

Programação para Ciência de Dados

Análise de Séries Temporais - Fundamentos e Manipulação

Arthur Casals

11 de Novembro de 2025

Agenda

- ▶ Introdução: O que são Séries Temporais?
- ▶ Bloco 1: DatetimeIndex e Manipulação Básica
- ▶ Bloco 2: Resample, Rolling Windows e Agregações
- ▶ Bloco 3: Decomposição e Análise de Componentes
- ▶ Bloco 4: Visualizações e Boas Práticas

O que são Séries Temporais?

Série Temporal = sequência de dados ordenados no tempo

Características principais:

- ▶ Dados coletados em intervalos regulares
- ▶ Ordem importa (temporal)
- ▶ Cada observação tem timestamp
- ▶ Passado influencia presente/futuro

Diferença de dados tradicionais:

| Dados Tradicionais | Séries Temporais |
|---------------------------|-------------------------|
| Linhas independentes | Linhas dependentes |
| Ordem não importa | Ordem CRÍTICA |
| Sem timestamp | Com timestamp |
| i.i.d. assumido | Autocorrelação presente |

Exemplos de Séries Temporais

Estão por toda parte!

Finanças:

- ▶ Preços de ações (segundo a segundo)
- ▶ Taxa de câmbio (diária)
- ▶ Volume de negociação

Negócios:

- ▶ Vendas (diária, mensal)
- ▶ Tráfego de website (horária)
- ▶ Demanda de produtos

Ciência:

- ▶ Temperatura (horária)
- ▶ Qualidade do ar (nossa aula de hoje!)
- ▶ Dados de sensores IoT

Por que Séries Temporais São Especiais?

Características únicas que exigem tratamento especial:

1. Dependência Temporal

- ▶ Valor de hoje depende de ontem
- ▶ Autocorrelação presente
- ▶ Não podemos embaralhar os dados!

2. Tendência

- ▶ Movimento de longo prazo (subindo/descendo)
- ▶ Crescimento ou declínio ao longo do tempo

Por que Séries Temporais São Especiais?

Características únicas que exigem tratamento especial:

3. Sazonalidade

- ▶ Padrões que se repetem periodicamente
- ▶ Diário, semanal, mensal, anual

4. Ciclicidade

- ▶ Flutuações de longo prazo (sem período fixo)
- ▶ Ciclos econômicos, por exemplo

Componentes de uma Série Temporal

Decomposição clássica:

$$Y_t = T_t + S_t + C_t + I_t$$

Onde:

- ▶ Y_t : Valor observado no tempo t
- ▶ T_t : Tendência (Trend) - direção de longo prazo
- ▶ S_t : Sazonalidade (Seasonal) - padrão repetitivo
- ▶ C_t : Ciclo (Cycle) - flutuações de longo prazo
- ▶ I_t : Irregular/Resíduo - ruído aleatório

Modelo alternativo (multiplicativo):

$$Y_t = T_t \times S_t \times C_t \times I_t$$

Objetivo: Separar componentes para entender melhor

Visualização dos Componentes

Exemplo: Vendas mensais

Tendência:

- ▶ Crescimento constante ao longo dos anos
- ▶ Linha suave subindo

Sazonalidade:

- ▶ Pico em dezembro (Natal)
- ▶ Queda em janeiro/fevereiro
- ▶ Repete todo ano

Visualização dos Componentes

Exemplo: Vendas mensais

Ciclo:

- ▶ Recessão econômica reduz vendas
- ▶ Duração: 3-5 anos
- ▶ Não tem período fixo

Irregular:

- ▶ Eventos aleatórios (promoção inesperada)
- ▶ Não previsível

Aplicações de Análise de Séries Temporais

O que podemos fazer:

1. Descrição e Exploração

- ▶ Entender padrões históricos
- ▶ Identificar anomalias
- ▶ Visualizar tendências

2. Previsão (Forecasting)

- ▶ Prever valores futuros
- ▶ ARIMA, Prophet, LSTM
- ▶ (Não veremos modelagem hoje - foco em manipulação)

O que podemos fazer:

3. Detecção de Anomalias

- ▶ Identificar comportamentos estranhos
- ▶ Fraude, falhas de sistema

4. Análise de Causa e Efeito

- ▶ Impacto de eventos
- ▶ Correlação temporal entre séries

Dataset de Hoje: Air Quality

Dados de qualidade do ar em uma cidade italiana

Origem:

- ▶ UCI Machine Learning Repository
- ▶ Sensor multivariado de qualidade do ar
- ▶ Localização: cidade industrial italiana
- ▶ Período: março 2004 - fevereiro 2005

Frequência:

- ▶ Medições HORÁRIAS
- ▶ 9358 observações (1 ano de dados)
- ▶ Timestamp completo: data + hora

Dataset de Hoje: Air Quality

Dados de qualidade do ar em uma cidade italiana

Variáveis principais:

- ▶ CO (Carbon Monoxide) - mg/m^3
- ▶ NO2 (Nitrogen Dioxide) - $\mu\text{g/m}^3$
- ▶ Temperature ($^{\circ}\text{C}$)
- ▶ Relative Humidity (%)
- ▶ Absolute Humidity

Por que Air Quality é Bom Exemplo?

Dataset ideal para aprender séries temporais:

Características interessantes:

- ▶ **Frequência horária:** Alta granularidade
- ▶ **Sazonalidade clara:**
 - ▶ Diária: poluição varia durante o dia
 - ▶ Semanal: diferença fim de semana vs dias úteis
 - ▶ Anual: diferença inverno vs verão
- ▶ **Tendências:** Mudanças ao longo do ano
- ▶ **Múltiplas variáveis:** Correlações temporais
- ▶ **Missing values:** Desafio realista

Relevância prática:

- ▶ Problema de saúde pública
- ▶ Previsão para alertas
- ▶ Identificar fontes de poluição

Desafios em Séries Temporais

Problemas comuns:

1. Missing Values

- ▶ Sensores falham
- ▶ Não pode usar métodos tradicionais
- ▶ Precisa interpolação temporal

2. Outliers

- ▶ Eventos extremos vs erros de medição
- ▶ Difícil distinguir

3. Frequências Irregulares

- ▶ Dados não coletados em intervalos exatos
- ▶ Precisa resample/regularização

Desafios em Séries Temporais

Problemas comuns:

4. Múltiplas Frequências

- ▶ Combinar dados diários + mensais
- ▶ Alinhamento temporal

5. Estacionariedade

- ▶ Propriedades estatísticas mudam no tempo
- ▶ Muitos modelos assumem estacionariedade

Pandas para Séries Temporais

DatetimeIndex:

- ▶ Índice temporal especializado
- ▶ Operações temporais nativas
- ▶ Acesso por datas intuitivo

Funções temporais:

- ▶ `resample()`: Mudar frequência
- ▶ `rolling()`: Médias móveis
- ▶ `shift()`: Defasagens (lags)
- ▶ `diff()`: Diferenças temporais
- ▶ `pct_change()`: Variações percentuais

Acessores temporais:

- ▶ `.dt.year`, `.dt.month`, `.dt.day`
- ▶ `.dt.hour`, `.dt.dayofweek`
- ▶ `.dt.quarter`, `.dt.dayofyear`

Objetivos de Aprendizado

Ao final desta aula, você será capaz de:

Conhecimento:

- ▶ Entender componentes de séries temporais
- ▶ Conhecer desafios específicos
- ▶ Interpretar padrões temporais
- ▶ Distinguir tendência vs sazonalidade

Habilidades:

- ▶ Criar e manipular DatetimeIndex
- ▶ Fazer resample (mudar frequência)
- ▶ Calcular médias móveis (rolling)
- ▶ Decompor séries em componentes
- ▶ Tratar missing values em séries
- ▶ Criar visualizações temporais eficazes
- ▶ Extrair features temporais

Setup: Imports e Configurações

</> 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 datetime import datetime, timedelta
7 import zipfile
8 import io
9 import requests
```

Setup: Imports e Configurações

</> Python

```
1
2 # Configuracoes
3 pd.set_option('display.max_columns', None)
4 pd.set_option('display.max_rows', 100)
5
6 plt.style.use('seaborn-v0_8-darkgrid')
7 plt.rcParams['figure.figsize'] = (14, 6)
8 plt.rcParams['font.size'] = 11
9
10 # Para plots melhores
11 sns.set_palette("husl")
12
13 print("Setup completo!")
14
```

Carregar Air Quality Dataset

</> Python

```
1 # Carregar dados
2 zip_url = "https://archive.ics.uci.edu/ml/machine-learning-databases
           /00360/AirQualityUCI.zip"
3 r = requests.get(zip_url)
4 z = zipfile.ZipFile(io.BytesIO(r.content))
5
6 # Carregar CSV
7 with z.open('AirQualityUCI.csv') as f:
8     df = pd.read_csv(f, sep=';', decimal=',', na_values=-200)
```

Carregar Air Quality Dataset

</> Python

```
1 # Combinar 'Date' e 'Time' em uma única coluna
2 df['timestamp'] = pd.to_datetime(
3     df['Date'] + ' ' + df['Time'],
4     format='%d/%m/%Y %H.%M.%S',
5     dayfirst=True,
6     errors='coerce'
7 )
8
9 print("=== DATASET CARREGADO ===")
10 print(f"Shape: {df.shape}")
11 print(f"\nPrimeiras linhas:")
12 print(df.head())
13
14 print(f"\nInfo:")
15 df.info()
```

Exploração Inicial

Python

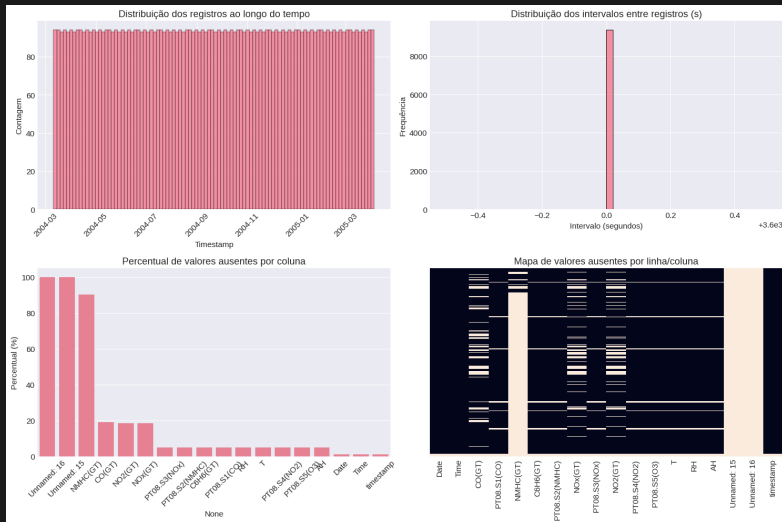
```
1 # Verificar range temporal
2 print("=== RANGE TEMPORAL ===")
3 print(f"Inicio: {df['timestamp'].min()}")
4 print(f"Fim: {df['timestamp'].max()}")
5 print(f"Duracao: {df['timestamp'].max() - df['timestamp'].min()}")
6
7 # Verificar frequencia
8 print("\n=== FREQUENCIA ===")
9 time_diff = df['timestamp'].diff().mode()[0]
10 print(f"Intervalo mais comum: {time_diff}")
```

Exploração Inicial

</> Python

