

Programação para Ciência de Dados

Pré-processamento de Dados - Transformação e Normalização

Arthur Casals

6 de Novembro de 2025

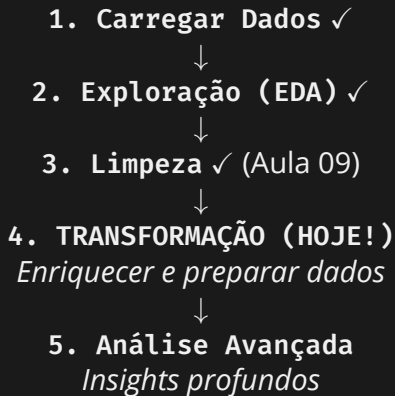
- ▶ Aula prática de amanhã (07/11): **TEM ENTREGA**

Agenda

- ▶ Introdução: Por que Transformar Dados?
- ▶ Bloco 0: Preview de Machine Learning (Contexto)
- ▶ Bloco 1: Feature Engineering
- ▶ Bloco 2: Encoding de Variáveis Categóricas
- ▶ Bloco 3: Normalização e Padronização
- ▶ Bloco 4: Pipeline Completo e Análise Final

A Próxima Etapa: Transformação

Onde estamos no pipeline de Ciência de Dados:



Por que Transformar Dados?

Dados limpos \neq Dados prontos para análise avançada

Problemas que transformação resolve:

- ▶ **Limitação em análises numéricas**

- ▶ Categorias (texto) não podem ser correlacionadas
- ▶ Precisa converter para números
- ▶ "Alto", "Médio", "Baixo" \rightarrow 2, 1, 0

- ▶ **Escalas muito diferentes dificultam comparação**

- ▶ Idade: 20-80 vs Salário: 1000-100000
- ▶ Impossível comparar "importância" visualmente
- ▶ Gráficos ficam distorcidos

Por que Transformar Dados?

Dados limpos \neq Dados prontos para análise avançada

Problemas que transformação resolve:

- ▶ **Informação "escondida" nos dados**

- ▶ Insights não são óbvios nas variáveis originais
- ▶ Precisa criar novas variáveis derivadas
- ▶ Ex: não temos "gasto médio mensal", mas podemos calcular

- ▶ **Distribuições problemáticas**

- ▶ Muito assimétrica dificulta visualização
- ▶ Transformações podem normalizar

Transformação: Impacto na Análise

Exemplo: Analisar Cancelamento de Clientes (Churn)

Sem transformação adequada:

- ▶ Análises superficiais e limitadas
- ▶ Não consegue correlacionar categorias
- ▶ Comparações injustas (escalas diferentes)
- ▶ Padrões importantes não são visíveis
- ▶ Visualizações ruins

Com transformação adequada:

- ▶ Análises profundas e quantitativas
- ▶ Correlações entre todas as variáveis
- ▶ Comparações justas e interpretáveis
- ▶ Padrões ocultos revelados
- ▶ Visualizações claras e informativas

Transformação: Impacto na Análise

Nota Importante

Transformação = Preparar dados para extrair insights máximos

Caso Real: Netflix

Sistema de Análise de Preferências

Desafio inicial:

- ▶ Entender quais fatores influenciam preferências
- ▶ Dados: gêneros (texto), duração (minutos), ano, rating
- ▶ Escalas diferentes: ano (1900-2024) vs rating (0-10)
- ▶ Categorias: mais de 100 gêneros diferentes

Transformações aplicadas:

- ▶ **Encoding:** Gêneros → Representação numérica
- ▶ **Feature Engineering:** "Anos desde lançamento", "Era do filme"
- ▶ **Normalização:** Todas as variáveis comparáveis
- ▶ **Agregação:** Combinações de gêneros populares

Sistema de Análise de Preferências

Resultado:

- ▶ Análises revelam padrões claros de preferência
- ▶ Base para sistema de recomendação (Machine Learning)
- ▶ Economia de \$1 bilhão/ano em retenção

Transformação vs Machine Learning

Relacionamento entre conceitos:

Transformação de Dados (Esta Aula):

- ▶ Técnicas para preparar e enriquecer dados
- ▶ Útil para *qualquer* análise quantitativa
- ▶ Análise exploratória avançada
- ▶ Visualizações melhores
- ▶ Cálculo de correlações e estatísticas

Machine Learning (Cursos Futuros):

- ▶ Usa *exatamente as mesmas transformações*
- ▶ Mas para construir modelos preditivos
- ▶ Objetivo: fazer previsões automáticas
- ▶ Requer dados transformados da mesma forma

Transformação vs Machine Learning

Nota Importante

Hoje: aprenderemos transformações para análise.
Bonus: são as mesmas usadas em ML!

Tipos de Transformações

1. Feature Engineering (CRIAÇÃO)

- ▶ Criar novas variáveis a partir das existentes
- ▶ Extrair informação "escondida"
- ▶ Combinar features para revelar padrões
- ▶ Ex: $\text{TotalCharges} / \text{tenure} = \text{"Gasto médio mensal"}$

2. Encoding (CONVERSÃO)

- ▶ Texto \rightarrow números
- ▶ Necessário para análises quantitativas
- ▶ Ex: "Sim"/"Não" \rightarrow 1/0

Tipos de Transformações

3. Scaling (NORMALIZAÇÃO)

- ▶ Colocar variáveis em escalas comparáveis
- ▶ Facilita interpretação e visualização
- ▶ Ex: Idade e Salário na mesma escala $[0, 1]$

4. Mathematical Transform (CORREÇÃO)

- ▶ Log, sqrt, potência
- ▶ Corrigir distribuições assimétricas
- ▶ Ex: $\log(\text{salário})$ deixa distribuição mais normal

Dataset de Hoje: Telco Customer Churn

Dados de clientes de empresa de telecomunicações

O que é Churn?

- ▶ Churn = Cancelamento (cliente sai da empresa)
- ▶ Métrica crítica para negócios de assinatura
- ▶ Muito caro adquirir novo cliente vs manter existente

Objetivo da Análise:

- ▶ **Entender:** Quais fatores estão associados ao cancelamento?
- ▶ **Comparar:** Como clientes que cancelam diferem dos que ficam?
- ▶ **Identificar:** Quais características são mais comuns em cada grupo?
- ▶ **Revelar:** Padrões não óbvios através de transformações

Dados disponíveis:

- ▶ 7000 clientes
- ▶ 20 variáveis (demográficas, serviços, contratos, pagamentos)
- ▶ Target: Churn (Yes/No)

Telco Churn: Variáveis do Dataset

Variáveis Demográficas:

- ▶ gender: Gênero (Male/Female)
- ▶ SeniorCitizen: Idoso? (0/1)
- ▶ Partner: Tem parceiro? (Yes/No)
- ▶ Dependents: Tem dependentes? (Yes/No)

Serviços Contratados:

- ▶ PhoneService, MultipleLines: Telefone
- ▶ InternetService: DSL, Fiber optic, No
- ▶ OnlineSecurity, OnlineBackup, DeviceProtection, TechSupport: Serviços adicionais
- ▶ StreamingTV, StreamingMovies: Streaming

Contrato e Pagamento:

- ▶ Contract: Month-to-month, One year, Two year
- ▶ PaperlessBilling: Sem papel? (Yes/No)
- ▶ PaymentMethod: Método de pagamento (4 tipos)

Telco Churn: Variáveis do Dataset (cont.)

Variáveis Numéricas:

- ▶ tenure: Meses como cliente (0-72)
- ▶ MonthlyCharges: Valor mensal (R\$ 18-118)
- ▶ TotalCharges: Total gasto (R\$ 18-8684)

Target (o que queremos analisar):

- ▶ Churn: Cancelou? (Yes/No)
- ▶ Aproximadamente 27% de churn
- ▶ Esta é nossa variável de interesse

Tipos de variáveis presentes:

- ▶ ✓ Numéricas contínuas (charges, tenure)
- ▶ ✓ Categóricas binárias (Yes/No)
- ▶ ✓ Categóricas ordinais (Contract tem ordem)
- ▶ ✓ Categóricas nominais (PaymentMethod sem ordem)

Por que Este Dataset é Bom Exemplo?

Requer todas as transformações que vamos aprender:

1. Feature Engineering

- ▶ Podemos criar: gasto médio mensal, razão charges/tenure
- ▶ Agregações: total de serviços contratados
- ▶ Flags: cliente premium? valor alto?

2. Encoding

- ▶ Binárias: gender, Partner, Dependents → 0/1
- ▶ Ordinais: Contract (ordem de compromisso)
- ▶ Nominais: PaymentMethod, InternetService

Por que Este Dataset é Bom Exemplo?

Requer todas as transformações que vamos aprender:

3. Normalização

- ▶ tenure (0-72) vs TotalCharges (18-8684)
- ▶ Escalas muito diferentes precisam normalização

4. Fácil de Entender

- ▶ Telecomunicações é familiar
- ▶ Variáveis têm significado claro
- ▶ Resultados são interpretáveis

Objetivos de Aprendizado

Ao final desta aula, você será capaz de:

Conhecimento:

- ▶ Identificar quando e por que transformar dados
- ▶ Compreender diferentes tipos de encoding
- ▶ Conhecer técnicas de normalização e suas diferenças
- ▶ Entender o conceito de feature engineering

Habilidades:

- ▶ Criar features úteis a partir de dados existentes
- ▶ Aplicar encoding apropriado para cada tipo de variável
- ▶ Normalizar dados com Min-Max, Standardization, Robust Scaling
- ▶ Construir pipeline completo de transformação
- ▶ Validar transformações através de análise exploratória

Objetivos de Aprendizado (cont.)

Atitudes:

- ▶ Pensar criativamente sobre features
- ▶ Documentar decisões de transformação
- ▶ Validar resultados sistematicamente

Pandas: Manipulação e transformação

- ▶ `pd.get_dummies()`: One-Hot Encoding
- ▶ `map()`, `replace()`: Label Encoding
- ▶ `apply()`, `transform()`: Operações customizadas

Scikit-learn: Transformações padronizadas

- ▶ `MinMaxScaler()`: Normalização [0, 1]
- ▶ `StandardScaler()`: Padronização (Z-score)
- ▶ `RobustScaler()`: Scaling robusto a outliers
- ▶ `LabelEncoder()`: Encoding de labels

NumPy: Operações matemáticas

- ▶ `np.log()`, `np.sqrt()`: Transformações matemáticas
- ▶ Operações vetorizadas eficientes

Visualização:

- ▶ Matplotlib, Seaborn
- ▶ Validar transformações visualmente

Ordem das Operações

Pipeline lógico de transformação:

1. Feature Engineering

Criar novas variáveis



2. Encoding

Categorias → números



3. Normalização/Scaling

Ajustar escalas



4. Validação

Verificar resultado

Ordem das Operações

Atenção

A ordem importa! Encoding antes de scaling, por exemplo.

Setup: Carregar Dados

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 from sklearn.preprocessing import MinMaxScaler, StandardScaler,
    RobustScaler
7
8 # Configuracoes
9 pd.set_option('display.max_columns', None)
10 pd.set_option('display.float_format', '{:.2f}'.format)
11 plt.style.use('seaborn-v0_8-darkgrid')
```

Setup: Carregar Dados (cont.)

</> Python

```
1 # URL
2 url = 'https://raw.githubusercontent.com/IBM/telco-customer-churn-on-icp4d
      /master/data/Telco-Customer-Churn.csv'
3
4 # Carregar
5 df_original = pd.read_csv(url)
6 df = df_original.copy()
7 print(f"Dataset carregado: {df.shape}")
8 print(f"Linhas: {df.shape[0]:,} | Colunas: {df.shape[1]}")
9 # Dataset carregado: (7043, 21)
10 # Linhas: 7,043 | Colunas: 21
```

Exploração Inicial

</> Python

```
1 # Primeiras linhas
2 print("=== PRIMEIRAS LINHAS ===")
3 display(df.head())
4
5 # Informacoes gerais
6 print("\n=== INFO ===")
7 df.info()
8
9 # Estatisticas descritivas
10 print("\n=== ESTATISTICAS ===")
11 display(df.describe())
```

Exploração Inicial (cont.)

</> Python

```
1 # Distribuicao do target
2 print("\n=== DISTRIBUICAO CHURN ===")
3 print(df['Churn'].value_counts())
4 print(df['Churn'].value_counts(normalize=True))
5 # Churn: No - 0.73, Yes - 0.27
```

Identificar Tipos de Variáveis

Python

```
1 # Separar variaveis por tipo
2 numeric_cols = df.select_dtypes(include=['int64', 'float64']).columns.
    tolist()
3 categorical_cols = df.select_dtypes(include=['object']).columns.tolist()
4 # Remover target da lista
5 if 'Churn' in categorical_cols:
6     categorical_cols.remove('Churn')
7 print(f"Numericas ({len(numeric_cols)}): {numeric_cols}")
8 print(f"Categoricas ({len(categorical_cols)}): {categorical_cols}")
9 # Analisar cardinalidade das categoricas
10 print("\n=== CARDINALIDADE ===")
11 for col in categorical_cols:
12     n_unique = df[col].nunique()
13     print(f"{col}: {n_unique} valores unicos")
```

Identificar Tipos de Variáveis (cont.)

</> Python

