretacao:
6 # - RangeIndex: indices numericos
7 # - Data columns: numero de colunas
8 # - Non-Null Count: quantos NAO sao nulos
9 # - Dtype: tipo de cada coluna
10 # - memory usage: memoria RAM usada
```

Info: Visão Geral do Dataset (cont.)

O que observar:

- ▶ Colunas com Non-Null Count < total → têm dados faltantes!
- ▶ Tipos incorretos (idade como object?)
- ▶ Uso excessivo de memória

Describe: Estatísticas Descritivas

</> Python

```
1 # Estatísticas das numericas
2 print("=== ESTATISTICAS NUMERICAS ===")
3 print(df.describe())
4
5 # Estatísticas das categoricas
6 print("\n=== ESTATISTICAS CATEGORICAS ===")
7 print(df.describe(include=['object', 'category']))
8
9 # Transpor para melhor visualizacao
10 print("\n=== TRANSPOSTO ===")
11 print(df.describe().T)
```

Describe: Estatísticas Descritivas

O que buscar:

- ▶ Min/max estranhos (negativos, zeros, muito altos)
- ▶ Média muito diferente da mediana (assimetria)
- ▶ Desvio padrão muito alto (outliers?)

Tipos de Dados: Verificação

</> Python

```
1 # Ver tipos de cada coluna
2 print("=== TIPOS DE DADOS ===")
3 print(df.dtypes)
4
5 # Contar tipos
6 print("\n=== CONTAGEM POR TIPO ===")
7 print(df.dtypes.value_counts())
8
9 # Separar por tipo
10 numericas = df.select_dtypes(include=['number']).columns
11 categoricas = df.select_dtypes(include=['object']).columns
12
13 print(f"\nNumericas ({len(numericas)}): {list(numericas)}")
14 print(f"Categoricas ({len(categoricas)}): {list(categoricas)}")
```

Valores Únicos: Primeira Análise

</> Python

```
1 # Contar valores unicos em cada coluna
2 print("=== VALORES UNICOS POR COLUNA ===")
3 for col in df.columns:
4     n_unique = df[col].nunique()
5     pct_unique = (n_unique / len(df)) * 100
6     print(f"{col:20s}: {n_unique:4d} unicos ({pct_unique:5.1f}%)")
7
8 # Valores unicos das categoricas
9 print("\n=== CATEGORIAS UNICAS ===")
10 for col in categoricas:
11     print(f"\n{col}:")
12     print(df[col].value_counts())
```


Checklist de Auditoria Inicial

Antes de começar a limpar, responda:

1. ☐ Quantas linhas e colunas tenho?
2. ☐ Quais são numéricas? Quais são categóricas?
3. ☐ Os tipos de dados estão corretos?
4. ☐ Há valores faltantes? Em quais colunas?
5. ☐ Quais são os ranges das variáveis numéricas?
6. ☐ Há valores impossíveis (negativos, zeros, extremos)?
7. ☐ As categóricas têm quantas categorias?
8. ☐ Há inconsistências visíveis (espaços, maiúsculas)?
9. ☐ Há duplicatas aparentes?
10. ☐ Os dados fazem sentido no contexto do negócio?

Atenção

Só avance quando puder responder todas essas perguntas!

Exercício Guiado: Auditoria Completa

Exercício Prático

Execute a auditoria inicial do Housing dataset:

1. Carregue os dados e crie cópia de segurança
2. Execute todos os comandos de inspeção:

```
1 df.head()
2 df.tail()
3 df.sample(10)
4 df.info()
5 df.describe()
6 df.dtypes
```

3. Anote TODOS os problemas que identificar
4. Classifique por tipo (missing, invalid, inconsistent)
5. Estime severidade (crítico, moderado, leve)

Transição para Próximo Bloco

Recapitulando Introdução:

Aprendemos:

- ▶ Por que limpeza é crítica (80% do tempo!)
- ▶ Tipos de problemas de qualidade
- ▶ Filosofia e princípios de limpeza
- ▶ Workflow sistemático em 4 fases
- ▶ Como fazer auditoria inicial completa

Agora vamos para o Bloco 1:

- ▶ Identificação detalhada de dados faltantes
- ▶ Classificação: MCAR, MAR, MNAR
- ▶ Análise de padrões de missing
- ▶ Decisão: quando remover vs imputar

Bloco 1

Identificando e Entendendo Dados Faltantes

O que São Dados Faltantes?

Missing data = ausência de valor onde deveria existir um

Como aparecem em diferentes contextos:

- ▶ **Pandas/NumPy:** NaN (Not a Number)
- ▶ **Python nativo:** None
- ▶ **Bancos SQL:** NULL
- ▶ **Excel:** Células vazias
- ▶ **CSV:** (string vazia)
- ▶ **Códigos especiais:** -999, 9999, "N/A", "unknown", "?"

Atenção

Cuidado: nem tudo que parece missing é missing! Zero pode ser válido.

Por que Valores Ficam Faltando?

Razões técnicas:

- ▶ Sensor/equipamento quebrado
- ▶ Falha na transmissão de dados
- ▶ Erro no sistema de coleta
- ▶ Problema de integração entre sistemas
- ▶ Corrupção de arquivo

Razões humanas:

- ▶ Esquecimento ao preencher formulário
- ▶ Erro de digitação
- ▶ Informação não disponível no momento
- ▶ Pessoa não quis responder (privacidade)
- ▶ Campo não aplicável ao caso

Exemplos Reais de Missing Data

Saúde:

- ▶ Paciente não fez exame de sangue → colesterol faltando
- ▶ Sensor de frequência cardíaca desconectou

E-commerce:

- ▶ Cliente não preencheu telefone no cadastro
- ▶ Produto sem categoria definida

Financeiro:

- ▶ Renda não declarada em formulário
- ▶ Valor de transação perdido em falha de rede

Imóveis (nosso caso):

- ▶ Preço não definido (negociação em andamento)
- ▶ Área não medida ainda
- ▶ Idade desconhecida (imóvel muito antigo)

Identificar Missing: Comandos Básicos

</> Python

```
1 # Verificar se ha valores faltantes
2 print("=== HA VALORES FALTANTES? ===")
3 print(df.isnull().any())
4 # True = tem missing, False = sem missing
5
6 # Contar missing por coluna
7 print("\n=== CONTAGEM DE MISSING ===")
8 print(df.isnull().sum())
9
10 # Total de missing no dataset inteiro
11 total_missing = df.isnull().sum().sum()
12 print(f"\nTotal de valores faltantes: {total_missing}")
```


Percentual de Missing

Python

```
1 # Calcular percentual de missing
2 print("=== PERCENTUAL DE MISSING ===")
3 missing_count = df.isnull().sum()
4 missing_pct = (missing_count / len(df)) * 100
5 # Criar DataFrame resumo
6 missing_summary = pd.DataFrame({
7     'Missing_Count': missing_count,
8     'Percentage': missing_pct
9 })
10 # Ordenar por percentual (decrecente)
11 missing_summary = missing_summary.sort_values(
12     'Percentage', ascending=False
13 )
14 # Mostrar apenas colunas com missing
15 print(missing_summary[missing_summary['Missing_Count'] > 0])
```

Interpretando Percentuais de Missing

Guia de decisão baseado em %:

0-5%: PROBLEMA PEQUENO

- ▶ ✓ Pode remover linhas com segurança
- ▶ ✓ Ou imputar com método simples
- ▶ Impacto mínimo nas análises
- ▶ Decisão rápida

5-20%: PROBLEMA MODERADO

- ▶ Remover pode perder muita informação
- ▶ Imputação deve ser cuidadosa
- ▶ Analise padrão de missing (é aleatório?)
- ▶ Considere imputação por grupos
- ▶ Documente estratégia escolhida

Interpretando Percentuais (continuação)

20-40%: PROBLEMA SÉRIO

- ▶ ✗ Muito missing para ignorar
- ▶ ✗ Imputação pode distorcer resultados
- ▶ Opções:
 - ▶ Coletar mais dados
 - ▶ Usar modelos que lidam com missing (alguns ML)
 - ▶ Criar flag "was_missing"+ imputação
 - ▶ Considerar remover a variável

40%: PROBLEMA CRÍTICO

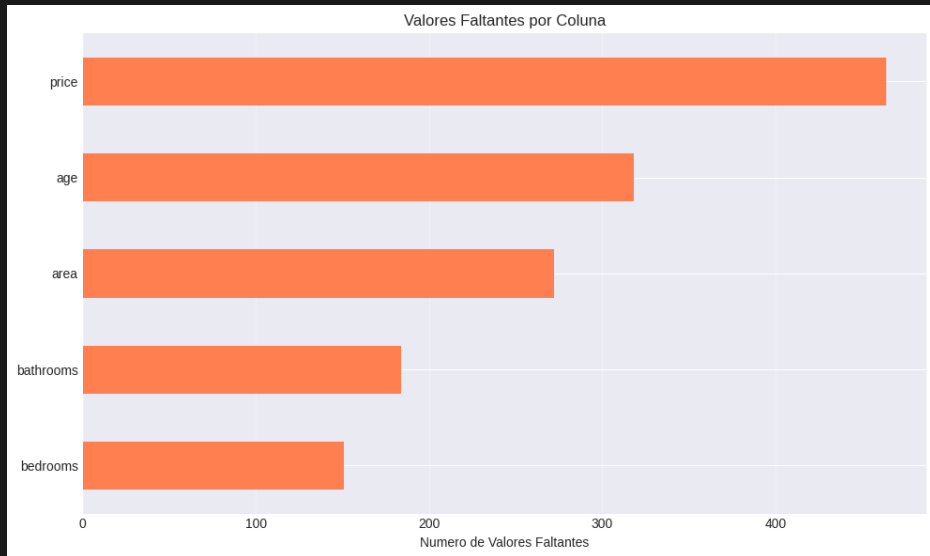
- ▶ ✗ Variável praticamente inútil
- ▶ ✗ Mais missing que dados!
- ▶ Recomendação: remover a coluna inteira
- ▶ Exceção: se missing é a informação (ex: "não respondeu")

Visualizar Missing: Gráfico de Barras

</> Python