```
1 # Valores faltantes
2 print("\n=== VALORES FALTANTES ===")
3 missing = df.isnull().sum()
4 missing_pct = (missing / len(df)) * 100
5 missing_df = pd.DataFrame({
6     'Missing': missing,
7     'Percent': missing_pct
8 }).sort_values('Percent', ascending=False)
9 print(missing_df[missing_df['Missing'] > 0])
```


Exploração Inicial (visual)



Fluxo de Análise de Séries Temporais

Processo sistemático:

FASE 1: Preparação

- ▶ Carregar dados
- ▶ Converter para DatetimeIndex
- ▶ Ordenar temporalmente
- ▶ Tratar missing values

FASE 2: Exploração

- ▶ Visualizar série completa
- ▶ Identificar padrões visuais
- ▶ Estatísticas descritivas
- ▶ Detectar anomalias

Fluxo de Análise de Séries Temporais

Processo sistemático:

FASE 3: Análise

- ▶ Decomposição
- ▶ Agregações temporais
- ▶ Médias móveis
- ▶ Correlações temporais

FASE 4: Feature Engineering

- ▶ Extrair componentes temporais
- ▶ Criar lags
- ▶ Calcular diferenças

Ferramentas e Bibliotecas

Manipulação:

- ▶ **Pandas:** Fundação principal
- ▶ **NumPy:** Operações numéricas
- ▶ **Matplotlib/Seaborn:** Visualizações

Análise avançada:

- ▶ **statsmodels:** Decomposição, ARIMA
- ▶ **scipy:** Análise de sinais
- ▶ **Prophet:** Forecasting por Facebook

Deep Learning:

- ▶ **TensorFlow/Keras:** LSTM, GRU
- ▶ **PyTorch:** Redes neurais recorrentes

Foco de hoje: Apenas Pandas + básicos de statsmodels

Tipos de Análise: Univariada vs Multivariada

Duas abordagens:

Univariada:

- ▶ Uma única série temporal
- ▶ Analisa apenas ela mesma
- ▶ Exemplo: apenas temperatura
- ▶ Previsão baseada em história própria

Multivariada:

- ▶ Múltiplas séries relacionadas
- ▶ Analisa relações entre elas
- ▶ Exemplo: temperatura + umidade + CO
- ▶ Uma pode ajudar a prever outra

Tipos de Análise: Univariada vs Multivariada

Air Quality:

- ▶ Dataset é multivariado
- ▶ Hoje: focaremos principalmente em univariadas
- ▶ Veremos correlações entre variáveis

Frequências Comuns

Diferentes granularidades temporais:

| Frequência | Código | Exemplo |
|--------------|--------------|------------------------|
| Anual | 'A' ou 'Y' | Vendas anuais |
| Trimestral | 'Q' | Relatórios financeiros |
| Mensal | 'M' | Faturamento mensal |
| Semanal | 'W' | Tráfego semanal |
| Diária | 'D' | Temperatura diária |
| Horária | 'H' | Consumo de energia |
| Minuto | 'T' ou 'min' | Dados de sensores |
| Segundo | 'S' | Trading de ações |
| Milissegundo | 'L' ou 'ms' | Dados científicos |

Frequências Comuns

Diferentes granularidades temporais:

Air Quality: Dados horários ('H')

Escolha depende:

- ▶ Natureza do fenômeno
- ▶ Disponibilidade de dados
- ▶ Objetivo da análise

Conceitos Importantes

Terminologia essencial:

Lag (Defasagem):

- ▶ Valor de t períodos atrás
- ▶ lag=1: ontem, lag=7: semana passada

Diferenciação:

- ▶ $Y_t - Y_{t-1}$
- ▶ Remove tendência
- ▶ Torna série estacionária

Conceitos Importantes

Terminologia essencial:

Autocorrelação:

- ▶ Correlação da série consigo mesma
- ▶ Em diferentes lags
- ▶ Mede dependência temporal

Estacionariedade:

- ▶ Propriedades estatísticas constantes
- ▶ Média e variância não mudam no tempo
- ▶ Importante para modelagem

Caso de Uso Motivador

Empresa de monitoramento ambiental:

Problema:

- ▶ Sensores geram dados horários
- ▶ 8760 observações por ano por sensor
- ▶ Precisa detectar padrões e anomalias
- ▶ Alertar quando qualidade do ar está ruim

Desafios:

- ▶ Volume massivo de dados
- ▶ Ruído nos dados
- ▶ Missing values ocasionais
- ▶ Múltiplas escalas (hora, dia, semana, mês)

Empresa de monitoramento ambiental:

Solução com séries temporais:

- ▶ Agregação temporal (horária → diária)
- ▶ Médias móveis para suavizar
- ▶ Decomposição para identificar tendências
- ▶ Detecção de anomalias via desvios

Expectativas Realistas

O que aprenderemos hoje:

SIM:

- ▶ ✓ Manipular séries temporais no Pandas
- ▶ ✓ Visualizar e explorar padrões
- ▶ ✓ Fazer agregações temporais
- ▶ ✓ Calcular médias móveis
- ▶ ✓ Decompor séries
- ▶ ✓ Tratar dados temporais

NÃO (ficará para o futuro):

- ▶ ✗ Modelagem preditiva (ARIMA, Prophet)
- ▶ ✗ Deep Learning para séries (LSTM)
- ▶ ✗ Análise estatística profunda
- ▶ ✗ Testes de estacionariedade formais

Nota Importante

Hoje é sobre MANIPULAÇÃO e EXPLORAÇÃO, não modelagem!

Bloco 1: DatetimeIndex

- ▶ Criar DatetimeIndex
- ▶ Acessar e filtrar por datas
- ▶ Extrair componentes temporais
- ▶ Operações básicas

Bloco 2: Agregações e Rolling

- ▶ Resample (mudar frequência)
- ▶ Rolling windows (médias móveis)
- ▶ Shifting e lags
- ▶ Diferenciação

Blocos 3 e 4: Decomposição + Práticas

- ▶ Decompor em componentes
- ▶ Visualizações avançadas
- ▶ Boas práticas
- ▶ Pipeline completo

Transição para Bloco 1

Fundação estabelecida!

Entendemos:

- ▶ O que são séries temporais
- ▶ Por que são especiais
- ▶ Componentes principais
- ▶ Dataset Air Quality
- ▶ Desafios comuns

Próximo: Bloco 1 - DatetimeIndex

- ▶ Coração da manipulação temporal no Pandas
- ▶ Converter strings para datas
- ▶ Definir índice temporal
- ▶ Acessar dados por período
- ▶ Extrair informações temporais

Bloco 1

DatetimeIndex e Manipulação Básica

O que é DatetimeIndex?

Índice especializado para séries temporais

Definição:

- ▶ Tipo especial de índice no Pandas
- ▶ Contém objetos datetime
- ▶ Otimizado para operações temporais
- ▶ Permite slicing por datas

Por que usar DatetimeIndex?

- ▶ Acesso intuitivo: `df['2024-01-01']`
- ▶ Funções temporais: `resample()`, `rolling()`
- ▶ Alinhamento automático por tempo
- ▶ Operações vetorizadas rápidas
- ▶ Componentes temporais: `.dt.year`, `.dt.month`

O que é DatetimeIndex?

Nota Importante

Sem DatetimeIndex, você não tem série temporal - só dados com datas!

DataFrame vs TimeSeries

DataFrame comum (índice numérico):

| Index | Date | Value |
|-------|------------|-------|
| 0 | 2024-01-01 | 10 |
| 1 | 2024-01-02 | 15 |
| 2 | 2024-01-03 | 12 |

TimeSeries (DatetimeIndex):

| Timestamp | Value |
|------------|-------|
| 2024-01-01 | 10 |
| 2024-01-02 | 15 |
| 2024-01-03 | 12 |

Diferença principal:

- ▶ TimeSeries: data É o índice
- ▶ Permite: `df['2024-01']` → todos de janeiro
- ▶ Habilita: `resample`, `rolling`, etc.

Criar DatetimeIndex: Método 1 - parse_dates

</> Python

```
1 # Metodo 1: parse_dates no pd.read_csv
2 # Quando data e hora estao em colunas separadas
3 # (ver slides anteriores)
4 print("=== TIPO DA COLUNA ===")
5 print(type(df['timestamp'][0]))
6 # <class 'pandas._libs.tslibs.timestamps.Timestamp'>
7 # Verificar
8 print(df['timestamp'].head())
9 # Ainda NÃO é indice! É apenas uma coluna
10 print(f"\nIndice atual: {type(df.index)}")
11 # Indice atual: <class 'pandas.core.indexes.range.RangeIndex'>
```

Criar DatetimeIndex: Método 2 - pd.to_datetime

Python

```
1 # Metodo 2: converter coluna existente
2 # Se ja tem coluna de texto com datas
3 df['timestamp'] = pd.to_datetime(df['timestamp_string'])
4
5 # Ou especificar formato (mais rapido!)
6 df['timestamp'] = pd.to_datetime(df['timestamp_string'],
7                                 format='%Y-%m-%d %H:%M:%S')
8 # Formatos comuns:
9 # '%Y-%m-%d' $ \rightarrow$ 2024-01-15
10 # '%d/%m/%Y' $ \rightarrow$ 15/01/2024
11 # '%Y-%m-%d %H:%M:%S' $ \rightarrow$ 2024-01-15 10:30:00
12 print("=== CONVERSAO ===")
13 print(df['timestamp'].dtype)
14 # datetime64[ns]
15 print(df['timestamp'].head())
```

Definir DatetimeIndex

</> Python

```
1 # Transformar coluna em indice
2
3 # Opcao 1: set_index
4 df = df.set_index('timestamp')
5
6 print("=== AGORA É DATETIMEINDEX ===")
7 print(type(df.index))
8 # <class 'pandas.core.indexes.datetimes.DatetimeIndex'>
9
10 print(f"\nIndice:\n{df.index}")
```


Definir DatetimeIndex (cont.)

</> Python

```
1 # Opcao 2: Direto no read_csv
2 df = pd.read_csv('data.csv',
3                 parse_dates=['timestamp'],
4                 index_col='timestamp')
5
6 # Verificar propriedades do indice
7 print(f"\nFrequencia: {df.index.freq}")
8 print(f"Inicio: {df.index.min()}")
9 print(f"Fim: {df.index.max()}")
10 # Frequencia: None
11 # Inicio: 2004-03-10 18:00:00
12 # Fim: 2005-04-04 14:00:00
```

Verificar e Ordenar DatetimeIndex

</> Python

```
1 # SEMPRE verificar se esta ordenado
2 print("=== ORDENACAO ===")
3 print(f"Está ordenado? {df.index.is_monotonic_increasing}")
4 # === ORDENACAO ===
5 # Está ordenado? False
6
7 # Se nao estiver, ordenar
8 df = df.sort_index()
9 # Verificar novamente
10 print(f"Agora esta ordenado? {df.index.is_monotonic_increasing}")
11 # Por que isso importa?
12 # - resample() e rolling() assumem ordem temporal
13 # - Slicing por datas precisa de ordem
14 # - Performance melhor
```

Verificar e Ordenar DatetimeIndex (cont.)

</> Python

```
1 # Checar duplicatas
2 print(df.index.duplicated().sum())
3 # Remover duplicatas de timestamp
4 df = df[~df.index.duplicated(keep='first')]
5 print(f"\nShape apos remover duplicatas: {df.shape}")
6 # Checar NaNs
7 print(df.index.hasnans)
8 # Remover NaNs
9 df = df.loc[df.index.dropna()]
```

Propriedades do DatetimeIndex

Range e frequência:

- ▶ `.min()`, `.max()`: Limites
- ▶ `.freq`: Frequência (se regular)
- ▶ `.inferred_freq`: Frequência inferida

Propriedades:

- ▶ `.is_monotonic_increasing`: Ordenado?
- ▶ `.is_unique`: Sem duplicatas?
- ▶ `.has_duplicates`: Tem duplicatas?

Componentes:

- ▶ `.year`, `.month`, `.day`
- ▶ `.hour`, `.minute`, `.second`
- ▶ `.dayofweek`, `.dayofyear`
- ▶ `.quarter`, `.weekofyear`

Acessar Componentes Temporais

</> Python

```
1 # Extrair componentes do DatetimeIndex
2
3 # Do índice direto
4 print("=== COMPONENTES DO INDICE ===")
5 print(f"Anos unicos: {df.index.year.unique()}")
6 print(f"Meses unicos: {df.index.month.unique()}")
7
8 # === COMPONENTES DO INDICE ===
9 # Anos unicos: Index([2004, 2005], dtype='int32', name='timestamp')
10 # Meses unicos: Index([3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 1, 2], dtype='
    int32', name='timestamp')
```

Acessar Componentes Temporais (cont.)

Python

```
1 # Criar colunas com componentes
2 df['year'] = df.index.year
3 df['month'] = df.index.month
4 df['day'] = df.index.day
5 df['hour'] = df.index.hour
6 df['dayofweek'] = df.index.dayofweek # 0=Segunda, 6=Domingo
7 df['quarter'] = df.index.quarter
8
9 print("\n=== NOVAS COLUNAS ===")
10 print(df[['year', 'month', 'day', 'hour', 'dayofweek']].head())
11
12 # Nome do dia da semana
13 df['weekday_name'] = df.index.day_name()
14 print(f"\nDias da semana: {df['weekday_name'].unique()}")
```

Acessar Componentes Temporais (cont.)

</> Python

```
1 # === NOVAS COLUNAS ===
2 #                               year  month  day  hour  dayofweek
3 # timestamp
4 # 2004-03-10 18:00:00  2004      3   10   18      2
5 # 2004-03-10 19:00:00  2004      3   10   19      2
6 # 2004-03-10 20:00:00  2004      3   10   20      2
7 # 2004-03-10 21:00:00  2004      3   10   21      2
8 # 2004-03-10 22:00:00  2004      3   10   22      2
9 #
10 # Dias da semana: ['Wednesday' 'Thursday' 'Friday' 'Saturday' 'Sunday' '
    Monday' 'Tuesday']
```

Slicing por Datas: Básico

</> Python

```
1 # Um dia específico
2 jan_15 = df.loc['2004-03-15']
3 print(f"Dados de 15/03/2004: {len(jan_15)} registros")
4 # Um mês inteiro
5 march_2004 = df.loc['2004-03']
6 print(f"Dados de março/2004: {len(march_2004)} registros")
7 # Um ano inteiro
8 year_2004 = df.loc['2004']
9 print(f"Dados de 2004: {len(year_2004)} registros")
10 # Range de datas
11 spring = df.loc['2004-03':'2004-05']
12 print(f"Primavera 2004: {len(spring)} registros")
13 # Até certa data
14 early = df.loc[:'2004-06']
15 print(f"Até junho/2004: {len(early)} registros")
```


Slicing por Datas: Básico (cont.)

</> Python

```
1 # Dados de 15/03/2004: 24 registros
2 # Dados de marco/2004: 510 registros
3 # Dados de 2004: 7110 registros
4 # Primavera 2004: 1974 registros
5 # Até junho/2004: 2694 registros
```

Slicing por Datas: Avançado

</> Python

```
1 # Horário específico em todos os dias
2 morning_10am = df[df.index.hour == 10]
3 print(f"Todas as 10h da manhã: {len(morning_10am)} registros")
4 # Fins de semana (sabado=5, domingo=6)
5 weekends = df[df.index.dayofweek >= 5]
6 print(f"Fins de semana: {len(weekends)} registros")
7 # Dias uteis
8 weekdays = df[df.index.dayofweek < 5]
9 print(f"Dias uteis: {len(weekdays)} registros")
10 # Horário de pico (7-9h)
11 rush_hour = df[df.index.hour.isin([7, 8, 9])]
12 print(f"Horário de pico: {len(rush_hour)} registros")
13 # Dezembro (mes de festas)
14 december = df[df.index.month == 12]
15 print(f"Dezembro: {len(december)} registros")
```

Slicing por Datas: Avançado (cont.)

</> Python

```
1 # Todas as 10h da manha: 390 registros
2 # Fins de semana: 2688 registros
3 # Dias uteis: 6669 registros
4 # Horario de pico: 1170 registros
5 # Dezembro: 744 registros
```

Filtrar por Múltiplas Condições

</> Python

```
1 # Dias uteis + horario comercial (9-18h)
2 business_hours = df[
3     (df.index.dayofweek < 5) &   # Segunda a sexta
4     (df.index.hour >= 9) &       # Depois das 9h
5     (df.index.hour < 18)         # Antes das 18h
6 ]
7 print(f"Horario comercial: {len(business_hours)} registros")
8 # Verao (junho-agosto) + dias quentes
9 summer_hot = df[
10     (df.index.month.isin([6, 7, 8])) &
11     (df['T'] > 25) # Temperatura > 25C
12 ]
```

Filtrar por Múltiplas Condições

</> Python

```
1 print(f"Dias quentes de verao: {len(summer_hot)} registros")
2 # Madrugada (0-6h) em fins de semana
3 weekend_nights = df[
4     (df.index.dayofweek >= 5) &
5     (df.index.hour < 6)
6 ]
7 print(f"Madrugadas de fim de semana: {len(weekend_nights)}")
```

Operações com Timestamps

</> Python