```
1 # === CARDINALIDADE ===  
2 # customerID: 7043 valores unicos  
3 # gender: 2 valores unicos  
4 # Partner: 2 valores unicos  
5 # Dependents: 2 valores unicos  
6 # PhoneService: 2 valores unicos  
7 # MultipleLines: 3 valores unicos  
8 # InternetService: 3 valores unicos  
9 # OnlineSecurity: 3 valores unicos  
10 # OnlineBackup: 3 valores unicos  
11 # DeviceProtection: 3 valores unicos  
12 # ...  
13 # PaymentMethod: 4 valores unicos  
14 # TotalCharges: 6531 valores unicos
```

Fluxo de Transformação

Processo sistemático que seguiremos:

ETAPA 1: Feature Engineering (Bloco 1)

- ▶ Criar variáveis derivadas úteis
- ▶ Agregações e combinações
- ▶ Flags e indicadores
- ▶ Validar utilidade (correlação com Churn)

ETAPA 2: Encoding (Bloco 2)

- ▶ Identificar tipo de cada categórica
- ▶ Aplicar encoding apropriado
- ▶ Binárias → 0/1
- ▶ Ordinais → Label Encoding
- ▶ Nominais → One-Hot Encoding

Fluxo de Transformação (cont.)

ETAPA 3: Normalização (Bloco 3)

- ▶ Identificar variáveis com escalas muito diferentes
- ▶ Escolher técnica apropriada
- ▶ Aplicar scaling
- ▶ Validar resultado

ETAPA 4: Pipeline Final (Bloco 4)

- ▶ Integrar todas as transformações
- ▶ Criar função reutilizável
- ▶ Documentar decisões
- ▶ Análise exploratória avançada com dados transformados

Nota Importante

Cada etapa será validada antes de prosseguir!

Considerações Importantes

1. Sempre preserve os dados originais

- ▶ Trabalhe em cópias
- ▶ Mantenha dataset original intacto
- ▶ Facilita comparações antes/depois

2. Documente suas decisões

- ▶ Por que criou cada feature?
- ▶ Por que escolheu aquele encoding?
- ▶ Raciocínio claro e reproduzível

3. Valide cada transformação

- ▶ Verificar se faz sentido
- ▶ Analisar correlações
- ▶ Visualizar resultados
- ▶ Estatísticas descritivas

Considerações Importantes

4. Pense no significado do negócio

- ▶ Features devem fazer sentido prático
- ▶ Consulte especialistas do domínio
- ▶ Interpretabilidade é importante

5. Comece simples, evolua gradualmente

- ▶ Transformações óbvias primeiro
- ▶ Teste uma de cada vez
- ▶ Adicione complexidade progressivamente

6. Cuidado com "data leakage"

- ▶ Não use informação "do futuro"
- ▶ Importante quando eventualmente usar ML
- ▶ Boa prática: sempre separar dados para validação

Considerações Importantes

Atenção

Transformação é ciência + arte + conhecimento de domínio!

Expectativas Realistas

O que você conseguirá fazer hoje:

- ▶ ✓ Criar features úteis sistematicamente
- ▶ ✓ Encodar qualquer tipo de variável categórica
- ▶ ✓ Normalizar dados para análises comparativas
- ▶ ✓ Construir pipeline completo reproduzível
- ▶ ✓ Validar transformações através de EDA

O que NÃO faremos hoje (mas você aprenderá depois):

- ▶ ✗ Construir modelos de Machine Learning
- ▶ ✗ Fazer previsões automáticas
- ▶ ✗ Otimizar hiperparâmetros
- ▶ ✗ Avaliar performance preditiva

Nota Importante

Foco: Transformação para análise exploratória avançada.
Bônus: Mesmas técnicas são usadas em ML!

Caso de Uso Motivador

Pergunta de Negócio:

"Quais fatores estão mais fortemente associados ao cancelamento de clientes?"

Para responder, precisamos:

1. **Enriquecer dados** com features derivadas
 - ▶ Ex: "Gasto médio mensal", "Total de serviços"
2. **Converter categorias** para análise quantitativa
 - ▶ Calcular correlações com Churn
3. **Normalizar escalas** para comparação justa
 - ▶ Identificar quais variáveis têm maior "peso"
4. **Visualizar padrões** de forma clara
 - ▶ Comparar clientes que cancelam vs que ficam

Ao final: Teremos insights acionáveis para reduzir churn!

Estrutura da Aula

Bloco 0: Preview de Machine Learning

- ▶ Contexto: o que é ML e por que importa
- ▶ Relação entre transformação e ML

Bloco 1: Feature Engineering

- ▶ Criar variáveis derivadas úteis

Bloco 2: Encoding

- ▶ Converter categorias em números

Bloco 3: Normalização

- ▶ Ajustar escalas para comparabilidade

Bloco 4: Pipeline Final

- ▶ Integrar tudo + análise final

Bloco 0: Preview de Machine Learning

Por que falar de ML numa aula de transformação?

- ▶ As transformações que aprenderemos são *essenciais* para ML
- ▶ Mas também são úteis para análise sem ML
- ▶ Importante entender o contexto mais amplo
- ▶ Ajuda a entender *por que* fazemos cada transformação

Estrutura deste bloco:

1. O que é Machine Learning? (conceitual)
2. Tipos de problemas de ML
3. Por que ML precisa de dados transformados?
4. Pipeline típico de Ciência de Dados
5. Transição para transformações (foco da aula)

Bloco 0: Preview de Machine Learning

Nota Importante

NÃO iremos aprender ML hoje! Esta parte da aula é para contexto somente.

O que é Machine Learning?

Definição simples:

"Ensinar computadores a aprender padrões nos dados e fazer previsões automaticamente"

Diferença fundamental:

Programação tradicional:

- ▶ Você escreve as regras explicitamente
- ▶ Ex: SE idade > 65 ENTÃO "idoso"

Machine Learning:

- ▶ Computador descobre as regras automaticamente
- ▶ Você fornece exemplos (dados)
- ▶ Algoritmo aprende padrões
- ▶ Ex: Mostrar 1000 casos de churn, algoritmo aprende padrão

Machine Learning: Analogia

Aprender a identificar spam:

Abordagem tradicional (regras):

- ▶ SE contém "ganhe dinheiro" ENTÃO spam
- ▶ SE contém "clique aqui" ENTÃO spam
- ▶ SE remetente desconhecido ENTÃO spam
- ▶ Problema: spammers se adaptam, regras ficam desatualizadas

Abordagem Machine Learning:

- ▶ Dar 10.000 exemplos de emails (spam e não-spam)
- ▶ Algoritmo aprende padrões automaticamente
- ▶ Consegue identificar spam novo que nunca viu
- ▶ Se adapta com novos dados

Machine Learning: Analogia

Nota Importante

ML = Aprender com exemplos ao invés de programar regras

Exemplos Cotidianos de ML

Recomendações:

- ▶ Netflix: "Você pode gostar deste filme"
- ▶ Amazon: "Clientes que compraram X também compraram Y"
- ▶ Spotify: Playlists personalizadas

Reconhecimento:

- ▶ Reconhecimento facial (desbloquear celular)
- ▶ Assistentes de voz (Alexa, Siri)
- ▶ Tradução automática (Google Translate)

Previsões:

- ▶ Previsão do tempo
- ▶ Detecção de fraude em cartão de crédito
- ▶ Diagnóstico médico assistido
- ▶ Risco de crédito em bancos

Tipos de Problemas de ML

1. Classificação (prever categoria)

- ▶ Exemplo: Email é spam ou não-spam?
- ▶ Exemplo: Cliente vai cancelar ou não? (Churn)
- ▶ Exemplo: Doença está presente ou ausente?
- ▶ Resposta: categoria/classe

2. Regressão (prever número)

- ▶ Exemplo: Qual será o preço desta casa?
- ▶ Exemplo: Quantas vendas teremos no próximo mês?
- ▶ Exemplo: Qual a temperatura amanhã?
- ▶ Resposta: valor numérico contínuo

3. Clustering (agrupar similares)

- ▶ Exemplo: Segmentar clientes em grupos
- ▶ Exemplo: Agrupar notícias por tema
- ▶ Sem "resposta correta"prévia

Problema de Classificação: Detalhes

Exemplo: Prever Churn (classificação binária)

Dados de entrada (features):

- ▶ tenure = 12 meses
- ▶ MonthlyCharges = R\$ 85
- ▶ Contract = "Month-to-month"
- ▶ InternetService = "Fiber optic"
- ▶ ... (outras variáveis)

Saída desejada (target):

- ▶ Churn = "Yes" ou "No"
- ▶ Ou: probabilidade de churn (0-100%)

Como funciona:

1. Treinar com exemplos históricos (já sabemos quem cancelou)
2. Algoritmo aprende padrões
3. Usar para prever novos clientes

Problema de Regressão: Detalhes

Exemplo: Prever Preço de Casa

Dados de entrada (features):

- ▶ Área = 120 m²
- ▶ Quartos = 3
- ▶ Banheiros = 2
- ▶ Localização = "Centro"
- ▶ Idade = 10 anos

Saída desejada (target):

- ▶ Preço = R\$ 450.000
- ▶ Valor numérico contínuo

Como funciona:

1. Treinar com casas já vendidas (preço conhecido)
2. Algoritmo aprende relação features → preço
3. Usar para estimar preço de casas novas

Por que ML Precisa de Dados Transformados?

Problema 1: Algoritmos trabalham com números

- ▶ ML usa matemática (álgebra linear, cálculo)
- ▶ Não consegue processar texto diretamente
- ▶ "Month-to-month" precisa virar número
- ▶ **Solução: Encoding**

Problema 2: Escalas diferentes confundem

- ▶ tenure: 0-72 vs TotalCharges: 0-8684
- ▶ Algoritmo pode dar peso errado
- ▶ Variável com números maiores parece "mais importante"
- ▶ **Solução: Normalização**

Problema 3: Informação está escondida

- ▶ Padrões não são óbvios nas variáveis originais
- ▶ Combinações podem ser mais úteis
- ▶ **Solução: Feature Engineering**

Exemplo: Por que Encoding Importa

Variável: Contract (tipo de contrato)

- ▶ Valores: "Month-to-month", "One year", "Two year"

Problema:

- ▶ Algoritmo não consegue processar texto
- ▶ Não dá para calcular: "Month-to-month" + "One year" = ?

Solução: Encoding

- ▶ Month-to-month → 0
- ▶ One year → 1
- ▶ Two year → 2

Agora o algoritmo consegue:

- ▶ Processar matematicamente
- ▶ Identificar padrão: contratos mais longos = menos churn
- ▶ Fazer previsões baseadas nisso

Exemplo: Por que Normalização Importa

Duas variáveis importantes:

- ▶ tenure: 0-72 (meses)
- ▶ TotalCharges: 18-8684 (reais)

Problema:

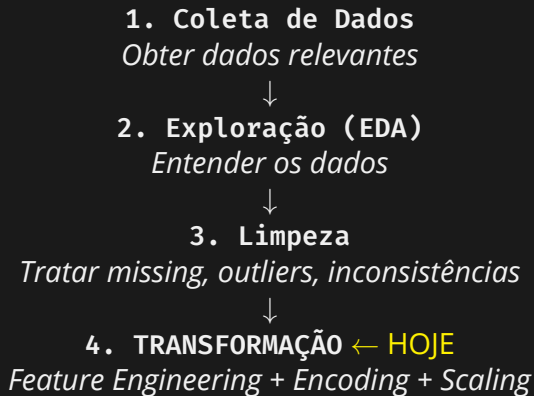
- ▶ TotalCharges tem números muito maiores
- ▶ Algoritmos baseados em distância ficam confusos
- ▶ Podem dar peso excessivo a TotalCharges
- ▶ Não porque é mais importante, mas porque números são maiores!

Solução: Normalização

- ▶ Colocar ambas na mesma escala [0, 1]
- ▶ tenure: 0-72 \rightarrow 0.0-1.0
- ▶ TotalCharges: 18-8684 \rightarrow 0.0-1.0
- ▶ Agora são comparáveis de forma justa

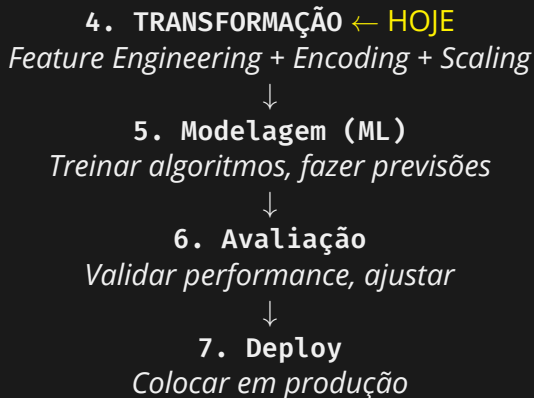
Pipeline Típico de Ciência de Dados com ML

Processo completo:



Pipeline Típico de Ciência de Dados com ML

Processo completo:



Onde Estamos no Curso?

Aula	Tema	Status
1-2	Python Básico	✓ Completo
3-4	NumPy	✓ Completo
5-6	Pandas	✓ Completo
7-8	EDA	✓ Completo
9	Limpeza	✓ Completo
10	Transformação	← HOJE
11	Séries Temporais	Futuro
12	Dados Multidimensionais	Futuro

Sobre Machine Learning:

- ▶ NÃO faz parte deste curso (introdutório)
- ▶ Será visto em cursos avançados
- ▶ Mas: as técnicas de hoje são fundamentais para ML
- ▶ Você já estará preparado quando chegar lá!

Relação: Análise vs Machine Learning

Mesmas técnicas, objetivos diferentes:

Análise Exploratória (nosso foco hoje):

- ▶ Calcular correlações entre variáveis
- ▶ Comparar grupos (clientes que cancelam vs não cancelam)
- ▶ Visualizações mais informativas
- ▶ Identificar fatores associados ao target
- ▶ Gerar insights acionáveis

Machine Learning (cursos futuros):

- ▶ Usar *exatamente as mesmas transformações*
- ▶ Mas para treinar modelos preditivos
- ▶ Objetivo: fazer previsões automáticas
- ▶ Aplicar em dados novos

Exemplo Concreto: Duas Abordagens

Problema: Entender Churn

Abordagem 1: Análise (hoje)

1. Transformar dados (encoding, normalização, features)
2. Calcular correlações: qual variável mais correlacionada com Churn?
3. Visualizar: comparar perfis de quem cancela vs quem fica
4. Insight: "Clientes com contrato mensal têm 3x mais churn"
5. Ação: Oferecer incentivos para contratos anuais

Abordagem 2: ML (futuro)

1. Transformar dados (*da mesma forma*)
2. Treinar modelo de classificação
3. Para cada cliente novo: prever probabilidade de churn
4. Sistema automático: se $\text{prob} > 70\%$, acionar retenção
5. Prevenção proativa

O que Aprenderemos Hoje

Foco: Transformações para análise exploratória

Feature Engineering:

- ▶ Criar variáveis como "gasto médio mensal"
- ▶ Objetivo: revelar padrões não óbvios
- ▶ Validação: correlação com Churn, visualizações

Encoding:

- ▶ Converter categorias em números
- ▶ Objetivo: permitir análises quantitativas
- ▶ Validação: verificar se faz sentido

Normalização:

- ▶ Colocar variáveis em escalas comparáveis
- ▶ Objetivo: comparações justas, visualizações melhores
- ▶ Validação: estatísticas descritivas, gráficos

O que NÃO Faremos Hoje

Deixando claro (para não gerar expectativas erradas):

NÃO faremos:

- ▶ ✗ Treinar modelos de Machine Learning
- ▶ ✗ Fazer previsões automáticas
- ▶ ✗ Usar algoritmos de classificação/regressão
- ▶ ✗ Avaliar métricas como accuracy, precision, recall
- ▶ ✗ Split train/test para validação de modelo
- ▶ ✗ Otimização de hiperparâmetros

Por quê?

- ▶ Este é um curso introdutório
- ▶ ML requer curso dedicado (muitos conceitos novos)
- ▶ Hoje: fundação essencial (transformação de dados)
- ▶ Futuro: você usará essas mesmas técnicas em ML

Contexto

O que sabemos:

- ▶ ✓ ML aprende padrões em dados automaticamente
- ▶ ✓ ML precisa de dados bem transformados
- ▶ ✓ As transformações também melhoram análises sem ML
- ▶ ✓ Vamos focar em análise exploratória

Próximos passos:

- ▶ **Bloco 1:** Feature Engineering (criar variáveis úteis)
- ▶ **Bloco 2:** Encoding (categorias → números)
- ▶ **Bloco 3:** Normalização (ajustar escalas)
- ▶ **Bloco 4:** Pipeline final + análise avançada

Método de validação:

- ▶ Após cada transformação: análise exploratória
- ▶ Visualizações, correlações, estatísticas
- ▶ Verificar se faz sentido para o negócio

Bloco 1

Feature Engineering

O que vamos aprender

Conceito:

- ▶ O que é feature engineering
- ▶ Por que criar novas variáveis
- ▶ Tipos de features que podemos criar

Prática:

- ▶ 12 técnicas diferentes com código
- ▶ Aplicadas ao dataset Telco Churn
- ▶ Validação através de análise exploratória

Objetivo:

- ▶ Revelar padrões escondidos nos dados
- ▶ Criar variáveis mais interpretáveis
- ▶ Melhorar análises e visualizações

O que é Feature Engineering?

Definição:

"Processo de criar novas variáveis (features) a partir das variáveis existentes no dataset"

Por que fazer isso?

- ▶ Informação útil está "escondida" nos dados brutos
- ▶ Combinações de variáveis revelam padrões
- ▶ Algumas relações não são óbvias
- ▶ Features derivadas podem ser mais interpretáveis

Exemplo simples:

- ▶ Temos: TotalCharges e tenure
- ▶ Criamos: $\text{AvgMonthlySpend} = \text{TotalCharges} / \text{tenure}$
- ▶ Por quê? Gasto médio mensal é mais interpretável
- ▶ Facilita comparações entre clientes

Por que Feature Engineering Importa?

Impacto na qualidade da análise:

Sem feature engineering:

- ▶ Análise limitada às variáveis originais
- ▶ Padrões importantes podem passar despercebidos
- ▶ Comparações são difíceis
- ▶ Visualizações menos informativas

Com feature engineering:

- ▶ Padrões ocultos são revelados
- ▶ Comparações mais justas e interpretáveis
- ▶ Visualizações mais ricas
- ▶ Insights mais profundos
- ▶ Recomendações mais embasadas

Por que Feature Engineering Importa?

Nota Importante

Boas features = Diferença entre análise superficial e insights profundos

Tipos de Feature Engineering

12 técnicas que aprenderemos:

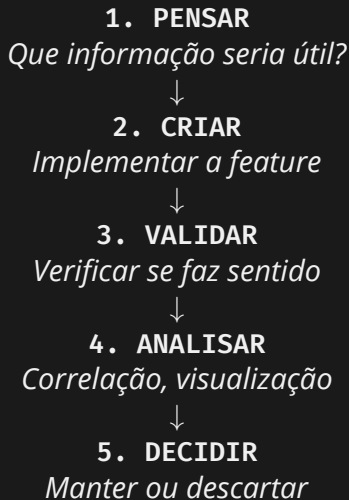
1. Operações aritméticas simples (+, -, *, /)
2. Diferenças entre variáveis
3. Razões e proporções
4. Binning (discretização)
5. Agregações (somas, contagens)
6. Flags booleanas (indicadores 0/1)
7. Interações entre variáveis
8. Extração de informação de strings
9. Features temporais (conceitual)
10. Percentis e rankings
11. Z-score para detecção de anomalias
12. Features baseadas em conhecimento do domínio

Tipos de Feature Engineering

Para cada técnica:

- ▶ Exemplo prático com código
- ▶ Validação através de análise exploratória
- ▶ Interpretação do resultado

Feature Engineering: Método sistemático



Feature 1: Operações Aritméticas Simples

Gasto médio mensal

</> Python

```
1 # Feature: Average Monthly Spend
2 # Raciocínio: Clientes com gasto médio alto podem ter
3 # comportamento diferente
4 df['TotalCharges'] = pd.to_numeric(df['TotalCharges'], errors='coerce')
5 df['AvgMonthlySpend'] = df['TotalCharges'] / df['tenure']
6 # Tratar casos especiais
7 # Se tenure = 0, evitar divisão por zero
8 df['AvgMonthlySpend'] = df['AvgMonthlySpend'].replace([np.inf, -np.inf],
9 np.nan)
9 # Para tenure = 0, usar MonthlyCharges como aproximação
10 df.loc[df['tenure'] == 0, 'AvgMonthlySpend'] = df.loc[df['tenure'] == 0, '
    MonthlyCharges']
11 print("=== FEATURE CRIADA: AvgMonthlySpend ===")
12 print(df['AvgMonthlySpend'].describe())
```

Por que AvgMonthlySpend é Útil?

Validação através de análise:

</> Python

```
1 # Comparar entre clientes que cancelam e que ficam
2 print("=== COMPARACAO POR CHURN ===")
3 print(df.groupby('Churn')['AvgMonthlySpend'].describe())
4
5 # Visualizar distribuicao
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 fig, ax = plt.subplots(1, 2, figsize=(14, 5))
10
11 # Boxplot
12 sns.boxplot(data=df, x='Churn', y='AvgMonthlySpend', ax=ax[0])
13 ax[0].set_title('Gasto Médio Mensal por Churn')
```

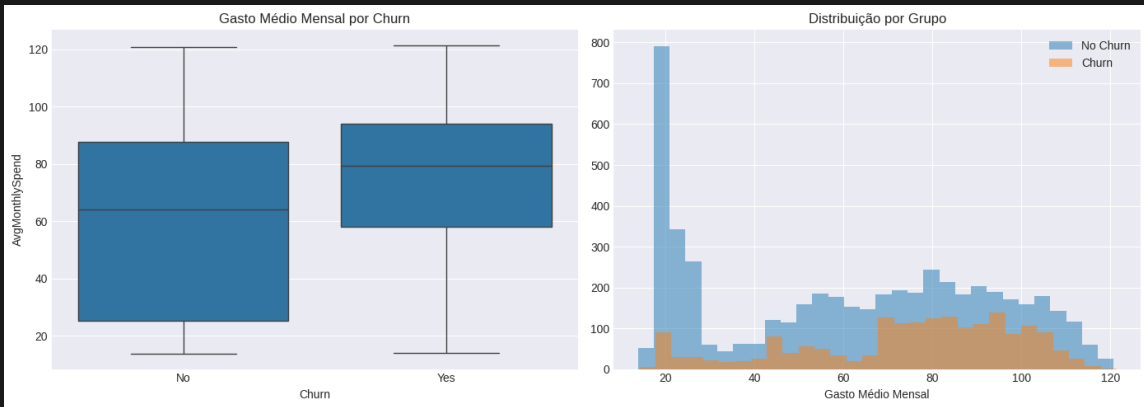
Por que AvgMonthlySpend é Útil? (cont.)

Validação através de análise:

</> Python

```
1 # Histograma
2 df[df['Churn']=='No']['AvgMonthlySpend'].hist(bins=30, alpha=0.5, label='
    No Churn', ax=ax[1])
3 df[df['Churn']=='Yes']['AvgMonthlySpend'].hist(bins=30, alpha=0.5, label='
    Churn', ax=ax[1])
4 ax[1].set_xlabel('Gasto Médio Mensal')
5 ax[1].set_title('Distribuição por Grupo')
6 ax[1].legend()
7
8 plt.tight_layout()
9 plt.show()
```

Por que AvgMonthlySpend é Útil? (cont.)



Feature 2: Diferença entre Variáveis

Variação no gasto

</> Python

```
1 # Feature: Charge Difference
2 # Raciocínio: Se MonthlyCharges >> AvgMonthlySpend,
3 # pode indicar aumento recente (cliente insatisfeito?)
4 df['ChargeDiff'] = df['MonthlyCharges'] - df['AvgMonthlySpend']
5 print("=== FEATURE CRIADA: ChargeDiff ===")
6 print(df['ChargeDiff'].describe())
7 # Interpretar
8 print("\n=== INTERPRETACAO ===")
9 print("ChargeDiff > 0: Gasto atual maior que média (aumento recente)")
10 print("ChargeDiff < 0: Gasto atual menor que média (diminuição)")
11 print("ChargeDiff 0: Gasto estável")
12 # Analisar por churn
13 print("\n=== MEDIA POR CHURN ===")
14 print(df.groupby('Churn')['ChargeDiff'].mean())
```

Feature 3: Razões/Proporções

Proporção do tempo como cliente

</> Python

```
1 # Feature: Tenure Ratio
2 # Raciocínio: Normalizar tenure em relação ao máximo possível
3 # Facilita comparações
4
5 max_tenure = df['tenure'].max()
6 df['TenureRatio'] = df['tenure'] / max_tenure
7
8 print("=== FEATURE CRIADA: TenureRatio ===")
9 print(f"Range: [{df['TenureRatio'].min():.2f}, {df['TenureRatio'].max():.2f}]" )
```

Feature 3: Razões/Proporções (cont.)

Proporção do tempo como cliente

</> Python

```
1 # Criar categorias
2 df['TenureCategory'] = pd.cut(df['TenureRatio'],
3                               bins=[0, 0.33, 0.66, 1.0],
4                               labels=['Novo', 'Médio', 'Veterano'])
5
6 print("\n=== DISTRIBUICAO POR CATEGORIA ===")
7 print(pd.crosstab(df['TenureCategory'], df['Churn'], normalize='index'))
```

Feature 4: Binning (Discretização)

Transformar numérica contínua em categórica ordinal

</> Python

```
1 # Feature: MonthlyCharges em faixas
2 # Raciocínio: Facilita análise por grupos de preço
3
4 # Método 1: Bins com larguras iguais
5 df['ChargesBin_Equal'] = pd.cut(df['MonthlyCharges'],
6                                 bins=5,
7                                 labels=['Muito Baixo', 'Baixo', 'Médio',
8                                       'Alto', 'Muito Alto'])
```

Feature 4: Binning (Discretização) (cont.)

Transformar numérica contínua em categórica ordinal

</> Python