```
1 import matplotlib.pyplot as plt
2
3 # Criar grafico de missing
4 missing_data = df.isnull().sum()
5 missing_data = missing_data[missing_data > 0].sort_values()
6
7 fig, ax = plt.subplots(figsize=(10, 6))
8 missing_data.plot(kind='barh', ax=ax, color='coral')
9
10 ax.set_xlabel('Numero de Valores Faltantes')
11 ax.set_title('Valores Faltantes por Coluna')
12 ax.grid(True, alpha=0.3, axis='x')
13
14 plt.tight_layout()
15 plt.show()
```

Visualizar Missing: Gráfico de Barras

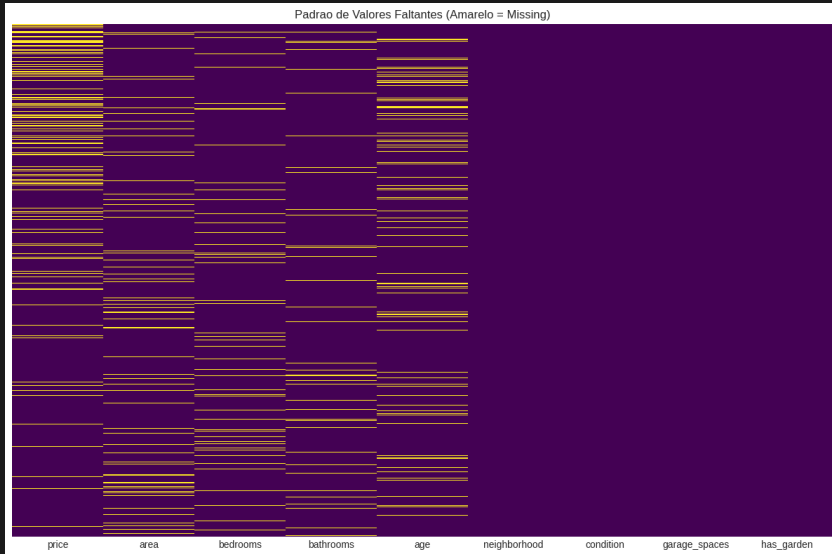


Visualizar Missing: Heatmap

Python

```
1 import seaborn as sns
2 # Heatmap de valores faltantes
3 fig, ax = plt.subplots(figsize=(12, 8))
4 # Criar matriz: True = missing, False = presente
5 missing_matrix = df.isnull()
6 # Plotar (amarelo = missing, roxo = presente)
7 sns.heatmap(missing_matrix,
8             yticklabels=False,
9             cbar=False,
10             cmap='viridis',
11             ax=ax)
12 ax.set_title('Padrao de Valores Faltantes (Amarelo = Missing)')
13 plt.tight_layout()
14 plt.show()
```

Visualizar Missing: Heatmap



Interpretando o Heatmap de Missing

O que procurar no heatmap:

Padrões Verticais (colunas):

- ▶ Linhas amarelas contínuas → muitas linhas com missing
- ▶ Linhas esparsas → missing aleatório

Padrões Horizontais (linhas):

- ▶ Colunas amarelas → variável com muito missing
- ▶ Colunas roxas → variável completa

Padrões Correlacionados:

- ▶ Duas colunas amarelas nos mesmos lugares
- ▶ Significa: quando uma falta, outra também falta
- ▶ Pode indicar problema sistemático

Missing por Linha (Casos Completos)

</> Python

```
1 # Quantas linhas estao completas (sem missing)?
2 complete_rows = df.notnull().all(axis=1).sum()
3 incomplete_rows = len(df) - complete_rows
4
5 print(f"Linhas completas: {complete_rows} ({complete_rows/len(df)*100:.1f
6      }%)")
7
8 print(f"Linhas com missing: {incomplete_rows} ({incomplete_rows/len(df)
9      *100:.1f}%)")
10
11 # Contar quantos missing por linha
12 missing_per_row = df.isnull().sum(axis=1)
13
14 print("\n=== DISTRIBUICAO DE MISSING POR LINHA ===")
15 print(missing_per_row.value_counts().sort_index())
```

Decisão: Remover vs Imputar

Árvore de decisão:

Quanto missing há?



< 5% → Pode remover linhas

5-20% → Imputação recomendada

20% → Análise profunda necessária



É variável crítica?



SIM (ex: target) → Remover linhas

NÃO (ex: feature) → Imputar

Tipos de Missing: Introdução

Nem todo missing é igual!

Três categorias fundamentais:

1. MCAR - Missing Completely At Random

- ▶ Ausência é 100% aleatória
- ▶ Não relacionada a NADA (nem observado, nem não observado)

2. MAR - Missing At Random

- ▶ Ausência relacionada a variáveis observadas
- ▶ Podemos explicar o missing com outras colunas

3. MNAR - Missing Not At Random

- ▶ Ausência relacionada ao próprio valor faltante
- ▶ Missing porque o valor é alto/baixo/específico

MCAR: Missing Completely At Random

Características:

- ▶ Probabilidade de missing é igual para todos
- ▶ Ausência não depende de nenhuma variável
- ▶ É o tipo mais "inofensivo" de missing
- ▶ Raro na prática

Exemplo clássico:

- ▶ Técnico derrubou café no equipamento
- ▶ Equipamento quebrou aleatoriamente
- ▶ Sorteio aleatório de quem responde pergunta
- ▶ Bug de software afetando amostras aleatórias

Implicação:

- ▶ ✓ Remover linhas não introduz viés
- ▶ ✓ Qualquer método de imputação funciona bem

MCAR: Exemplo Prático

Cenário: Medição de Pressão Arterial

Situação:

- ▶ Estudo com 1000 pacientes
- ▶ Medidor de pressão quebrou em dia específico
- ▶ 50 medições perdidas (5%)
- ▶ Não há padrão: jovens, idosos, saudáveis, doentes

ID	Idade	Peso	Pressão	Dia
1	25	70	120	01/10
2	30	75	NaN	15/10 ← quebrou
3	35	80	130	01/10
4	40	85	NaN	15/10 ← quebrou
5	45	90	125	02/10

MCAR: Exemplo Prático

Análise:

- ▶ Missing não relacionado a idade, peso ou pressão real
- ▶ É MCAR!

MAR: Missing At Random

Características:

- ▶ Ausência relacionada a outras variáveis observadas
- ▶ Podemos "explicar" o missing com dados que temos
- ▶ Mais comum que MCAR
- ▶ Requer tratamento mais sofisticado

Exemplos:

- ▶ Homens não respondem pergunta sobre gravidez
- ▶ Jovens pulam pergunta sobre aposentadoria
- ▶ Pessoas sem carro deixam vazio "marca do carro"
- ▶ Imóveis sem garagem → vagas = NaN

Implicação:

- ▶ Remover linhas pode introduzir viés
- ▶ ✓ Imputação por grupos funciona bem
- ▶ ✓ Podemos usar outras variáveis para prever

MAR: Exemplo Prático

Cenário: Pesquisa de Saúde

Situação:

- ▶ Pergunta sobre "quantos filhos tem?"
- ▶ Muitos jovens (< 25 anos) não responderam
- ▶ Idosos responderam quase todos

ID	Idade	Sexo	Filhos
1	22	F	NaN ← jovem
2	24	M	NaN ← jovem
3	45	F	2
4	50	M	3
5	23	F	NaN ← jovem
6	48	F	1

MAR: Exemplo Prático

Análise:

- ▶ Missing relacionado à IDADE (observada)
- ▶ Jovens ainda não têm filhos ou não quiseram responder
- ▶ É MAR!

MNAR: Missing Not At Random

Características:

- ▶ Ausência relacionada ao próprio valor que falta
- ▶ Valor está faltando PORQUE é alto, baixo ou específico
- ▶ Tipo mais problemático
- ▶ Introduz viés sistemático

Exemplos:

- ▶ Pessoas com salário muito alto não revelam
- ▶ Deprimidos não respondem sobre humor
- ▶ Muito endividados omitem dívidas
- ▶ Alunos com notas baixas faltam à prova
- ▶ Imóveis muito caros sem preço (em negociação)

Implicação:

- ▶ ✗ Difícil de detectar
- ▶ ✗ Qualquer imputação pode ser enviesada
- ▶ ✗ Pode exigir coletar mais dados

MNAR: Exemplo Prático

Cenário: Pesquisa Salarial

Situação:

- ▶ Pergunta: "Qual seu salário mensal?"
- ▶ Executivos e CEOs não responderam
- ▶ Porque salário é MUITO alto (privacidade)

ID	Cargo	Salário
1	Analista	R\$ 5.000
2	Gerente	R\$ 15.000
3	Diretor	NaN ← muito alto!
4	CEO	NaN ← muito alto!
5	Coordenador	R\$ 10.000
6	VP	NaN ← muito alto!

MNAR: Exemplo Prático

Análise:

- ▶ Missing porque o salário é alto
- ▶ Não podemos ver o salário, mas sabemos que falta porque é alto
- ▶ É MNAR!

Comparação: MCAR vs MAR vs MNAR

Aspecto	MCAR	MAR	MNAR
Aleatoriedade	Completamente aleatório	Relacionado a observados	Relacionado ao próprio valor
Frequência	Raro	Comum	Moderado
Detectabilidade	Fácil	Moderada	Difícil
Viés ao remover	Nenhum	Possível	Sério
Imputação simples	Funciona bem	Funciona razoavelmente	Problemática
Tratamento	Qualquer método	Por grupos	Muito cuidadoso

Investigar Tipo de Missing: Análise Exploratória

</> Python

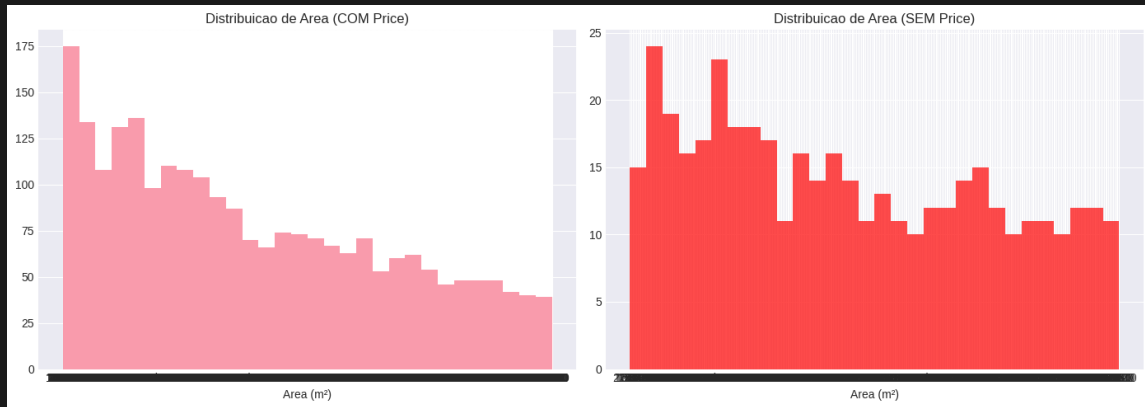
```
1 # Comparar características de quem tem vs nao tem missing
2 # Exemplo: price faltante
3 # Grupo COM price
4 with_price = df[df['price'].notnull()]
5 # Grupo SEM price
6 without_price = df[df['price'].isnull()]
7 print("=== COM PRECO ===")
8 print(with_price[['area', 'bedrooms', 'age']].describe())
9 print("\n=== SEM PRECO ===")
10 print(without_price[['area', 'bedrooms', 'age']].describe())
11
12 # Sao muito diferentes? Pode ser MAR ou MNAR
13 # Sao similares? Provavelmente MCAR
```

Teste Visual: Comparar Distribuições

</> Python