```
1 from datetime import timedelta
2 # Adicionar/subtrair tempo
3 tomorrow = df.index + timedelta(days=1)
4 print(f"Amanha: {tomorrow[:5]}")
5 week_ago = df.index - timedelta(days=7)
6 print(f"Semana passada: {week_ago[:5]}")
7 # Diferença entre timestamps
8 time_diff = df.index.to_series().diff()
9 print(f"\nDiferença de tempo:")
10 print(time_diff.head())
11 # Tempo decorrido desde início
12 elapsed = df.index - df.index.min()
13 elapsed_hours = elapsed.total_seconds() / 3600
14 print(f"\nHoras desde início:")
15 print(f"Min: {elapsed_hours.min():.1f}h")
16 print(f"Max: {elapsed_hours.max():.1f}h")
```

Criar DatetimeIndex do Zero

</> Python

```
1 # Diario
2 daily = pd.date_range(start='2024-01-01',
3                        end='2024-12-31',
4                        freq='D')
5 print(f"Diario: {len(daily)} dias")
6 # Horário
7 hourly = pd.date_range(start='2024-01-01',
8                        periods=24*7, # 1 semana
9                        freq='H')
10 print(f"Horario: {len(hourly)} horas")
```

Criar DatetimeIndex do Zero (cont.)

</> Python

```
1 # Mensal (fim do mes)
2 monthly = pd.date_range(start='2024-01-01',
3                           periods=12,
4                           freq='M')
5 print(f"Mensal: {len(monthly)} meses")
6 # Criar DataFrame com indice temporal
7 df_new = pd.DataFrame({
8     'value': np.random.randn(len(hourly))
9 }, index=hourly)
```


Frequências do Pandas

| Código | Descrição | Exemplo | Alias |
|--------|---------------|------------|--------------|
| D | Dia | 2024-01-01 | day |
| H | Hora | 10:00 | hour |
| T, min | Minuto | 10:30 | minute |
| S | Segundo | 10:30:45 | second |
| L, ms | Milissegundo | - | - |
| W | Semana | Domingos | week |
| M | Fim de mês | 2024-01-31 | month |
| MS | Início de mês | 2024-01-01 | - |
| Q | Trimestre | 2024-03-31 | quarter |
| A, Y | Ano | 2024-12-31 | year |
| B | Dias úteis | Seg-Sex | business day |

Modificadores:

- ▶ '2H' = a cada 2 horas
- ▶ '15T' = a cada 15 minutos
- ▶ '3D' = a cada 3 dias

Detectar e Preencher Gaps Temporais

</> Python

```
1 # Identificar lacunas na serie temporal
2
3 # Frequencia esperada
4 expected_freq = '1H' # Horaria
5
6 # Criar range completo esperado
7 full_range = pd.date_range(start=df.index.min(),
8                             end=df.index.max(),
9                             freq=expected_freq)
10
11 print(f"Timestamps esperados: {len(full_range)}")
12 print(f"Timestamps reais: {len(df)}")
13 print(f"Missing: {len(full_range) - len(df)}")
```

Detectar e Preencher Gaps Temporais (cont.)

</> Python

```
1 # Reindexar para preencher gaps
2 df_complete = df.reindex(full_range)
3
4 print(f"\nApos reindex: {df_complete.shape}")
5 print(f"NaN introduzidos: {df_complete.isnull().sum().sum()}")
6
7 # Agora podemos tratar missing adequadamente
```

Tratar Missing Values: Forward Fill

</> Python

```
1 # Forward fill: propagar ultimo valor válido
2
3 # Metodo 1: ffill (forward fill)
4 df_filled = df_complete.fillna(method='ffill')
5
6 print("=== FORWARD FILL ===")
7 print(f"Missing antes: {df_complete['CO(GT)'].isnull().sum()}")
8 print(f"Missing depois: {df_filled['CO(GT)'].isnull().sum()}")
9
10 # Limitar quantos valores propagar
11 df_filled = df_complete.fillna(method='ffill', limit=3)
12 # Propaga no maximo 3 valores
```

Tratar Missing Values: Forward Fill (cont.)

Python

```
1 # Exemplo visual
2 print("\nAntes:")
3 print(df_complete['CO(GT)'][100:110])
4
5 print("\nDepois:")
6 print(df_filled['CO(GT)'][100:110])
7
8 # Util quando: valor tende a permanecer constante
```

Tratar Missing Values: Interpolação

</> Python

```
1 # Interpolacao linear (padrao)
2 df_interp = df_complete.interpolate(method='linear')
3
4 print("=== INTERPOLACAO ===")
5 print(f"Missing antes: {df_complete['CO(GT)'].isnull().sum()}")
6 print(f"Missing depois: {df_interp['CO(GT)'].isnull().sum()}")
7
8 # Outros metodos
9 # method='time': considera distancia temporal real
10 df_interp = df_complete.interpolate(method='time')
11
12 # method='spline': curva suave
13 df_interp = df_complete.interpolate(method='spline', order=2)
```

Tratar Missing Values: Interpolação (cont.)

</> Python

```
1 # Comparar
2 print("\nOriginal (com gaps):")
3 print(df_complete['CO(GT)'][100:110])
4 print("\nInterpolado:")
5 print(df_interp['CO(GT)'][100:110])
```

Escolher Método de Imputação Temporal

Guia de decisão:

Forward Fill (ffill):

- ▶ Use: valores mudam raramente
- ▶ Exemplo: status de sensor (ligado/desligado)
- ▶ Exemplo: categoria (tipo de dia)
- ▶ Limite: poucos gaps pequenos

Backward Fill (bfill):

- ▶ Use: raramente (só em casos especiais)
- ▶ Exemplo: preencher início da série

Escolher Método de Imputação Temporal

Guia de decisão:

Interpolação Linear:

- ▶ Use: valores contínuos que mudam suavemente
- ▶ Exemplo: temperatura, umidade, poluição
- ▶ Melhor para séries temporais típicas

NUNCA usar:

- ▶ Média global (ignora contexto temporal!)

Estatísticas por Período

</> Python

```
1 # Por hora do dia
2 hourly_stats = df.groupby(df.index.hour)['CO(GT)'].agg([
3     'mean', 'std', 'min', 'max', 'count'
4 ])
5
6 print("=== ESTATISTICAS POR HORA ===")
7 print(hourly_stats)
8
9 # Por dia da semana
10 weekday_stats = df.groupby(df.index.dayofweek)['CO(GT)'].mean()
11 weekday_stats.index = ['Seg', 'Ter', 'Qua', 'Qui', 'Sex', 'Sab', 'Dom']
12
13 print("\n=== MEDIA POR DIA DA SEMANA ===")
14 print(weekday_stats)
```

Estatísticas por Período (cont.)

</> Python

```
1 # Por mes
2 monthly_stats = df.groupby(df.index.month)['CO(GT)'].mean()
3 monthly_stats.index = ['Jan', 'Fev', 'Mar', 'Abr', 'Mai', 'Jun',
4                        'Jul', 'Ago', 'Set', 'Out', 'Nov', 'Dez']
5
6 print("\n=== MEDIA POR MES ===")
7 print(monthly_stats)
```

Comparar Períodos

</> Python

```
1 # Comparar dias uteis vs fins de semana
2 weekday_co = df[df.index.dayofweek < 5]['CO(GT)'].mean()
3 weekend_co = df[df.index.dayofweek >= 5]['CO(GT)'].mean()
4
5 print("=== DIAS UTEIS VS FINS DE SEMANA ===")
6 print(f"Dias uteis: {weekday_co:.2f} mg/m³")
7 print(f"Fins de semana: {weekend_co:.2f} mg/m³")
8 print(f"Diferenca: {weekday_co - weekend_co:.2f} mg/m³")
9 print(f"Variacao: {((weekday_co/weekend_co - 1) * 100):.1f}%")
```

Comparar Períodos

</> Python

```
1 # Comparar horarios de pico vs nao-pico
2 peak_hours = df[df.index.hour.isin([7,8,9,17,18,19])]['CO(GT)'].mean()
3 off_peak = df[~df.index.hour.isin([7,8,9,17,18,19])]['CO(GT)'].mean()
4
5 print("\n=== PICO VS NAO-PICO ===")
6 print(f"Horario de pico: {peak_hours:.2f} mg/m³")
7 print(f"Fora de pico: {off_peak:.2f} mg/m³")
```

Criar Features Temporais

</> Python

```
1 # 1. Parte do dia
2 def get_part_of_day(hour):
3     if 6 <= hour < 12:
4         return 'Morning'
5     elif 12 <= hour < 18:
6         return 'Afternoon'
7     elif 18 <= hour < 22:
8         return 'Evening'
9     else:
10         return 'Night'
11
12 df['part_of_day'] = df.index.hour.map(get_part_of_day)
```

Criar Features Temporais

Python

```
1 # 2. Fim de semana
2 df['is_weekend'] = (df.index.dayofweek >= 5).astype(int)
3
4 # 3. Horário de pico
5 df['is_rush_hour'] = df.index.hour.isin([7,8,9,17,18,19]).astype(int)
6
7 # 4. Estacao do ano (hemisferio norte)
8 df['season'] = df.index.month % 12 // 3 + 1
9 season_map = {1: 'Winter', 2: 'Spring', 3: 'Summer', 4: 'Fall'}
10 df['season_name'] = df['season'].map(season_map)
11
12 print(df[['part_of_day', 'is_weekend', 'season_name']].head(10))
```

Features Cíclicas

Python

```
1 # Representar componentes temporais como cíclicas
2 # (preserva a natureza circular do tempo em análises!)
3
4 import numpy as np
5
6 # Hora do dia (0-23) como seno/cosseno
7 df['hour_sin'] = np.sin(2 * np.pi * df.index.hour / 24)
8 df['hour_cos'] = np.cos(2 * np.pi * df.index.hour / 24)
9
10 # Dia da semana (0-6) como seno/cosseno
11 df['dow_sin'] = np.sin(2 * np.pi * df.index.dayofweek / 7)
12 df['dow_cos'] = np.cos(2 * np.pi * df.index.dayofweek / 7)
```


Features Cíclicas

</> Python

```
1 # Mes do ano (1-12) como seno/cosseno
2 df['month_sin'] = np.sin(2 * np.pi * df.index.month / 12)
3 df['month_cos'] = np.cos(2 * np.pi * df.index.month / 12)
4
5 print("=== FEATURES CICLICAS ===")
6 print(df[['hour_sin', 'hour_cos', 'dow_sin', 'month_sin']].head())
7
8 # Por que?
9 # - Preserva distancia real: 23h e 0h estao proximas
10 # - Dezembro (12) e Janeiro (1) sao continuos
11 # - Util para correlacoes e visualizacoes temporais
```

Por que Features Cíclicas?

O problema com encoding ordinal:

Hora do dia (0-23):

- ▶ Ordinal: 0, 1, 2, ..., 22, 23
- ▶ Problema: 23h e 0h codificadas como muito distantes!
- ▶ Realidade: 23h e 0h são 1 hora de diferença

Solução cíclica:

- ▶ $\sin(2\pi \times \text{hour}/24)$ e $\cos(2\pi \times \text{hour}/24)$
- ▶ 23h: $\sin \approx -0.26$, $\cos \approx 0.97$
- ▶ 0h: $\sin \approx 0.0$, $\cos \approx 1.0$
- ▶ Distância euclidiana pequena! ✓

Por que Features Cíclicas?

Quando usar:

- ▶ Hora do dia, dia da semana, mês
- ▶ Qualquer variável cíclica/periódica
- ▶ Análises de correlação temporal
- ▶ Cálculos de distância e similaridade

💡 Nota Importante

Representação cíclica é essencial para análises que dependem de distâncias/similaridades temporais!

Timezone Awareness

</> Python

```
1 # Trabalhar com fusos horarios
2
3 # Verificar se tem timezone
4 print(f"Timezone aware? {df.index.tz is not None}")
5
6 # Adicionar timezone (se nao tiver)
7 df_tz = df.copy()
8 df_tz.index = df_tz.index.tz_localize('UTC')
9
10 print(f"Agora é timezone aware: {df_tz.index.tz}")
11
12 # Converter para outro fuso
13 df_sp = df_tz.copy()
14 df_sp.index = df_sp.index.tz_convert('America/Sao_Paulo')
```

Timezone Awareness (cont.)

</> Python

```
1 print(f"\nUTC: {df_tz.index[0]}")
2 print(f"Sao Paulo: {df_sp.index[0]}")
3
4 # Importante para:
5 # - Dados de multiplas localizacoes
6 # - Horário de verao
7 # - Coordenacao entre sistemas
```

Encontrar Valores Extremos

</> Python

```
1 # Maior valor de CO
2 max_co_idx = df['CO(GT)'].idxmax()
3 max_co_val = df['CO(GT)'].max()
4
5 print("=== PIOR QUALIDADE DO AR (CO) ===")
6 print(f"Quando: {max_co_idx}")
7 print(f"Valor: {max_co_val:.2f} mg/m³")
8 print(f"Hora: {max_co_idx.hour}h")
9 print(f"Dia da semana: {max_co_idx.day_name()}")
```

Encontrar Valores Extremos

</> Python

```
1 # Melhor valor
2 min_co_idx = df['CO(GT)'].idxmin()
3 min_co_val = df['CO(GT)'].min()
4
5 print("\n=== MELHOR QUALIDADE DO AR (CO) ===")
6 print(f"Quando: {min_co_idx}")
7 print(f"Valor: {min_co_val:.2f} mg/m³")
8
9 # Top 10 piores momentos
10 worst_10 = df.nlargest(10, 'CO(GT)')
11 print("\n=== TOP 10 PIORES ===")
12 print(worst_10[['CO(GT)', 'T', 'RH']].to_string())
```

Exercício: Análise Temporal Básica

Exercício Prático

Explore o Air Quality dataset temporalmente:

1. Configure DatetimeIndex:

```
1 # - Carregar com parse_dates
2 # - Definir como indice
3 # - Ordenar e verificar
```

2. Extraia componentes temporais:

- ▶ year, month, day, hour, dayofweek
- ▶ Crie coluna 'season'

Exercício: Análise Temporal Básica

Exercício Prático

Explore o Air Quality dataset temporalmente (cont.):

3. Análises:

- ▶ Média de CO por hora do dia
- ▶ Média de temperatura por mês
- ▶ Comparar dias úteis vs fim de semana

4. Identifique:

- ▶ 10 momentos com pior qualidade do ar
- ▶ Padrões (quando ocorrem?)

Bloco 2

Resample, Rolling Windows e Agregações

O que é Resample?

Resample = mudar a frequência temporal dos dados

Dois tipos principais:

1. Downsampling (reduzir frequência):

- ▶ Horária → Diária
- ▶ Diária → Mensal
- ▶ Minuto → Hora
- ▶ Precisa agregar (média, soma, max, etc.)

2. Upsampling (aumentar frequência):

- ▶ Diária → Horária
- ▶ Mensal → Diária
- ▶ Hora → Minuto
- ▶ Cria gaps (precisa preencher)

O que é Resample?

Resample = mudar a frequência temporal dos dados

Por que fazer?

- ▶ Reduzir ruído (*downsampling*)
- ▶ Alinhar séries com frequências diferentes
- ▶ Visualizar tendências
- ▶ Reduzir volume de dados

Downsampling: Conceito

Exemplo: Horária → Diária

Original (horária - 24 valores por dia):

| Timestamp | CO |
|------------------|-----|
| 2004-03-10 00:00 | 2.6 |
| 2004-03-10 01:00 | 2.0 |
| ... | ... |
| 2004-03-10 23:00 | 2.3 |

Após resample diário (média):

| Date | CO |
|------------|-----|
| 2004-03-10 | 2.2 |

24 valores → 1 valor (média dos 24)

Resample: Sintaxe Básica

</> Python

```
1 # Sintaxe geral
2 # df_resampled = df.resample('frequencia').funcao_agregacao()
3
4 # Exemplo 1: Horaria -> Diaria (media)
5 daily_avg = df['CO(GT)'].resample('D').mean()
6
7 print("=== RESAMPLE DIARIO ===")
8 print(f"Original: {len(df)} registros (horarios)")
9 print(f"Resampled: {len(daily_avg)} registros (diarios)")
10 print(f"\nPrimeiros valores:")
11 print(daily_avg.head())
```

Resample: Sintaxe Básica

</> Python

```
1 # Exemplo 2: Horaria -> Diaria (max)
2 daily_max = df['CO(GT)'].resample('D').max()
3
4 print("\n=== MAX DIARIO ===")
5 print(daily_max.head())
6
7 # Frequencias: 'D'=dia, 'W'=semana, 'M'=mes, 'Q'=trimestre, 'A'=ano
```

Resample: Múltiplas Agregações

</> Python

```
1 # Uma coluna, multiplas funcoes
2 daily_stats = df['CO(GT)'].resample('D').agg([
3     'mean', 'std', 'min', 'max', 'count'
4 ])
5
6 print("=== ESTATISTICAS DIARIAS DE CO ===")
7 print(daily_stats.head())
8
9 # Multiplas colunas, diferentes funcoes
10 daily_multi = df.resample('D').agg({
11     'CO(GT)': ['mean', 'max'],
12     'NO2(GT)': ['mean', 'max'],
13     'T': ['mean', 'min', 'max'],
14     'RH': 'mean'
15 })
```


Resample: Múltiplas Agregações (cont.)

</> Python

```
1 print("\n=== MULTIPLAS COLUNAS ===")
2 print(daily_multi.head())
3
4 # Renomear colunas se necessario
5 daily_multi.columns = ['_'.join(col).strip() for col in daily_multi.
6                        columns]
```

Resample para Diferentes Períodos

</> Python

```
1 # Diário (media)
2 daily = df['CO(GT)'].resample('D').mean()
3 print(f"Diário: {len(daily)} registros")
4 # Semanal (media)
5 weekly = df['CO(GT)'].resample('W').mean()
6 print(f"Semanal: {len(weekly)} registros")
7 # Mensal (media)
8 monthly = df['CO(GT)'].resample('M').mean()
9 print(f"Mensal: {len(monthly)} registros")
10 # A cada 6 horas
11 six_hourly = df['CO(GT)'].resample('6H').mean()
12 print(f"6 em 6 horas: {len(six_hourly)} registros")
13 # A cada 3 dias
14 three_daily = df['CO(GT)'].resample('3D').mean()
15 print(f"A cada 3 dias: {len(three_daily)} registros")
```

Escolher Função de Agregação

Qual função usar ao fazer downsample?

mean() - Média:

- ▶ Use: valores contínuos, visão geral
- ▶ Exemplo: temperatura média diária

sum() - Soma:

- ▶ Use: totais, acumulados
- ▶ Exemplo: vendas totais do mês

max() / min():

- ▶ Use: valores extremos
- ▶ Exemplo: pico de poluição, temperatura mínima

Escolher Função de Agregação

Qual função usar ao fazer downsample?

median() - Mediana:

- ▶ Use: quando há outliers
- ▶ Mais robusta que média

count():

- ▶ Use: número de observações válidas
- ▶ Detectar missing values

Upsampling: Aumentar Frequência

</> Python

```
1 # Criar dados mensais primeiro (downsample)
2 monthly = df['CO(GT)'].resample('M').mean()
3
4 print("=== DADOS MENSAIS ===")
5 print(f"Registros: {len(monthly)}")
6 print(monthly.head())
7
8 # Agora upsampling: mensal -> diario
9 daily_from_monthly = monthly.resample('D').asfreq()
10
11 print("\n=== APOS UPSAMPLING ===")
12 print(f"Registros: {len(daily_from_monthly)}")
13 print(daily_from_monthly.head(10))
```

Upsampling: Aumentar Frequência (cont.)

</> Python

```
1 # Nota: Crie NaN! Precisa preencher
2
3 # Preencher com forward fill
4 daily_filled = monthly.resample('D').ffill()
5
6 print("\n=== COM FORWARD FILL ===")
7 print(daily_filled.head(10))
```

Upsampling com Interpolação

Python

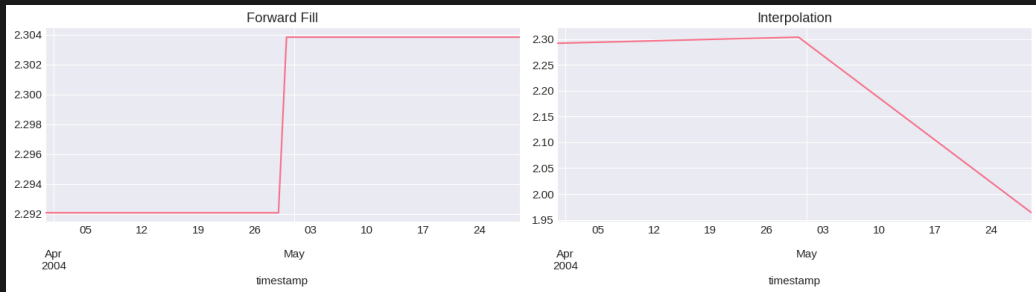
```
1 # Upsampling com interpolacao (melhor!)
2
3 # Mensal -> Diario com interpolacao
4 daily_interp = monthly.resample('D').interpolate(method='linear')
5
6 print("=== UPSAMPLING COM INTERPOLACAO ===")
7 print(daily_interp.head(40)) # Primeiro mes
8
9 # Comparar metodos
10 import matplotlib.pyplot as plt
11
12 fig, axes = plt.subplots(1, 2, figsize=(14, 4))
```

Upsampling com Interpolação (cont.)

Python

```
1 # Forward fill
2 monthly.resample('D').ffill().head(60).plot(
3     ax=axes[0], title='Forward Fill'
4 )
5
6 # Interpolacao
7 monthly.resample('D').interpolate().head(60).plot(
8     ax=axes[1], title='Interpolation'
9 )
10
11 plt.tight_layout()
12 plt.show()
```


Upsampling com Interpolação (cont.)



Rolling Windows: Conceito

Rolling = janela deslizante que calcula estatística

Como funciona:

[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

Janela de tamanho 3:

[1, 2, 3] → média = 2.0

[2, 3, 4] → média = 3.0

[3, 4, 5] → média = 4.0

...

Resultado:

[NaN, NaN, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0, 9.0]

Rolling Windows: Conceito

Rolling = janela deslizante que calcula estatística

Usos:

- ▶ Suavizar ruído (média móvel)
- ▶ Detectar tendências
- ▶ Calcular volatilidade (desvio móvel)
- ▶ Features para ML

Rolling: Média Móvel Simples

</> Python

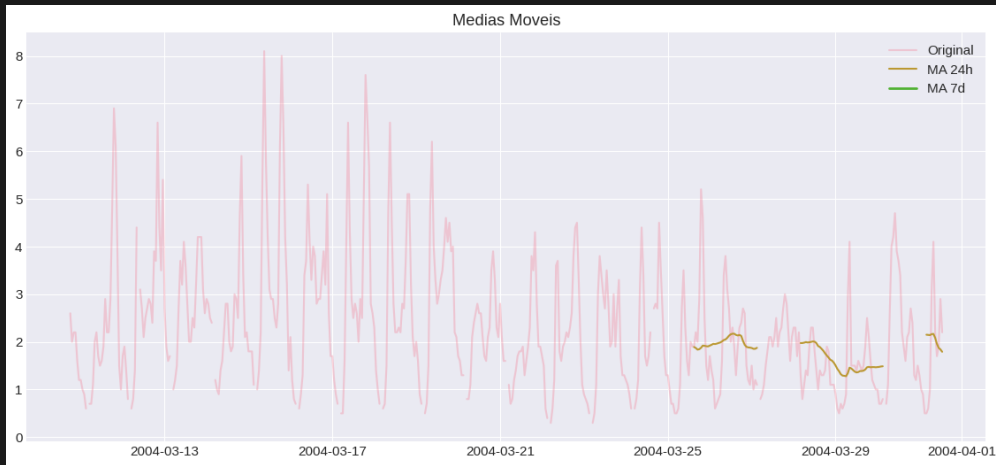
```
1 # Janela de 24 horas (suavizar ruído diário)
2 df['CO_MA24'] = df['CO(GT)'].rolling(window=24).mean()
3
4 print("=== MEDIA MOVEL 24H ===")
5 print(df[['CO(GT)', 'CO_MA24']].head(30))
6
7 # Primeiras 23 serão NaN (janela incompleta)
8
9 # Janela de 7 dias (168 horas)
10 df['CO_MA168'] = df['CO(GT)'].rolling(window=168).mean()
```

Rolling: Média Móvel Simples (cont.)

</> Python

```
1 # Visualizar
2 import matplotlib.pyplot as plt
3
4 plt.figure(figsize=(14, 6))
5 plt.plot(df.index[:500], df['CO(GT)'][:500],
6          alpha=0.3, label='Original')
7 plt.plot(df.index[:500], df['CO_MA24'][:500],
8          label='MA 24h')
9 plt.plot(df.index[:500], df['CO_MA168'][:500],
10          label='MA 7d', linewidth=2)
11 plt.legend()
12 plt.title('Medias Moveis')
13 plt.show()
```

Rolling: Média Móvel Simples (cont.)



Problemas no gráfico: dados faltants (NaN)

Original:

- ▶ Não possui NaN iniciais - série completa desde o primeiro timestamp
- ▶ Pode ter gaps esparsos devido a falhas no sensor (-200 convertido para NaN)

MA 24h:

- ▶ Primeiras 23 observações são NaN (janela incompleta)
- ▶ Requer 24 valores consecutivos para calcular a primeira média
- ▶ Visível apenas a partir da 24^a hora do dataset

MA 7d:

- ▶ Primeiras 167 observações são NaN (janela incompleta)
- ▶ Requer 168 valores (7 dias × 24 horas) para calcular a primeira média
- ▶ Visível apenas após 7 dias completos de dados
- ▶ Por isso não aparece no gráfico que mostra apenas [:500] (primeiras 21 dias)

Solução: usar `min_periods` (próximos slides)

Escolher Tamanho da Janela

Window size (tamanho da janela):

Janela pequena:

- ▶ Mais sensível a mudanças
- ▶ Segue dados mais de perto
- ▶ Mais ruído
- ▶ Exemplo: window=5

Janela grande:

- ▶ Mais suave
- ▶ Menos sensível
- ▶ Menos ruído
- ▶ Mais lag (atraso)
- ▶ Exemplo: window=100

Escolher Tamanho da Janela

Window size (tamanho da janela):

Regras práticas:

- ▶ Dados horários: 24 (1 dia) ou 168 (1 semana)
- ▶ Dados diários: 7 (1 semana) ou 30 (1 mês)
- ▶ Dados mensais: 12 (1 ano)
- ▶ Experimente! Visualize diferentes tamanhos

Rolling: Outras Estatísticas

Python

```
1 # Desvio padrao movel (volatilidade)
2 df['CO_std24'] = df['CO(GT)'].rolling(window=24).std()
3
4 # Maximo movel
5 df['CO_max24'] = df['CO(GT)'].rolling(window=24).max()
6
7 # Minimo movel
8 df['CO_min24'] = df['CO(GT)'].rolling(window=24).min()
9
10 # Mediana movel (mais robusta)
11 df['CO_median24'] = df['CO(GT)'].rolling(window=24).median()
12
13 # Soma movel
14 df['CO_sum24'] = df['CO(GT)'].rolling(window=24).sum()
```

Rolling: Outras Estatísticas (cont.)

</> Python

```
1 print("=== ROLLING STATISTICS ===")
2 print(df[['CO(GT)', 'CO_MA24', 'CO_std24',
3           'CO_max24', 'CO_min24']].tail(10))
4
5 # Interpretacao:
6 # - std alta = alta variabilidade
7 # - max/min = range de variacao
```

Rolling: center=True

</> Python

```
1 # Centrar a janela (ao inves de olhar pra tras)
2
3 # Padrao: olha para tras
4 df['CO_MA_backward'] = df['CO(GT)'].rolling(window=24).mean()
5
6 # Centrado: usa valores antes E depois
7 df['CO_MA_centered'] = df['CO(GT)'].rolling(window=24, center=True).mean()
8
9 print("=== BACKWARD VS CENTERED ===")
10 print(df[['CO(GT)', 'CO_MA_backward', 'CO_MA_centered']].iloc[20:30])
```

Rolling: center=True (cont.)

</> Python