```
1 # Método 2: Quantis (mesma quantidade em cada bin)
2 df['ChargesBin_Quantile'] = pd.qcut(df['MonthlyCharges'],
3                                     q=5,
4                                     labels=['Q1', 'Q2', 'Q3', 'Q4', 'Q5']
5                                     ])
6
7 print("=== BINNING CRIADO ===")
8 print("\nDistribuição - Bins Iguais:")
9 print(df['ChargesBin_Equal'].value_counts().sort_index())
```

Feature 4: Binning (Discretização) (cont.)

Python

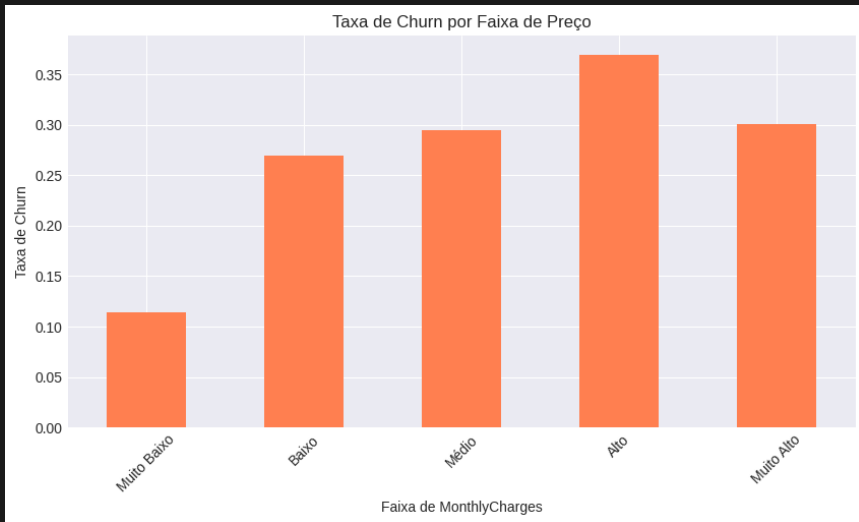
```
1 print("\nDistribuição - Quantis:")
2 print(df['ChargesBin_Quantile'].value_counts().sort_index())
3
4 # Analisar churn por faixa de preço
5 print("\n=== CHURN RATE POR FAIXA ===")
6 churn_by_bin = df.groupby('ChargesBin_Equal', observed=False)['Churn'].
    apply(
7     lambda x: (x == 'Yes').mean()
8 )
9 print(churn_by_bin)
```

Feature 4: Binning (Discretização) (cont.)

</> Python

```
1 # Visualizar
2 churn_by_bin.plot(kind='bar', color='coral', figsize=(10, 5))
3 plt.title('Taxa de Churn por Faixa de Preço')
4 plt.ylabel('Taxa de Churn')
5 plt.xlabel('Faixa de MonthlyCharges')
6 plt.xticks(rotation=45)
7 plt.show()
```

Feature 4: Binning (Discretização) (cont.)



Por que Binning é Útil?

Vantagens da discretização:

1. Interpretabilidade

- ▶ "Clientes de preço alto" é mais claro que "R\$ 89,47"
- ▶ Facilita comunicação com não-técnicos
- ▶ Categorias são intuitivas

2. Análise por grupos

- ▶ Comparar comportamento entre faixas
- ▶ Identificar thresholds importantes
- ▶ Ex: "Churn dispara acima de R\$ 70/mês"

3. Lidar com não-linearidade

- ▶ Relação pode não ser linear
- ▶ Ex: churn alto em valores muito baixos E muito altos
- ▶ Bins capturam melhor esses padrões

Feature 5: Agregação - Contar Serviços

Total de serviços adicionais contratados

</> Python

```
1 # Feature: Total Services
2 # Raciocínio: Clientes com mais serviços podem estar
3 # mais "engajados" (menos propensos a cancelar?)
4
5 # Lista de colunas de serviços
6 service_cols = [
7     'PhoneService', 'MultipleLines',
8     'InternetService', 'OnlineSecurity',
9     'OnlineBackup', 'DeviceProtection',
10    'TechSupport', 'StreamingTV', 'StreamingMovies'
11 ]
```

Feature 5: Agregação - Contar Serviços (cont.)

Total de serviços adicionais contratados

</> Python

```
1 # Contar quantos serviços tem "Yes" ou não é "No"
2 df['TotalServices'] = 0
3
4 for col in service_cols:
5     # Adiciona 1 se não for "No" ou "No internet service" ou "No phone
    service"
6     df['TotalServices'] += (~df[col].isin(['No', 'No internet service',
7                                             'No phone service'])).astype(
    int)
```

Feature 5: Agregação - Contar Serviços (cont.)

</> Python

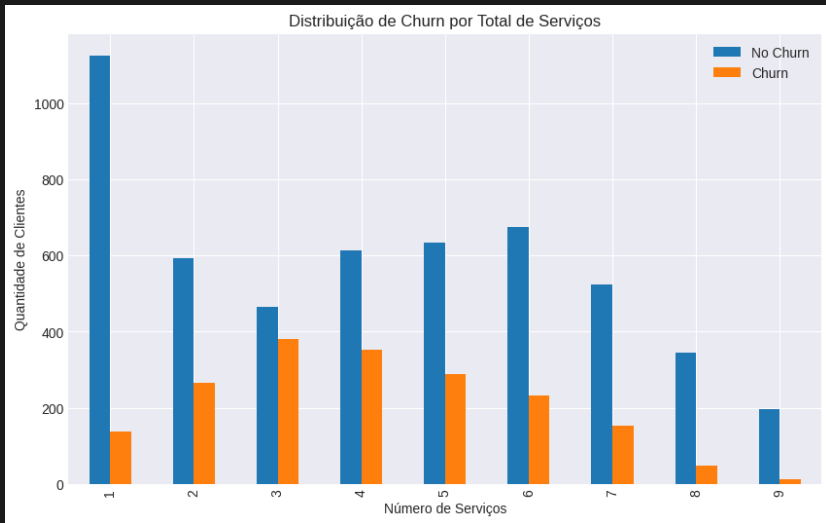
```
1 print("=== FEATURE CRIADA: TotalServices ===")
2 print(df['TotalServices'].describe())
3
4 print("\n=== DISTRIBUICAO ===")
5 print(df['TotalServices'].value_counts().sort_index())
6
7 # Analisar relação com Churn
8 print("\n=== MEDIA DE SERVICOS POR CHURN ===")
9 print(df.groupby('Churn')['TotalServices'].mean())
```

Feature 5: Agregação - Contar Serviços (cont.)

</> Python

```
1 # Visualizar
2 fig, ax = plt.subplots(figsize=(10, 6))
3 df.groupby(['TotalServices', 'Churn']).size().unstack().plot(kind='bar',
4     ax=ax)
5 ax.set_title('Distribuição de Churn por Total de Serviços')
6 ax.set_xlabel('Número de Serviços')
7 ax.set_ylabel('Quantidade de Clientes')
8 plt.legend(['No Churn', 'Churn'])
9 plt.show()
```

Feature 5: Agregação - Contar Serviços (cont.)



Agregação: Versão Otimizada

Método mais eficiente com apply:

</> Python

```
1 # Alternativa: usar apply (mais eficiente)
2 def count_services(row):
3     """Conta serviços ativos excluindo 'No' e variações"""
4     count = 0
5     for col in service_cols:
6         value = str(row[col]).lower()
7         if value not in ['no', 'no internet service', 'no phone service']:
8             count += 1
9     return count
10
11 df['TotalServices'] = df.apply(count_services, axis=1)
```

Agregação: Versão Otimizada (cont.)

Método mais eficiente com apply:

</> Python

```
1 # Ou ainda mais simples com sum:
2 df['TotalServices'] = df[service_cols].apply(
3     lambda x: (~x.isin(['No', 'No internet service',
4                          'No phone service'])).sum(),
5     axis=1
6 )
7
8 print("Total de serviços calculado!")
9 print(df['TotalServices'].value_counts())
```


Feature 6: Flags Booleanas

Indicadores binários úteis

</> Python

```
1 # Feature: Has Internet
2 df['HasInternet'] = (df['InternetService'] != 'No').astype(int)
3 # Feature: Has Phone
4 df['HasPhone'] = (df['PhoneService'] == 'Yes').astype(int)
5 # Feature: Senior with Dependents (combinação interessante)
6 df['SeniorWithDependents'] = (
7     (df['SeniorCitizen'] == 1) &
8     (df['Dependents'] == 'Yes')
9 ).astype(int)
10 # Feature: Premium Customer (muitos serviços + contrato longo)
11 df['IsPremium'] = (
12     (df['TotalServices'] >= 6) &
13     (df['Contract'].isin(['One year', 'Two year']))
14 ).astype(int)
```

Flags: Por que São Úteis?

Validando as flags criadas:

</> Python

```
1 # Analisar cada flag
2 flags = ['HasInternet', 'HasPhone', 'SeniorWithDependents', 'IsPremium']
3
4 print("=== ANALISE DE FLAGS ===\n")
5 for flag in flags:
6     print(f"\n--- {flag} ---")
7
8     # Distribuição
9     print(f"Total com flag=1: {df[flag].sum()} ({df[flag].mean()*100:.1f}%)\n")
```

Flags: Por que São Úteis?

Validando as flags criadas:

</> Python

```
1  # Taxa de churn
2  churn_rate = df.groupby(flag)['Churn'].apply(
3      lambda x: (x == 'Yes').mean()
4  )
5  print("Taxa de Churn:")
6  print(churn_rate)
7
8  # Diferença
9  diff = churn_rate[1] - churn_rate[0]
10 print(f"Diferença: {diff*100:.1f} pontos percentuais")
```

Interações entre Variáveis

O que são interações?

- ▶ Efeito combinado de duas variáveis
- ▶ Uma variável modifica o efeito da outra
- ▶ Relação não é simplesmente aditiva

Exemplo conceitual:

- ▶ Variável A: Preço alto
- ▶ Variável B: Contrato curto
- ▶ Isoladamente: Ambos aumentam churn
- ▶ **Interação:** Combinação (preço alto + contrato curto) aumenta MUITO mais
- ▶ Efeito não é $A + B$, é $A * B$

Como criar:

- ▶ Multiplicação de variáveis
- ▶ Combinação de flags
- ▶ Categorias compostas

Quando combinar variáveis (interações)?

Sinais de que interação pode ser útil:

1. Conhecimento do domínio sugere

- ▶ "Idosos com dependentes" podem ter padrão específico
- ▶ "Preço alto + contrato curto" é combinação arriscada
- ▶ Especialista identifica combinações relevantes

2. Análise exploratória revela

- ▶ Crosstabs mostram padrões interessantes
- ▶ Ex: Churn é 80% quando A=high E B=yes
- ▶ Mas apenas 20% quando só A=high ou só B=yes

3. Hipóteses de negócio

- ▶ "Clientes jovens com múltiplas linhas têm perfil específico"
- ▶ "Serviço de fibra + streaming indica uso intenso"
- ▶ Testar hipóteses com interações

Quando NÃO combinar variáveis?

Evite interações desnecessárias:

1. Sem justificativa clara

- ▶ Não combine aleatoriamente
- ▶ "Gênero * PaymentMethod" provavelmente não faz sentido
- ▶ Precisa haver raciocínio de negócio

2. Variáveis já correlacionadas

- ▶ TotalCharges * tenure é redundante (TC já é função de tenure)
- ▶ MonthlyCharges * AvgMonthlySpend é circular

3. Explosão combinatória

- ▶ Cuidado com variáveis de alta cardinalidade
- ▶ PaymentMethod (4) * InternetService (3) = 12 combinações
- ▶ Pode criar muitas features sem valor

Feature 7: Interações entre Variáveis

Criando interações úteis:

</> Python

```
1 # Interação 1: Contrato Curto + Preço Alto
2 # Hipótese: Combinação perigosa para churn
3 df['ShortContract_HighPrice'] = (
4     (df['Contract'] == 'Month-to-month') &
5     (df['MonthlyCharges'] > df['MonthlyCharges'].quantile(0.75))
6 ).astype(int)
7
8 # Interação 2: Sem segurança + Idoso
9 # Hipótese: Idosos sem proteção podem estar insatisfeitos
10 df['Senior_NoSecurity'] = (
11     (df['SeniorCitizen'] == 1) &
12     (df['OnlineSecurity'] == 'No')
13 ).astype(int)
```

Feature 7: Interações entre Variáveis (cont.)

Criando interações úteis:

</> Python

```
1 # Interação 3: Fibra + Streaming (usuário power)
2 df['FiberStreamer'] = (
3     (df['InternetService'] == 'Fiber optic') &
4     ((df['StreamingTV'] == 'Yes') | (df['StreamingMovies'] == 'Yes'))
5 ).astype(int)
```


Feature 7: Interações entre Variáveis (cont.)

Validando as interações:

</> Python

```
1 # Analisar impacto das interações
2 interactions = ['ShortContract_HighPrice', 'Senior_NoSecurity', '
   FiberStreamer']
3
4 print("=== ANALISE DE INTERACOES ===\n")
5 for interaction in interactions:
6     print(f"\n--- {interaction} ---")
7
8     # Crosstab com Churn
9     ct = pd.crosstab(df[interaction], df['Churn'], normalize='index')
10    print(ct)
```

Feature 7: Interações entre Variáveis (cont.)

Validando as interações:

</> Python

```
1 # Taxa de churn
2 churn_rate = df[df[interaction] == 1]['Churn'].apply(
3     lambda x: 1 if x == 'Yes' else 0
4 ).mean()
5
6 overall_rate = (df['Churn'] == 'Yes').mean()
7
8 print(f"\nChurn Rate com {interaction}=1: {churn_rate*100:.1f}%")
9 print(f"Churn Rate geral: {overall_rate*100:.1f}%")
10 print(f"Diferença: {(churn_rate - overall_rate)*100:.1f} p.p.")
```

Feature 7: Interações entre Variáveis (cont.)

Validando as interações:

</> Python

```
1 # ...  
2 # Churn Rate com FiberStreamer=1: 39.9%  
3 # Churn Rate geral: 26.5%  
4 # Diferença: 13.4 p.p.
```

Feature 8: Extrair de Strings

Extrair informação útil de textos:

</> Python

```
1 # Feature: Payment Type Category
2 # Extrair categoria principal do método de pagamento
3 def categorize_payment(payment_method):
4     """Categorizar método de pagamento"""
5     if pd.isna(payment_method):
6         return 'Unknown'
7     payment = str(payment_method).lower()
8     if 'electronic' in payment or 'bank' in payment:
9         return 'Automatic'
10    elif 'credit' in payment or 'check' in payment:
11        return 'Manual'
12    else:
13        return 'Other'
```

Feature 8: Extrair de Strings (cont.)

Extrair informação útil de textos:

Python

```
1 df['PaymentCategory'] = df['PaymentMethod'].apply(categorize_payment)
2 print("=== CATEGORIAS DE PAGAMENTO ===")
3 print(df['PaymentCategory'].value_counts())
```

Feature 9: Features Temporais (conceitual)

Conceito: Extrair informação de timestamps

</> Python

```
1 # EXEMPLO CONCEITUAL (dataset não tem datas reais)
2 # ESTE CÓDIGO NÃO VAI FUNCIONAR NO COLAB
3 # Se tivéssemos uma coluna 'SignupDate':
4
5 # df['SignupDate'] = pd.to_datetime(df['SignupDate'])
6
7 # Extrair componentes temporais
8 # df['SignupYear'] = df['SignupDate'].dt.year
9 # df['SignupMonth'] = df['SignupDate'].dt.month
10 # df['SignupDayOfWeek'] = df['SignupDate'].dt.dayofweek
11 # df['SignupQuarter'] = df['SignupDate'].dt.quarter
```

Feature 9: Features Temporais (conceitual) (cont.)

Conceito: Extrair informação de timestamps

</> Python

```
1 # EXEMPLO CONCEITUAL (dataset não tem datas reais)
2 # ESTE CÓDIGO NÃO VAI FUNCIONAR NO COLAB
3
4 # Features derivadas
5 # df['DaysSinceSignup'] = (pd.Timestamp.now() - df['SignupDate']).dt.days
6 # df['IsWeekend'] = df['SignupDayOfWeek'].isin([5, 6]).astype(int)
7 # df['IsEndOfMonth'] = (df['SignupDate'].dt.day > 25).astype(int)
8
9 # Sazonalidade
10 # df['IsSummerSignup'] = df['SignupMonth'].isin([12, 1, 2]).astype(int)
```

Feature 9: Features Temporais (conceitual) (cont.)

Usando 'tenure' como proxy temporal:

</> Python

```
1 # Já que temos tenure (meses), podemos extrair padrões temporais
2 # Converter tenure em anos / Categorizar por período
3 df['TenureYears'] = (df['tenure'] / 12).astype(int)
4 def categorize_tenure(months):
5     if months <= 6:
6         return 'Very New'
7     elif months <= 12:
8         return 'New'
9     elif months <= 24:
10        return 'Established'
11    elif months <= 48:
12        return 'Long-term'
13    else:
14        return 'Veteran'
```


Feature 9: Features Temporais (conceitual) (cont.)

Usando 'tenure' como proxy temporal:

</> Python

```
1 df['TenurePeriod'] = df['tenure'].apply(categorize_tenure)
2
3 print("=== DISTRIBUICAO POR PERIODO ===")
4 print(df['TenurePeriod'].value_counts())
```

Features Temporais: Importância

Por que features temporais importam?

Padrões temporais revelam insights:

- ▶ Sazonalidade: churn pode aumentar em certos meses
- ▶ Dia da semana: comportamento diferente fim de semana
- ▶ Tempo decorrido: "primeiros 3 meses" são críticos
- ▶ Tendências: mudanças ao longo do tempo

Exemplos de aplicação:

- ▶ E-commerce: compras aumentam em novembro/dezembro
- ▶ Telecom: cancelamentos após fim de promoção
- ▶ Bancos: transações suspeitas à noite
- ▶ Saúde: internações aumentam no inverno

Features Temporais: Importância

Nota Importante

Sempre que houver timestamp, extraia componentes temporais!

Feature 10: Percentis e Rankings

Posição relativa no dataset:

</> Python

```
1 # Feature: Rank de TotalCharges
2 # Útil para identificar clientes extremos
3
4 # Percentil (0-100)
5 df['TotalCharges_Percentile'] = df['TotalCharges'].rank(pct=True) * 100
6
7 # Rank absoluto
8 df['TotalCharges_Rank'] = df['TotalCharges'].rank(ascending=False)
9
10 # Categorizar por percentil
11 df['ChargesLevel'] = pd.cut(df['TotalCharges_Percentile'],
12                             bins=[0, 25, 50, 75, 100],
13                             labels=['Bottom 25%', 'Q2', 'Q3', 'Top 25%'])
```

Feature 10: Percentis e Rankings (cont.)

Posição relativa no dataset:

</> Python

```
1 print("=== DISTRIBUICAO POR NIVEL ===")
2 print(df['ChargesLevel'].value_counts())
3
4 # Analisar churn por nível
5 print("\n=== CHURN POR NIVEL ===")
6 print(pd.crosstab(df['ChargesLevel'], df['Churn'], normalize='index'))
```

Feature 11: Z-Score para Detecção de Anomalias

Identificar valores extremos (outliers):

</> Python

```
1 # Feature: Z-score de MonthlyCharges
2 # Identifica clientes com gastos muito atípicos
3
4 from scipy import stats
5
6 # Calcular Z-score
7 df['MonthlyCharges_Zscore'] = stats.zscore(df['MonthlyCharges'])
8
9 # Flag para outliers ( $|Z| > 3$ )
10 df['IsOutlier_Charges'] = (abs(df['MonthlyCharges_Zscore']) > 3).astype(
    int)
```

Feature 11: Z-Score para Detecção de Anomalias

Identificar valores extremos (outliers):

</> Python

```
1 print("=== OUTLIERS IDENTIFICADOS ===")
2 print(f"Total de outliers: {df['IsOutlier_Charges'].sum()}")
3 print(f"Percentual: {df['IsOutlier_Charges'].mean()*100:.2f}%")
4
5 # Analisar outliers
6 print("\n=== CHURN ENTRE OUTLIERS ===")
7 print(pd.crosstab(df['IsOutlier_Charges'], df['Churn'], normalize='index')
      )
```

Feature 12: Conhecimento do Domínio

Features que requerem expertise do negócio:

</> Python

```
1 # Feature: High Risk Profile
2 # Baseado em conhecimento de que certos perfis cancelam mais
3
4 df['HighRiskProfile'] = (
5     (df['Contract'] == 'Month-to-month') & # Sem compromisso
6     (df['tenure'] < 12) & # Cliente novo
7     (df['TotalServices'] <= 2) & # Poucos serviços
8     (df['PaymentMethod'] == 'Electronic check') # Método menos estável
9 ).astype(int)
10
11 # Feature: Value Customer
12 # Perfil de cliente valioso (menos propenso a sair)
```


Feature 12: Conhecimento do Domínio (cont.)

Features que requerem expertise do negócio:

</> Python