```
1 import matplotlib.pyplot as plt
2 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
3 # Distribuicao de area: COM price
4 axes[0].hist(with_price['area'].dropna(),
5              bins=30, alpha=0.7, label='Com Price')
6 axes[0].set_title('Distribuicao de Area (COM Price)')
7 axes[0].set_xlabel('Area (m²)')
8 # Distribuicao de area: SEM price
9 axes[1].hist(without_price['area'].dropna(),
10             bins=30, alpha=0.7, color='red',
11             label='Sem Price')
12 axes[1].set_title('Distribuicao de Area (SEM Price)')
13 axes[1].set_xlabel('Area (m²)')
14 plt.tight_layout()
15 plt.show()
```

Teste Visual: Comparar Distribuições

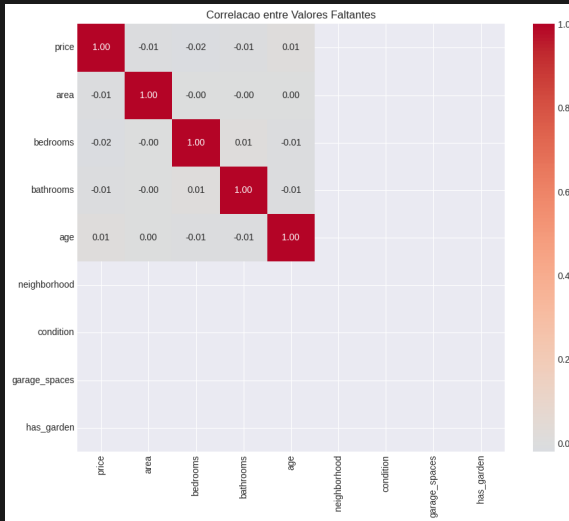


Correlação entre Missings

</> Python

```
1 # Verificar se missings ocorrem juntos
2 # Criar matriz de correlacao de missing
3 missing_matrix = df.isnull().astype(int)
4 missing_corr = missing_matrix.corr()
5 # Plotar heatmap
6 import seaborn as sns
7 fig, ax = plt.subplots(figsize=(10, 8))
8 sns.heatmap(missing_corr, annot=True, fmt='.2f',
9             cmap='coolwarm', center=0,
10            square=True, ax=ax)
11 ax.set_title('Correlacao entre Valores Faltantes')
12 plt.tight_layout()
13 plt.show()
14 # Correlacao alta? Missings ocorrem juntos
```

Correlação entre Missings



Como Identificar o Tipo na Prática

Checklist de investigação:

1. Consulte especialistas do domínio

- ▶ Como os dados foram coletados?
- ▶ Há razões conhecidas para missing?
- ▶ Processo de coleta introduz viés?

2. Análise exploratória

- ▶ Compare grupos com/sem missing
- ▶ Distribuições são similares? (MCAR)
- ▶ Diferem em variáveis observadas? (MAR)

3. Use lógica e contexto

- ▶ Missing faz sentido dado o valor?
- ▶ Há incentivo para não responder?
- ▶ É questão sensível?

Decisão de Tratamento por Tipo

Se identificou como MCAR:

- ▶ ✓ Pode remover linhas se $<5\%$
- ▶ ✓ Imputação simples (média/mediana) funciona
- ▶ ✓ Sem grandes preocupações

Se identificou como MAR:

- ▶ Cuidado ao remover (pode introduzir viés)
- ▶ ✓ Use imputação por grupos
- ▶ ✓ Considere variáveis relacionadas
- ▶ ✓ Modelos preditivos podem ajudar

Se identificou como MNAR:

- ▶ ✗ Muito difícil de tratar bem
- ▶ Qualquer imputação pode ser enviesada
- ▶ Considere: criar flag, coletar mais dados, ou remover variável

Caso de Uso: Housing Prices Missing

Análise do nosso dataset:

price missing (12%):

- ▶ Hipótese: Imóveis em negociação ainda
- ▶ Pode ser MNAR se preços muito altos estão faltando
- ▶ Investigar: comparar area/bedrooms de missing vs não-missing

area missing (8%):

- ▶ Hipótese: Não foi medido ainda
- ▶ Provavelmente MCAR (aleatório)
- ▶ Pode imputar com mediana ou por bairro

age missing (15%):

- ▶ Hipótese: Imóveis muito antigos (sem registro)
- ▶ Pode ser MNAR!
- ▶ Cuidado: age faltante pode significar "muito velho"

Criar Relatório de Missing

</> Python

```
1 def generate_missing_report(df):
2     """Gera relatorio completo de missing"""
3     missing_count = df.isnull().sum()
4     missing_pct = (missing_count / len(df)) * 100
5     report = pd.DataFrame({
6         'Column': df.columns,
7         'Missing_Count': missing_count.values,
8         'Missing_Pct': missing_pct.values,
9         'Dtype': df.dtypes.values
10    })
11    report = report[report['Missing_Count'] > 0]
12    report = report.sort_values('Missing_Pct', ascending=False)
13    return report
14 # Usar
15 print(generate_missing_report(df))
```

Documentar Classificação de Missing

</> Python

```
1 # Criar dicionario documentando tipo de missing
2 missing_classification = {
3     'price': {
4         'type': 'MNAR',
5         'reason': 'Imoveis caros em negociacao',
6         'evidence': 'Area maior em missing',
7         'strategy': 'Imputar por grupo + flag'
8     },
9     'area': {
10        'type': 'MCAR',
11        'reason': 'Nao medido ainda',
12        'evidence': 'Distribuicao similar',
13        'strategy': 'Imputar com mediana'
14    },
```

Documentar Classificação de Missing (cont.)

</> Python

```
1     'age': {  
2         'type': 'MNAR',  
3         'reason': 'Imoveis muito antigos sem registro',  
4         'evidence': 'Logica do negocio',  
5         'strategy': 'Valor alto fixo + flag'  
6     }  
7 }
```


Estratégia Geral de Tratamento

Remoção (Deletion):

- ▶ Listwise: remover linha inteira se tem qualquer missing
- ▶ Pairwise: usar linhas completas para cada análise
- ▶ Quando: MCAR + pequeno % + muitos dados disponíveis

Imputação (Imputation):

- ▶ Simples: média, mediana, moda
- ▶ Por grupos: considerar outras variáveis
- ▶ Avançada: modelos preditivos
- ▶ Quando: MAR + % moderado + poucos dados

Manter como está:

- ▶ Alguns modelos aceitam missing (XGBoost, LightGBM)
- ▶ Missing é a informação (flag binária)
- ▶ Quando: MNAR + missing é relevante

Exercício Completo: Análise de Missing

Exercício Prático

Investigue completamente o missing no Housing:

1. Identifique todas as colunas com missing
2. Para cada uma:

```
1 # a) Calcule percentual
2 # b) Compare grupo com/sem missing
3 # c) Visualize distribuicoes
4 # d) Verifique correlacao com outras vars
```

3. Classifique como MCAR, MAR ou MNAR
4. Justifique sua classificação
5. Documente estratégia recomendada
6. Crie relatório consolidado

Resumo do Bloco 1

O que aprendemos:

Identificação:

- ▶ Como detectar missing: `isnull()`, `sum()`
- ▶ Calcular percentuais e criar relatórios
- ▶ Visualizar com gráficos e heatmaps

Classificação:

- ▶ MCAR: completamente aleatório
- ▶ MAR: relacionado a observados
- ▶ MNAR: relacionado ao próprio valor

Investigação:

- ▶ Comparar grupos com/sem missing
- ▶ Analisar correlações entre missings
- ▶ Consultar especialistas do domínio

Transição para Próximo Bloco

Agora que identificamos e classificamos missing:

Bloco 2: Estratégias de Imputação

- ▶ Remoção de linhas (quando e como)
- ▶ Imputação simples (média, mediana, moda)
- ▶ Imputação por grupos (groupby + transform)
- ▶ Forward fill e backward fill
- ▶ Interpolação
- ▶ Valor constante
- ▶ Criar flags de imputação
- ▶ Casos de uso práticos

Nota Importante

Aprenderemos COMO preencher os valores faltantes de forma inteligente!

Bloco 2

Estratégias de Imputação

O que é Imputação?

Imputação = preencher valores faltantes com estimativas

Por que imputar?

- ▶ Muitos algoritmos de ML não aceitam NaN
- ▶ Remover todas as linhas pode perder muita informação
- ▶ Informação parcial é melhor que nenhuma informação
- ▶ Manter tamanho do dataset

Quando NÃO imputar?

- ▶ > 50% da coluna é missing (mais missing que dados!)
- ▶ Missing é a informação (ex: não respondeu = desinteressado)
- ▶ Dados são MNAR com forte viés
- ▶ Você tem dados suficientes após remoção

O que é Imputação?

Atenção

Imputação é sempre uma estimativa - há incerteza!

Métodos de Imputação: Panorama

Nível 1: Remoção

- ▶ Remover linhas (dropna)
- ▶ Remover colunas

Nível 2: Imputação Simples

- ▶ Média, mediana, moda
- ▶ Valor constante

Nível 3: Imputação Avançada

- ▶ Forward fill / Backward fill
- ▶ Imputação por grupos
- ▶ Interpolação

Nível 4: Modelagem (não veremos hoje)

- ▶ Regressão, KNN, Random Forest

Estratégia 1: Remoção de Linhas (Listwise)

Quando usar:

- ▶ ✓ Você tem MUITOS dados (perder 5% é OK)
- ▶ ✓ Missing é pequeno ($< 5\%$)
- ▶ ✓ Dados são MCAR
- ▶ ✓ Variável crítica está missing (ex: target)

Vantagens:

- ▶ Muito simples
- ▶ Rápido
- ▶ Sem estimativas (dados reais apenas)

Desvantagens:

- ▶ Perde informação
- ▶ Reduz tamanho do dataset
- ▶ Pode introduzir viés se não for MCAR
- ▶ Desperdiça dados parcialmente completos

Remoção: Implementação Básica

</> Python