```
1 # Quando usar cada um?
2
3 # Backward (padrao): PREVISAO
4 # - Usa apenas passado
5 # - Pode usar em producao
6 # - Mais comum
7
8 # Centered: ANALISE/SUAVIZACAO
9 # - Usa passado + futuro
10 # - NAO pode usar para previsao!
11 # - Suavizacao mais precisa
```

Rolling: min_periods

</> Python

```
1 # Define quantos valores válidos (não-NaN) são necessários no mínimo para
  # calcular a estatística da janela móvel.
2
3 # Padrao: precisa janela completa
4 df['CO_MA24_strict'] = df['CO(GT)'].rolling(window=24).mean()
5 # Primeiras 23 = NaN
6
7 # Com min_periods: aceita janela parcial
8 df['CO_MA24_partial'] = df['CO(GT)'].rolling(
9     window=24,
10    min_periods=1 # Minimo de 1 valor valido
11 ).mean()
12 # Todas tem valores!
13 print("=== MIN_PERIODS ===")
14 print(df[['CO(GT)', 'CO_MA24_strict', 'CO_MA24_partial']].head(25))
```

Rolling: min_periods (cont.)

</> Python

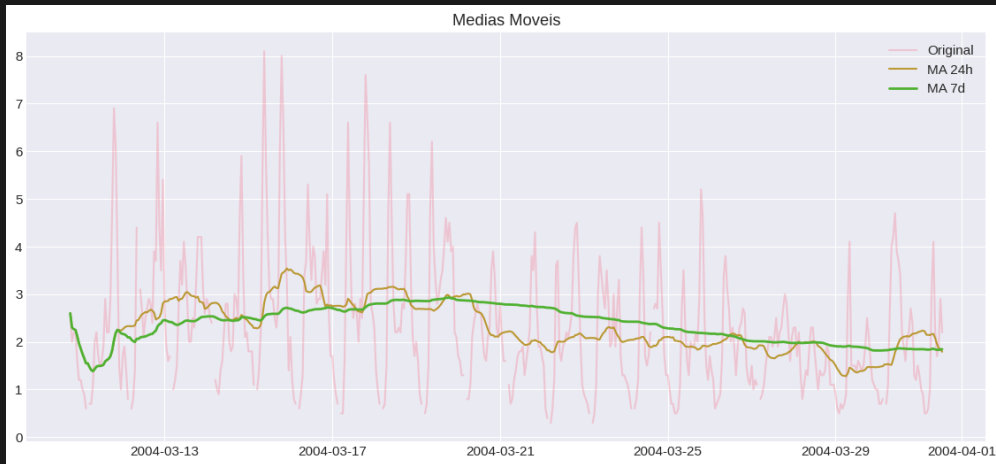
```
1 # Quando usar min_periods?  
2 # - Para evitar NaN no inicio  
3 # - Quando poucos dados faltantes OK  
4 # - Cuidado: janelas pequenas no inicio sao menos confiaveis
```

Rolling: Média Móvel Simples (cont.)

Python

```
1 # Janela de 24 horas com min_periods
2 df['CO_MA24'] = df['CO(GT)'].rolling(window=24, min_periods=1).mean()
3
4 # Janela de 7 dias (168 horas) com min_periods
5 df['CO_MA168'] = df['CO(GT)'].rolling(window=168, min_periods=1).mean()
6
7 # Visualizar
8 ...
```


Rolling: Média Móvel Simples (cont.)



Expanding: Estatísticas Cumulativas

Python

```
1 # Expanding = janela que cresce desde o inicio
2
3 # Media cumulativa
4 df['CO_cumulative_mean'] = df['CO(GT)'].expanding().mean()
5
6 # No tempo t, usa TODOS os valores de 0 ate t
7
8 print("=== EXPANDING MEAN ===")
9 print(df[['CO(GT)', 'CO_cumulative_mean']].head(20))
10
11 # Outras estatisticas cumulativas
12 df['CO_cumulative_max'] = df['CO(GT)'].expanding().max()
13 df['CO_cumulative_min'] = df['CO(GT)'].expanding().min()
14 df['CO_cumulative_std'] = df['CO(GT)'].expanding().std()
```

Expanding: Estatísticas Cumulativas

</> Python

```
1 # Uso:  
2 # - Ver tendencia de longo prazo  
3 # - Comparar valor atual com historico  
4 # - Detectar mudancas de regime
```

Rolling vs Expanding vs Resample

Três técnicas, três propósitos:

| Método | O que faz | Janela | Uso |
|---------------|------------------|-----------------------|----------------------|
| Resample | Muda frequência | Período fixo | Agregar por dia/-mês |
| Rolling | Média móvel | Tamanho fixo, desliza | Suavizar, features |
| Expanding | Acumulativo | Cresce desde início | Tendência geral |

Rolling vs Expanding vs Resample

Exemplo com 10 valores:

- ▶ **Resample('2D')**: 10 valores → 5 valores (média a cada 2 dias)
- ▶ **Rolling(3)**: Janela [1,2,3], [2,3,4], ... (10 valores)
- ▶ **Expanding()**: [1], [1,2], [1,2,3], ... (10 valores)

Shift: Defasagens (Lags)

</> Python

```
1 # Shift: mover valores no tempo
2
3 # Lag 1 (valor de 1 hora atras)
4 df['CO_lag1'] = df['CO(GT)'].shift(1)
5
6 # Lag 24 (valor de 24 horas atras - ontem)
7 df['CO_lag24'] = df['CO(GT)'].shift(24)
8
9 # Lag 168 (7 dias atras)
10 df['CO_lag168'] = df['CO(GT)'].shift(168)
11
12 print("=== LAGS ===")
13 print(df[['CO(GT)', 'CO_lag1', 'CO_lag24', 'CO_lag168']].head(30))
```

Shift: Defasagens (Lags) (cont.)

Python

```
1 # Uso:
2 # - Autocorrelacao: correlacao com proprio passado
3 # - Features para ML: valor de ontem ajuda prever hoje
4 # - Comparacoes temporais
5
6 # Shift negativo = lead (valores futuros)
7 df['CO_lead1'] = df['CO(GT)'].shift(-1) # Proxima hora
8 # CUIDADO: data leakage se usar para previsao!
```

Diff: Diferenças Temporais

Python

```
1 # Diff: diferenca entre periodos consecutivos
2
3 # Diferenca de 1 periodo (mudanca horaria)
4 df['CO_diff1'] = df['CO(GT)'].diff(1)
5 # CO_diff1[t] = CO[t] - CO[t-1]
6
7 # Diferenca de 24 periodos (mudanca dia-a-dia)
8 df['CO_diff24'] = df['CO(GT)'].diff(24)
9 # CO_diff24[t] = CO[t] - CO[t-24]
10
11 print("=== DIFFERENCES ===")
12 print(df[['CO(GT)', 'CO_diff1', 'CO_diff24']].head(30))
```


Diff: Diferenças Temporais (cont.)

</> Python

```
1 # Interpretacao:
2 # - diff > 0: aumentou
3 # - diff < 0: diminuiu
4 # - diff = 0: estavel
5
6 # Uso:
7 # - Tornar serie estacionaria
8 # - Detectar mudancas abruptas
9 # - Velocidade de mudanca
```

pct_change: Variações Percentuais

</> Python

```
1 # Variacao percentual de 1 periodo
2 df['CO_pct_change1'] = df['CO(GT)'].pct_change(1)
3 # pct_change[t] = (CO[t] - CO[t-1]) / CO[t-1]
4
5 # Variacao percentual dia-a-dia
6 df['CO_pct_change24'] = df['CO(GT)'].pct_change(24)
7
8 print("=== PERCENT CHANGES ===")
9 print(df[['CO(GT)', 'CO_pct_change1', 'CO_pct_change24']].head(30))
```

pct_change: Variações Percentuais (cont.)

</> Python

```
1 # Interpretacao:
2 # - 0.1 = +10% de aumento
3 # - -0.05 = -5% de queda
4 # - 0.0 = sem mudanca
5
6 # Converter para percentual
7 df['CO_pct_change24_formatted'] = df['CO_pct_change24'] * 100
8
9 print(f"\nVariacao percentual media dia-a-dia: {df['CO_pct_change24'].mean()
      *100:.2f}%")
```

Autocorrelação: Correlação com o Passado

Python

```
1 # Medir correlacao com valores defasados
2
3 # Criar lags
4 for lag in [1, 6, 12, 24, 168]:
5     df[f'CO_lag{lag}'] = df['CO(GT)'].shift(lag)
6
7 # Calcular correlacoes
8 correlations = []
9 for lag in [1, 6, 12, 24, 168]:
10     corr = df['CO(GT)'].corr(df[f'CO_lag{lag}'])
11     correlations.append((lag, corr))
12     print(f"Lag {lag:3d}h: correlacao = {corr:.3f}")
```

Autocorrelação: Correlação com o Passado (cont.)

</> Python

```
1 # Interpretacao:
2 # - Lag 1h: correlacao alta (0.9+) = muito autocorrelacionado
3 # - Lag 24h: correlacao alta = padrao diario forte
4 # - Lag 168h (7d): correlacao = padrao semanal
5
6 # Usar para:
7 # - Identificar periodicidade
8 # - Escolher lags para features
```

Plotar Autocorrelação

</> Python

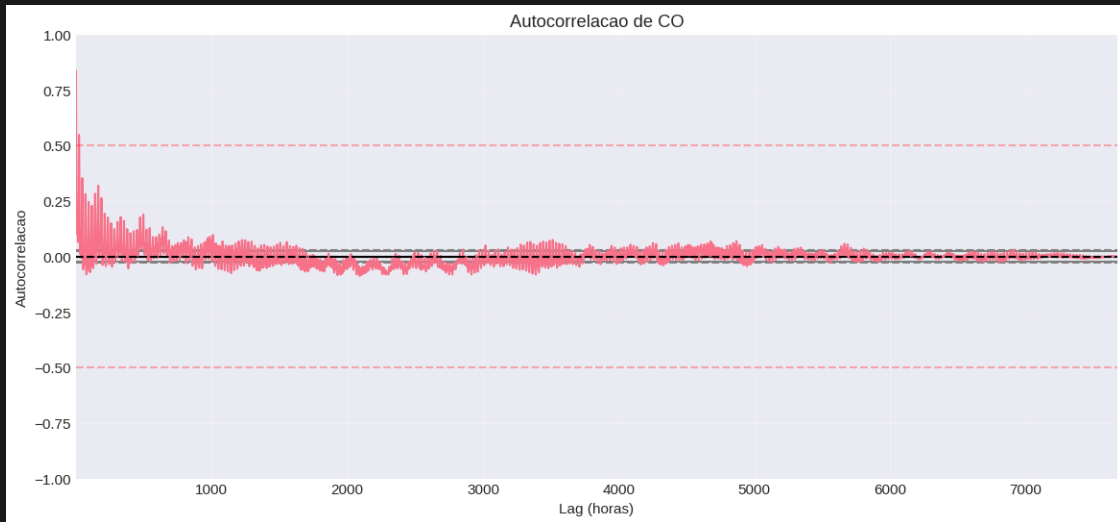
```
1 from pandas.plotting import autocorrelation_plot
2
3 # Plot de autocorrelacao
4 plt.figure(figsize=(14, 6))
5 autocorrelation_plot(df['CO(GT)'].dropna())
6 plt.title('Autocorrelacao de CO')
7 plt.xlabel('Lag (horas)')
8 plt.ylabel('Autocorrelacao')
9 plt.axhline(y=0, color='k', linestyle='--')
10 plt.axhline(y=0.5, color='r', linestyle='--', alpha=0.3)
11 plt.axhline(y=-0.5, color='r', linestyle='--', alpha=0.3)
12 plt.grid(True, alpha=0.3)
13 plt.show()
```

Plotar Autocorrelação (cont.)

</> Python

```
1 # Linhas horizontais: limites de significancia
2 # Se autocorr > linha: correlacao significativa
3 # Picos em multiplos de 24: sazonalidade diaria
4 # Picos em multiplos de 168: sazonalidade semanal
```

Plotar Autocorrelação (cont.)



Criar Features com Rolling e Lags

</> Python

```
1 # Criar features uteis para modelagem
2
3 # 1. Lags importantes
4 df['CO_lag1'] = df['CO(GT)'].shift(1)
5 df['CO_lag24'] = df['CO(GT)'].shift(24)
6 df['CO_lag168'] = df['CO(GT)'].shift(168)
7
8 # 2. Medias moveis
9 df['CO_MA24'] = df['CO(GT)'].rolling(24).mean()
10 df['CO_MA168'] = df['CO(GT)'].rolling(168).mean()
11
12 # 3. Diferenca entre valor e media
13 df['CO_vs_MA24'] = df['CO(GT)'] - df['CO_MA24']
```

Criar Features com Rolling e Lags (cont.)

Python