```
1 df['ValueCustomer'] = (  
2     (df['Contract'].isin(['One year', 'Two year'])) & # Contrato longo  
3     (df['TotalServices'] >= 4) & # Muitos serviços  
4     (df['tenure'] > 24) & # Cliente antigo  
5     (df['MonthlyCharges'] > df['MonthlyCharges'].median()) # Gasto alto  
6 ).astype(int)
```

Feature 12: Conhecimento do Domínio (cont.)

Validando features baseadas em domínio:

</> Python

```
1 # Analisar profiles criados
2 profiles = ['HighRiskProfile', 'ValueCustomer']
3
4 print("=== ANALISE DE PROFILES ===\n")
5 for profile in profiles:
6     print(f"\n--- {profile} ---")
7
8     # Quantidade
9     count = df[profile].sum()
10    pct = df[profile].mean() * 100
11    print(f"Total: {count} ({pct:.1f}%)")
```

Feature 12: Conhecimento do Domínio (cont.)

Validando features baseadas em domínio:

</> Python

```
1  # Churn rate
2  churn_rate = df[df[profile] == 1]['Churn'].apply(
3      lambda x: 1 if x == 'Yes' else 0
4  ).mean()
5
6  print(f"Churn Rate: {churn_rate*100:.1f}%")
7
8  # Comparar com baseline
9  baseline = (df['Churn'] == 'Yes').mean()
10 print(f"Baseline: {baseline*100:.1f}%")
11 print(f"Diferença: {(churn_rate - baseline)*100:.1f} p.p.")
```

Conhecimento do Domínio: Essencial

Por que expertise importa:

Dados sozinhos não bastam:

- ▶ Padrões estatísticos \neq causação
- ▶ Correlação espúria é comum
- ▶ Contexto de negócio valida hipóteses

Como obter conhecimento do domínio:

- ▶ Conversar com especialistas
- ▶ Estudar literatura do setor
- ▶ Analisar casos históricos
- ▶ Testar hipóteses de negócio
- ▶ Iterar com stakeholders

Conhecimento do Domínio: Essencial

Por que expertise importa:

Melhores features vêm de:

- ▶ 50% conhecimento dos dados
- ▶ 50% conhecimento do negócio

💡 Nota Importante

Cientista de dados efetivo = Técnica + Domínio!

Validar Features: Correlação com Target

Quais features são mais úteis?

</> Python

```
1 # Converter Churn para binário
2 df['Churn_binary'] = (df['Churn'] == 'Yes').astype(int)
3
4 # Selecionar apenas features numéricas
5 numeric_features = df.select_dtypes(include=['int64', 'float64']).columns.
    tolist()
6
7 # Remover colunas que não são features
8 exclude = ['Churn_binary', 'customerID']
9 numeric_features = [f for f in numeric_features if f not in exclude]
```

Validar Features: Correlação com Target (cont.)

Quais features são mais úteis?

</> Python

```
1 # Calcular correlações com Churn
2 correlations = df[numeric_features + ['Churn_binary']].corr()['Churn_binary'].drop('Churn_binary')
3
4 # Ordenar por valor absoluto
5 correlations_abs = correlations.abs().sort_values(ascending=False)
6
7 print("=== TOP 15 FEATURES MAIS CORRELACIONADAS ===")
8 print(correlations_abs.head(15))
```

Visualizar Importância das Features

</> Python

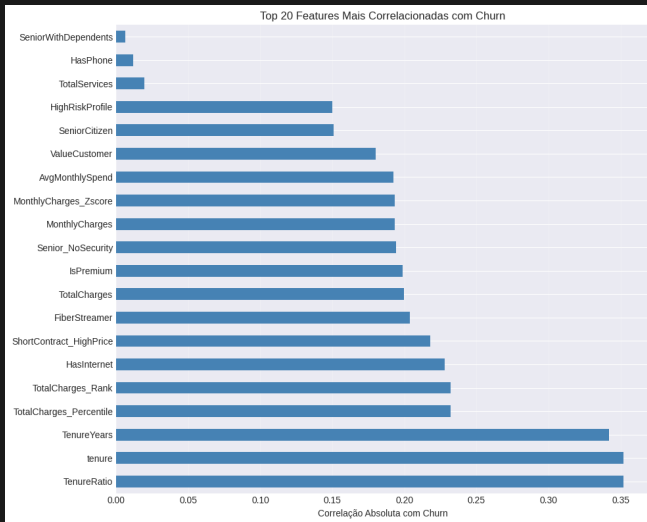
```
1 # Plotar correlações
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 fig, ax = plt.subplots(figsize=(10, 8))
6
7 # Top 20 features
8 top_features = correlations_abs.head(20)
9
10 # Barplot horizontal
11 top_features.plot(kind='barh', ax=ax, color='steelblue')
12 ax.set_xlabel('Correlação Absoluta com Churn')
13 ax.set_title('Top 20 Features Mais Correlacionadas com Churn')
14 ax.grid(True, alpha=0.3, axis='x')
```


Visualizar Importância das Features (cont.)

</> Python

```
1 plt.tight_layout()
2 plt.show()
3
4 # Insights
5 print("\n=== INSIGHTS ===")
6 print("- Features criadas aparecem no top?")
7 print("- São melhores que features originais?")
8 print("- Faz sentido do ponto de vista de negócio?")
```

Visualizar Importância das Features (cont.)



Boas Práticas de Feature Engineering

1. Documente TUDO

- ▶ Raciocínio de cada feature
- ▶ Como foi calculada
- ▶ Por que acredita que ajuda
- ▶ Código comentado

2. Mantenha simplicidade

- ▶ Comece com features óbvias
- ▶ Adicione complexidade gradualmente
- ▶ Teste impacto incremental
- ▶ Descarte features que não agregam

3. Validação contínua

- ▶ Verifique correlação com target
- ▶ Análise visual (distribuições, boxplots)
- ▶ Faz sentido para o negócio?
- ▶ Não adicione features "porque sim"

4. Atenção à interpretabilidade

- ▶ Features devem ser explicáveis
- ▶ Stakeholders precisam entender
- ▶ Muito complexo = difícil validar
- ▶ Preferir clareza a sofisticação

5. Cuidado com leakage

- ▶ Não use informação "do futuro"
- ▶ Features devem ser calculáveis em produção
- ▶ Ex: não use "cancelou_mes_seguinte" para prever churn atual

6. Itere e refine

- ▶ Feature engineering é iterativo
- ▶ Primeira versão raramente é perfeita
- ▶ Aprenda com análises
- ▶ Descarte o que não funciona

Resumo: Features Criadas

Total de features novas no dataset:

Python

```
1 # Listar todas as features criadas
2 created_features = [
3     'AvgMonthlySpend', 'ChargeDiff', 'TenureRatio',
4     'ChargesBin_Equal', 'ChargesBin_Quantile',
5     'TotalServices', 'HasInternet', 'HasPhone',
6     'SeniorWithDependents', 'IsPremium',
7     'ShortContract_HighPrice', 'Senior_NoSecurity', 'FiberStreamer',
8     'PaymentCategory', 'TenureYears', 'TenurePeriod',
9     'TotalCharges_Percentile', 'ChargesLevel',
10    'MonthlyCharges_Zscore', 'IsOutlier_Charges',
11    'HighRiskProfile', 'ValueCustomer'
12 ]
```

Resumo: Features Criadas (cont.)

Total de features novas no dataset:

</> Python

```
1 print(f"=== TOTAL DE FEATURES CRIADAS: {len(created_features)} ===")
2 print("\nColunas no dataset:")
3 print(f"Original: {len(df_original.columns)}")
4 print(f"Atual: {len(df.columns)}")
5 print(f"Features adicionadas: {len(df.columns) - len(df_original.columns)}")
6 # === TOTAL DE FEATURES CRIADAS: 22 ===
7 #
8 # Colunas no dataset:
9 # Original: 21
10 # Atual: 46
11 # Features adicionadas: 25
```

Bloco 2

Encoding de Variáveis Categóricas

Bloco 2: Encoding de Variáveis Categóricas

Conceito:

- ▶ Por que precisamos converter categorias em números
- ▶ Tipos de variáveis categóricas
- ▶ Estratégias diferentes para cada tipo

Técnicas:

- ▶ Label Encoding (variáveis ordinais)
- ▶ One-Hot Encoding (variáveis nominais)
- ▶ Binary Encoding (variáveis binárias)
- ▶ Estratégias para alta cardinalidade

Objetivo:

- ▶ Permitir análises quantitativas
- ▶ Calcular correlações
- ▶ Criar visualizações numéricas
- ▶ Preparar dados para comparações

Por que Encoding Importa?

Limitações de variáveis categóricas (texto):

O que NÃO podemos fazer com texto:

- ▶ ✗ Calcular correlações
- ▶ ✗ Fazer operações matemáticas
- ▶ ✗ Criar certos tipos de visualização
- ▶ ✗ Comparar quantitativamente
- ▶ ✗ Usar em análises estatísticas avançadas

O que PODEMOS fazer após encoding:

- ▶ ✓ Correlacionar com outras variáveis
- ▶ ✓ Análises estatísticas completas
- ▶ ✓ Heatmaps de correlação
- ▶ ✓ Comparações numéricas diretas
- ▶ ✓ Normalização e padronização

Por que Encoding Importa?

Nota Importante

Encoding = Traduzir categorias para linguagem numérica

Tipos de Variáveis Categóricas

Importante: Tipo determina estratégia de encoding!

1. Nominais (sem ordem natural):

- ▶ Exemplos: cor, cidade, método de pagamento
- ▶ Valores são apenas "diferentes", não "maiores" ou "menores"
- ▶ Estratégia: One-Hot Encoding

2. Ordinais (com ordem natural):

- ▶ Exemplos: educação (fundamental < médio < superior)
- ▶ Valores têm hierarquia clara
- ▶ Estratégia: Label Encoding (preservar ordem)

3. Binárias (apenas 2 valores):

- ▶ Exemplos: Yes/No, True/False, Male/Female
- ▶ Caso especial (mais simples)
- ▶ Estratégia: 0/1 direto

Identificar Tipo no Dataset Telco Churn

Classificando nossas variáveis:

Binárias (Yes/No):

- ▶ gender, Partner, Dependents, PhoneService
- ▶ PaperlessBilling, SeniorCitizen

Ordinais (têm ordem):

- ▶ **Contract:** Month-to-month < One year < Two year
- ▶ (Ordem de compromisso/duração)

Nominais (sem ordem):

- ▶ **PaymentMethod:** Electronic check, Mailed check, Bank transfer, Credit card
- ▶ **InternetService:** DSL, Fiber optic, No
- ▶ **MultipleLines, OnlineSecurity, etc:** Yes, No, No service

Label Encoding: Para Ordinais

Quando usar:

- ▶ Variável tem ordem natural clara
- ▶ Queremos preservar essa ordem
- ▶ Exemplos: níveis de educação, tamanho (S, M, L), satisfação (baixa, média, alta)

Como funciona:

- ▶ Atribuir números inteiros seguindo a ordem
- ▶ 0, 1, 2, 3, ...
- ▶ Ordem numérica = ordem categórica

Label Encoding: Para Ordinais

Exemplo - Contract:

Original	Encoded
Month-to-month	0
One year	1
Two year	2

Interpretação: $2 > 1 > 0$ reflete compromisso crescente

Label Encoding: Implementação

Método 1: Mapeamento manual (recomendado para ordinais)

</> Python

```
1 # Criar dicionário de mapeamento
2 contract_mapping = {
3     'Month-to-month': 0,
4     'One year': 1,
5     'Two year': 2
6 }
7
8 # Aplicar mapeamento
9 df['Contract_encoded'] = df['Contract'].map(contract_mapping)
```


Label Encoding: Implementação (cont.)

Método 1: Mapeamento manual (recomendado para ordinais)

</> Python

```
1 # Verificar resultado
2 print("=== LABEL ENCODING: Contract ===")
3 print("\nValores originais e encodados:")
4 print(df[['Contract', 'Contract_encoded']].drop_duplicates().sort_values('
    Contract_encoded'))
5 print("\nDistribuição:")
6 print(df['Contract_encoded'].value_counts().sort_index())
```

Label Encoding: Implementação (cont.)

Método 1: Mapeamento manual (recomendado para ordinais)

</> Python

```
1 # === LABEL ENCODING: Contract ===
2 #
3 # Valores originais e encodados:
4 #      Contract      Contract_encoded
5 # 0      Month-to-month      0
6 # 1          One year        1
7 # 11         Two year        2
8 #
9 # Distribuição:
10 # Contract_encoded
11 # 0      3875
12 # 1      1473
13 # 2      1695
```

Label Encoding - `.cat.codes`

Encoding Stateless Nativo do Pandas:

- ▶ **Stateless:** Não mantém objeto com mapeamento persistente
- ▶ **Operação ad-hoc:** Executa encoding no DataFrame atual
- ▶ **Ordem automática:** Por padrão usa ordem alfabética/lexicográfica
- ▶ **Type-safe:** Cria um dtype Categorical eficiente

💡 Nota Importante

Ponto-chave: `.cat.codes` é uma *propriedade* de objetos Categorical do pandas, não um transformer reutilizável.

Label Encoding - `.cat.codes`: Comportamento

Vantagens:

- ▶ Performance otimizada
- ▶ Memória eficiente
- ▶ Integração natural com pandas
- ▶ Suporta ordem customizada
- ▶ Dtype explícito (ordered)

Desvantagens:

- ▶ Não persiste mapeamento
- ▶ Inconsistência sem cuidado
- ▶ Requer categories explícitas
- ▶ Não integra com sklearn
- ▶ Difícil de serializar

⚠ Atenção

Armadilha comum: Usar `astype('category')` sem definir categories pode gerar encodings diferentes em train vs test!

Label Encoding - .cat.codes: Implementação

Método 2: Mapeamento automático (usando cat.codes)

</> Python

```
1 # Definir ordem explícita
2 df['Contract'] = pd.Categorical(
3     df['Contract'],
4     categories=['Month-to-month', 'One year', 'Two year'],
5     ordered=True
6 )
7 df['Contract_encoded'] = df['Contract'].cat.codes
8 # 0, 1, 2 na ordem definida
```

Label Encoding - Sklearn LabelEncoder

Encoding Stateful e Reutilizável:

- ▶ **Stateful:** Mantém objeto com mapeamento interno (classes_)
- ▶ **Fit/Transform pattern:** Aprende no treino, aplica no teste
- ▶ **Ordem alfabética:** Sempre ordena classes alfabeticamente
- ▶ **Serializável:** Pode ser salvo e carregado (joblib)

💡 Nota Importante

Ponto-chave: LabelEncoder é um *transformer sklearn* que armazena estado e pode ser reutilizado em novos dados.

Label Encoding - Sklearn LabelEncoder: Comportamento

Vantagens:

- ▶ Consistência garantida
- ▶ Fit/transform separation
- ▶ Erro explícito em unknowns
- ▶ Fácil serialização
- ▶ Inverse transform nativo
- ▶ Ideal para target encoding

Limitações:

- ▶ Apenas 1D arrays
- ▶ Ordem alfabética fixa
- ▶ Não integra diretamente em ColumnTransformer
- ▶ Menos controle sobre ordem
- ▶ Sem ordered categorical

💡 Nota Importante

Para features em pipelines, usar OrdinalEncoder (2D). LabelEncoder é ideal para target (y).

Label Encoding - Sklearn LabelEncoder: Implementação

Método 3: LabelEncoder do sklearn

</> Python

```
1 from sklearn.preprocessing import LabelEncoder
2
3 # Criar encoder
4 le = LabelEncoder()
5
6 # Fit e transform
7 df['Contract_encoded_sklearn'] = le.fit_transform(df['Contract'])
8
9 # Ver mapeamento criado
10 print("=== SKLEARN LABEL ENCODER ===")
11 print("Classes encontradas:", le.classes_)
12 print("Mapeamento:")
13 for i, label in enumerate(le.classes_):
14     print(f" {label} -> {i}")
```


Label Encoding - Sklearn LabelEncoder (cont.)

Método 3: LabelEncoder do sklearn

</> Python

```
1 # === SKLEARN LABEL ENCODER ===  
2 # Classes encontradas: ['Month-to-month' 'One year' 'Two year']  
3 # Mapeamento:  
4 #   Month-to-month -> 0  
5 #   One year -> 1  
6 #   Two year -> 2  
7  
8 # ATENÇÃO: sklearn mapeia alfabeticamente!  
9 # Pode não preservar ordem desejada para ordinais  
10 # Por isso, mapeamento manual é preferível para ordinais
```

Comparação: Estado e Persistência

cat.codes (Stateless)

- ▶ Não mantém objeto
- ▶ Mapeamento implícito
- ▶ Requer categories externas
- ▶ Serialização manual
- ▶ Recriação necessária

Implicação: Desenvolvedor é responsável por manter consistência entre treino e produção.

LabelEncoder (Stateful)

- ▶ Mantém objeto encoder
- ▶ Mapeamento em classes_
- ▶ Autocontido
- ▶ joblib.dump() direto
- ▶ Reutilização simples

Implicação: Encoder garante consistência automaticamente via transform().

Comparação: Controle de Ordem

cat.codes: Flexibilidade

Opções de ordenação:

- ▶ Alfabética
- ▶ Por frequência
- ▶ Semântica customizada
- ▶ Qualquer ordem desejada

Exemplo de uso:

- ▶ Ordenar por importância
- ▶ Ordenar por risco
- ▶ Ordem de negócio

Decisão: Se ordem customizada é crítica → cat.codes. Se consistência é crítica → LabelEncoder/OrdinalEncoder.

LabelEncoder: Consistência

Ordem fixa:

- ▶ Sempre alfabética
- ▶ Determinística
- ▶ Não customizável
- ▶ Previsível

Trade-off:

- ▶ Menos flexibilidade
- ▶ Mais consistência
- ▶ Menos erros humanos

Comparação: Performance e Memória

Aspectos de eficiência computacional:

cat.codes

- ▶ Dtype Categorical eficiente
- ▶ Economia de memória
- ▶ Operações otimizadas
- ▶ Cache-friendly
- ▶ Ideal para datasets grandes

Vantagem: 50-90% menos memória com muitas categorias repetidas

LabelEncoder

- ▶ Retorna array numpy padrão
- ▶ Memória convencional
- ▶ Performance adequada
- ▶ Overhead mínimo
- ▶ Compatibilidade total

Vantagem: Simplicidade e integração universal

Nota Importante

Na prática, diferença de performance é marginal.

Tabela Comparativa Resumida

Aspecto	cat.codes	LabelEncoder
Estado	Stateless	Stateful
Consistência train/test	Manual	Automática
Valores desconhecidos	-1 (silencioso)	ValueError
Ordem	Customizável	Alfabética fixa
Serialização	Manual	joblib.dump()
Pipeline sklearn	Não	Sim (1D)
Inverse transform	Manual	Nativo
Arrays 2D	Não	Não (use OrdinalEncoder)
Memória	Eficiente	Padrão
Dtype	Categorical	numpy array

💡 Nota Importante

Verde indica vantagem clara. Nenhuma ferramenta é superior em todos aspectos.

Validar Label Encoding

Verificar se encoding faz sentido:

</> Python

```
1 # Analisar relação entre encoded e Churn
2 print("=== VALIDACAO: Contract_encoded vs Churn ===")
3
4 # Média de churn por nível de contract
5 churn_by_contract = df.groupby('Contract_encoded')['Churn'].apply(
6     lambda x: (x == 'Yes').mean()
7 )
8 print("\nTaxa de Churn por Contract (encoded):")
9 print(churn_by_contract)
10
11 # Correlação com Churn_binary
12 from scipy.stats import pearsonr
13 corr, pval = pearsonr(df['Contract_encoded'], df['Churn_binary'])
14 print(f"\nCorrelação com Churn: {corr:.3f} (p={pval:.4f})")
```

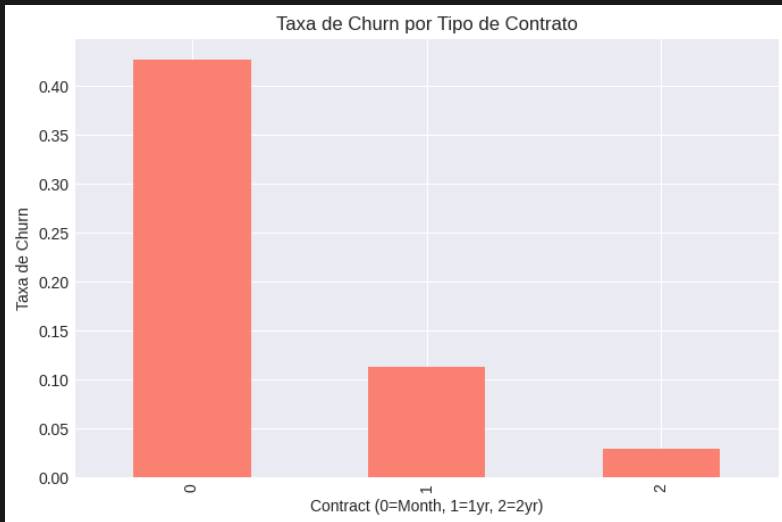
Validar Label Encoding (cont.)

Verificar se encoding faz sentido:

</> Python