```
1 # Contar linhas antes
2 print(f"Linhas originais: {len(df)}")
3
4 # Remover todas as linhas com QUALQUER NaN
5 df_clean = df.dropna()
6 print(f"Apos dropna(): {len(df_clean)}")
7 print(f"Perdidos: {len(df) - len(df_clean)} linhas")
8
9 # Percentual perdido
10 pct_lost = ((len(df) - len(df_clean)) / len(df)) * 100
11 print(f"Percentual perdido: {pct_lost:.1f}%")
12
13 # CUIDADO: Se perder > 20%, reconsidere!
14 if pct_lost > 20:
15     print("ALERTA: Perdendo muitos dados!")
```

Remoção: Opções Avançadas

</> Python

```
1 # Remover apenas se TODAS as colunas forem NaN
2 df_clean = df.dropna(how='all')
3 # Use: linhas completamente vazias
4
5 # Remover se NaN em colunas específicas
6 df_clean = df.dropna(subset=['price', 'area'])
7 # Use: quando certas colunas são críticas
8
9 # Remover se missing em N ou mais colunas
10 # (manter apenas se tem pelo menos 5 valores)
11 df_clean = df.dropna(thresh=5)
12
13 # Remover linhas, mas manter todas as colunas
14 df_clean = df.dropna(axis=0) # axis=0 = linhas
```

Remoção: Caso de Uso - Análise de Satisfação de Clientes

Dados:

- ▶ 100.000 respostas de pesquisa
- ▶ Pergunta crítica: "Nota de 0-10 para o produto"
- ▶ 3% não responderam esta pergunta
- ▶ Outras perguntas têm <1% missing

Decisão:

- ▶ Remover 3.000 linhas sem nota (variável crítica)
- ▶ Ainda restam 97.000 respostas (excelente!)
- ▶ Missing é pequeno e não introduz viés significativo
- ▶ Melhor que tentar "adivinhar" a nota

Resultado:

- ▶ ✓ Dataset limpo com target válido
- ▶ ✓ Análise confiável

Estratégia 2: Imputação Simples - Conceito

Preencher com estatística central

1. Média (mean):

- ▶ Soma de todos / quantidade
- ▶ Para: dados numéricos simétricos
- ▶ Sensível a outliers

2. Mediana (median):

- ▶ Valor do meio
- ▶ Para: dados numéricos com outliers
- ▶ Mais robusta (RECOMENDADA!)

3. Moda (mode):

- ▶ Valor mais frequente
- ▶ Para: dados categóricos
- ▶ Ou numéricos discretos

Imputação com Média

</> Python

```
1 # Calcular media ANTES de preencher
2 df['area'] = pd.to_numeric(df['area'], errors='coerce')
3 mean_area = df['area'].mean()
4 print(f"Media de area: {mean_area:.2f} m²")
5 # Verificar quantos missing
6 n_missing = df['area'].isnull().sum()
7 print(f"Valores faltantes: {n_missing}")
8 # Preencher NaN com a media
9 df['area'] = df['area'].fillna(mean_area)
10 # Verificar se funcionou
11 print(f"Missing apos imputacao: {df['area'].isnull().sum()}")
12 # Deve ser 0!
13 print(f"\n{n_missing} valores imputados com media")
```

Imputação com Mediana (Recomendado)

</> Python

```
1 # Mediana é mais robusta que media!
2 median_price = df['price'].median()
3 print(f"Mediana de price: R$ {median_price:,.2f}")
4
5 # Comparar com media
6 mean_price = df['price'].mean()
7 print(f"Media de price: R$ {mean_price:,.2f}")
8
9 # Se muito diferente, ha outliers!
10 diff_pct = abs(median_price - mean_price) / median_price * 100
11 if diff_pct > 10:
12     print(f"{diff_pct:.1f}% diferenca! Use MEDIANA.")
13
14 # Imputar com mediana
15 df['price'] = df['price'].fillna(median_price)
```

Média vs Mediana: Quando Usar?

Use MÉDIA quando:

- ▶ Distribuição é simétrica (sino)
- ▶ Sem outliers significativos
- ▶ Quer refletir o "total"
- ▶ Exemplo: altura, temperatura, notas

Use MEDIANA quando:

- ▶ Distribuição é assimétrica
- ▶ Há outliers
- ▶ Quer o valor "típico"
- ▶ Exemplo: preços, salários, renda
- ▶ **EM CASO DE DÚVIDA: USE MEDIANA!**

Teste prático:

- ▶ Se média e mediana são similares → qualquer uma OK
- ▶ Se diferem muito → use mediana

Imputação com Moda

</> Python

```
1 # Moda para variaveis CATEGORICAS
2 mode_neighborhood = df['neighborhood'].mode()[0]
3 print(f"Bairro mais comum: {mode_neighborhood}")
4
5 # Contar frequencias
6 print("\nDistribuicao de bairros:")
7 print(df['neighborhood'].value_counts())
8
9 # Imputar com moda
10 n_missing = df['neighborhood'].isnull().sum()
11 df['neighborhood'] = df['neighborhood'].fillna(mode_neighborhood)
12
13 print(f"\n\cmark {n_missing} valores imputados com moda")
14
15 # NOTA: mode() retorna Series, por isso [0]
```

Imputar Múltiplas Colunas de Uma Vez

</> Python

```
1 # Imputar todas as numericas com mediana
2 numeric_cols = df.select_dtypes(include=['number']).columns
3
4 for col in numeric_cols:
5     if df[col].isnull().sum() > 0:
6         median_val = df[col].median()
7         n_missing = df[col].isnull().sum()
8         df[col] = df[col].fillna(median_val)
9         print(f"{col}: {n_missing} valores  $\rightarrow$  mediana {
median_val:.2f}")
```

Imputar Múltiplas Colunas de Uma Vez (cont.)

</> Python

```
1 # Imputar todas as categoricas com moda
2 cat_cols = df.select_dtypes(include=['object']).columns
3
4 for col in cat_cols:
5     if df[col].isnull().sum() > 0:
6         mode_val = df[col].mode()[0]
7         n_missing = df[col].isnull().sum()
8         df[col] = df[col].fillna(mode_val)
9         print(f"{col}: {n_missing} valores  $\rightarrow$  moda '{mode_val}'")
```

Imputação Simples: Prós e Contras

Vantagens:

- ▶ ✓ Muito fácil de implementar
- ▶ ✓ Rápido (mesmo para milhões de linhas)
- ▶ ✓ Funciona razoavelmente para MCAR
- ▶ ✓ Mantém tamanho do dataset
- ▶ ✓ Não requer outras variáveis

Desvantagens:

- ▶ ✗ Reduz variabilidade (todos ficam iguais)
- ▶ ✗ Não considera relações entre variáveis
- ▶ ✗ Pode distorcer distribuição
- ▶ ✗ Subestima desvio padrão
- ▶ ✗ Não funciona bem para MAR ou MNAR

Problema da Imputação Simples: Exemplo

Dados originais (sem missing):

10, 15, 20, 25, 30, 35, 40
Média = 25 | Desvio = 10.6

Dados com missing:

10, NaN, 20, NaN, 30, NaN, 40

Após imputar com média (25):

10, 25, 20, 25, 30, 25, 40
Média = 25 (OK) | Desvio = 8.9 (X menor!)

Atenção

Imputação reduz artificialmente a variabilidade!

Estratégia 3: Imputação por Grupos

Ideia:

- ▶ Não usar média global
- ▶ Usar média do GRUPO a que pertence
- ▶ Mais sofisticado e realista

Quando usar:

- ▶ Dados são MAR (relacionados a outras variáveis)
- ▶ Há agrupamento natural nos dados
- ▶ Valores variam significativamente entre grupos

Exemplos:

- ▶ Área de imóvel: média POR BAIRRO
- ▶ Salário: média POR CARGO
- ▶ Gasto: média POR FAIXA ETÁRIA
- ▶ Preço: mediana POR CATEGORIA

Por que Imputação por Grupos é Melhor?

Cenário:

- ▶ **Centro:** Apartamentos pequenos (média 60 m²)
- ▶ **Subúrbio:** Casas médias (média 120 m²)
- ▶ **Condomínio:** Casas grandes (média 200 m²)

Missing: área de imóvel no Centro

Imputação global (ruim):

- ▶ Média geral: $(60+120+200)/3 = 127 \text{ m}^2$
- ▶ Imputar 127 m² no Centro
- ▶ ✗ Ruim!

Imputação por grupo (melhor):

- ▶ Média do Centro: 60 m²
- ▶ Imputar 60 m²
- ▶ ✓ Muito mais realista!

Imputação por Grupos: Implementação

</> Python

