

# Geração de Datas Sintéticas com NumPy & Pandas para Dashboards Complexos em Plotly/Dash

A criação de **datasets temporais artificiais de alta qualidade** é um passo essencial quando precisamos prototipar ou demonstrar painéis interativos antes de ter acesso a dados de produção. A seguir, apresento um guia detalhado — do momento em que se define a granularidade das séries até a orquestração de múltiplas frequências num mesmo dashboard — combinando NumPy, Pandas, Plotly e Dash.

## 1. Fundamentos de Datas em NumPy e Pandas

### 1.1 Frequências Básicas e `date_range`

O ponto de partida para a maioria dos fluxos é a função `pd.date_range`, capaz de gerar índices temporais com frequências diárias ('D'), horárias ('H') ou anuais ('A'), obedecendo inclusive calendários de negócios ('B'). Ao especificar apenas o `periods`, Pandas calcula a data final automaticamente, o que facilita loops de simulação quando o horizonte ainda é incerto.

### 1.2 Frequências Personalizadas com Offsets

Para granularidades incomuns, como "cada 15 minutos em dias úteis", combinamos offsets: `freq='15T'` com um filtro por `Timestamp.weekday`. A classe `pandas.tseries.offsets.CustomBusinessDay` permite excluir feriados nacionais, úteis para séries financeiras.

### 1.3 Índices de Período (`PeriodIndex`)

Quando análises envolvem fechamentos mensais ou trimestres, a classe `PeriodIndex` fornece coerência semântica e otimização de memória. Converter datas diárias para períodos (`to_period('M')`) simplifica agregações antes da exportação para dashboards.

## 2. Estratégias de Simulação de Séries Temporais

### 2.1 Componentes Clássicos: Tendência, Sazonalidade e Ruído

Uma série realista combina:

- **Tendência:** vetor linear `np.linspace`, exponencial com `np.exp`, ou polinômios de grau `n`.
- **Sazonalidade:** função seno/componente de calendário (`np.sin(2π t / período)`) que respeita periodicidade anual, semanal ou horária.

- **Ruído:** processo gaussiano via `np.random.normal`; ajustes em `scale` aumentam volatilidade.

Exemplo em poucas linhas:

```
t = np.arange(n)
trend = 50 + 0.1 * t                # crescimento suave
season = 10 * np.sin(2*np.pi*t/365) # ciclo anual
noise = np.random.normal(0, 2, n)   # ruído
series = trend + season + noise
```

## 2.2 Correlação entre Séries de Frequências Diferentes

Ao gerar **vendas diárias**, **usuários semanais** e **receita mensal**, manter coerência estatística é crucial. Uma abordagem prática calcula as métricas agregadas a partir da série de maior granularidade e depois adiciona ruído proporcional:

```
daily_sales = base + noise
weekly_users = daily_sales.resample('W').mean() * 5 + np.random.normal(0, 20, len_weeks)
monthly_revenue = daily_sales.resample('M').sum() * 50 + np.random.normal(0, 2000, len_mo)
```

Essa técnica garante que um pico promocional diário repercuta na semana e no mês subsequentes.

## 3. Unificando Granularidades num Único DataFrame

### 3.1 Alinhamento Temporal

Para combinar séries, iniciamos com índices alinhados no **tipo Timestamp**. Se algo estiver em `PeriodIndex`, convertemos com `.to_timestamp(how='end')` antes de `merge` ou `join`.