```
1 # Visualizar
2 import matplotlib.pyplot as plt
3 churn_by_contract.plot(kind='bar', color='salmon', figsize=(8, 5))
4 plt.title('Taxa de Churn por Tipo de Contrato')
5 plt.xlabel('Contract (0=Month, 1=1yr, 2=2yr)')
6 plt.ylabel('Taxa de Churn')
7 plt.show()
```

Validar Label Encoding (cont.)



One-Hot Encoding: Para Nominais

Quando usar:

- ▶ Variável NÃO tem ordem natural
- ▶ Categorias são simplesmente "diferentes"
- ▶ Exemplos: cor, cidade, método de pagamento

Como funciona:

- ▶ Criar uma coluna binária para CADA categoria
- ▶ 1 se categoria está presente, 0 caso contrário
- ▶ Também chamado de "dummy variables"

One-Hot Encoding: Para Nominais

Exemplo - InternetService (3 categorias):

Original	InternetService_DSL	InternetService_Fiber	InternetService_No
DSL	1	0	0
Fiber optic	0	1	0
No	0	0	1

Resultado: 3 colunas binárias ao invés de 1 categórica

One-Hot: Por que NÃO usar Label Encoding?

Exemplo: PaymentMethod

- ▶ Categorias: Electronic check, Mailed check, Bank transfer, Credit card
- ▶ Se usarmos Label: 0, 1, 2, 3

O que está errado?

- ▶ Números implicam ordem: $3 > 2 > 1 > 0$
- ▶ Mas não há ordem real entre métodos de pagamento!
- ▶ Credit card (3) não é "maior" que Electronic check (0)
- ▶ Análises podem interpretar erroneamente como ordinal
- ▶ Correlações ficam sem sentido

Solução: One-Hot

- ▶ Cada método vira coluna binária independente
- ▶ Sem implicação de ordem
- ▶ Cada coluna pode ser correlacionada separadamente

One-Hot Encoding: Implementação

Método: `pd.get_dummies()`

</> Python

```
1 # One-Hot Encoding de InternetService
2 print("=== ONE-HOT ENCODING: InternetService ===")
3
4 # Antes
5 print("Coluna original:")
6 print(df['InternetService'].value_counts())
7
8 # Aplicar One-Hot
9 internet_dummies = pd.get_dummies(df['InternetService'],
10                                   prefix='InternetService',
11                                   dtype=int)
```

One-Hot Encoding: Implementação

Método: `pd.get_dummies()`

</> Python

```
1 # Visualizar
2 print("\nColunas criadas:")
3 print(internet_dummies.columns.tolist())
4
5 print("\nExemplo de dados:")
6 print(internet_dummies.head())
7
8 # Adicionar ao dataframe
9 df = pd.concat([df, internet_dummies], axis=1)
10
11 print(f"\nShape antes: {df.shape[0]} x {df.shape[1] - len(internet_dummies
12     .columns)}")
12 print(f"Shape depois: {df.shape}")
```

One-Hot: Múltiplas Colunas de Uma Vez

</> Python

```
1 # Aplicar One-Hot a múltiplas colunas nominais
2 nominal_cols = ['InternetService', 'PaymentMethod']
3
4 # Método 1: Uma por vez (mais controle)
5 for col in nominal_cols:
6     dummies = pd.get_dummies(df[col], prefix=col, dtype=int)
7     df = pd.concat([df, dummies], axis=1)
8
9 # Método 2: Todas de uma vez (mais simples)
10 df_encoded = pd.get_dummies(df,
11                               columns=nominal_cols,
12                               prefix=nominal_cols,
13                               dtype=int,
14                               drop_first=False) # Veremos a seguir
```

One-Hot: Múltiplas Colunas de Uma Vez (cont.)

</> Python

```
1 print("=== COLUNAS CRIADAS ===")
2 new_cols = [col for col in df_encoded.columns if any(nom in col for nom in
    nominal_cols)]
3 print(f"Total: {len(new_cols)} novas colunas")
4 print(new_cols)
5
6 # === COLUNAS CRIADAS ===
7 # Total: 17 novas colunas
8 # ['InternetService_DSL', 'InternetService_Fiber optic',
9 #  'InternetService_No', 'InternetService_DSL',
10 #  'InternetService_Fiber optic', 'InternetService_No', ...]
```

One-Hot com `drop_first`

Problema:

- ▶ Se temos `InternetService` com 3 categorias
- ▶ One-Hot cria 3 colunas: `DSL`, `Fiber`, `No`
- ▶ Mas: se `DSL=0` e `Fiber=0`, então necessariamente `No=1`
- ▶ Informação redundante! (colinearidade perfeita)

Solução: `drop_first=True`

- ▶ Remove primeira categoria
- ▶ Serve como referência (baseline)
- ▶ N categorias $\rightarrow N-1$ colunas

Exemplo:

- ▶ `InternetService`: 3 categorias \rightarrow 2 colunas
- ▶ Se `Fiber=0` e `No=0`, então é `DSL` (implícito)
- ▶ `DSL` é a categoria de referência

One-Hot com drop_first: Implementação

</> Python

```
1 # Comparar com e sem drop_first
2 print("=== COMPARACAO: drop_first ===\n")
3
4 # Sem drop_first (3 colunas)
5 dummies_full = pd.get_dummies(df['InternetService'],
6                               prefix='Internet',
7                               dtype=int,
8                               drop_first=False)
9 print(f"Sem drop_first: {dummies_full.shape[1]} colunas")
10 print(dummies_full.columns.tolist())
11
12 # ['Internet_DSL', 'Internet_Fiber optic', 'Internet_No']
```

Quando usar drop_first?

One-Hot com drop_first: Implementação (cont.)

</> Python

```
1 # Com drop_first (2 colunas)
2 dummies_dropped = pd.get_dummies(df['InternetService'],
3                                   prefix='Internet',
4                                   dtype=int,
5                                   drop_first=True)
6 print(f"\nCom drop_first: {dummies_dropped.shape[1]} colunas")
7 print(dummies_dropped.columns.tolist())
8 # ['Internet_Fiber optic', 'Internet_No']
9 print("\nCategoria de referência (implícita): DSL")
```

One-Hot com drop_first: Implementação (cont.)

Quando usar drop_first?

- ▶ Para análise exploratória: tanto faz
- ▶ Para ML (futuro): geralmente sim (evita multicolinearidade)

One-Hot: Vantagens e Desvantagens

Vantagens:

- ▶ ✓ Não impõe ordem artificial
- ▶ ✓ Cada categoria é independente
- ▶ ✓ Interpretação clara (1 = tem, 0 = não tem)
- ▶ ✓ Funciona para qualquer tipo de categoria
- ▶ ✓ Permite correlações separadas

Desvantagens:

- ▶ ✗ Aumenta dimensionalidade (curse of dimensionality)
- ▶ ✗ N categorias \rightarrow N (ou N-1) colunas
- ▶ ✗ Problemático para alta cardinalidade (muitas categorias)
- ▶ ✗ Dataset fica "esparso" (muitos zeros)

One-Hot: Vantagens e Desvantagens

Atenção

Para variável com 100 categorias: One-Hot cria 100 colunas!
Veremos alternativas para alta cardinalidade.

Label vs One-Hot: Quando usar cada um

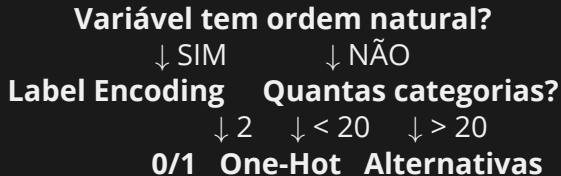
Aspecto	Label Encoding	One-Hot Encoding
Tipo de variável	Ordinal	Nominal
Ordem importa?	Sim	Não
Colunas criadas	1	N (ou N-1)
Espaço usado	Eficiente	Pode explodir
Interpretação	Ordem numérica	Presença/ausência
Exemplo	Contract, Educação, Tamanho	PaymentMethod, Cidade, Cor

💡 Nota Importante

Regra de ouro: Se tem ordem natural clara → Label
Se não tem ordem → One-Hot

Fluxograma: Escolher Encoding

Árvore de decisão:



Alternativas para alta cardinalidade (20):

- ▶ Frequency Encoding
- ▶ Target Encoding (cuidado com leakage!)
- ▶ Agrupar categorias raras
- ▶ Embeddings (deep learning - futuro)

Caso Prático: Pipeline de Encoding Telco

Aplicar encoding sistemático ao dataset:

</> Python

```
1 # Criar cópia para trabalhar
2 df_encoded = df.copy()
3
4 print("=== PIPELINE DE ENCODING ===\n")
5 # 1. ORDINAIS: Label Encoding
6 print("1. Encoding ORDINAIS (Label)...")
7 ordinal_mappings = {
8     'Contract': {
9         'Month-to-month': 0,
10        'One year': 1,
11        'Two year': 2
12    }
13 }
```


Caso Prático: Pipeline de Encoding Telco

Aplicar encoding sistemático ao dataset:

</> Python

```
1 for col, mapping in ordinal_mappings.items():  
2     df_encoded[f'{col}_encoded'] = df_encoded[col].map(mapping)  
3     print(f"        {col} encoded")
```

Caso Prático: Pipeline de Encoding (cont.)

</> Python

```
1 # 2. BINÁRIAS: Conversão 0/1
2 print("\n2. Encoding BINÁRIAS (0/1)...")
3 binary_cols = ['gender', 'Partner', 'Dependents',
4               'PhoneService', 'PaperlessBilling']
5
6 for col in binary_cols:
7     # Assumindo que 'Yes'/'Male' = 1, 'No'/'Female' = 0
8     df_encoded[f'{col}_encoded'] = df_encoded[col].map({
9         'Yes': 1, 'No': 0,
10        'Male': 1, 'Female': 0
11    }).fillna(df_encoded[col].map({
12        1: 1, 0: 0 # Se já for numérico (SeniorCitizen)
13    }))
14 print(f"        {col} encoded")
```

Caso Prático: Pipeline de Encoding (cont.)

</> Python

```
1 # 3. NOMINAIS: One-Hot
2 print("\n3. Encoding NOMINAIS (One-Hot)...")
3 nominal_cols = ['PaymentMethod', 'InternetService']
4
5 df_encoded = pd.get_dummies(df_encoded,
6                             columns=nominal_cols,
7                             prefix=nominal_cols,
8                             drop_first=True,
9                             dtype=int)
10 print(f"      {len(nominal_cols)} columnas expandidas")
```

Variáveis Binárias: Caso Especial

Tipos de binárias:

- ▶ Yes/No
- ▶ True/False
- ▶ Male/Female
- ▶ 1/0

Como encodar:

- ▶ Opção 1: 0 e 1 (mais comum e recomendado)
- ▶ Opção 2: -1 e 1 (às vezes usado)
- ▶ Opção 3: One-Hot (desnecessário e redundante!)

Recomendação:

- ▶ Use 0 para "negativo/ausente/feminino"
- ▶ Use 1 para "positivo/presente/masculino"
- ▶ Não precisa One-Hot (seria criar 2 colunas redundantes)
- ▶ Mais eficiente: 1 coluna ao invés de 2

Variáveis Binárias: Caso Especial

Nota Importante

Binárias: encoding mais simples possível (0/1 direto)

Encoding de Binárias: Implementação

</> Python

```
1 # Método 1: map() - mais explícito
2 df['gender_encoded'] = df['gender'].map({'Male': 1, 'Female': 0})
3
4 # Método 2: operador booleano - mais conciso
5 df['Partner_encoded'] = (df['Partner'] == 'Yes').astype(int)
6
7 # Método 3: replace() - similar ao map
8 df['Dependents_encoded'] = df['Dependents'].replace({'Yes': 1, 'No': 0})
```

Encoding de Binárias: Implementação

</> Python

```
1 # Verificar
2 print("=== ENCODING DE BINÁRIAS ===")
3 for col in ['gender', 'Partner', 'Dependents']:
4     encoded_col = f'{col}_encoded'
5     print(f"\n{col}:")
6     print(df[[col, encoded_col]].drop_duplicates())
7
8 # Validar: não deve haver valores além de 0 e 1
9 print("\nValores únicos (devem ser apenas 0 e 1):")
10 print(df[[c for c in df.columns if '_encoded' in c]].nunique())
```

Alta Cardinalidade: O Problema

O que é alta cardinalidade?

- ▶ Variável categórica com MUITAS categorias únicas
- ▶ Exemplos: CEP (milhares), UserID, ProductID

Por que é problema para One-Hot?

- ▶ Explosão dimensional
- ▶ 1000 categorias → 1000 colunas!
- ▶ Dataset fica enorme e esperso
- ▶ 99% dos valores são zeros
- ▶ Dificulta análise e visualização
- ▶ Consome muita memória

Alta Cardinalidade: O Problema

Threshold comum:

- ▶ < 10 categorias: One-Hot sem problemas
- ▶ 10-20 categorias: One-Hot ainda OK
- ▶ > 20 categorias: Considerar alternativas
- ▶ > 50 categorias: Definitivamente usar alternativas

Soluções para Alta Cardinalidade

Alternativas ao One-Hot:

1. Agrupar categorias raras

- ▶ Categorias com frequência $< 5\%$ \rightarrow "Other"
- ▶ Reduz número de categorias
- ▶ Mantém informação das principais

2. Frequency Encoding

- ▶ Substituir categoria por sua frequência
- ▶ Ex: "SP" aparece 30% das vezes $\rightarrow 0.30$
- ▶ Uma coluna numérica ao invés de N binárias

3. Target Encoding (avançado)

- ▶ Substituir categoria pela média do target naquela categoria
- ▶ Ex: Churn médio em "SP" = 0.25 $\rightarrow 0.25$
- ▶ **Cuidado:** risco de data leakage!
- ▶ Mais usado em ML

Solução 1: Agrupar Categorias Raras

</> Python

```
1 # Exemplo: Agrupar métodos de pagamento raros
2 def group_rare_categories(series, threshold=0.05, other_label='Other'):
3     """
4     Agrupa categorias com frequência < threshold em 'Other'
5
6     Parameters:
7     -----
8     series : pd.Series
9         Coluna categórica
10    threshold : float
11        Frequência mínima (default: 5%)
12    other_label : str
13        Label para categorias raras
14    """
```

Solução 1: Agrupar Categorias Raras (cont.)

</> Python

```
1  # Calcular frequências
2  freq = series.value_counts(normalize=True)
3
4  # Identificar categorias raras
5  rare_categories = freq[freq < threshold].index.tolist()
6
7  # Substituir raras por 'Other'
8  return series.replace(rare_categories, other_label)
```

Solução 1: Agrupar Categorias Raras (cont.)

</> Python

```
1 # Aplicar agrupamento
2 print("=== AGRUPAMENTO DE CATEGORIAS RARAS ===")
3
4 # Antes
5 print("\nANTES - PaymentMethod:")
6 print(df['PaymentMethod'].value_counts(normalize=True))
7
8 # Agrupar
9 df['PaymentMethod_grouped'] = group_rare_categories(
10     df['PaymentMethod'],
11     threshold=0.10 # Agrupar se < 10%
12 )
```

Solução 1: Agrupar Categorias Raras (cont.)

</> Python

```
1 # Depois
2 print("\nDEPOIS - PaymentMethod_grouped:")
3 print(df['PaymentMethod_grouped'].value_counts(normalize=True))
4
5 # Agora podemos aplicar One-Hot com menos colunas
6 dummies = pd.get_dummies(df['PaymentMethod_grouped'],
7                           prefix='Payment',
8                           drop_first=True)
9 print(f"\nColunas criadas: {len(dummies.columns)}")
```

Solução 2: Frequency Encoding

</> Python

```
1 # Frequency Encoding: substituir por frequência
2 def frequency_encoding(series):
3     """
4     Codifica categorias por sua frequência relativa
5
6     Returns:
7     -----
8     pd.Series com frequências
9     """
10    freq = series.value_counts(normalize=True)
11    return series.map(freq)
```

Vantagem: 1 coluna numérica ao invés de N binárias

Solução 2: Frequency Encoding

</> Python

```
1 # Aplicar
2 print("=== FREQUENCY ENCODING ===")
3
4 df['PaymentMethod_freq'] = frequency_encoding(df['PaymentMethod'])
5
6 print("\nMapeamento criado:")
7 print(df[['PaymentMethod', 'PaymentMethod_freq']].drop_duplicates().
8         sort_values('PaymentMethod_freq', ascending=False))
9
10 print("\nEstatísticas da feature encodada:")
11 print(df['PaymentMethod_freq'].describe())
```

Vantagem: 1 coluna numérica ao invés de N binárias

Comparar Estratégias de Encoding

Análise comparativa:

</> Python

```
1 # Comparar diferentes encodings de PaymentMethod
2 print("=== COMPARACAO DE ENCODINGS ===\n")
3
4 # 1. One-Hot (baseline)
5 onehot = pd.get_dummies(df['PaymentMethod'], prefix='PM')
6 print(f"One-Hot: {onehot.shape[1]} colunas")
7
8 # 2. One-Hot com agrupamento
9 grouped_onehot = pd.get_dummies(df['PaymentMethod_grouped'], prefix='PM')
10 print(f"One-Hot + Grouped: {grouped_onehot.shape[1]} colunas")
11
12 # 3. Frequency Encoding
13 print(f"Frequency: 1 coluna")
```

Comparar Estratégias de Encoding (cont.)

Análise comparativa:

</> Python

```
1 # 4. Label (se formos - não recomendado para nominal!)
2 label_enc = pd.factorize(df['PaymentMethod'])[0]
3 print(f"Label (NÃO RECOMENDADO para nominal): 1 coluna")
4
5 print("\n=== RECOMENDACAO ===")
6 print("Para PaymentMethod (nominal, 4 categorias):")
7 print(" - One-Hot é apropriado (apenas 4 colunas)")
```

Comparar Estratégias de Encoding (cont.)

Análise comparativa:

</> Python

```
1 # === COMPARACAO DE ENCODINGS ===
2 #
3 # One-Hot: 4 colunas
4 # One-Hot + Grouped: 4 colunas
5 # Frequency: 1 coluna
6 # Label (NÃO RECOMENDADO para nominal): 1 coluna
7 #
8 # === RECOMENDACAO ===
9 # Para PaymentMethod (nominal, 4 categorias):
10 # - One-Hot é apropriado (apenas 4 colunas)
11
```

Cuidados com Encoding

Armadilhas comuns:

1. Usar Label para nominais

- ▶ ✗ Impõe ordem artificial
- ▶ ✗ Correlações sem sentido
- ▶ Exemplo: PaymentMethod como 0,1,2,3 sugere ordem que não existe

2. Esquecer de encodar todas as categóricas

- ▶ Deixar alguma coluna object atrapalha análises
- ▶ Verificar: `df.select_dtypes('object')`

Armadilhas comuns:

3. Perder track das colunas originais

- ▶ Documentar qual encoding foi usado onde
- ▶ Manter mapeamento para interpretação posterior
- ▶ Comentar código claramente

4. Não validar resultado

- ▶ Sempre verificar distribuições
- ▶ Checar valores únicos
- ▶ Garantir que faz sentido para o negócio

Verificar Resultado do Encoding

Checklist de validação:

</> Python