```
1 # Imputar area com media POR BAIRRO
2 # (mais sofisticado que media geral!)
3 # ANTES: ver medias por bairro
4 print("Medias por bairro:")
5 print(df.groupby('neighborhood')['area'].mean())
6 # Imputar usando groupby + transform
7 df['area'] = df.groupby('neighborhood')['area'] \
8     .transform(lambda x: x.fillna(x.mean()))
9
10 # O que faz:
11 # 1. Agrupa por neighborhood
12 # 2. Para cada grupo, calcula mean de area
13 # 3. Preenche NaN com mean do GRUPO
14 # 4. Transform retorna shape original
```


Imputação por Grupos: Mediana

</> Python

```
1 # Usar MEDIANA por grupo (mais robusto)
2 # Imputar price com mediana por bairro
3 df['price'] = df.groupby('neighborhood')['price'] \
4     .transform(lambda x: x.fillna(x.median()))
5 # Verificar resultado
6 print("Medianas de price por bairro:")
7 print(df.groupby('neighborhood')['price'].median())
8 # Se algum grupo tem TODOS missing?
9 # Pandas usa mediana global nesses casos!
10 # Ou especifique fallback:
11 global_median = df['price'].median()
12 df['price'] = df.groupby('neighborhood')['price'] \
13     .transform(lambda x: x.fillna(x.median() if x.notna().any() else
14                               global_median))
```

Múltiplos Grupos: Hierárquico

</> Python

```
1 # Imputar considerando MÚLTIPLOS grupos
2 # Exemplo: por bairro E condicao
3
4 # Tentar nivel 1: bairro + condicao
5 df['price'] = df.groupby(['neighborhood', 'condition'])['price'] \
6     .transform(lambda x: x.fillna(x.median()))
7 # Se ainda tem missing, tentar nivel 2: so bairro
8 df['price'] = df.groupby('neighborhood')['price'] \
9     .transform(lambda x: x.fillna(x.median()))
10 # Se AINDA tem missing, usar global
11 df['price'] = df['price'].fillna(df['price'].median())
12 print(f"Missing restante: {df['price'].isnull().sum()}")
13 # Deve ser 0!
```

Estratégia 4: Valor Constante

Preencher com valor específico do domínio

Quando usar:

- ▶ Valor padrão faz sentido logicamente
- ▶ Missing significa "zero" ou "nenhum"
- ▶ Você conhece bem o domínio

Exemplos:

- ▶ `garage_spaces = 0` (sem garagem)
- ▶ `discount = 0` (sem desconto)
- ▶ `children = 0` (sem filhos)
- ▶ `response_time = 999` (não respondeu)
- ▶ `category = "Unknown"` (desconhecido)

Estratégia 4: Valor Constante

Atenção

CUIDADO: Valor deve fazer sentido! Não preencha idade com 0.

Imputação com Valor Constante

</> Python

```
1 # Exemplo 1: Vagas de garagem
2 # Se nao mencionou, provavelmente nao tem
3 df['garage_spaces'] = df['garage_spaces'].fillna(0)
4 # Exemplo 2: Categoria desconhecida
5 df['neighborhood'] = df['neighborhood'].fillna('Unknown')
6 # Exemplo 3: Estado medio (mais conservador)
7 df['condition'] = df['condition'].fillna('Average')
8 # Exemplo 4: Jardim = nao (se nao mencionou)
9 df['has_garden'] = df['has_garden'].fillna('No')
10 # Exemplo 5: Valor especifico do negocio
11 # Imóveis muito antigos (age missing) = 100 anos
12 df['age'] = df['age'].fillna(100)
```

Estratégia 5: Forward e Backward Fill

Propagação de valores (útil para séries temporais)

Forward Fill (ffill):

- ▶ Propaga último valor válido para frente
- ▶ 10, NaN, NaN, 20 \rightarrow 10, 10, 10, 20
- ▶ Assume: valor permanece até mudar

Backward Fill (bfill):

- ▶ Propaga próximo valor válido para trás
- ▶ 10, NaN, NaN, 20 \rightarrow 10, 20, 20, 20
- ▶ Assume: próximo valor já era válido

Quando usar:

- ▶ Dados ordenados (tempo, sequência)
- ▶ Valores mudam lentamente
- ▶ NÃO usar para dados não ordenados!

Forward e Backward Fill: Código

</> Python

```
1 # Forward fill: propaga ultimo valor valido
2 df_sorted = df.sort_values('age') # Ordenar primeiro!
3 df_sorted['price'] = df_sorted['price'].ffill()
4
5 # Backward fill: propaga proximo valor valido
6 df_sorted['price'] = df_sorted['price'].bfill()
7
8 # Exemplo visual:
9 # Original: 10, NaN, NaN, 20, NaN, 30, NaN
10 # ffill:    10, 10, 10, 20, 20, 30, 30
11 # bfill:    10, 20, 20, 20, 30, 30, NaN (ultimo fica)
12
13 # Combinar: ffill + bfill (preenche tudo)
14 df['price'] = df['price'].ffill().bfill()
```

Forward e Backward Fill: Código

Nota Importante

Só funciona bem com dados ORDENADOS (séries temporais)

Estratégia 6: Interpolação

Estimar valor entre dois pontos conhecidos

Interpolação Linear:

- ▶ Assume mudança linear entre pontos
- ▶ age: 10, NaN, 30 \rightarrow age: 10, 20, 30
- ▶ price: 100, NaN, NaN, 200 \rightarrow 100, 133, 166, 200

Quando usar:

- ▶ Dados ordenados e contínuos
- ▶ Séries temporais
- ▶ Gradiente suave esperado

Quando NÃO usar:

- ▶ Dados categóricos
- ▶ Dados não ordenados
- ▶ Mudanças abruptas

Interpolação: Implementação

</> Python

```
1 # Ordenar por variavel relevante
2 df_sorted = df.sort_values('age')
3
4 # Interpolacao linear
5 df_sorted['price'] = df_sorted['price'].interpolate(method='linear')
6
7 # Outros metodos:
8 # - 'polynomial': curva polinomial
9 # - 'spline': curva suave
10 # - 'nearest': valor mais proximo
```

Interpolação: Implementação (cont.)

</> Python

```
1 # Exemplo de resultado:
2 # age:      10,  15,  NaN,  25,  30
3 # price: 100, 110,  NaN, 130, 140
4 #
5 # Apos interpolacao:
6 # age:      10,  15,  20,  25,  30
7 # price: 100, 110, 120, 130, 140
```

Estratégia 7: Criar Flag de Imputação

SEMPRE uma boa prática!

Por que criar flag?

- ▶ Marca onde valores foram imputados
- ▶ Preserva informação sobre missing
- ▶ ML pode aprender que missing importa
- ▶ Permite análise de impacto

Como funciona:

- ▶ Criar coluna binária (0/1)
- ▶ 1 = foi imputado
- ▶ 0 = valor original
- ▶ Nome: `was_missing_[coluna]`

Exemplo:

- ▶ `price = 300000`, `price_was_missing = 0` (original)
- ▶ `price = 250000`, `price_was_missing = 1` (imputado)

Criar Flag: Implementação

</> Python

```
1 # ANTES de imputar, crie a flag!
2 # Flag indicando onde era missing
3 df['price_was_missing'] = df['price'].isnull().astype(int)
4 # Agora impute
5 df['price'] = df['price'].fillna(df['price'].median())
6 # Resultado:
7 # - price: valores completos (alguns imputados)
8 # - price_was_missing: 1 onde foi imputado, 0 original
9 # Verificar
10 print(f"Valores imputados: {df['price_was_missing'].sum()}")
11 # Analise: preco imputado vs original difere?
12 print(df.groupby('price_was_missing')['area'].mean())
```

Criar Flags para Múltiplas Colunas

</> Python

```
1 # Criar flags para todas as colunas com missing
2 cols_with_missing = df.columns[df.isnull().any()].tolist()
3 for col in cols_with_missing:
4     flag_name = f'{col}_was_missing'
5     df[flag_name] = df[col].isnull().astype(int)
6     print(f"\cmark  Flag criada: {flag_name}")
7 # Verificar flags criadas
8 flag_cols = [c for c in df.columns if '_was_missing' in c]
9 print(f"\nFlags de imputacao: {flag_cols}")
10 # Ver quantos foram imputados
11 print("\nTotal de imputacoes por coluna:")
12 print(df[flag_cols].sum())
```

Escolher Estratégia: Árvore de Decisão

PASSO 1: Tipo de variável?

- ▶ Numérica → média ou mediana
- ▶ Categórica → moda ou constante

PASSO 2: Há agrupamento natural?

- ▶ SIM → imputação por grupo
- ▶ NÃO → imputação global

PASSO 3: Dados têm ordem?

- ▶ Temporal → interpolação ou ffill
- ▶ Não ordenado → média/mediana

PASSO 4: Tipo de missing?

- ▶ MCAR → qualquer método
- ▶ MAR → por grupos
- ▶ MNAR → cuidado! Talvez constante + flag

Caso de Uso Completo: Housing Prices - Estratégia por variável

price (12% missing, MNAR):

- ▶ Imputar com mediana por bairro
- ▶ Criar flag price_was_missing
- ▶ Razão: preços variam muito por bairro

area (8% missing, MCAR):

- ▶ Imputar com mediana por bairro
- ▶ Razão: áreas também variam por bairro

age (15% missing, MNAR):

- ▶ Imputar com valor alto (100 anos) + flag
- ▶ Razão: missing pode significar "muito antigo"

neighborhood (5% missing):

- ▶ Imputar com moda ("Centro")
- ▶ Ou categoria "Unknown"

Pipeline de Imputação Completo

</> Python