### 3.2 Resampling e Upsampling

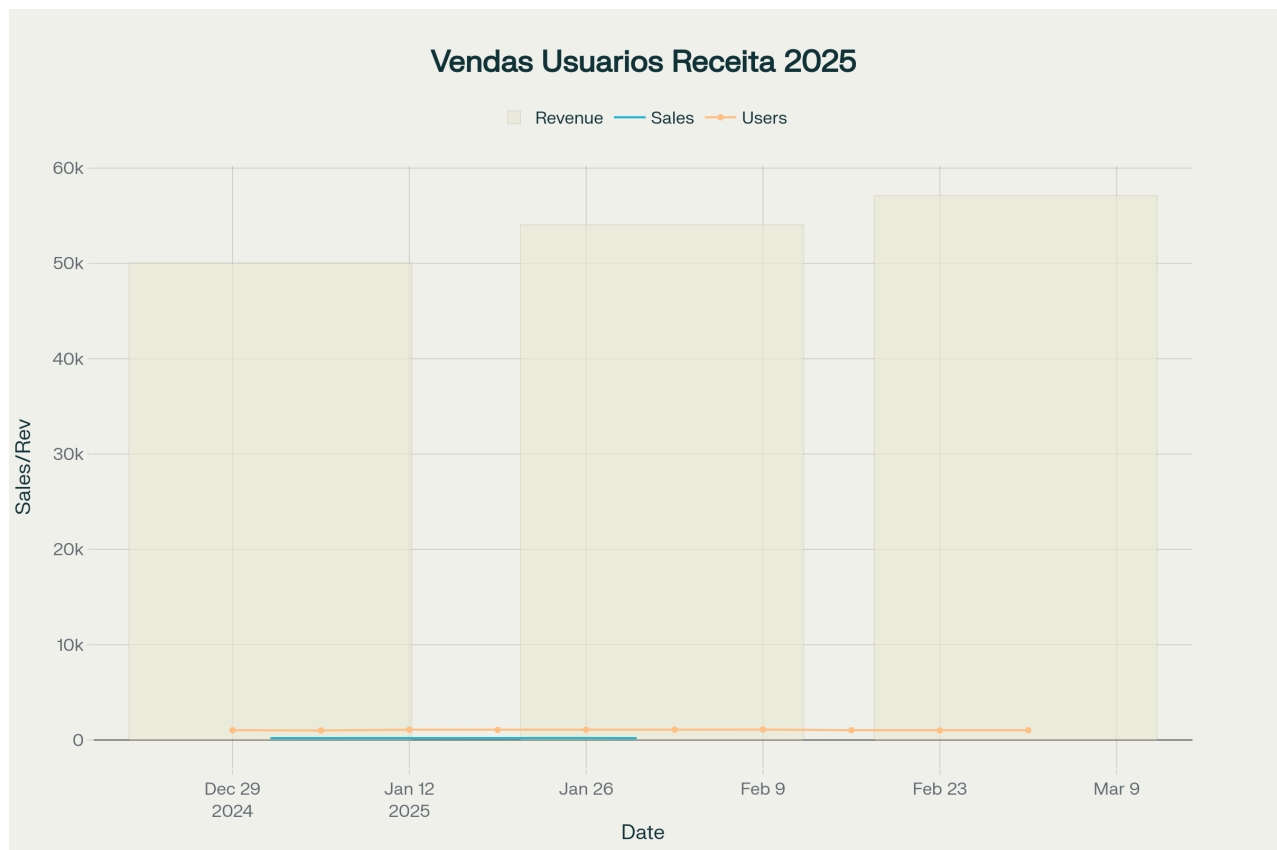
- **Downsample** (`resample('W').mean()`) reduz ruído.
- **Upsample** (`resample('D').ffill()`) preenche lacunas para visualizações contínuas.
- A combinação **asfreq + interpolação** gera transições suaves entre frequências.

### 3.3 Pipeline de Integração

1. Gerar cada série de forma independente.
2. Guardar em dicionários ou em arquivos Parquet quando o volume excede a RAM.
3. Carregar e padronizar colunas (`rename`, `astype`).
4. Executar `pd.concat(axis=1)` seguido de `sort_index`.

Após esses passos, o dataset está pronto para Plotly.