```
1 # Após encoding completo
2 print("=== VALIDACAO DO ENCODING ===\n")
3
4 # 1. Verificar se ainda há colunas object
5 object_cols = df_encoded.select_dtypes(include='object').columns.tolist()
6 # Remover colunas que devem ser object (IDs, etc)
7 object_cols = [c for c in object_cols if c not in ['customerID', 'Churn']]
8
9 print(f"1. Colunas object restantes: {len(object_cols)}")
10 if len(object_cols) > 0:
11     print(f"    ATENÇÃO: {object_cols}")
12 else:
13     print("    Todas as features são numéricas")
14 # 1. Colunas object restantes: 15 ...
```

Verificar Resultado do Encoding (cont.)

Checklist de validação:

Python

```
1 # 2. Verificar shape
2 print(f"\n2. Shape:")
3 print(f"    Original: {df.shape}")
4 print(f"    Encoded: {df_encoded.shape}")
5 print(f"    Colunas adicionadas: {df_encoded.shape[1] - df.shape[1]}")
```

Verificar Resultado do Encoding (cont.)

</> Python

```
1 # 3. Verificar valores únicos das encoded
2 print("\n3. Valores únicos por coluna encoded:")
3 encoded_cols = [c for c in df_encoded.columns if '_encoded' in c or
4                 any(nom in c for nom in ['PaymentMethod_', '
5                 InternetService_'])]
6
7 for col in encoded_cols[:10]: # Primeiras 10
8     n_unique = df_encoded[col].nunique() #ATENÇÃO
9     unique_vals = sorted(df_encoded[col].unique())
10    print(f"    {col}: {n_unique} valores    $\rightarrow$    {unique_vals}")
```


Verificar Resultado do Encoding (cont.)

</> Python

```
1 # 4. Correlação entre features encodadas e target
2 print("\n4. Correlação com Churn:")
3 correlations = df_encoded[encoded_cols].corrwith(df_encoded['Churn_binary'
4 ])
5 top_corr = correlations.abs().sort_values(ascending=False).head(10)
6 print(top_corr)
7 print("\n Validação completa!")
```

Função Reutilizável: Pipeline de Encoding

</> Python

```
1 def encode_telco_churn(df):
2     """
3     Pipeline completo de encoding para dataset Telco Churn
4
5     Parameters:
6     -----
7     df : DataFrame
8         Dataset original
9
10    Returns:
11    -----
12    DataFrame encodado
13    """
14    df_enc = df.copy()
```

Função Reutilizável: Pipeline (cont.)

</> Python

```
1  # 1. Ordinais (Label)
2  ordinal_maps = {
3      'Contract': {'Month-to-month': 0, 'One year': 1, 'Two year': 2}
4  }
5  for col, mapping in ordinal_maps.items():
6      df_enc[f'{col}_encoded'] = df_enc[col].map(mapping)
```

Função Reutilizável: Pipeline (cont.)

</> Python

```
1  # 2. Binárias (0/1)
2  binary_cols = ['gender', 'Partner', 'Dependents',
3                'PhoneService', 'PaperlessBilling']
4  for col in binary_cols:
5      df_enc[f'{col}_encoded'] = (df_enc[col] == 'Yes').astype(int) \
6          if df_enc[col].dtype == 'object' else df_enc[col]
7
8  # SeniorCitizen já é 0/1
9  df_enc['SeniorCitizen_encoded'] = df_enc['SeniorCitizen']
```

Função Reutilizável: Pipeline (cont.)

</> Python

```
1  # 3. Nominais (One-Hot)
2  nominal_cols = ['PaymentMethod', 'InternetService']
3  df_enc = pd.get_dummies(df_enc, columns=nominal_cols,
4                          prefix=nominal_cols, drop_first=True, dtype=
5  int)
6
7  # 4. Target binário
8  df_enc['Churn_binary'] = (df_enc['Churn'] == 'Yes').astype(int)
9
10 return df_enc
```

Usar Pipeline de Encoding

</> Python

```
1 # Aplicar pipeline completo
2 print("=== APLICANDO PIPELINE DE ENCODING ===\n")
3
4 # Encoding
5 df_final = encode_telco_churn(df)
6
7 print(f"Shape original: {df.shape}")
8 print(f"Shape final: {df_final.shape}")
9
10 # Listar features numéricas
11 numeric_features = df_final.select_dtypes(include=['int64', 'float64']).
    columns
12 print(f"\nTotal de features numéricas: {len(numeric_features)}")
```

Usar Pipeline de Encoding (cont.)

</> Python

```
1 # Verificar se há categorical restante (exceto IDs e target original)
2 cat_cols = df_final.select_dtypes('object').columns
3 cat_cols = [c for c in cat_cols if c not in ['customerID', 'Churn']]
4 print(f"Colunas categóricas restantes: {len(cat_cols)}")
5
6 if len(cat_cols) == 0:
7     print("\n Encoding completo! Todos os dados são numéricos.")
8 else:
9     print(f"\n Atenção: ainda há colunas categóricas: {cat_cols}")
```

Análise Após Encoding

Exploração com dados encodados:

</> Python

```
1 # Matriz de correlação
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Selecionar features relevantes
6 features_to_analyze = [
7     'tenure', 'MonthlyCharges', 'TotalCharges',
8     'Contract_encoded', 'TotalServices',
9     'InternetService_Fiber optic', 'InternetService_No',
10    'PaymentMethod_Electronic check',
11    'Churn_binary'
12 ]
```

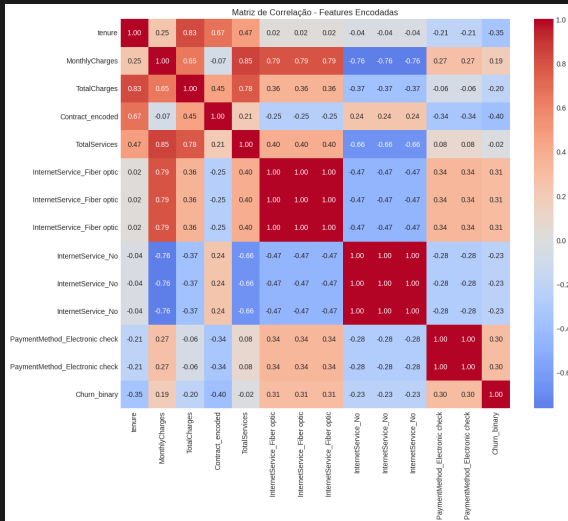

Análise Após Encoding (cont.)

Exploração com dados encodados:

</> Python

```
1 # Heatmap
2 fig, ax = plt.subplots(figsize=(12, 10))
3 corr_matrix = df_final[features_to_analyze].corr()
4 sns.heatmap(corr_matrix, annot=True, fmt='.2f', cmap='coolwarm',
5             center=0, square=True, ax=ax)
6 ax.set_title('Matriz de Correlação - Features Encodadas')
7 plt.tight_layout()
8 plt.show()
```

Análise Após Encoding (cont.)



Resumo do Bloco 2

O que aprendemos:

Conceitos:

- ▶ Por que encoding é necessário
- ▶ Tipos de variáveis categóricas
- ▶ Escolher estratégia apropriada

Técnicas:

- ▶ Label Encoding (ordinais)
- ▶ One-Hot Encoding (nominais)
- ▶ Binary Encoding (binárias)
- ▶ Soluções para alta cardinalidade

O que aprendemos:

Habilidades:

- ▶ Identificar tipo de variável
- ▶ Aplicar encoding correto
- ▶ Criar pipeline reproduzível
- ▶ Validar resultados
- ▶ Analisar dados encodados

Bloco 3

Normalização e Padronização

Introdução

Problema:

- ▶ Variáveis numéricas em escalas muito diferentes
- ▶ Dificulta comparações e visualizações
- ▶ Como tornar variáveis comparáveis?

Técnicas:

- ▶ Min-Max Scaling (normalização $[0,1]$)
- ▶ Standardization (Z-score)
- ▶ Robust Scaling (resistente a outliers)
- ▶ Transformações matemáticas (log, sqrt)

Objetivo:

- ▶ Comparações justas entre variáveis
- ▶ Visualizações mais claras
- ▶ Identificar importância relativa
- ▶ Facilitar interpretação

O Problema das Escalas Diferentes

Exemplo concreto - Dataset Telco Churn:

Variável	Mínimo	Máximo
tenure (meses)	0	72
MonthlyCharges (R\$)	18	118
TotalCharges (R\$)	18	8684

Problemas:

- ▶ TotalCharges varia em 8000+ unidades
- ▶ tenure varia apenas em 72 unidades
- ▶ Em gráficos: escala de TotalCharges domina visualmente
- ▶ Difícil comparar: "tenure=36 é mais importante que TotalCharges=5000?"
- ▶ **Escala diferente \neq Importância diferente**

O Problema das Escalas Diferentes

Exemplo concreto - Dataset Telco Churn:

Solução:

- ▶ Colocar todas em escala comum
- ▶ Preservar relações dentro de cada variável
- ▶ Permitir comparações justas

Benefícios da Normalização

Para análise exploratória:

1. Visualizações mais claras

- ▶ Múltiplas variáveis no mesmo gráfico
- ▶ Escalas compatíveis
- ▶ Comparações visuais diretas

2. Comparabilidade

- ▶ Identificar qual variável varia mais
- ▶ Comparar desvios padrão
- ▶ Entender distribuições relativamente

Benefícios da Normalização

Para análise exploratória:

3. Interpretação de correlações

- ▶ Correlações mais significativas
- ▶ Coeficientes comparáveis
- ▶ Heatmaps mais informativos

4. Preparação para análises avançadas

- ▶ Clustering (agrupamento)
- ▶ PCA (análise de componentes principais)
- ▶ E futuramente: Machine Learning

Três Técnicas Principais

1. Min-Max Scaling (Normalização)

- ▶ Escala para $[0, 1]$
- ▶
$$X_{scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$
- ▶ **Use quando:** Quer todos os valores entre 0 e 1
- ▶ **Sensível a:** Outliers

2. Standardization (Z-score, Padronização)

- ▶ Média = 0, Desvio Padrão = 1
- ▶ $X_{scaled} = \frac{X - \mu}{\sigma}$
- ▶ **Use quando:** Quer centralizar em zero
- ▶ **Sensível a:** Outliers moderados

3. Robust Scaling

- ▶ Usa mediana e IQR
- ▶ $X_{scaled} = \frac{X - median}{IQR}$
- ▶ **Use quando:** Tem MUITOS outliers
- ▶ **Resistente a:** Valores extremos

Min-Max Scaling: Normalização para [0, 1]

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Exemplo manual:

- ▶ Dados: [10, 20, 30, 40, 50]
- ▶ Min = 10, Max = 50, Range = 40
- ▶ $10 \rightarrow (10-10)/40 = 0.0$
- ▶ $30 \rightarrow (30-10)/40 = 0.5$
- ▶ $50 \rightarrow (50-10)/40 = 1.0$
- ▶ Resultado: [0.0, 0.25, 0.5, 0.75, 1.0]

Propriedades:

- ▶ Menor valor $\rightarrow 0$
- ▶ Maior valor $\rightarrow 1$
- ▶ Preserva forma da distribuição
- ▶ Interpretação: "posição relativa no range"

Min-Max: Implementação com Sklearn

</> Python

```
1 from sklearn.preprocessing import MinMaxScaler
2 import numpy as np
3
4 # Selecionar colunas numéricas para normalizar
5 numeric_cols = ['tenure', 'MonthlyCharges', 'TotalCharges']
6
7 # Criar scaler
8 scaler = MinMaxScaler()
9 scaled_array = scaler.fit_transform(df[numeric_cols])
10
11 # Gerar novas colunas para features normalizadas
12 normalized_cols = [col + '_minmax' for col in numeric_cols]
13 for idx, col in enumerate(normalized_cols):
14     df[col] = scaled_array[:, idx]
```

Min-Max: Implementação com Sklearn (cont.)

</> Python

```
1 # Aplicar fit and transform
2 df[normalized_cols] = scaler.fit_transform(df[numeric_cols])
3
4 # Alternativa: atribuição multicolumna
5 # scaled_df = pd.DataFrame(scaler.fit_transform(df[numeric_cols]), columns
6 #                           =normalized_cols, index=df.index)
7
8 # Verificar resultado
9 print("=== MIN-MAX SCALING ===")
10 print("\nANTES:")
11 print(df[numeric_cols].describe())
12
13 print("\nDEPOIS:")
14 print(df[[c + '_minmax' for c in numeric_cols]].describe())
```


Min-Max: Verificar Transformação

Python

```
1 # Comparar valores originais e normalizados
2 print("=== COMPARACAO: Valores Originais vs Normalizados ===\n")
3
4 for col in numeric_cols:
5     orig_col = col
6     norm_col = col + '_minmax'
7
8     print(f"\n{col}:")
9     print(f"    Original - Min: {df[orig_col].min():.2f}, Max: {df[orig_col].max():.2f}")
10    print(f"    Normalized - Min: {df[norm_col].min():.2f}, Max: {df[norm_col].max():.2f}")
```

Min-Max: Verificar Transformação (cont.)

</> Python

```
1  # Exemplos
2  sample_idx = df.sample(3).index
3  print(f"  Exemplos:")
4  for idx in sample_idx:
5      orig_val = df.loc[idx, orig_col]
6      norm_val = df.loc[idx, norm_col]
7      print(f"      {orig_val:.2f} - {norm_val:.3f}")
```

Min-Max: Vantagens e Desvantagens

Vantagens:

- ▶ ✓ Range fixo e previsível $[0, 1]$
- ▶ ✓ Fácil de interpretar: 0 = mínimo, 1 = máximo
- ▶ ✓ Preserva exatamente a forma da distribuição
- ▶ ✓ Bom para visualizações
- ▶ ✓ Útil quando há limites naturais nos dados

Desvantagens:

- ▶ ✗ Muito sensível a outliers
- ▶ ✗ Um único valor extremo "comprime" todos os outros
- ▶ ✗ Dados novos podem ficar fora de $[0, 1]$
- ▶ ✗ Exemplo: se novo cliente tem tenure=100 (maior que max=72), valor normalizado seria > 1

Min-Max: Vantagens e Desvantagens

Quando usar:

- ▶ Dados sem outliers extremos
- ▶ Quando precisa de range específico
- ▶ Dados com limites naturais (ex: proporções, percentuais)

Standardization (Z-score): Detalhes

Padronização: média 0, desvio padrão 1

$$X_{std} = \frac{X - \mu}{\sigma}$$

onde μ = média, σ = desvio padrão

Exemplo manual:

- ▶ Dados: [10, 20, 30, 40, 50]
- ▶ Média = 30, Desvio Padrão = 15.8
- ▶ $10 \rightarrow (10-30)/15.8 = -1.27$
- ▶ $30 \rightarrow (30-30)/15.8 = 0.0$
- ▶ $50 \rightarrow (50-30)/15.8 = 1.27$
- ▶ Resultado: [-1.27, -0.63, 0.0, 0.63, 1.27]

Standardization (Z-score): Detalhes

Exemplo manual:

- ▶ Dados: [10, 20, 30, 40, 50]
- ▶ Média = 30, Desvio Padrão = 15.8
- ▶ $10 \rightarrow (10-30)/15.8 = -1.27$
- ▶ $30 \rightarrow (30-30)/15.8 = 0.0$
- ▶ $50 \rightarrow (50-30)/15.8 = 1.27$
- ▶ Resultado: [-1.27, -0.63, 0.0, 0.63, 1.27]

Interpretação Z-score:

- ▶ 0 = valor na média
- ▶ +1 = um desvio padrão acima da média
- ▶ -2 = dois desvios padrão abaixo da média
- ▶ $|Z| > 3$ = outlier extremo (regra geral)

Standardization: Implementação

</> Python

```
1 from sklearn.preprocessing import StandardScaler
2
3 # Criar scaler
4 scaler = StandardScaler()
5 std_cols = [c + '_std' for c in numeric_cols]
6
7 # Aplicar fit_transform
8 scaled_array = scaler.fit_transform(df[numeric_cols])
9 for i, col in enumerate(std_cols):
10     df[col] = scaled_array[:, i]
```

Standardization: Implementação (cont.)

</> Python

```
1 # Verificar resultado
2 print("=== STANDARDIZATION (Z-SCORE) ===")
3 print("\nANTES:")
4 print(df[numeric_cols].describe())
5
6 print("\nDEPOIS:")
7 print(df[std_cols].describe())
8
9 print("\n Média    0, Desvio Padrão    1")
```


Standardization: Interpretação

Python

```
1 # Interpretar Z-scores
2 print("=== INTERPRETACAO DE Z-SCORES ===\n")
3
4 # Exemplo: tenure
5 tenure_std = df['tenure_std']
6
7 print("Distribuição dos Z-scores (tenure):")
8 print(f"  Entre -1 e +1: {((tenure_std >= -1) & (tenure_std <= 1)).mean()
9       *100:.1f}%")
9 print(f"  Entre -2 e +2: {((tenure_std >= -2) & (tenure_std <= 2)).mean()
10      *100:.1f}%")
10 print(f"  Outliers (|Z| > 3): {(abs(tenure_std) > 3).sum()} clientes")
```

Standardization: Interpretação (cont.)

</> Python

```
1 # Identificar extremos
2 print("\n=== CLIENTES COM TENURE EXTREMO ===")
3 print("Muito baixo (Z < -2):")
4 print(df[tenure_std < -2][['tenure', 'tenure_std', 'Churn']].head())
5
6 print("\nMuito alto (Z > 2):")
7 print(df[tenure_std > 2][['tenure', 'tenure_std', 'Churn']].head())
```

Standardization: Vantagens e Desvantagens

Vantagens:

- ▶ ✓ Não está limitado a range específico
- ▶ ✓ Interpretação estatística clara (desvios padrão)
- ▶ ✓ Funciona bem com dados que seguem distribuição normal
- ▶ ✓ Facilita identificação de outliers
- ▶ ✓ Dados novos não "quebram" a escala
- ▶ ✓ Muito usado em análises estatísticas

Desvantagens:

- ▶ ✗ Ainda sensível a outliers (usa média e desvio padrão)
- ▶ ✗ Valores não estão em range fixo
- ▶ ✗ Menos intuitivo para não-estatísticos
- ▶ ✗ Requer distribuição aproximadamente normal para melhor efeito

Standardization: Vantagens e Desvantagens

Quando usar:

- ▶ Dados aproximadamente normais
- ▶ Quando interpretação estatística importa
- ▶ Outliers moderados (não extremos)
- ▶ Análises que assumem normalidade

Robust Scaling: Para Outliers

Escala resistente a valores extremos

$$X_{robust} = \frac{X - median}{IQR}$$

onde $IQR = Q3 - Q1$ (intervalo interquartil)

Por que é robusto?

- ▶ Usa mediana (não média)
- ▶ Usa IQR (não desvio padrão)
- ▶ Ambos são resistentes a outliers
- ▶ Outliers não "inflam" a escala

Robust Scaling: Para Outliers

Exemplo:

- ▶ Dados: [10, 20, 30, 40, 1000] (1000 é outlier)
- ▶ Mediana = 30, Q1 = 20, Q3 = 40, IQR = 20
- ▶ $10 \rightarrow (10-30)/20 = -1.0$
- ▶ $30 \rightarrow (30-30)/20 = 0.0$
- ▶ $1000 \rightarrow (1000-30)/20 = 48.5$
- ▶ Outlier não afeta escala dos valores normais!

Robust Scaling: Implementação

</> Python

```
1 from sklearn.preprocessing import RobustScaler
2
3 # Criar scaler
4 scaler = RobustScaler()
5
6 # Aplicar fit_transform
7 robust_cols = [col + '_robust' for col in numeric_cols]
8 scaled_df = pd.DataFrame(scaler.fit_transform(df[numeric_cols]), columns=
    robust_cols, index=df.index)
9 df = pd.concat([df, scaled_df], axis=1)
```

Robust Scaling: Implementação (cont.)

</> Python

```
1 # Verificar resultado
2 print("=== ROBUST SCALING ===")
3 print("\nANTES:")
4 print(df[numeric_cols].describe())
5
6 print("\nDEPOIS:")
7 robust_cols = [c + '_robust' for c in numeric_cols]
8 print(df[robust_cols].describe())
9
10 print("\n Mediana    0, IQR-based scaling")
```


Comparação: Min-Max vs Standard vs Robust

</> Python

```
1 print("=== COMPARACAO DAS TECNICAS ===\n")
2 scaling_methods = {
3     'Original': 'tenure',
4     'Min-Max': 'tenure_minmax',
5     'Standard': 'tenure_std',
6     'Robust': 'tenure_robust'
7 }
8 for name, col in scaling_methods.items():
9     if col in df.columns:
10         print(f"\n{name}:")
11         print(f"    Min: {df[col].min():.2f}")
12         print(f"    Median: {df[col].median():.2f}")
13         print(f"    Mean: {df[col].mean():.2f}")
14         print(f"    Max: {df[col].max():.2f}")
15         print(f"    Std: {df[col].std():.2f}")
```

Visualizar Efeito do Scaling

Python

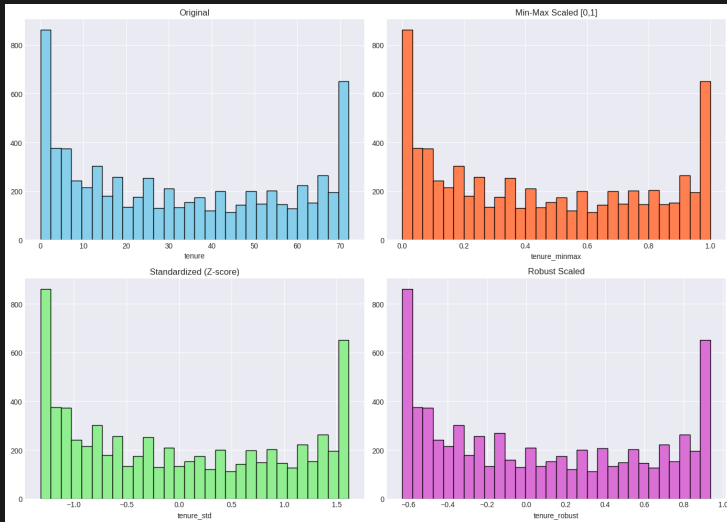
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 fig, axes = plt.subplots(2, 2, figsize=(14, 10))
4 # Original
5 df['tenure'].hist(bins=30, ax=axes[0,0], color='skyblue', edgecolor='black')
6 axes[0,0].set_title('Original')
7 axes[0,0].set_xlabel('tenure')
8
9 # Min-Max
10 df['tenure_minmax'].hist(bins=30, ax=axes[0,1], color='coral', edgecolor='black')
11 axes[0,1].set_title('Min-Max Scaled [0,1]')
12 axes[0,1].set_xlabel('tenure_minmax')
```

Visualizar Efeito do Scaling (cont.)