```
1 def impute_housing_data(df):
2     """Pipeline de imputacao completo"""
3     df = df.copy()
4
5     # 1. Criar flags ANTES de imputar
6     df['price_was_missing'] = df['price'].isnull().astype(int)
7     df['age_was_missing'] = df['age'].isnull().astype(int)
8
9     # 2. Imputar price por bairro
10    df['price'] = df.groupby('neighborhood')['price'] \
11        .transform(lambda x: x.fillna(x.median()))
12
13    # 3. Imputar area por bairro
14    df['area'] = df.groupby('neighborhood')['area'] \
15        .transform(lambda x: x.fillna(x.median()))
```

Pipeline de Imputação (cont.)

</> Python

```
1  # 4. Imputar age com valor alto (antigos)
2  df['age'] = df['age'].fillna(100)
3  # 5. Imputar bedrooms/bathrooms com mediana
4  df['bedrooms'] = df['bedrooms'].fillna(df['bedrooms'].median())
5  df['bathrooms'] = df['bathrooms'].fillna(df['bathrooms'].median())
6  # 6. Imputar categoricas com moda
7  df['neighborhood'] = df['neighborhood'].fillna(df['neighborhood'].mode
8  ()[0])
9  df['condition'] = df['condition'].fillna('Average')
10 # 7. Validacao final
11 assert df.isnull().sum().sum() == 0, "Ainda há missing!"
12 return df
13 # Usar pipeline
14 df_imputed = impute_housing_data(df)
```

Validar Imputação: Comparar Distribuições

Python

```
1 import matplotlib.pyplot as plt
2 # Comparar distribuicao antes/depois
3 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
4 # Antes: apenas valores originais
5 original_values = df_original['price'].dropna()
6 axes[0].hist(original_values, bins=30, alpha=0.7)
7 axes[0].set_title('Distribuicao Original (sem missing)')
8 axes[0].set_xlabel('Price')
9 # Depois: com valores imputados
10 axes[1].hist(df_imputed['price'], bins=30, alpha=0.7, color='green')
11 axes[1].set_title('Distribuicao Apos Imputacao')
12 axes[1].set_xlabel('Price')
13 plt.tight_layout()
14 plt.show()
```

Exercício Completo: Imputação Housing

Exercício Prático

Implemente imputação completa do Housing:

1. Para cada variável com missing:

- ▶ Decida qual estratégia usar
- ▶ Justifique a escolha

2. Implemente o pipeline:

```
1 # - Criar flags ANTES de imputar
2 # - Aplicar estrategias escolhidas
3 # - Validar: nenhum missing deve restar
4 assert df.isnull().sum().sum() == 0
```

3. Compare distribuições antes/depois

4. Documente todas as decisões

5. Gere relatório de imputação

Resumo do Bloco 2

Remoção:

- ▶ `dropna()` - quando $< 5\%$ e MCAR

Imputação Simples:

- ▶ Média, mediana, moda
- ▶ Rápido mas reduz variabilidade

Imputação Avançada:

- ▶ Por grupos (`groupby` + `transform`)
- ▶ Forward/backward fill (séries)
- ▶ Interpolação (dados ordenados)
- ▶ Valor constante (conhecimento do domínio)

Boa Prática:

- ▶ Sempre criar flags de imputação!

Transição para Próximo Bloco

Já sabemos lidar com missing!

Agora vamos para:

- ▶ Limpeza de strings (espaços, maiúsculas)
- ▶ Padronização de formatos
- ▶ Validação de ranges
- ▶ Detecção de valores impossíveis
- ▶ Tratamento de duplicatas
- ▶ Conversão de tipos
- ▶ Consistência lógica
- ▶ Pipeline completo de limpeza

💡 Nota Importante

Bloco 3: Deixar os dados impecáveis!

Bloco 3

Limpeza de Inconsistências e Validação

Além de Missing: Outros Problemas

Dados podem estar "presentes" mas ERRADOS:

1. Problemas de Formato

- ▶ Espaços extras: "Casa "vs "Casa"
- ▶ Maiúsculas/minúsculas: "SIM"vs "sim"vs "Sim"
- ▶ Inconsistências: "São Paulo"vs "Sao Paulo"vs "SP"
- ▶ Caracteres especiais: "R\$ 1.000"vs "1000"

2. Valores Impossíveis

- ▶ Idade = -5 ou 200 anos
- ▶ Área = 0 m² ou 10000 m²
- ▶ Data de nascimento no futuro
- ▶ Preço = R\$ 1 ou R\$ 999.999.999

3. Duplicatas

- ▶ Mesmo registro múltiplas vezes
- ▶ Inconsistências entre duplicatas

Limpeza de Strings: Por que Importa?

Strings sujas causam problemas:

Exemplo problemático:

"Centro"
"Centro"
"Centro "
"Centro "
"centro"
"CENTRO"

Limpeza de Strings: Por que Importa?

Para o computador:

- ▶ São 6 categorias diferentes!
- ▶ `value_counts()` mostra 6 valores únicos
- ▶ Análises ficam fragmentadas
- ▶ Gráficos com categorias duplicadas

Para o humano:

- ▶ É tudo a mesma coisa: Centro

Limpeza: Remover Espaços

</> Python

```
1 # Problema: espacos extras
2 print("ANTES:")
3 print(df['neighborhood'].unique())
4 # [' Centro', 'Centro ', ' Centro ', 'Suburbio', ...]
5
6 # Solucao: strip() remove espacos das pontas
7 df['neighborhood'] = df['neighborhood'].str.strip()
8
9 print("\nDEPOIS:")
10 print(df['neighborhood'].unique())
11 # ['Centro', 'Suburbio', ...]
12
13 # Verificar impacto
14 print(f"\nCategorias antes: {df_original['neighborhood'].nunique()}")
15 print(f"Categorias depois: {df['neighborhood'].nunique()}")
```

Limpeza: Padronizar Capitalização

</> Python

```
1 # Problema: capitalizacao inconsistente
2 print("ANTES:")
3 print(df['condition'].value_counts())
4 # Excelente      50
5 # EXCELENTE      30
6 # excelente      20
7 # Bom            40
8 # BOM            15
9
10 # Solucao 1: Tudo minusculo (RECOMENDADO)
11 df['condition'] = df['condition'].str.lower()
12
13 # Solucao 2: Tudo maiusculo
14 df['condition'] = df['condition'].str.upper()
```

Limpeza: Padronizar Capitalização (cont.)

</> Python

```
1 # Solucao 3: Primeira letra maiuscula
2 # df['condition'] = df['condition'].str.title()
3
4 print("\nDEPOIS:")
5 print(df['condition'].value_counts())
6 # excelente    100
7 # bom          55
```

Qual Capitalização Usar?

Guia de decisão:

Use **.lower() (minúsculo)**:

- ▶ ✓ RECOMENDADO para análise
- ▶ ✓ Mais fácil de digitar e buscar
- ▶ ✓ Menos propenso a erros
- ▶ Exemplo: "centro", "excelente", "bom"

Use **.upper() (maiúsculo)**:

- ▶ Códigos e siglas (ISO, UF)
- ▶ Exemplo: "SP", "RJ", "USD", "BRL"

Use **.title() (capitalizado)**:

- ▶ Para apresentação final
- ▶ Nomes próprios
- ▶ Exemplo: "São Paulo", "Excelente"

Limpeza: Substituir Variações

</> Python

```
1 # Problema: multiplas formas de escrever o mesmo
2 print("ANTES:")
3 print(df['neighborhood'].unique())
4 # ['Sao Paulo', 'São Paulo', 'Sao Paulo (SP)',
5 #  'S. Paulo', 'SP']
6
7 # Solucao: Substituir variacoes
8 df['neighborhood'] = df['neighborhood'] \
9     .str.replace('Sao Paulo', 'São Paulo') \
10    .str.replace('S. Paulo', 'São Paulo') \
11    .str.replace(' (SP)', '') \
12    .str.replace('SP', 'São Paulo') \
13    .str.strip()
14
15 print("\nDEPOIS:")
16 print(df['neighborhood'].unique())
```

Limpeza: Remover Caracteres Especiais

</> Python

```
1 # Problema: caracteres indesejados
2 print("ANTES:")
3 print(df['price'].head())
4 # 'R$ 300.000', 'R$ 250,000', '150000', '$200k'
5
6 # Se price esta como string, converter!
7 # Remover caracteres nao numericos
8 df['price'] = df['price'].str.replace('R$', '') \
9     .str.replace('$', '') \
10    .str.replace('.', '') \
11    .str.replace(',', '.') \
12    .str.replace('k', '000') \
13    .str.strip()
```


Limpeza: Remover Caracteres Especiais (cont.)

</> Python

```
1 # Converter para numero
2 df['price'] = pd.to_numeric(df['price'], errors='coerce')
3
4 print("\nDEPOIS:")
5 print(df['price'].head())
6 # 300000, 250000, 150000, 200000
```

Padronização com Mapeamento

</> Python

```
1 # Problema: categorias nao padronizadas
2 print("ANTES:")
3 print(df['condition'].unique())
4 # ['Excelente', 'Muito Bom', 'Bom', 'Regular',
5 #  'Ruim', 'Pessimo', 'Ok']
6 # Solucao: Mapear para valores padrao
7 condition_map = {
8     'excelente': 'excellent',
9     'muito bom': 'good',
10    'bom': 'good',
11    'ok': 'average',
12    'regular': 'average',
13    'ruim': 'poor',
14    'pessimo': 'poor'
15 }
```

Padronização com Mapeamento (cont.)

</> Python

```
1 df['condition'] = df['condition'].str.lower().map(condition_map)
2
3 print("\nDEPOIS:")
4 print(df['condition'].unique())
5 # ['excellent', 'good', 'average', 'poor']
```

Aplicar Limpeza em Todas as Categóricas

</> Python