A ilustração abaixo exhibe um resultado típico, combinando 30 pontos diários, 10 semanais e 3 mensais, com uso de dois eixos-y:



Frequências temporais múltiplas

## 4. Visualização Avançada em Plotly

### 4.1 Gráficos Line-Bar Mix

Plotly permite empilhar **barras semitransparentes** (receita) com **linhas** (vendas e usuários) usando `go.Bar` e `go.Scatter`, enquanto o parâmetro `secondary_y` de `make_subplots` define escalas independentes.

```
from plotly.subplots import make_subplots
import plotly.graph_objects as go

fig = make_subplots(specs=[[{"secondary_y": True}]]
fig.add_trace(go.Scatter(x=df.index, y=df.sales,
                        name='Vendas', line=dict(color='royalblue')), secondary_y=False)
fig.add_trace(go.Scatter(x=df.index, y=df.users,
                        name='Usuários', mode='lines+markers',
                        line=dict(color='darkorange')), secondary_y=True)
fig.add_trace(go.Bar(x=df.index, y=df.revenue,
                    name='Receita', marker_color='green', opacity=0.3), secondary_y=False)
```

## 4.2 Interatividade

- **Hover templates** revelam múltiplas métricas no mesmo tooltip.
- **Rangoslider** (`xaxis_rangeslider_visible=True`) facilita zoom temporal.
- **Dropdown menus** podem alternar escalas logarítmicas ou períodos de agregação.

## 5. Dashboards em Dash

### 5.1 Estrutura Básica

Um aplicativo Dash segue três blocos: **layout**, **callbacks** e **server**. O layout contém componentes `Graph`, `DatePickerRange`, `Dropdown`. Callbacks ligam os widgets aos filtros Pandas:

```
@app.callback(
    Output('graph', 'figure'),
    Input('period-dropdown', 'value'),
    Input('date-range', 'start_date'),
    Input('date-range', 'end_date'))
def update(period, start, end):
    mask = (df.index >= start) & (df.index <= end)
    subset = df.loc[mask]
    if period != 'D':
        subset = subset.resample(period).agg({'sales': 'sum', 'users': 'mean', 'revenue': 'sum'})
    return build_plotly_figure(subset)
```

### 5.2 Timezones e Localização

Para ambientes globais, converta sempre para **UTC** logo após a geração (`tz_localize('UTC')`) e apenas então exiba em `'America/Sao_Paulo'` para inspeções locais, evitando discrepâncias em horário de verão.

### 5.3 Desempenho

Séries sintéticas podem escalar a milhões de linhas. Empregue **caching**, **DataTables lazy loading** e, se necessário, **WebGL-rendered traces** (`Scattergl`) para alta fluidez.

## 6. Boas Práticas, Armadilhas e Cenários Especiais

Desafio	Estratégia de Mitigação
Ruído excessivo degrada insights	Ajustar <code>scale</code> de <code>np.random.normal</code> e aplicar filtros <code>ewm</code>
Séries com granularidades muito díspares	Normalizar escalas ou usar subplots independentes
“Buracos” por feriados regionais	<code>CustomBusinessDay(holidays=lista)</code> evita datas inexistentes
Dependências de performance em notebooks	Persistir em Parquet ou Feather e importar via Dask

Outros cenários incluem **calendários fiscais 4-4-5**, uso de **frequências "Q-APR"** para empresas que encerram o exercício em abril, e a simulação de **dados irregularmente espaçados** (IoT), onde `pd.to_timedelta(np.random.exponential(scale))` define intervalos estocásticos.

## Conclusão

A combinação de **NumPy para aleatoriedade controlada** e **Pandas para gerenciamento temporal** viabiliza a geração de datasets sintéticos robustos, possibilitando prototipar dashboards Plotly/Dash que espelham fielmente as complexidades dos ambientes de produção. Seguindo as estratégias descritas — modelagem de tendência e sazonalidade, consolidação de múltiplas granularidades e boas práticas de performance — desenvolvedores podem antecipar desafios de visualização, validar design de interações e demonstrar valor analítico antes mesmo de o ETL oficial estar disponível.