</> Python

```
1 # Standardized
2 df['tenure_std'].hist(bins=30, ax=axes[1,0], color='lightgreen', edgecolor=
   = 'black')
3 axes[1,0].set_title('Standardized (Z-score)')
4 axes[1,0].set_xlabel('tenure_std')
5
6 # Robust
7 df['tenure_robust'].hist(bins=30, ax=axes[1,1], color='orchid', edgecolor=
   = 'black')
8 axes[1,1].set_title('Robust Scaled')
9 axes[1,1].set_xlabel('tenure_robust')
10
11 plt.tight_layout()
12 plt.show()
```

Visualizar Efeito do Scaling (cont.)



Quando NÃO Normalizar

Situações onde scaling não é necessário ou até prejudicial:

1. Variáveis já estão na mesma escala

- ▶ Ex: proporções (todas entre 0 e 1)
- ▶ Ex: percentuais (todas entre 0 e 100)
- ▶ Scaling seria redundante

2. Escala original tem significado importante

- ▶ Ex: temperatura em Celsius tem interpretação clara
- ▶ Ex: valor em reais é mais interpretável que normalizado
- ▶ Para apresentações de negócio: manter original

Quando NÃO Normalizar

Situações onde scaling não é necessário ou até prejudicial:

3. Análises que não dependem de escala

- ▶ Árvores de decisão (splits são ordinais)
- ▶ Correlações (já são normalizadas $[-1, 1]$)
- ▶ Análises descritivas simples

4. Interpretabilidade é prioridade

- ▶ Stakeholders não-técnicos
- ▶ Valores originais mais fáceis de explicar
- ▶ Relatórios de negócio

Boas Práticas: Separar Dados para Validação

Conceito importante:

Problema:

- ▶ Quando normalizamos, calculamos parâmetros (min, max, média, etc)
- ▶ Esses parâmetros são calculados de TODO o dataset
- ▶ Se no futuro tivermos dados novos, como normalizar?
- ▶ Precisamos usar os MESMOS parâmetros!

Solução: Fit vs Transform

- ▶ `fit()`: Aprende parâmetros (min, max, média, etc)
- ▶ `transform()`: Aplica transformação usando parâmetros aprendidos
- ▶ `fit_transform()`: Faz ambos de uma vez

Boas Práticas: Separar Dados para Validação

Boa prática:

- ▶ Separar dados em conjunto principal e validação
- ▶ Fit apenas no conjunto principal
- ▶ Transform em ambos usando os mesmos parâmetros
- ▶ Simula como seria com dados futuros/novos

Implementar Separação de Dados

</> Python

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 # 1. Separar dados (80% principal, 20% validação)
4 df_main, df_val = train_test_split(df, test_size=0.2, random_state=42)
5 print(f"Conjunto principal: {len(df_main)} registros")
6 print(f"Conjunto validação: {len(df_val)} registros")
7
8 # 2. Fit scaler APENAS no conjunto principal
9 scaler = StandardScaler()
10 scaler.fit(df_main[numeric_cols])
11 print("\n== PARAMETROS APRENDIDOS (do conjunto principal) ==")
12 print(f"Médias: {scaler.mean_}")
13 print(f"Desvios: {scaler.scale_}")
```

Aplicar Transformação aos Dois Conjuntos

</> Python

```
1 # 3. Transform em AMBOS os conjuntos (usando mesmos parâmetros)
2 df_main_scaled = scaler.transform(df_main[numeric_cols])
3 df_val_scaled = scaler.transform(df_val[numeric_cols])
4 # Converter de volta para DataFrame
5 df_main[std_cols] = df_main_scaled
6 df_val[std_cols] = df_val_scaled
7
8 # 4. Verificar que validação usa mesmos parâmetros
9 print("=== VERIFICACAO ===")
10 print("\nConjunto Principal:")
11 print(f" Média tenure_std: {df_main['tenure_std'].mean():.3f} ( 0)")
12 print(f" Std tenure_std: {df_main['tenure_std'].std():.3f} ( 1)")
13 print("\nConjunto Validação:")
14 print(f" Média tenure_std: {df_val['tenure_std'].mean():.3f} (pode 0)")
15 print(f" Std tenure_std: {df_val['tenure_std'].std():.3f} (pode 1)")
16 print("\n Validação usa parâmetros do principal (correto!)")
```

Por que Separar é Importante?

Simula cenário real:

Analogia:

- ▶ Conjunto principal = dados históricos que temos hoje
- ▶ Conjunto validação = dados futuros que chegarão amanhã
- ▶ Quando dados novos chegarem, não podemos "re-calcular" min/max/média
- ▶ Precisamos usar os parâmetros que já tínhamos

Exemplo prático:

- ▶ Hoje: normalizamos tenure com $\text{min}=0$, $\text{max}=72$
- ▶ Amanhã: chega cliente com $\text{tenure}=80$ (novo máximo!)
- ▶ Usamos: $(80-0)/72 = 1.11$ (>1 é OK!)
- ▶ NÃO usamos: re-calcular com novo $\text{max}=80$
- ▶ Por quê? Consistência e comparabilidade

Por que Separar é Importante?

Relevância para ML (futuro):

- ▶ Em ML, isso é CRÍTICO (evita "data leakage")
- ▶ Mesma boa prática que usamos aqui
- ▶ Você já estará familiarizado quando chegar lá!

Transformações Matemáticas

Para corrigir distribuições assimétricas:

</> Python

```
1 import numpy as np
2 # 1. Log Transform (para assimetria à direita)
3 # Útil para dados exponenciais ou com cauda longa
4 df['TotalCharges_log'] = np.log1p(df['TotalCharges'])
5 # log1p = log(1+x), evita log(0)
6 # 2. Square Root (menos agressivo que log)
7 df['TotalCharges_sqrt'] = np.sqrt(df['TotalCharges'])
8 # 3. Inverso (para assimetria à esquerda)
9 # Cuidado: inverte ordem!
10 df['tenure_inverse'] = 1 / (df['tenure'] + 1) # +1 evita divisão por 0
11 print("=== TRANSFORMACOES MATEMATICAS ===")
12 print("\nOriginal vs Transformada:")
13 print(df[['TotalCharges', 'TotalCharges_log', 'TotalCharges_sqrt']].
        describe())
```

Quando Usar Transformações Matemáticas

Python

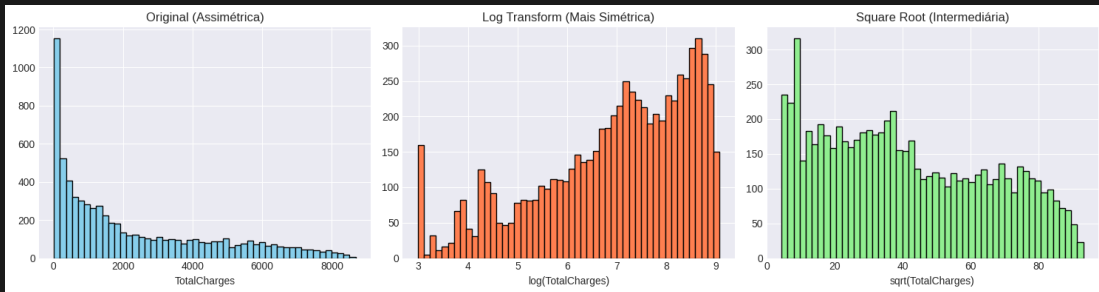
```
1 # Visualizar efeito da transformação
2 import matplotlib.pyplot as plt
3
4 fig, axes = plt.subplots(1, 3, figsize=(15, 4))
5
6 # Original (assimétrica)
7 df['TotalCharges'].hist(bins=50, ax=axes[0], color='skyblue', edgecolor='
    black')
8 axes[0].set_title('Original (Assimétrica)')
9 axes[0].set_xlabel('TotalCharges')
```

Quando Usar Transformações Matemáticas (cont.)

Python

```
1 # Log (mais simétrica)
2 df['TotalCharges_log'].hist(bins=50, ax=axes[1], color='coral', edgecolor=
   'black')
3 axes[1].set_title('Log Transform (Mais Simétrica)')
4 axes[1].set_xlabel('log(TotalCharges)')
5
6 # Sqrt (intermediária)
7 df['TotalCharges_sqrt'].hist(bins=50, ax=axes[2], color='lightgreen',
   edgecolor='black')
8 axes[2].set_title('Square Root (Intermediária)')
9 axes[2].set_xlabel('sqrt(TotalCharges)')
10
11 plt.tight_layout()
12 plt.show()
```

Quando Usar Transformações Matemáticas (cont.)



Transformações: Quando Usar Cada Uma

Guia de decisão:

Log Transform:

- ▶ Distribuição muito assimétrica à direita
- ▶ Dados exponenciais (renda, população)
- ▶ Range muito grande (vários orders de magnitude)
- ▶ **Cuidado:** Requer valores > 0

Square Root:

- ▶ Assimetria moderada à direita
- ▶ Menos agressivo que log
- ▶ Dados de contagem (número de eventos)
- ▶ Aceita zero

Transformações: Quando Usar Cada Uma

Box-Cox (avançado):

- ▶ Encontra transformação ótima automaticamente
- ▶ Família de transformações parametrizada
- ▶ Requer valores > 0
- ▶ Mais usado em análises estatísticas formais

Pipeline Completo de Normalização

</> Python

```
1 def normalize_features(df, numeric_cols, method='standard',
2                       fit_data=None, return_scaler=False):
3     """
4     Parâmetros:
5     -----
6     df : DataFrame
7         Dados a normalizar
8     numeric_cols : list
9         Colunas numéricas
10    method : str
11        'minmax', 'standard', ou 'robust'
```

Pipeline Completo (cont.)

</> Python

```
1 fit_data : DataFrame, optional
2     Se fornecido, fit nos fit_data e transform no df
3 return_scaler : bool
4     Se True, retorna (df_normalized, scaler)
5 """
6 from sklearn.preprocessing import MinMaxScaler, StandardScaler,
7 RobustScaler
8
9 # Escolher scaler
10 scalers = {
11     'minmax': MinMaxScaler(),
12     'standard': StandardScaler(),
13     'robust': RobustScaler()
14 }
15 scaler = scalers[method]
```

Pipeline Completo (cont.)

</> Python

```
1  # Fit e transform
2  if fit_data is not None:
3      # Fit em fit_data, transform em df
4      scaler.fit(fit_data[numeric_cols])
5  else:
6      # Fit e transform no mesmo df
7      scaler.fit(df[numeric_cols])
8
9  # Transform
10 df_normalized = df.copy()
11 df_normalized[numeric_cols] = scaler.transform(df[numeric_cols])
12
13 if return_scaler:
14     return df_normalized, scaler
15 return df_normalized
```

Pipeline Completo (cont.)

Python

```
1 # Uso
2 df_normalized = normalize_features(df, numeric_cols, method='standard')
3 print("Normalização aplicada!")
```

Análise Exploratória com Dados Normalizados

</> Python

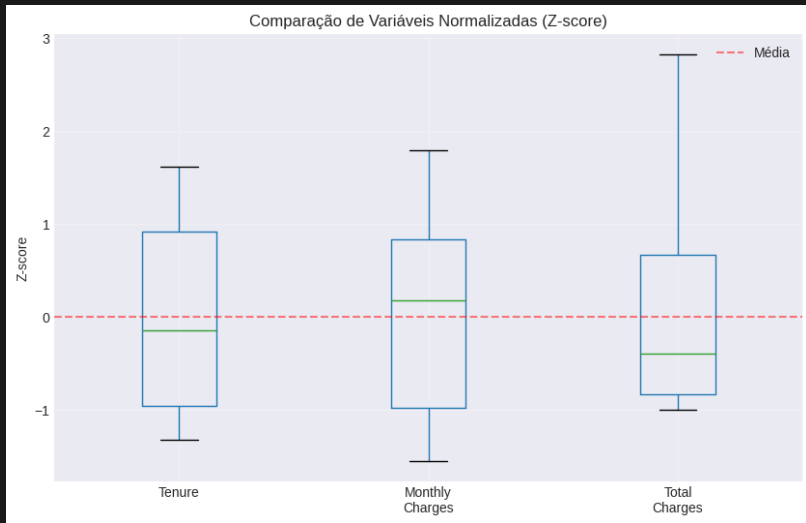
```
1 # Comparar múltiplas variáveis no mesmo gráfico
2 import matplotlib.pyplot as plt
3
4 # Selecionar variáveis normalizadas
5 normalized_features = ['tenure_std', 'MonthlyCharges_std', '
    TotalCharges_std']
```

Análise Exploratória com Dados Normalizados (cont.)

</> Python

```
1 # Boxplot comparativo
2 fig, ax = plt.subplots(figsize=(10, 6))
3 df[normalized_features].boxplot(ax=ax)
4 ax.set_title('Comparação de Variáveis Normalizadas (Z-score)')
5 ax.set_ylabel('Z-score')
6 ax.set_xticklabels(['Tenure', 'Monthly\nCharges', 'Total\nCharges'])
7 ax.axhline(y=0, color='r', linestyle='--', alpha=0.5, label='Média')
8 ax.grid(True, alpha=0.3)
9 plt.legend()
10 plt.show()
11 print(" Agora podemos comparar variáveis diretamente!")
12 print("Todas na mesma escala (desvios padrão)")
```


Análise Exploratória com Dados Normalizados (cont.)



Heatmap com Variáveis Normalizadas

</> Python

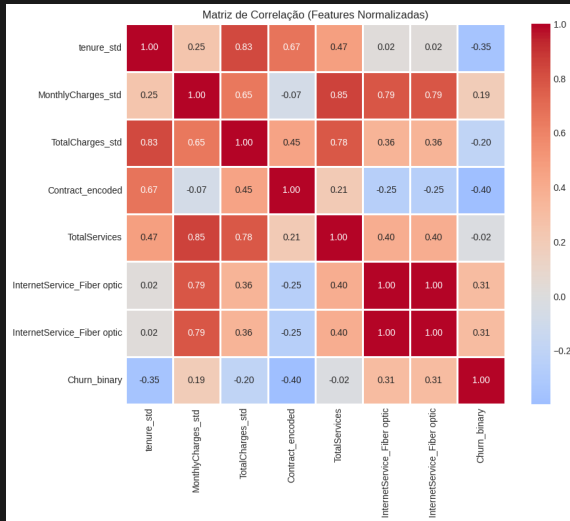
```
1 # Matriz de correlação com features normalizadas
2 import seaborn as sns
3
4 features_for_heatmap = [
5     'tenure_std', 'MonthlyCharges_std', 'TotalCharges_std',
6     'Contract_encoded', 'TotalServices',
7     'InternetService_Fiber optic',
8     'Churn_binary'
9 ]
10
11 # Calcular correlações
12 corr_matrix = df[features_for_heatmap].corr()
```

Heatmap com Variáveis Normalizadas (cont.)

</> Python

```
1 # Heatmap
2 fig, ax = plt.subplots(figsize=(10, 8))
3 sns.heatmap(corr_matrix, annot=True, fmt='.2f', cmap='coolwarm',
4             center=0, square=True, linewidths=1, ax=ax)
5 ax.set_title('Matriz de Correlação (Features Normalizadas)')
6 plt.tight_layout()
7 plt.show()
8
9 print("Correlações mais interpretáveis com dados normalizados")
```

Heatmap com Variáveis Normalizadas (cont.)



Resumo do Bloco 3

O que aprendemos:

Conceitos:

- ▶ Por que normalizar (comparabilidade, visualização)
- ▶ Diferenças entre técnicas de scaling
- ▶ Quando usar cada técnica
- ▶ Quando NÃO normalizar

Técnicas:

- ▶ Min-Max Scaling (range fixo $[0,1]$)
- ▶ Standardization (Z-score, média=0, std=1)
- ▶ Robust Scaling (resistente a outliers)
- ▶ Transformações matemáticas (log, sqrt)

Boas Práticas:

- ▶ Separar dados para validação
- ▶ Fit vs Transform (parâmetros consistentes)
- ▶ Validar resultado visualmente
- ▶ Documentar escolhas

Checklist de Normalização

Antes de normalizar, pergunte:

1. ☐ Variáveis estão em escalas muito diferentes?
 - ▶ Se não: talvez não precise normalizar
2. ☐ Qual objetivo da análise?
 - ▶ Comparação visual → Min-Max
 - ▶ Análise estatística → Standardization
 - ▶ Muitos outliers → Robust
3. ☐ Há outliers extremos?
 - ▶ Sim → Robust ou tratar outliers primeiro
 - ▶ Não → Min-Max ou Standard
4. ☐ Interpretabilidade é crítica?
 - ▶ Sim → Considerar manter valores originais
 - ▶ Para stakeholders: valores originais
 - ▶ Para análise técnica: pode normalizar
5. ☐ Tenho dados que simulariam dados futuros?
 - ▶ Sim → Separar e validar processo
 - ▶ Sempre boa prática!

Bloco 4

Pipeline Completo

O que vamos fazer

Integração:

- ▶ Juntar Feature Engineering + Encoding + Normalização
- ▶ Criar função pipeline end-to-end
- ▶ Aplicar ao dataset Telco Churn

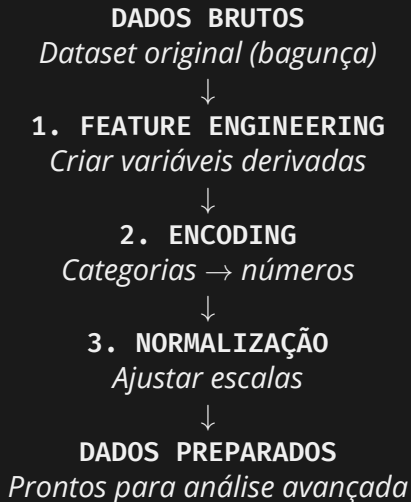
Análise Avançada:

- ▶ Explorar dados completamente transformados
- ▶ Identificar insights que não eram visíveis antes
- ▶ Correlações, visualizações, comparações

Fechamento:

- ▶ Preview: conexão com Machine Learning
- ▶ Checklist de boas práticas
- ▶ Próximos passos no curso

Pipeline Completo: Estrutura



Pipeline: Função Completa

</> Python