```
1 # Pipeline de limpeza para todas as categoricas
2 categorical_cols = df.select_dtypes(include=['object']).columns
3 for col in categorical_cols:
4     print(f"Limpendo {col}...")
5     # 1. Remover espacos
6     df[col] = df[col].str.strip()
7     # 2. Minusculo
8     df[col] = df[col].str.lower()
9     # 3. Remover espacos duplos internos
10    df[col] = df[col].str.replace('  ', ' ')
11    print(f"{col}: {df[col].nunique()} categorias unicas")
12 print("\nLimpeza concluida!")
```

Validação de Dados

1. Range (Intervalo):

- ▶ Valor está dentro do esperado?
- ▶ Ex: idade entre 0-150, área entre 20-1000

2. Tipo:

- ▶ Variável tem o tipo correto?
- ▶ Ex: price é numérico, neighborhood é texto

3. Consistência Lógica:

- ▶ Relações entre variáveis fazem sentido?
- ▶ Ex: $\text{bathrooms} \leq \text{bedrooms} + 2$

4. Formato:

- ▶ Dados seguem padrão esperado?
- ▶ Ex: data em formato DD/MM/YYYY, CEP com 8 dígitos

Definir Regras de Validação

Para Housing Prices - regras de negócio:

Ranges aceitáveis:

- ▶ **price:** R\$ 50.000 a R\$ 10.000.000
- ▶ **area:** 20 m² a 1.000 m²
- ▶ **bedrooms:** 0 a 10
- ▶ **bathrooms:** 0 a 10
- ▶ **age:** 0 a 150 anos
- ▶ **garage_spaces:** 0 a 10

Como definir regras?

- ▶ Consulte especialistas do domínio
- ▶ Analise distribuição (percentis 1% e 99%)
- ▶ Use senso comum e conhecimento do mundo real
- ▶ Seja conservador (não muito restrito)
- ▶ Documente as regras e suas fontes

Validação: Identificar Valores Impossíveis

</> Python

```
1 # Converter todas as colunas de interesse para numérico
2 cols = ['age', 'price', 'area', 'bedrooms', 'bathrooms']
3 df[cols] = df[cols].apply(pd.to_numeric, errors='coerce')
4
5 # Verificar valores impossíveis
6 print("=== VALORES NEGATIVOS ===")
7 for col in cols:
8     negative = df[df[col] < 0]
9     if len(negative) > 0:
10         print(f"{col}: {len(negative)} valores negativos")
11         print(negative[[col]].head())
```

Validação: Identificar Valores Impossíveis (cont.)

</> Python

```
1 print("\n=== VALORES EXTREMOS ===")
2 print(f"Bedrooms > 20: {len(df[df['bedrooms'] > 20])}")
3 print(f"Age > 200: {len(df[df['age'] > 200])}")
4 print(f"Area > 10000: {len(df[df['area'] > 10000])}")
5
6 print("\n=== ZEROS SUSPEITOS ===")
7 print(f"Area = 0: {len(df[df['area'] == 0])}")
8 print(f"Price = 0: {len(df[df['price'] == 0])}")
```


Validação: Aplicar Regras

</> Python

```
1 # Definir e aplicar regras de validacao
2
3 rules = {
4     'price': (50_000, 10_000_000),
5     'area': (20, 1000),
6     'bedrooms': (0, 10),
7     'bathrooms': (0, 10),
8     'age': (0, 150),
9     'garage_spaces': (0, 10)
10 }
```

Validação: Aplicar Regras (cont.)

</> Python

```
1 print("=== VALIDACAO DE RANGES ===")
2 for col, (min_val, max_val) in rules.items():
3     invalid = df[(df[col] < min_val) | (df[col] > max_val)]
4
5     if len(invalid) > 0:
6         print(f"\n{col}: {len(invalid)} valores invalidos")
7         print(f"    Range esperado: [{min_val}, {max_val}]")
8         print(f"    Valores encontrados:")
9         print(f"        Min: {invalid[col].min()}")
10        print(f"        Max: {invalid[col].max()}")
```

O que Fazer com Valores Inválidos?

1. Corrigir manualmente (MELHOR):

- ▶ Se erro óbvio: 3000 → 300 m²
- ▶ Se tem acesso à fonte original
- ▶ Se poucos casos (< 10)

2. Tratar como missing (RECOMENDADO):

- ▶ Substituir por NaN
- ▶ Aplicar estratégias de imputação
- ▶ Mais conservador e seguro

O que Fazer com Valores Inválidos?

3. Clipar nos limites:

- ▶ $\text{age} = -5 \rightarrow 0$
- ▶ $\text{area} = 10000 \rightarrow 1000$
- ▶ Menos recomendado (distorce)

4. Remover registro:

- ▶ Apenas se muitos valores inválidos na linha
- ▶ Última opção

Tratar Inválidos como Missing

</> Python

```
1 # Substituir valores fora do range por NaN
2 # Depois aplicar imputacao
3
4 # Exemplo: bedrooms
5 print(f"ANTES: {df['bedrooms'].isnull().sum()} missing")
6
7 # Valores invalidos virao NaN
8 df.loc[df['bedrooms'] < 0, 'bedrooms'] = np.nan
9 df.loc[df['bedrooms'] > 10, 'bedrooms'] = np.nan
```

Tratar Inválidos como Missing

</> Python

```
1 print(f"DEPOIS: {df['bedrooms'].isnull().sum()} missing")
2
3 # Ou usando between() (mais elegante)
4 df.loc[~df['bedrooms'].between(0, 10), 'bedrooms'] = np.nan
5
6 # Agora imputar
7 df['bedrooms'] = df['bedrooms'].fillna(df['bedrooms'].median())
8
9 print(f"FINAL: {df['bedrooms'].isnull().sum()} missing")
```

Identificar Duplicatas

</> Python

```
1 # Verificar duplicatas completas (todas as colunas iguais)
2 print("=== DUPLICATAS COMPLETAS ===")
3 n_duplicates = df.duplicated().sum()
4 print(f"Total de duplicatas: {n_duplicates}")
5
6 # Ver as duplicatas
7 if n_duplicates > 0:
8     duplicates = df[df.duplicated(keep=False)]
9     print("\nExemplos de duplicatas:")
10    print(duplicates.sort_values(by=df.columns[0]).head(10))
```

Identificar Duplicatas (cont.)

</> Python

```
1 # Duplicatas em colunas especificas
2 # (ex: mesmo endereco)
3 print("\n=== DUPLICATAS PARCIAIS ===")
4 dup_address = df.duplicated(
5     subset=['neighborhood', 'area', 'bedrooms'],
6     keep=False
7 )
8 print(f"Duplicatas de endereco: {dup_address.sum()}")
```


Remover Duplicatas

Python

```
1 print(f"Linhas ANTES: {len(df)}")
2
3 # Remover duplicatas completas
4 df_clean = df.drop_duplicates()
5 print(f"Linhas DEPOIS (completas): {len(df_clean)}")
6 print(f"Removidos: {len(df) - len(df_clean)}")
7
8 # Manter primeira ocorrencia (padrao)
9 df_clean = df.drop_duplicates(keep='first')
10
11 # Manter ultima ocorrencia
12 df_clean = df.drop_duplicates(keep='last')
```

Remover Duplicatas (cont.)

</> Python

```
1 # Remover TODAS (incluindo originais)
2 # df_clean = df.drop_duplicates(keep=False)
3
4 # Duplicatas em colunas especificas
5 df_clean = df.drop_duplicates(
6     subset=['neighborhood', 'area', 'bedrooms']
7 )
8 print(f"Linhas finais: {len(df_clean)}")
```

Conversão de Tipos

</> Python

```
1 # Verificar tipos atuais
2 print("=== TIPOS ATUAIS ===")
3 print(df.dtypes)
4
5 # Converter para tipos corretos
6 print("\n=== CONVERTENDO ===")
7
8 # Numericas inteiras
9 df['bedrooms'] = df['bedrooms'].astype('int')
10 df['bathrooms'] = df['bathrooms'].astype('int')
11 df['age'] = df['age'].astype('int')
12
13 # Numericas decimais
14 df['price'] = df['price'].astype('float')
15 df['area'] = df['area'].astype('float')
```

Conversão de Tipos (cont.)

</> Python

```
1 # Categorias (economiza memoria!)
2 df['neighborhood'] = df['neighborhood'].astype('category')
3 df['condition'] = df['condition'].astype('category')
4
5 print("\n=== TIPOS FINAIS ===")
6 print(df.dtypes)
7
8 # CUIDADO: NaN vai gerar erro
```

Conversão Segura com to_numeric

Python

```
1 # Se coluna tem valores nao numericos misturados
2
3 # Opcao 1: Converter, invalidos viram NaN
4 df['area'] = pd.to_numeric(df['area'], errors='coerce')
5 # "100"  $\\rightarrow$  100
6 # "abc"  $\\rightarrow$  NaN
7 # "50.5" $\\rightarrow$  50.5
8
9 # Opcao 2: Ignorar erros (mantem original)
10 df['area'] = pd.to_numeric(df['area'], errors='ignore')
```

Conversão Segura com to_numeric (cont.)

Python

```
1 # Opcao 3: Levantar erro (aborta se tiver invalido)
2 try:
3     df['area'] = pd.to_numeric(df['area'], errors='raise')
4 except ValueError as e:
5     print(f"Erro: {e}")
6     print("Ha valores nao numericos!")
7
8 # RECOMENDADO: usar 'coerce' e depois imputar NaNs
```

Tratamento de Datas

Python

```
1 # Converter strings para datetime
2 # Formato padrao (YYYY-MM-DD)
3 df['sale_date'] = pd.to_datetime(df['sale_date'])
4 # Formato brasileiro (DD/MM/YYYY)
5 df['sale_date'] = pd.to_datetime(df['sale_date'],
6                                 format='%d/%m/%Y')
7 # Erros viram NaT (Not a Time)
8 df['sale_date'] = pd.to_datetime(df['sale_date'],
9                                 errors='coerce')
10 # Extrair componentes
11 df['sale_year'] = df['sale_date'].dt.year
12 df['sale_month'] = df['sale_date'].dt.month
13 df['sale_day'] = df['sale_date'].dt.day
14 df['sale_weekday'] = df['sale_date'].dt.dayofweek
15 print(df[['sale_date', 'sale_year', 'sale_month']].head())
```

Validação de Consistência Lógica

</> Python

```
1 # Verificar relacoes logicas entre variaveis
2
3 print("=== VALIDACAO DE LOGICA ===")
4
5 # Regra 1: Banheiros nao pode ser >> quartos
6 invalid_bath = df[df['bathrooms'] > df['bedrooms'] + 2]
7 print(f"Banheiros > Quartos+2: {len(invalid_bath)}")
8 if len(invalid_bath) > 0:
9     print(invalid_bath[['bedrooms', 'bathrooms']].head())
```


Validação de Consistência Lógica (cont.)

</> Python