```
1 # 4. Volatilidade (desvio movel)
2 df['CO_volatility'] = df['CO(GT)'].rolling(24).std()
3
4 # 5. Range (max - min)
5 df['CO_range24'] = (df['CO(GT)'].rolling(24).max() -
6                     df['CO(GT)'].rolling(24).min())
7
8 print("=== FEATURES CRIADAS ===")
9 print(df[['CO(GT)', 'CO_lag24', 'CO_MA24',
10          'CO_vs_MA24', 'CO_volatility']].tail())
```

Features para Forecasting com Dados Históricos

Features mais úteis para previsão (criadas com valores passados):

1. Lags (valores passados):

- ▶ lag1, lag24, lag168 (1h, 1d, 1sem atrás)
- ▶ Use autocorrelação para escolher

2. Médias móveis:

- ▶ MA7, MA30 (tendência de curto/longo prazo)
- ▶ Diferença entre valor atual e MA

3. Volatilidade:

- ▶ Desvio padrão móvel
- ▶ Range móvel (max - min)

Features para Forecasting com Dados Históricos

Features mais úteis para previsão (criadas com valores passados):

4. Diferenças:

- ▶ diff1, diff24 (taxa de mudança)
- ▶ pct_change (variação percentual)

5. Componentes temporais:

- ▶ hour, dayofweek, month
- ▶ is_weekend, is_holiday

Nota Importante

Estas features são úteis para entender padrões e relações temporais nos dados!

Ewm: Exponential Weighted Moving Average

</> Python

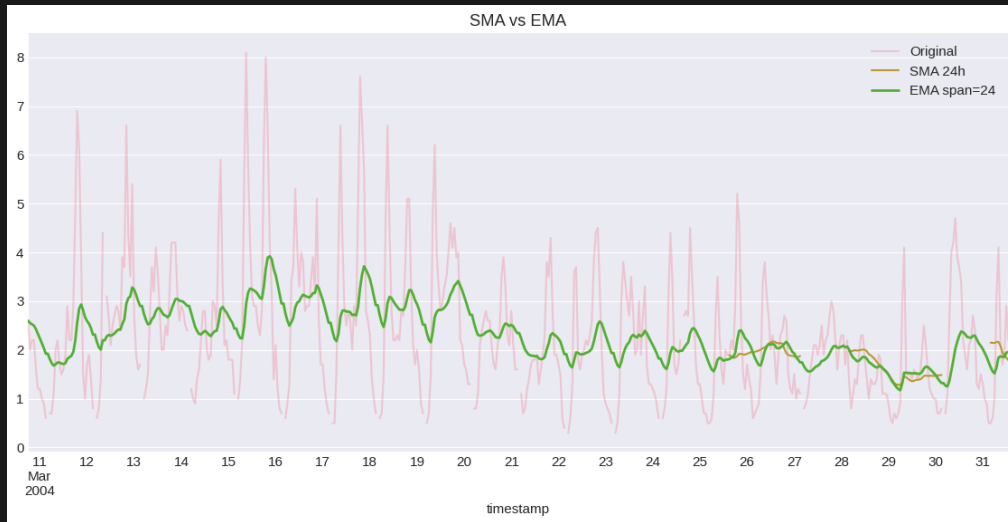
```
1 # EWM: da mais peso aos valores recentes
2
3 # Media movel simples (SMA)
4 df['CO_SMA24'] = df['CO(GT)'].rolling(24).mean()
5
6 # Media movel exponencial (EMA/EWMA)
7 df['CO_EMA'] = df['CO(GT)'].ewm(span=24, adjust=False).mean()
8
9 # Comparar
10 print("=== SMA vs EMA ===")
11 print(df[['CO(GT)', 'CO_SMA24', 'CO_EMA']].tail(20))
```

Ewm: Exponential Weighted Moving Average (cont.)

</> Python

```
1 # Visualizar
2 plt.figure(figsize=(14, 6))
3 df['CO(GT)'][:500].plot(alpha=0.3, label='Original')
4 df['CO_SMA24'][:500].plot(label='SMA 24h')
5 df['CO_EMA'][:500].plot(label='EMA span=24', linewidth=2)
6 plt.legend()
7 plt.title('SMA vs EMA')
8 plt.show()
9
10 # EMA: mais responsiva a mudancas recentes
11 # SMA: trata todos valores igualmente
```

Ewm: Exponential Weighted Moving Average (cont.)



SMA vs EMA: Diferenças

Simple Moving Average (SMA):

- ▶ Peso igual para todos os valores na janela
- ▶ Fórmula: $(x_1 + x_2 + \dots + x_n)/n$
- ▶ Mais suave
- ▶ Lag maior

Exponential Moving Average (EMA):

- ▶ Peso decai exponencialmente
- ▶ Mais peso aos valores recentes
- ▶ Menos lag
- ▶ Mais responsiva

SMA vs EMA: Diferenças

Quando usar cada um:

- ▶ **SMA:** Análise de tendências de longo prazo
- ▶ **EMA:** Trading, detecção rápida de mudanças
- ▶ **Ambos:** Cruzamento de médias (sinais de compra/venda)

Aplicar Função Customizada em Rolling

</> Python

```
1 # Aplicar funcao customizada
2
3 # Funcao: range (max - min)
4 def window_range(window):
5     return window.max() - window.min()
6
7 # Aplicar
8 df['CO_range_custom'] = df['CO(GT)'].rolling(24).apply(window_range)
9
10 # Funcao: quantil 75%
11 df['CO_q75'] = df['CO(GT)'].rolling(24).quantile(0.75)
```

Aplicar Função Customizada em Rolling (cont.)

</> Python

```
1 # Funcao: contagem de valores acima da media
2 def count_above_mean(window):
3     return (window > window.mean()).sum()
4
5 df['CO_count_above'] = df['CO(GT)'].rolling(24).apply(count_above_mean)
6
7 print("=== FUNCOES CUSTOMIZADAS ===")
8 print(df[['CO(GT)', 'CO_range_custom',
9           'CO_q75', 'CO_count_above']].tail())
10
11 # Usar .apply() para qualquer logica customizada!
```

Comparar Diferentes Janelas

Python

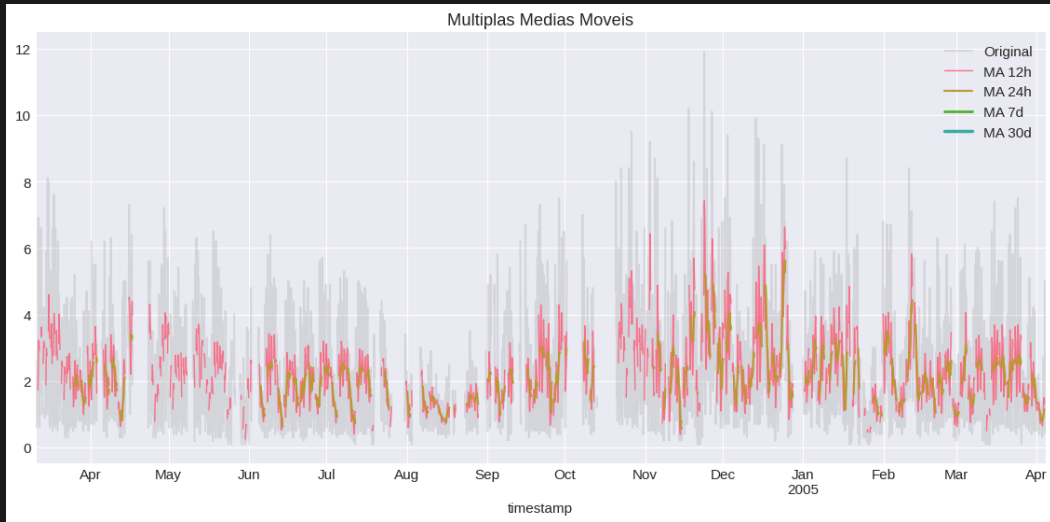
```
1 # Múltiplas janelas para capturar diferentes escalas
2
3 # Curto prazo (12h)
4 df['CO_MA12'] = df['CO(GT)'].rolling(12).mean()
5
6 # Médio prazo (24h)
7 df['CO_MA24'] = df['CO(GT)'].rolling(24).mean()
8
9 # Longo prazo (168h = 7 dias)
10 df['CO_MA168'] = df['CO(GT)'].rolling(168).mean()
11
12 # Muito longo prazo (720h = 30 dias)
13 df['CO_MA720'] = df['CO(GT)'].rolling(720).mean()
```

Comparar Diferentes Janelas (cont.)

</> Python

```
1 # Visualizar todas
2 plt.figure(figsize=(14, 6))
3 df['CO(GT)'].plot(alpha=0.2, label='Original', color='gray')
4 df['CO_MA12'].plot(label='MA 12h', linewidth=1)
5 df['CO_MA24'].plot(label='MA 24h', linewidth=1.5)
6 df['CO_MA168'].plot(label='MA 7d', linewidth=2)
7 df['CO_MA720'].plot(label='MA 30d', linewidth=2.5)
8 plt.legend()
9 plt.title('Múltiplas Médias Móveis')
10 plt.show()
```

Comparar Diferentes Janelas (cont.)



Detector Anomalias com Rolling

Python

```
1 # Detectar valores anormais usando rolling statistics
2
3 # Media e desvio movel
4 df['CO_rolling_mean'] = df['CO(GT)'].rolling(24).mean()
5 df['CO_rolling_std'] = df['CO(GT)'].rolling(24).std()
6
7 # Z-score movel
8 df['CO_zscore'] = (df['CO(GT)'] - df['CO_rolling_mean']) / df['
    CO_rolling_std']
9
10 # Anomalias: |z-score| > 3
11 df['is_anomaly'] = (df['CO_zscore'].abs() > 3).astype(int)
```

Detector Anomalias com Rolling (cont.)

</> Python

```
1 print("=== ANOMALIAS DETECTADAS ===")
2 anomalies = df[df['is_anomaly'] == 1]
3 print(f"Total de anomalias: {len(anomalies)}")
4 print(f"Percentual: {(len(anomalies)/len(df))*100:.2f}%")
5
6 print("\nExemplos:")
7 print(anomalies[['CO(GT)', 'CO_rolling_mean',
8                  'CO_zscore']].head())
9
10 # Visualizar
11 # ...
```


Exercício: Análise com Rolling e Resample

Exercício Prático

Explore agregações temporais no Air Quality:

1. Resample: diaria e semanal; comparar com dados horários
2. Rolling windows:
 - ▶ Calcular MA de 24h e 168h
 - ▶ Plotar original vs médias móveis
 - ▶ Calcular desvio padrão móvel
3. Features:
 - ▶ Criar lags (1h, 24h, 168h)
 - ▶ Criar diferenças (diff)
 - ▶ Calcular autocorrelação

Exercício: Análise com Rolling e Resample (cont.)

Exercício Prático

Explore agregações temporais no Air Quality:

4. Detecção de anomalias:

- ▶ Use z-score móvel
- ▶ Identifique top 10 anomalias
- ▶ Quando ocorreram?

Bloco 3

Decomposição de Séries Temporais

O que é Decomposição?

Decomposição = separar série em componentes

Modelo aditivo:

$$Y_t = T_t + S_t + R_t$$

Modelo multiplicativo:

$$Y_t = T_t \times S_t \times R_t$$

Componentes:

- ▶ T_t (**Trend**): Tendência de longo prazo
- ▶ S_t (**Seasonal**): Padrão sazonal repetitivo
- ▶ R_t (**Residual**): Ruído/irregular

O que é Decomposição?

Decomposição = separar série em componentes

Modelo aditivo:

$$Y_t = T_t + S_t + R_t$$

Modelo multiplicativo:

$$Y_t = T_t \times S_t \times R_t$$

Quando usar cada modelo:

- ▶ **Aditivo:** Amplitude sazonal constante
- ▶ **Multiplicativo:** Amplitude sazonal cresce com tendência

Aditivo vs Multiplicativo

Exemplo visual:

Modelo Aditivo:

- ▶ Vendas oscilam ± 100 unidades em jan/2020
- ▶ Vendas oscilam ± 100 unidades em dez/2024
- ▶ Amplitude sazonal CONSTANTE
- ▶ Variação independente do nível

Modelo Multiplicativo:

- ▶ Vendas oscilam $\pm 10\%$ em jan/2020 (1000 ± 100)
- ▶ Vendas oscilam $\pm 10\%$ em dez/2024 (5000 ± 500)
- ▶ Amplitude sazonal CRESCE
- ▶ Variação proporcional ao nível

Aditivo vs Multiplicativo

Regra prática:

- ▶ Se em dúvida, comece com aditivo
- ▶ Se amplitude cresce, use multiplicativo
- ▶ Ou aplique log e use aditivo

Decomposição: Preparação dos Dados

</> Python

```
1 # Preparar dados para decomposicao
2
3 # 1. Resample para diario (decomposicao funciona melhor em freq regular)
4 daily = df['CO(GT)'].resample('D').mean()
5
6 # 2. Remover NaN
7 daily = daily.dropna()
8
9 print("=== DADOS PARA DECOMPOSICAO ===")
10 print(f"Frequencia: Diaria")
11 print(f"Registros: {len(daily)}")
12 print(f"Periodo: {daily.index.min()} a {daily.index.max()}")
13 print(f"Missing: {daily.isnull().sum()}")
```


Decomposição: Preparação dos Dados

Python

```
1 # Verificar
2 print(f"\nAmostra:")
3 print(daily.head())
4
5 # IMPORTANTE: Decomposicao precisa de:
6 # - Frequencia regular (sem gaps)
7 # - Sem missing values
8 # - Período suficiente (pelo menos 2 ciclos sazonais)
```

Decomposição: statsmodels

</> Python

```
1 from statsmodels.tsa.seasonal import seasonal_decompose
2
3 # Decomposicao aditiva
4 decomposition = seasonal_decompose(
5     daily,
6     model='additive',
7     period=7, # Sazonalidade semanal (7 dias)
8     extrapolate_trend='freq'
9 )
10
11 # Extrair componentes
12 trend = decomposition.trend
13 seasonal = decomposition.seasonal
14 residual = decomposition.resid
```

Decomposição: statsmodels (cont.)

</> Python

```
1 print("=== COMPONENTES ===")
2 print(f"Trend: {trend.shape}")
3 print(f"Seasonal: {seasonal.shape}")
4 print(f"Residual: {residual.shape}")
5
6 # Verificar decomposicao
7 # Original = Trend + Seasonal + Residual
8 reconstructed = trend + seasonal + residual
9 print(f"\nErro de reconstrucao: {(daily - reconstructed).abs().max():.10f}")
10
```

Visualizar Decomposição

Python

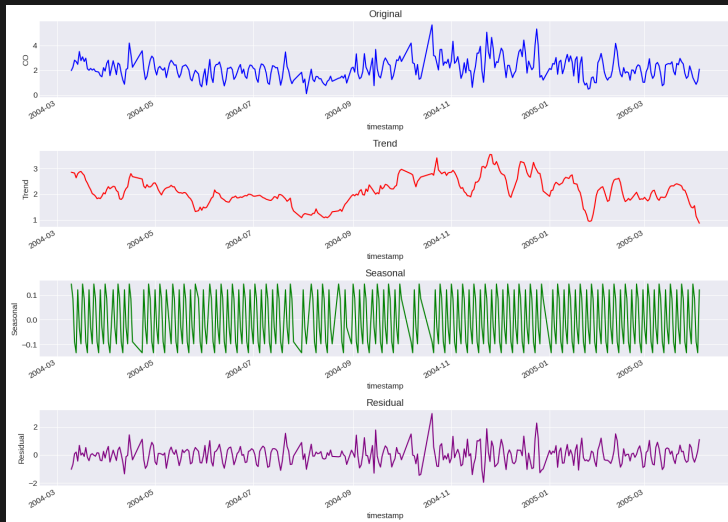
```
1 import matplotlib.pyplot as plt
2
3 # Plot de decomposicao
4 fig = decomposition.plot()
5 fig.set_size_inches(14, 10)
6 plt.tight_layout()
7 plt.show()
8
9 # Ou manualmente com mais controle
10 fig, axes = plt.subplots(4, 1, figsize=(14, 10))
11
12 # Original
13 daily.plot(ax=axes[0], title='Original', color='blue')
14 axes[0].set_ylabel('C0')
```

Visualizar Decomposição (cont.)

Python