```
1 def transform_telco_churn(df, fit_scaler=True, scaler=None):
2     """
3     Pipeline completo de transformação para Telco Churn
4
5     Parameters:
6     -----
7     df : DataFrame
8         Dataset bruto
9     fit_scaler : bool
10         Se True, fit novo scaler. Se False, usa scaler fornecido
11     scaler : sklearn scaler, optional
12         Scaler pré-treinado (para dados novos)
```

Pipeline: Função Completa (cont.)

</> Python

```
1 Returns:
2 -----
3 df_transformed, scaler
4 """
5 from sklearn.preprocessing import StandardScaler
6
7 df_transformed = df.copy()
8
9 # ===== 1. FEATURE ENGINEERING =====
10 print("1. Feature Engineering...")
```

Pipeline: Função Completa (cont.)

</> Python

```
1  # Operações aritméticas
2  df_transformed['AvgMonthlySpend'] = (
3      df_transformed['TotalCharges'] / df_transformed['tenure']
4  )
5  df_transformed['AvgMonthlySpend'].replace([np.inf, -np.inf], np.nan,
6  inplace=True)
7  df_transformed.loc[df_transformed['tenure'] == 0, 'AvgMonthlySpend'] = \
    df_transformed.loc[df_transformed['tenure'] == 0, 'MonthlyCharges']
```

Pipeline: Função Completa (cont.)

</> Python

```
1  # Agregações
2  service_cols = ['PhoneService', 'MultipleLines', 'InternetService',
3                  'OnlineSecurity', 'OnlineBackup', 'DeviceProtection',
4                  'TechSupport', 'StreamingTV', 'StreamingMovies']
5  df_transformed['TotalServices'] = 0
6  for col in service_cols:
7      df_transformed['TotalServices'] += (
8          ~df_transformed[col].isin(['No', 'No internet service',
9                                     'No phone service'])
10         ).astype(int)
```

Pipeline: Função Completa (cont.)

</> Python

```
1  # Flags
2  df_transformed['IsPremium'] = (
3      (df_transformed['TotalServices'] >= 6) &
4      (df_transformed['Contract'].isin(['One year', 'Two year'])))
5  ).astype(int)
```

Pipeline: Função Completa (cont.)

</> Python

```
1  # ===== 2. ENCODING =====
2  print("2. Encoding...")
3
4  # Ordinais (Label)
5  contract_map = {'Month-to-month': 0, 'One year': 1, 'Two year': 2}
6  df_transformed['Contract_encoded'] = df_transformed['Contract'].map(
    contract_map)
7
8  # Binárias (0/1)
9  binary_cols = ['gender', 'Partner', 'Dependents',
10                'PhoneService', 'PaperlessBilling']
11  for col in binary_cols:
12      df_transformed[f'{col}_encoded'] = (
13          df_transformed[col] == 'Yes'
14      ).astype(int)
```


Pipeline: Função Completa (cont.)

</> Python

```
1 df_transformed['SeniorCitizen_encoded'] = df_transformed['  
SeniorCitizen']  
2  
3 # Nominais (One-Hot)  
4 nominal_cols = ['PaymentMethod', 'InternetService']  
5 df_transformed = pd.get_dummies(df_transformed, columns=nominal_cols,  
6                                prefix=nominal_cols, drop_first=True,  
dtype=int)
```

Pipeline: Função Completa (cont.)

</> Python

```
1      # ===== 3. NORMALIZAÇÃO =====
2      print("3. Normalização...")
3
4      # Colunas numéricas para normalizar
5      numeric_to_scale = ['tenure', 'MonthlyCharges', 'TotalCharges',
6                          'AvgMonthlySpend', 'TotalServices']
7
8      if fit_scaler:
9          # Criar e treinar novo scaler
10         scaler = StandardScaler()
11         df_transformed[numeric_to_scale] = scaler.fit_transform(
12             df_transformed[numeric_to_scale]
13         )
```

Pipeline: Função Completa (cont.)

</> Python

```
1  else:
2      # Usar scaler existente
3      if scaler is None:
4          raise ValueError("Forneça scaler quando fit_scaler=False")
5      df_transformed[numeric_to_scale] = scaler.transform(
6          df_transformed[numeric_to_scale]
7      )
8
9      # Target binário
10     df_transformed['Churn_binary'] = (df_transformed['Churn'] == 'Yes').
11     astype(int)
12
13     print(" Pipeline completo!")
14     return df_transformed, scaler
```

Aplicar Pipeline ao Dataset

</> Python

```
1 # Aplicar pipeline completo
2 print("=== APLICANDO PIPELINE COMPLETO ===\n")
3
4 # Transformar
5 df_final, trained_scaler = transform_telco_churn(df, fit_scaler=True)
6
7 print(f"\nShape original: {df.shape}")
8 print(f"Shape final: {df_final.shape}")
9 print(f"Colunas adicionadas: {df_final.shape[1] - df.shape[1]}")
```

Aplicar Pipeline ao Dataset (cont.)

</> Python

```
1 # Verificar tipos
2 print("\n=== TIPOS DE DADOS ===")
3 print(f"Numéricas: {len(df_final.select_dtypes(include='number').columns)}")
4 print(f"Categóricas: {len(df_final.select_dtypes(include='object').columns)}")
5
6 # Listar features finais
7 numeric_features = df_final.select_dtypes(include='number').columns
8 print(f"\nTotal de features numéricas: {len(numeric_features)}")
9 print("\n Dataset completamente transformado!")
```

Validar Pipeline

Checklist de validação:

</> Python

```
1 # Validações pós-pipeline
2 print("=== VALIDACAO DO PIPELINE ===\n")
3
4 # 1. Verificar missing values
5 print("1. Missing values:")
6 missing = df_final.isnull().sum()
7 print(f"    Total: {missing.sum()}")
8 if missing.sum() > 0:
9     print("    Colunas com missing:")
10    print(missing[missing > 0])
```

Validar Pipeline (cont.)

</> Python

```
1 # 2. Verificar ranges das normalizadas
2 print("\n2. Ranges das variáveis normalizadas:")
3 for col in ['tenure', 'MonthlyCharges', 'TotalCharges']:
4     print(f"    {col}: [{df_final[col].min():.2f}, {df_final[col].max():.2f}
5         ]")
6     print(f"        Média: {df_final[col].mean():.2f}, Std: {df_final[col].
7         std():.2f}")
8
9 # 3. Verificar encoding
10 print("\n3. Encoding:")
11 encoded_cols = [c for c in df_final.columns if '_encoded' in c or
12                 'PaymentMethod_' in c or 'InternetService_' in c]
13 print(f"    Colunas encodadas: {len(encoded_cols)}")
```

Validar Pipeline (cont.)

</> Python

```
1 # 4. Verificar features criadas
2 print("\n4. Features criadas:")
3 created = ['AvgMonthlySpend', 'TotalServices', 'IsPremium']
4 for feat in created:
5     if feat in df_final.columns:
6         print(f"    {feat}")
7     else:
8         print(f"    {feat} - FALTANDO!")
```


Validar Pipeline (cont.)

</> Python

```
1 # 5. Verificar se há categóricas restantes
2 print("\n5. Categóricas restantes (exceto IDs e target):")
3 cat_cols = df_final.select_dtypes('object').columns
4 cat_cols = [c for c in cat_cols if c not in ['customerID', 'Churn']]
5 print(f"    {len(cat_cols)} colunas")
6 if len(cat_cols) > 0:
7     print(f"    {cat_cols}")
```

Análise Exploratória Avançada

Agora podemos fazer análises que não eram possíveis antes:

</> Python

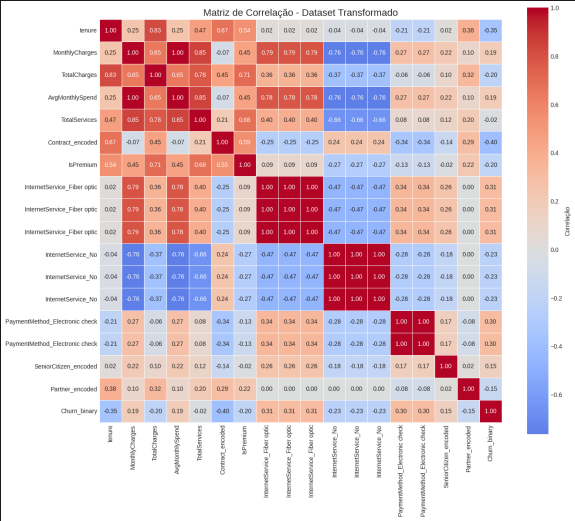
```
1 # Matriz de correlação completa
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Selecionar features relevantes
6 analysis_features = [
7     'tenure', 'MonthlyCharges', 'TotalCharges',
8     'AvgMonthlySpend', 'TotalServices',
9     'Contract_encoded', 'IsPremium',
10    'InternetService_Fiber optic', 'InternetService_No',
11    'PaymentMethod_Electronic check',
12    'SeniorCitizen_encoded', 'Partner_encoded',
13    'Churn_binary'
14 ]
```

Análise Exploratória Avançada (cont.)

</> Python

```
1 # Calcular correlações
2 corr_matrix = df_final[analysis_features].corr()
3
4 # Heatmap
5 fig, ax = plt.subplots(figsize=(14, 12))
6 sns.heatmap(corr_matrix, annot=True, fmt='.2f', cmap='coolwarm',
7             center=0, square=True, linewidths=0.5,
8             cbar_kws={'label': 'Correlação'}, ax=ax)
9 ax.set_title('Matriz de Correlação - Dataset Transformado', fontsize=16)
10 plt.tight_layout()
11 plt.show()
```

Análise Exploratória Avançada (cont.)



Insights dos Dados Transformados

</> Python

```
1 # Top correlações com Churn
2 print("=== TOP CORRELACOES COM CHURN ===\n")
3
4 correlations = corr_matrix['Churn_binary'].drop('Churn_binary')
5 correlations_abs = correlations.abs().sort_values(ascending=False)
6
7 print("Positivas (aumentam churn):")
8 positive = correlations[correlations > 0].sort_values(ascending=False)
9 for feat, corr in positive.head(5).items():
10     print(f"    {feat}: {corr:.3f}")
```

Insights dos Dados Transformados (cont.)

</> Python

```
1 print("\nNegativas (reduzem churn):")
2 negative = correlations[correlations < 0].sort_values()
3 for feat, corr in negative.head(5).items():
4     print(f"    {feat}: {corr:.3f}")
5
6 # Insights de negócio
7 print("\n=== INSIGHTS DE NEGOCIO ===")
8 print("    Contratos longos (Contract_encoded) reduzem churn")
9 print("    Mais serviços (TotalServices) reduz churn")
10 print("    Fibra óptica aumenta churn (provavelmente por preço)")
11 print("    Clientes premium têm menos churn")
```

Comparação: Antes vs Depois das Transformações

Python

```
1 # Visualizar impacto das transformações
2 fig, axes = plt.subplots(2, 2, figsize=(14, 10))
3
4 # 1. Antes: variáveis em escalas diferentes
5 axes[0,0].scatter(df['tenure'], df['TotalCharges'], alpha=0.3)
6 axes[0,0].set_xlabel('tenure (0-72)')
7 axes[0,0].set_ylabel('TotalCharges (0-8684)')
8 axes[0,0].set_title('ANTES: Escalas Diferentes')
9
10 # 2. Depois: variáveis normalizadas
11 axes[0,1].scatter(df_final['tenure'], df_final['TotalCharges'], alpha=0.3)
12 axes[0,1].set_xlabel('tenure (normalizado)')
13 axes[0,1].set_ylabel('TotalCharges (normalizado)')
14 axes[0,1].set_title('DEPOIS: Escalas Comparáveis')
```

Comparação: Antes vs Depois das Transformações (cont.)

</> Python

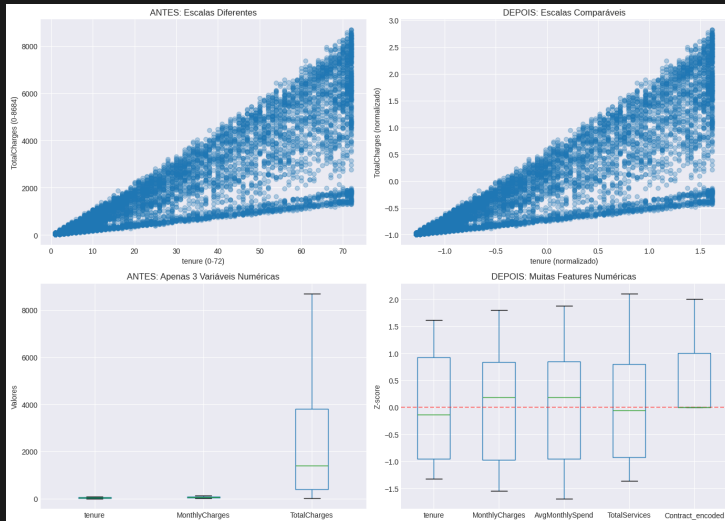
```
1 # 3. Antes: limitado a numéricas originais
2 original_numeric = ['tenure', 'MonthlyCharges', 'TotalCharges']
3 df[original_numeric].boxplot(ax=axes[1,0])
4 axes[1,0].set_title('ANTES: Apenas 3 Variáveis Numéricas')
5 axes[1,0].set_ylabel('Valores')
```


Comparação: Antes vs Depois das Transformações (cont.)

</> Python

```
1 # 4. Depois: muitas features numéricas
2 analysis_subset = ['tenure', 'MonthlyCharges', 'AvgMonthlySpend',
3                    'TotalServices', 'Contract_encoded']
4 df_final[analysis_subset].boxplot(ax=axes[1,1])
5 axes[1,1].set_title('DEPOIS: Muitas Features Numéricas')
6 axes[1,1].set_ylabel('Z-score')
7 axes[1,1].axhline(y=0, color='r', linestyle='--', alpha=0.5)
8 plt.tight_layout()
9 plt.show()
10 print("\n=== GANHOS ===")
11 print("  Variáveis comparáveis visualmente")
12 print("  Features criadas revelam padrões")
13 print("  Todas as variáveis são analisáveis numericamente")
14 print("  Pronto para análises avançadas e ML (futuro)")
```

Comparação: Antes vs Depois das Transformações (cont.)



Conexão com Machine Learning

Como tudo isso se conecta ao ML:

O que fizemos hoje:

- ▶ Feature Engineering: criar variáveis úteis
- ▶ Encoding: converter categorias em números
- ▶ Normalização: colocar em escalas comparáveis
- ▶ **Objetivo:** análise exploratória avançada

Em Machine Learning (cursos futuros):

- ▶ Usaremos *exatamente o mesmo pipeline*
- ▶ Mas ao invés de análise manual...
- ▶ Treinaremos algoritmos para aprender padrões automaticamente
- ▶ Ex: algoritmo aprende que "Contract_encoded alto + TotalServices alto = baixo churn"
- ▶ Depois aplica em clientes novos para prever risco de churn

Nota Importante

Transformação é pré-requisito tanto para análise quanto para ML!
Você já domina a base essencial.

Checklist Final de Transformação

Antes de considerar dados prontos:

1. ☐ **Feature Engineering aplicado**

- ▶ Features úteis criadas e validadas
- ▶ Correlação com target verificada
- ▶ Faz sentido para o negócio

2. ☐ **Encoding completo**

- ▶ Nenhuma coluna categórica restante (exceto IDs)
- ▶ Estratégia apropriada para cada tipo
- ▶ Ordinais preservam ordem, nominais viram One-Hot

3. ☐ **Normalização aplicada (se necessário)**

- ▶ Variáveis em escalas comparáveis
- ▶ Técnica apropriada escolhida
- ▶ Parâmetros salvos para dados futuros

Checklist Final de Transformação

Antes de considerar dados prontos:

4. ☐ **Validações passam**

- ▶ Missing values tratados
- ▶ Ranges fazem sentido
- ▶ Apenas dados numéricos

5. ☐ **Pipeline documentado**

- ▶ Código organizado e comentado
- ▶ Reproduzível em novos dados
- ▶ Decisões justificadas

Recap: Aula 10 Completa

O que cobrimos:

Bloco 0: Preview de ML

- ▶ Contexto: o que é ML e por que transformações importam

Bloco 1: Feature Engineering

- ▶ 12 técnicas para criar variáveis derivadas
- ▶ Validação através de correlações e visualizações

Bloco 2: Encoding

- ▶ Label, One-Hot, Binary encoding
- ▶ Estratégias para diferentes tipos de variáveis

Recap: Aula 10 Completa

O que cobrimos:

Bloco 3: Normalização

- ▶ Min-Max, Standardization, Robust Scaling
- ▶ Comparabilidade e visualização

Bloco 4: Pipeline

- ▶ Integração completa
- ▶ Análise exploratória avançada

Habilidades Adquiridas

Você agora é capaz de:

Técnicas:

- ▶ ✓ Criar features úteis sistematicamente
- ▶ ✓ Encodar qualquer tipo de variável categórica
- ▶ ✓ Normalizar dados apropriadamente
- ▶ ✓ Construir pipeline reproduzível

Análise:

- ▶ ✓ Validar transformações através de EDA
- ▶ ✓ Identificar padrões em dados transformados
- ▶ ✓ Interpretar correlações e visualizações
- ▶ ✓ Extrair insights acionáveis

Habilidades Adquiridas

Boas Práticas:

- ▶ ✓ Documentar decisões
- ▶ ✓ Separar dados para validação
- ▶ ✓ Pensar em conhecimento de domínio
- ▶ ✓ Criar código reproduzível

Aplicação Prática

Como usar essas habilidades:

Em projetos de análise:

- ▶ Preparar dados para análises quantitativas
- ▶ Criar visualizações mais informativas
- ▶ Identificar fatores associados a outcomes
- ▶ Gerar recomendações baseadas em dados

No mercado de trabalho:

- ▶ Analista de Dados: preparação de relatórios
- ▶ Cientista de Dados: base para ML
- ▶ Business Intelligence: dashboards avançados
- ▶ Pesquisa: análises estatísticas robustas

Próximos passos:

- ▶ Praticar com diferentes datasets
- ▶ Explorar domínios novos
- ▶ Combinar com técnicas de visualização avançada
- ▶ Preparação para cursos de ML

Preview: Aula 11 - Séries Temporais

Tema: Séries Temporais

- ▶ Dados que variam ao longo do tempo
- ▶ Exemplos: vendas mensais, temperatura diária, preços de ações

O que aprenderemos:

- ▶ Manipulação de timestamps com Pandas
- ▶ Resampling e agregações temporais
- ▶ Decomposição de séries (tendência, sazonalidade)
- ▶ Visualizações temporais
- ▶ Rolling windows e médias móveis

Conexão com hoje:

- ▶ Usaremos as mesmas técnicas de transformação
- ▶ Mas adaptadas para dados temporais
- ▶ Feature engineering temporal (lag features, diferenças)

Exercícios sugeridos:

1. Aplicar pipeline a novo dataset

- ▶ Escolha dataset do Kaggle ou UCI ML Repository
- ▶ Aplique todas as transformações aprendidas
- ▶ Compare antes/depois

2. Criar features criativas

- ▶ Para dataset Telco: invente 5 features novas
- ▶ Valide com correlação e visualização
- ▶ Justifique cada uma

Exercícios sugeridos:

3. Experimentar diferentes scalings

- ▶ Compare visualmente Min-Max vs Standard vs Robust
- ▶ Em dataset com outliers
- ▶ Documente diferenças

4. Refinar pipeline

- ▶ Adicionar tratamento de erros
- ▶ Logging de transformações
- ▶ Testes automatizados

Exercício Prático

Tempo: 60 minutos

Entrega: via Moodle (notebook)

Tarefas:

1. (atualizado durante a aula)

Notebook: Disponível no Moodle

Obrigado!