```
1 # Regra 2: Preço/m² deve estar em range razoavel
2 df['price_per_sqm'] = df['price'] / df['area']
3 median_price_sqm = df['price_per_sqm'].median()
4
5 outliers = df[
6     (df['price_per_sqm'] < median_price_sqm * 0.1) |
7     (df['price_per_sqm'] > median_price_sqm * 10)
8 ]
9 print(f"\nPreço/m² anormal: {len(outliers)}")
```

Validação: Criar Função Reutilizável

</> Python

```
1 def validate_housing_data(df):
2     """Valida regras de negocio"""
3     issues = []
4     # Verificar negativos
5     for col in ['price', 'area', 'age']:
6         n_negative = (df[col] < 0).sum()
7         if n_negative > 0:
8             issues.append(f"{col}: {n_negative} negativos")
9     # Verificar ranges
10    if (df['bedrooms'] > 20).any():
11        issues.append("Bedrooms > 20 encontrados")
12    if (df['area'] > 10000).any():
13        issues.append("Area > 10000 encontrada")
```

Validação: Criar Função Reutilizável (cont.)

</> Python

```
1  # ... outras validacoes
2
3  return issues
4
5  # Usar
6  problems = validate_housing_data(df)
7  if problems:
8      print("PROBLEMAS ENCONTRADOS:")
9      for p in problems:
10         print(f"    - {p}")
```

Exercício: Limpeza e Validação Completa

Exercício Prático

Limpe e valide completamente o Housing dataset:

1. Limpeza de strings: `strip()`, `lower()`, `replace()`
2. Validação de ranges:
 - ▶ Defina regras para cada variável
 - ▶ Identifique valores inválidos
 - ▶ Trate como missing e impute
3. Duplicatas:
 - ▶ Identifique e remova
4. Conversão de tipos corretos
5. Validação de consistência lógica
6. Documente TUDO!

Resumo do Bloco 3

Técnicas de limpeza aprendidas:

Strings:

- ▶ strip(), lower(), upper(), replace()
- ▶ Padronização de categóricas
- ▶ Mapeamento de valores

Validação:

- ▶ Definir regras de negócio
- ▶ Identificar valores impossíveis
- ▶ Tratar inválidos (missing ou corrigir)

Outros:

- ▶ Remover duplicatas
- ▶ Converter tipos corretamente
- ▶ Tratar datas
- ▶ Validar consistência lógica

Bloco 4

Pipeline Completo e Documentação

Pipeline Completo: Estrutura

</> Python

```
1 def clean_housing_pipeline(df):
2     """
3     Pipeline completo de limpeza:
4     1. Copiar dados originais
5     2. Limpar strings
6     3. Validar e tratar invalidos
7     4. Criar flags de missing
8     5. Imputar valores faltantes
9     6. Remover duplicatas
10    7. Converter tipos
11    8. Validar resultado
12    """
13    df_clean = df.copy()
14    # ... implementacao
15    return df_clean
```

Pipeline: Implementação (cont.)

</> Python

```
1 def clean_housing_pipeline(df):
2     df_clean = df.copy()
3     print("Pipeline de limpeza iniciado...")
4
5     # PASSO 1: Limpar strings
6     print("\n[1/7] Limpando strings...")
7     cat_cols = df_clean.select_dtypes(include=['object']).columns
8     for col in cat_cols:
9         df_clean[col] = df_clean[col].str.strip().str.lower()
```


Pipeline: Implementação (cont.)

</> Python

```
1  # PASSO 2: Validar ranges e marcar invalidos
2  print("[2/7] Validando ranges...")
3  df_clean.loc[~df_clean['bedrooms'].between(0, 10), 'bedrooms'] = np.
   nan
4  df_clean.loc[~df_clean['area'].between(20, 1000), 'area'] = np.nan
5  df_clean.loc[~df_clean['age'].between(0, 150), 'age'] = np.nan
```

Pipeline: Implementação (cont.)

Python

```
1  # PASSO 3: Criar flags de imputacao
2  print("[3/7] Criando flags...")
3  df_clean['price_was_missing'] = df_clean['price'].isnull().astype(int)
4  df_clean['age_was_missing'] = df_clean['age'].isnull().astype(int)
5
6  # PASSO 4: Imputar valores faltantes
7  print("[4/7] Imputando valores...")
8  # Por grupo
9  df_clean['price'] = df_clean.groupby('neighborhood')['price'] \
10     .transform(lambda x: x.fillna(x.median()))
11  df_clean['area'] = df_clean.groupby('neighborhood')['area'] \
12     .transform(lambda x: x.fillna(x.median()))
```

Pipeline: Implementação (cont.)

</> Python

```
1  # Global
2  df_clean['age'] = df_clean['age'].fillna(100)
3  df_clean['bedrooms'] = df_clean['bedrooms'].fillna(df_clean['bedrooms']
4  ].median())
5
6  # PASSO 5: Remover duplicatas
7  print("[5/7] Removendo duplicatas...")
8  n_before = len(df_clean)
9  df_clean = df_clean.drop_duplicates()
10 n_removed = n_before - len(df_clean)
    print(f"  Removidos: {n_removed} duplicatas")
```

Pipeline: Implementação (cont.)

</> Python

```
1  # PASSO 6: Converter tipos - Cuidado com conversão numérica
2  print("[6/7] Convertendo tipos...")
3  df_clean['bedrooms'] = df_clean['bedrooms'].astype('int')
4  df_clean['bathrooms'] = df_clean['bathrooms'].astype('int')
5  df_clean['neighborhood'] = df_clean['neighborhood'].astype('category')
6  df_clean['condition'] = df_clean['condition'].astype('category')
```

Pipeline: Implementação (cont.)

</> Python

```
1  # PASSO 7: Validacao final
2  print("[7/7] Validando resultado...")
3  assert df_clean.isnull().sum().sum() == 0, "Ainda ha missing!"
4
5  print("\n\cmark Pipeline concluido com sucesso!")
6  return df_clean
7
```

Gerar Relatório de Limpeza

</> Python

```
1 def generate_cleaning_report(df_original, df_clean):
2     """Relatorio completo de limpeza"""
3     report = {
4         'Linhas originais': len(df_original),
5         'Linhas finais': len(df_clean),
6         'Linhas removidas': len(df_original) - len(df_clean),
7         'Missing antes': df_original.isnull().sum().sum(),
8         'Missing depois': df_clean.isnull().sum().sum(),
9         'Duplicatas removidas': len(df_original) - len(df_original.
10 drop_duplicates()),
11         'Memoria antes (MB)': df_original.memory_usage(deep=True).sum() /
12         1024**2,
13         'Memoria depois (MB)': df_clean.memory_usage(deep=True).sum() /
14         1024**2
15     }
```

Gerar Relatório de Limpeza (cont.)

</> Python

```
1 print("="*50)
2 print("RELATÓRIO DE LIMPEZA")
3 print("="*50)
4 for key, value in report.items():
5     print(f"{key:25s}: {value:,.2f}")
6
7 return report
8
```

Checklist Final de Limpeza

Antes de considerar dados "limpos":

1. ☐ Valores faltantes identificados e tratados
2. ☐ Flags de imputação criadas
3. ☐ Strings padronizadas (strip, lower)
4. ☐ Tipos de dados corretos
5. ☐ Valores impossíveis removidos/corrigidos
6. ☐ Ranges validados
7. ☐ Duplicatas removidas
8. ☐ Consistência lógica verificada
9. ☐ Todas as decisões documentadas
10. ☐ Relatório de limpeza gerado
11. ☐ Pipeline reproduzível criado
12. ☐ Testes de validação passando

Checklist Final de Limpeza

Atenção

Só avance se TODOS os itens estiverem ✓

Live Coding

Vamos praticar juntos:

1. (atualizado em aula)

Recap: Aula 09 Completa

Jornada de hoje:

Bloco 1: Identificação de Missing

- ▶ Tipos: MCAR, MAR, MNAR
- ▶ Visualização e análise de padrões
- ▶ Decisão: remover vs imputar

Bloco 2: Estratégias de Imputação

- ▶ Simples: média, mediana, moda
- ▶ Avançada: por grupos, interpolação
- ▶ Flags de imputação

Recap: Aula 09 Completa

Jornada de hoje:

Bloco 3: Limpeza e Validação

- ▶ Padronização de strings
- ▶ Validação de ranges
- ▶ Remoção de duplicatas
- ▶ Conversão de tipos

Bloco 4: Pipeline

- ▶ Workflow completo e reproduzível

Erros Comuns - Evite!

1. Não fazer cópia dos dados

- ▶ ✗ `df_clean = df` (referência!)
- ▶ ✓ `df_clean = df.copy()`

2. Imputar sem análise

- ▶ ✗ `fillna(0)` em tudo
- ▶ ✓ Entender tipo de missing primeiro

3. Não documentar decisões

- ▶ ✗ Fazer mudanças silenciosamente
- ▶ ✓ Comentar e documentar cada passo

4. Ser muito agressivo

- ▶ ✗ Remover 50% dos dados
- ▶ ✓ Ser conservador

5. Não validar resultado

- ▶ ✗ Assumir que funcionou
- ▶ ✓ Sempre verificar com asserts

Quinta-feira (06/11) - Transformação e Normalização:

Feature Engineering:

- ▶ Criar novas variáveis úteis
- ▶ Extrair informação de datas
- ▶ Interações entre variáveis

Encoding de Categóricas:

- ▶ Label Encoding
- ▶ One-Hot Encoding
- ▶ Quando usar cada um

Normalização e Escalonamento:

- ▶ Min-Max Scaling
- ▶ Z-score Standardization
- ▶ Robust Scaling

Exercício Prático

Tempo: 60 minutos

Entrega: via Moodle (notebook)

Tarefas:

1. (atualizado durante a aula)

Notebook: Disponível no Moodle

Obrigado!