```
1 # Tendencia
2 trend.plot(ax=axes[1], title='Trend', color='red')
3 axes[1].set_ylabel('Trend')
4
5 # Sazonalidade
6 seasonal.plot(ax=axes[2], title='Seasonal', color='green')
7 axes[2].set_ylabel('Seasonal')
8
9 # Residual
10 residual.plot(ax=axes[3], title='Residual', color='purple')
11 axes[3].set_ylabel('Residual')
12
13 plt.tight_layout()
14 plt.show()
```

Visualizar Decomposição (cont.)



Interpretar Componentes

O que cada componente revela:

1. Tendência (Trend):

- ▶ Direção de longo prazo
- ▶ Crescimento ou declínio
- ▶ Ciclos muito longos
- ▶ Se está subindo/descendo/estável

2. Sazonalidade (Seasonal):

- ▶ Padrão que se repete
- ▶ Picos e vales regulares
- ▶ Período fixo (diário, semanal, anual)
- ▶ Quando esperar valores altos/baixos

O que cada componente revela:

3. Resíduo (Residual):

- ▶ O que não foi explicado
- ▶ Ruído aleatório
- ▶ Eventos irregulares
- ▶ Deve ser aleatório (sem padrão)

Analisar Tendência

</> Python

```
1 # Analise do componente de tendencia
2
3 print("=== ANALISE DE TENDENCIA ===")
4
5 # Direcao geral
6 start_value = trend.dropna().iloc[0]
7 end_value = trend.dropna().iloc[-1]
8 change = end_value - start_value
9 pct_change = (change / start_value) * 100
10
11 print(f"Valor inicial: {start_value:.2f}")
12 print(f"Valor final: {end_value:.2f}")
13 print(f"Mudanca: {change:.2f} ({pct_change:+.1f}%)")
```

Analisar Tendência (cont.)

</> Python

```
1 if pct_change > 5:
2     print("Tendencia de CRESCIMENTO")
3 elif pct_change < -5:
4     print("Tendencia de DECLINIO")
5 else:
6     print("Tendencia ESTAVEL")
7
8 # Taxa de crescimento media
9 growth_rate = trend.pct_change().mean() * 100
10 print(f"\nTaxa de crescimento diaria media: {growth_rate:.3f}%")
```

Analisar Tendência (cont.)

</> Python

```
1 # === ANALISE DE TENDENCIA ===  
2 # Valor inicial: 2.85  
3 # Valor final: 0.87  
4 # Mudanca: -1.99 (-69.6%)  
5 # Tendencia de DECLINIO  
6 # Taxa de crescimento diaria media: -0.083%
```

Analisar Sazonalidade

Python

```
1 # Analise do componente sazonal
2
3 print("=== ANALISE DE SAZONALIDADE ===")
4
5 # Amplitude sazonal
6 seasonal_amplitude = seasonal.max() - seasonal.min()
7 print(f"Amplitude sazonal: {seasonal_amplitude:.2f}")
8
9 # Percentual da variacao total
10 total_variation = daily.std()
11 seasonal_strength = (seasonal.std() / total_variation) * 100
12 print(f"Forca da sazonalidade: {seasonal_strength:.1f}% da variacao")
```

Analisar Sazonalidade (cont.)

</> Python

```
1 # Padrao semanal (se period=7)
2 # Media por dia da semana
3 seasonal_pattern = seasonal.groupby(seasonal.index.dayofweek).mean()
4 seasonal_pattern.index = ['Seg', 'Ter', 'Qua', 'Qui',
5                           'Sex', 'Sab', 'Dom']
6 print("\nPadrao semanal (media):")
7 print(seasonal_pattern)
8
9 # Dia com maior/menor valor sazonal
10 print(f"\nPior dia: {seasonal_pattern.idxmax()}")
11 print(f"Melhor dia: {seasonal_pattern.idxmin()}")
```

Analisar Resíduo

Python

```
1 # Analise do residuo
2
3 print("=== ANALISE DE RESIDUO ===")
4
5 # Estatisticas
6 print(f"Media: {residual.mean():.6f}") # Deve ser ~0
7 print(f"Desvio padrao: {residual.std():.3f}")
8 print(f"Min: {residual.min():.3f}")
9 print(f"Max: {residual.max():.3f}")
10
11 # Outliers no residuo (eventos incomuns)
12 residual_threshold = 3 * residual.std()
13 outliers = residual[residual.abs() > residual_threshold]
14 print(f"\nOutliers (|residual| > 3): {len(outliers)}")
```

Analisar Resíduo (cont.)

Python

```
1 if len(outliers) > 0:
2     print("\nMaiores residuos:")
3     top_outliers = residual.abs().nlargest(5)
4     for date, value in top_outliers.items():
5         print(f"  {date.date()}: {value:.3f}")
6
7 # Teste de aleatoriedade (autocorrelacao do residuo)
8 residual_autocorr = residual.autocorr(lag=1)
9 print(f"\nAutocorrelacao lag-1: {residual_autocorr:.3f}")
10 print("(Deve ser proximo de 0 se residuo é aleatorio)")
```

Comparar Aditivo vs Multiplicativo

</> Python

```
1 # Decompor com ambos os modelos
2
3 # Aditivo
4 decomp_add = seasonal_decompose(daily, model='additive',
5                                 period=7, extrapolate_trend='freq')
6
7 # Multiplicativo
8 decomp_mult = seasonal_decompose(daily, model='multiplicative',
9                                  period=7, extrapolate_trend='freq')
10
11 # Comparar resíduos (qual modelo se ajusta melhor?)
12 residual_add = decomp_add.resid
13 residual_mult = decomp_mult.resid
```


Comparar Aditivo vs Multiplicativo (cont.)

</> Python

```
1 print("== COMPARACAO DE MODELOS ==")
2 print(f"Aditivo - Residuo std: {residual_add.std():.3f}")
3 print(f"Multiplicativo - Residuo std: {residual_mult.std():.3f}")
4
5 # Menor desvio = melhor ajuste
6 if residual_add.std() < residual_mult.std():
7     print("\nModelo ADITIVO se ajusta melhor")
8 else:
9     print("\nModelo MULTIPLICATIVO se ajusta melhor") # É o caso aqui
```

Escolher Período de Sazonalidade

Python

```
1 # Como escolher o periodo correto?
2
3 # Para dados diarios:
4 # - period=7: sazonalidade SEMANAL
5 # - period=30: sazonalidade MENSAL
6 # - period=365: sazonalidade ANUAL
7
8 # Para dados horarios:
9 # - period=24: sazonalidade DIARIA
10 # - period=168: sazonalidade SEMANAL (7*24)
11
12 # Testar diferentes periodos
```

Escolher Período de Sazonalidade

</> Python

```
1 periods_to_test = [7, 14, 30]
2
3 for period in periods_to_test:
4     decomp = seasonal_decompose(daily, model='additive',
5                                 period=period,
6                                 extrapolate_trend='freq')
7     residual_std = decomp.resid.std()
8     print(f"Period={period:2d}: Residual std = {residual_std:.4f}")
9
10 print("\nMenor residual std = melhor periodo")
11
12 # Ou use conhecimento do dominio:
13 # - Qualidade do ar: padrao semanal (trabalho vs fim de semana)
```

Remover Tendência (Detrending)

</> Python

```
1 # Remover tendencia da serie original
2
3 # Opcao 1: Subtrair tendencia (aditivo)
4 detrended = daily - trend
5
6 print("=== SERIE SEM TENDENCIA ===")
7 print(detrended.head())
8
9 # Opcao 2: Dividir por tendencia (multiplicativo)
10 detrended_mult = daily / trend
11
12 # Visualizar
13 fig, axes = plt.subplots(2, 1, figsize=(14, 8))
```

Remover Tendência (Detrending)

</> Python

```
1 # Original
2 daily.plot(ax=axes[0], title='Original', color='blue')
3
4 # Sem tendencia (destaca sazonalidade)
5 detrended.plot(ax=axes[1], title='Detrended', color='green')
6
7 plt.tight_layout()
8 plt.show()
9
10 # Por que remover tendencia?
11 # - Analise de sazonalidade mais clara
12 # - Tornar serie estacionaria
13 # - Pre-processamento para modelagem
```

Remover Sazonalidade (Seasonal Adjustment)

</> Python

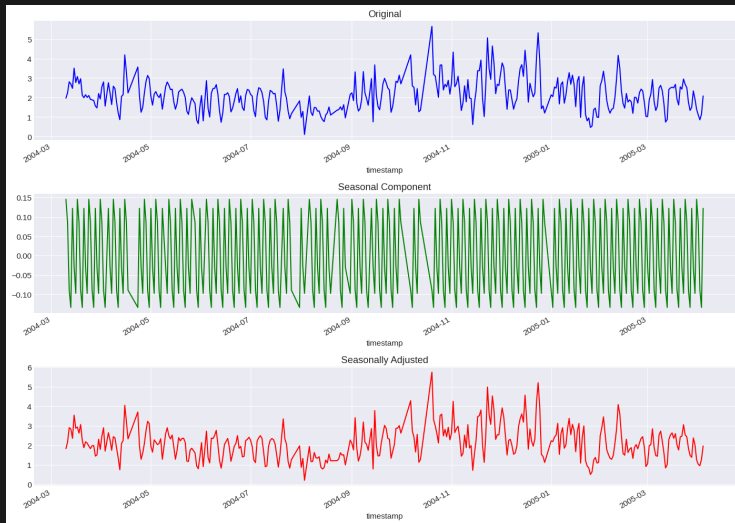
```
1 # Remover sazonalidade (dessazonalizar)
2
3 # Serie ajustada = Original - Sazonal (aditivo)
4 seasonally_adjusted = daily - seasonal
5
6 print("=== SERIE DESSAZONALIZADA ===")
7 print(seasonally_adjusted.head())
8
9 # Visualizar
10 fig, axes = plt.subplots(3, 1, figsize=(14, 10))
11
12 # Original
13 daily.plot(ax=axes[0], title='Original', color='blue')
```

Remover Sazonalidade (Seasonal Adjustment) (cont.)

</> Python

```
1 # Sazonalidade
2 seasonal.plot(ax=axes[1], title='Seasonal Component', color='green')
3
4 # Ajustada (sem sazonalidade)
5 seasonally_adjusted.plot(ax=axes[2],
6                           title='Seasonally Adjusted',
7                           color='red')
8 plt.tight_layout()
9 plt.show()
10 # Uso: ver tendencia sem oscilacoes sazonais
```

Remover Sazonalidade (Seasonal Adjustment) (cont.)



Aplicações da Decomposição

Por que decompor séries temporais?

1. Entendimento:

- ▶ Identificar padrões ocultos
- ▶ Separar efeitos de curto vs longo prazo
- ▶ Comunicar insights ao negócio

2. Análise:

- ▶ Medir força da sazonalidade
- ▶ Detectar mudanças de tendência
- ▶ Encontrar eventos anômalos (resíduo)

Aplicações da Decomposição

Por que decompor séries temporais?

3. Previsão:

- ▶ Prever cada componente separadamente
- ▶ Modelos específicos para cada parte
- ▶ Combinar previsões

4. Pré-processamento:

- ▶ Remover tendência → estacionariedade
- ▶ Ajuste sazonal → comparações justas
- ▶ Isolar ruído → detecção de anomalias

Limitações da Decomposição

Cuidados importantes:

1. Assume periodicidade constante:

- ▶ Sazonalidade deve ter período fixo
- ▶ Não funciona se período muda

2. Perda de dados nas pontas:

- ▶ Tendência usa média móvel
- ▶ NaN no início/fim da série
- ▶ Use `extrapolate_trend='freq'`

Limitações da Decomposição

Cuidados importantes:

3. Escolha do período é crítica:

- ▶ Período errado → decomposição sem sentido
- ▶ Precisa conhecimento do domínio

4. Múltiplas sazonalidades:

- ▶ Decomposição clássica captura apenas UMA
- ▶ Exemplo: diária E semanal E anual
- ▶ Precisa métodos avançados (STL, TBATS)

Bloco 4

Visualizações e Boas Práticas

Visualização 1: Line Plot Básico

</> Python

```
1 # Plot de linha: padrao para series temporais
2
3 plt.figure(figsize=(14, 6))
4 plt.plot(daily.index, daily.values, linewidth=1, color='steelblue')
5 plt.title('CO (daily average)', fontsize=16, fontweight='bold')
6 plt.xlabel('Date', fontsize=12)
7 plt.ylabel('CO (mg/m³)', fontsize=12)
8 plt.grid(True, alpha=0.3)
9 plt.tight_layout()
10 plt.show()
```

Visualização 1: Line Plot Básico (cont.)

</> Python

```
1 # Boas praticas:
2 # - Sempre rotular eixos com unidades
3 # - Título descritivo
4 # - Grid para facilitar leitura
5 # - Tamanho adequado (14x6 ou 16x6)
6 # - Linha nao muito grossa (linewidth=1 ou 1.5)
```

Visualização 2: Comparar Múltiplas Séries

Python

```
1 # Comparar varias series no mesmo plot
2
3 daily_multi = df[['CO(GT)', 'NO2(GT)', 'T']].resample('D').mean()
4
5 fig, axes = plt.subplots(3, 1, figsize=(14, 10), sharex=True)
6
7 # CO
8 axes[0].plot(daily_multi.index, daily_multi['CO(GT)'],
9             color='red', linewidth=1)
10 axes[0].set_ylabel('CO (mg/m³)')
11 axes[0].set_title('Carbon Monoxide')
12 axes[0].grid(True, alpha=0.3)
```


Visualização 2: Comparar Múltiplas Séries (cont.)

</> Python

```
1 # NO2
2 axes[1].plot(daily_multi.index, daily_multi['NO2(GT)'],
3             color='blue', linewidth=1)
4 axes[1].set_ylabel('NO2 (µg/m³)')
5 axes[1].set_title('Nitrogen Dioxide')
6 axes[1].grid(True, alpha=0.3)
7 # Temperatura
8 axes[2].plot(daily_multi.index, daily_multi['T'],
9             color='green', linewidth=1)
10 axes[2].set_ylabel('Temperature (°C)')
11 axes[2].set_title('Temperature')
12 axes[2].grid(True, alpha=0.3)
13 plt.tight_layout()
14 plt.show()
```

Visualização 3: Heatmap Temporal

</> Python

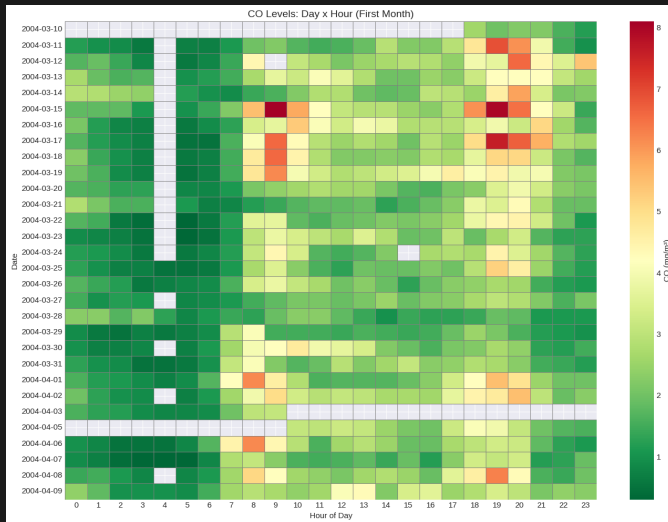
```
1 import seaborn as sns
2
3 # Heatmap: dia x hora
4
5 # Criar matriz: linhas=dias, colunas=horas
6 pivot = df.pivot_table(
7     values='CO(GT)',
8     index=df.index.date,
9     columns=df.index.hour,
10    aggfunc='mean'
11 )
```

Visualização 3: Heatmap Temporal (cont.)

</> Python

```
1 # Plot
2 plt.figure(figsize=(14, 10))
3 sns.heatmap(pivot[:30], # Primeiro mes
4             cmap='RdYlGn_r', # Vermelho=alto, Verde=baixo
5             cbar_kws={'label': 'CO (mg/m³)'},
6             linewidths=0.5,
7             linecolor='gray')
8 plt.title('CO Levels: Day x Hour (First Month)', fontsize=14)
9 plt.xlabel('Hour of Day')
10 plt.ylabel('Date')
11 plt.tight_layout()
12 plt.show()
13 # Revela: padroes diarios e semanais juntos!
```

Visualização 3: Heatmap Temporal (cont.)



Visualização 4: Box Plot por Período

Python

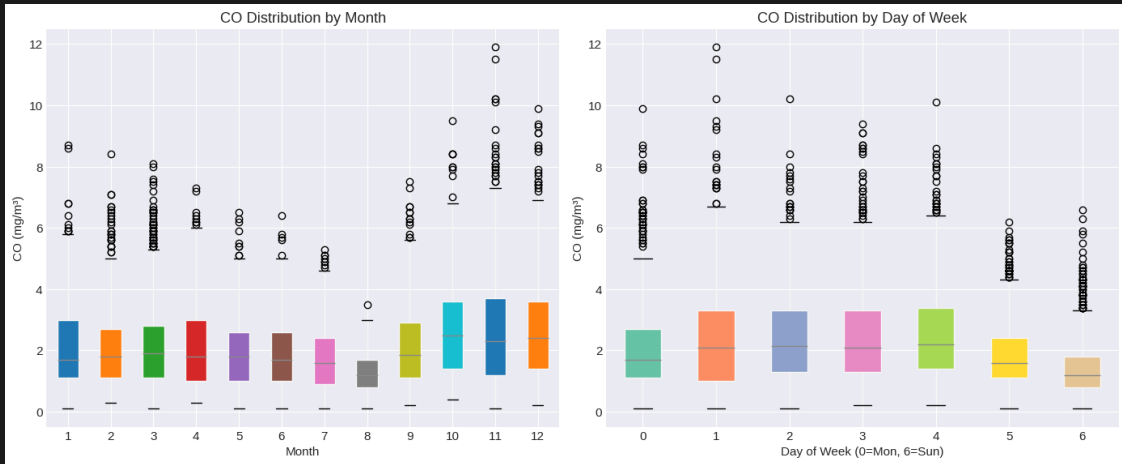
```
1 # Box plots para ver distribuicao por periodo
2
3 # Por mes
4 fig, axes = plt.subplots(1, 2, figsize=(14, 6))
5
6 # Box plot por mes
7 df['month'] = df.index.month
8 df.boxplot(column='CO(GT)', by='month', ax=axes[0])
9 axes[0].set_xlabel('Month')
10 axes[0].set_ylabel('CO (mg/m³)')
11 axes[0].set_title('CO Distribution by Month')
12 axes[0].get_figure().suptitle('') # Remove titulo auto
```

Visualização 4: Box Plot por Período (cont.)

Python

```
1 # Box plot por dia da semana
2 df['dayofweek'] = df.index.dayofweek
3 df.boxplot(column='CO(GT)', by='dayofweek', ax=axes[1])
4 axes[1].set_xlabel('Day of Week (0=Mon, 6=Sun)')
5 axes[1].set_ylabel('CO (mg/m³)')
6 axes[1].set_title('CO Distribution by Day of Week')
7 axes[1].get_figure().suptitle('')
8
9 plt.tight_layout()
10 plt.show()
```

Visualização 4: Box Plot por Período (cont.)



Visualização 5: Lag Scatter Plot

</> Python

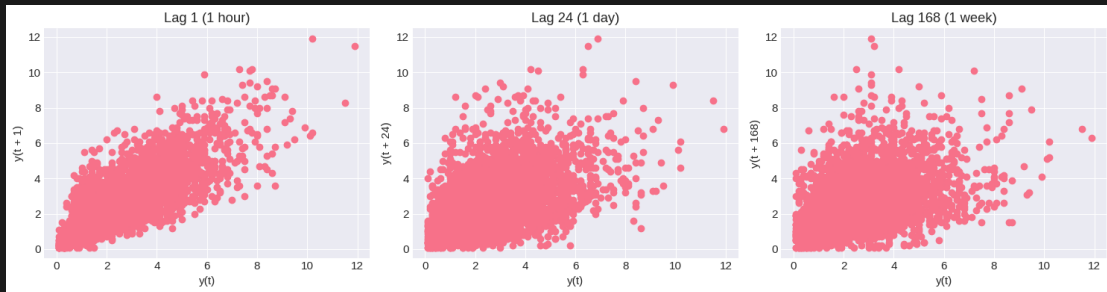
```
1 from pandas.plotting import lag_plot
2
3 # Lag plot: hoje vs ontem
4
5 fig, axes = plt.subplots(1, 3, figsize=(15, 4))
6
7 # Lag 1 (1 hora atras)
8 lag_plot(df['CO(GT)'], lag=1, ax=axes[0])
9 axes[0].set_title('Lag 1 (1 hour)')
10
11 # Lag 24 (1 dia atras)
12 lag_plot(df['CO(GT)'], lag=24, ax=axes[1])
13 axes[1].set_title('Lag 24 (1 day)')
```


Visualização 5: Lag Scatter Plot (cont.)

</> Python

```
1 # Lag 168 (1 semana atras)
2 lag_plot(df['CO(GT)'], lag=168, ax=axes[2])
3 axes[2].set_title('Lag 168 (1 week)')
4 plt.tight_layout()
5 plt.show()
6
7 # Interpretacao:
8 # - Pontos em linha diagonal = alta correlacao
9 # - Pontos dispersos = baixa correlacao
```

Visualização 5: Lag Scatter Plot (cont.)



Boas Práticas: Pipeline Completo

</> Python

```
1 def analyze_timeseries(df, col, freq='D'):
2     """Pipeline completo de analise temporal"""
3
4     # 1. Preparar
5     series = df[col].resample(freq).mean().dropna()
6
7     # 2. Estatisticas basicas
8     print("=== ESTATISTICAS ===")
9     print(series.describe())
10
11    # 3. Decomposicao
12    from statsmodels.tsa.seasonal import seasonal_decompose
13    decomp = seasonal_decompose(series, model='additive',
14                                period=7, extrapolate_trend='freq')
```

Boas Práticas: Pipeline Completo (cont.)

</> Python

```
1      # 4. Visualizacoes
2      fig, axes = plt.subplots(5, 1, figsize=(14, 12))
3      series.plot(ax=axes[0], title='Original')
4      decomp.trend.plot(ax=axes[1], title='Trend')
5      decomp.seasonal.plot(ax=axes[2], title='Seasonal')
6      decomp.resid.plot(ax=axes[3], title='Residual')
7      # Autocorrelacao
8      from pandas.plotting import autocorrelation_plot
9      autocorrelation_plot(series, ax=axes[4])
10     plt.tight_layout()
11     plt.show()
12     return decomp
13 # Usar
14 decomp = analyze_timeseries(df, 'CO(GT)')
```

Antes de analisar séries temporais: Checklist

1. ☐ **DatetimeIndex configurado**
 - ▶ Coluna datetime como índice
 - ▶ Ordenado temporalmente
2. ☐ **Dados limpos**
 - ▶ Missing values tratados
 - ▶ Outliers investigados
 - ▶ Duplicatas removidas
3. ☐ **Frequência regular**
 - ▶ Sem gaps temporais (ou preenchidos)
 - ▶ Frequência consistente
4. ☐ **Exploração inicial**
 - ▶ Visualizar série completa
 - ▶ Estatísticas descritivas
 - ▶ Identificar padrões visuais
5. ☐ **Componentes identificados**
 - ▶ Tendência analisada
 - ▶ Sazonalidade detectada
 - ▶ Período determinado

Erros Comuns - Evite!

1. Não definir `DatetimeIndex`

- ▶ ✗ Manter data como coluna comum
- ▶ ✓ Usar `set_index()`

2. Ignorar ordem temporal

- ▶ ✗ Shuffle ou embaralhar dados
- ▶ ✓ Sempre manter ordem cronológica

3. Validação em períodos aleatórios

- ▶ ✗ Embaralhar dados e testar em períodos passados
- ▶ ✓ Sempre validar em períodos futuros (cronológico)
- ▶ ✓ Exemplo: 80% inicial para análise, 20% final para validação

Erros Comuns - Evite!

4. Usar média global para missing

- ▶ ✗ `fillna(mean())`
- ▶ ✓ `interpolate()` ou `ffill()`

5. Ignorar sazonalidade

- ▶ ✗ Não considerar padrões periódicos
- ▶ ✓ Decompor e analisar

6. Escala de tempo errada

- ▶ ✗ Analisar horária quando padrão é mensal
- ▶ ✓ Resample para frequência adequada

Salvar Análise para Reutilização

Python

```
1 # Salvar objetos importantes
2
3 import joblib
4
5 # 1. Serie processada
6 daily.to_csv('co_daily.csv')
7
8 # 2. Componentes da decomposicao
9 decomposition.trend.to_csv('trend.csv')
10 decomposition.seasonal.to_csv('seasonal.csv')
```


Salvar Análise para Reutilização (cont.)

Python

```
1 # 3. Features criadas
2 features = df[['CO_MA24', 'CO_MA168', 'CO_lag24',
3               'hour_sin', 'hour_cos']].copy()
4 features.to_csv('temporal_features.csv')
5 # 4. Scaler (se usado)
6 # scaler = StandardScaler()
7 # scaler.fit(features)
8 # joblib.dump(scaler, 'scaler.pkl')
9 print("Análise salva!")
10 print("Para recarregar:")
11 print("daily = pd.read_csv('co_daily.csv', index_col=0, parse_dates=True)"
      )
```

Exercício Final: Análise Completa

Exercício Prático

Pipeline completo no Air Quality:

1. Setup:

```
1 # - Carregar dados
2 # - Configurar DatetimeIndex
3 # - Verificar e tratar missing
```

2. Exploração:

- ▶ Visualizar série completa
- ▶ Estatísticas por período
- ▶ Comparar dias úteis vs fim de semana

Exercício Prático

Pipeline completo no Air Quality:

3. Agregações:

- ▶ Resample para diário
- ▶ Calcular médias móveis (24h, 7d)
- ▶ Criar features temporais

4. Decomposição:

- ▶ Decompor série diária
- ▶ Analisar cada componente
- ▶ Visualizar decomposição

5. Insights:

- ▶ Quando poluição é pior?
- ▶ Há tendência de melhora/piora?
- ▶ Qual a força da sazonalidade?

Recap: Aula 11 Completa

Jornada de hoje:

Bloco 1: DatetimeIndex

- ▶ Criar e manipular índice temporal
- ▶ Filtrar e acessar por datas
- ▶ Extrair componentes temporais
- ▶ Criar features básicas

Bloco 2: Agregações

- ▶ Resample (mudar frequência)
- ▶ Rolling windows (médias móveis)
- ▶ Shift, diff, pct_change
- ▶ Autocorrelação

Recap: Aula 11 Completa

Jornada de hoje:

Bloco 3: Decomposição

- ▶ Separar tendência, sazonalidade, resíduo
- ▶ Analisar cada componente
- ▶ Aditivo vs multiplicativo

Bloco 4: Visualizações

- ▶ Plots temporais eficazes
- ▶ Heatmaps, box plots
- ▶ Boas práticas

Próximos Passos

Para aprofundar em séries temporais:

Modelagem (futuro):

- ▶ ARIMA, SARIMA (modelos clássicos)
- ▶ Prophet (Facebook)
- ▶ LSTM, GRU (deep learning)
- ▶ Exponential Smoothing

Análise avançada:

- ▶ Testes de estacionariedade
- ▶ Cointegração
- ▶ VAR (Vector Autoregression)
- ▶ Análise espectral

Próximos Passos

Para aprofundar em séries temporais:

Bibliotecas:

- ▶ statsmodels: análise estatística
- ▶ Prophet: forecasting fácil
- ▶ sktime: ML para séries
- ▶ pmdarima: auto-ARIMA

Exercício Prático

Tempo: 60 minutos

Entrega: via Moodle (notebook)

Tarefas:

1. (atualizado durante a aula)

Notebook: Disponível no Moodle

Obrigado!