

# **Previsão de vendas - Análise exploratória e modelagem**

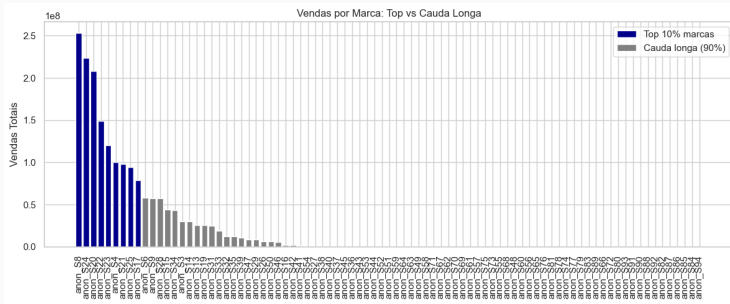
---

A precisão na previsão de vendas (*sell-out*) é um pilar estratégico para a eficiência operacional. Essa previsão alimenta diretamente o S&OP (Sales and Operations Planning), influenciando desde o plano de produção industrial até a estratégia de distribuição. O desafio é criar um modelo que sirva como um *input* confiável para esse processo decisório.

## Objetivo geral

Desenvolver uma solução de previsão de vendas que seja robusta, escalável e defensável. Aqui nós vamos apresentar uma solução completa, considerando os desafios práticos de uma implementação no mundo real, ou seja, construir um modelo preditivo robusto por segmento, validar estatisticamente e justificar o uso no negócio.

## Diagnóstico de cauda longa nas marcas



- Top 10% das marcas concentram a maior parte das vendas;
- Cauda longa representa 90% das marcas com baixo volume;
- **Hill**:  $0.61 \Rightarrow$  provável Pareto  $\Rightarrow$  variância infinita;
- Teste KS  $\Rightarrow$  confirma cauda longa tipo-Pareto.

# Justificativa estatística da cauda longa

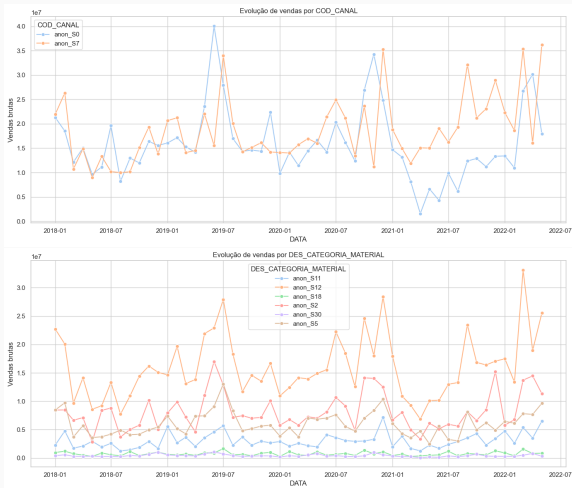
- **Distribuição de Pareto:** evidencia que poucas marcas concentram grande parte da receita;
- **Índice de Hill:** estima o peso da cauda  $\Rightarrow$  valor  $> 0.05$  indica *variância infinita*, dificultando previsibilidade;

## Valor para o negócio:

- Ajuda a identificar segmentos com alta instabilidade e risco;
- Pode direcionar foco preditivo para top marcas;
- Fundamenta decisão sobre agrupamentos ou tratamentos diferenciados.

- Top 10% marcas dominam o faturamento;
- Cauda longa sugere política de sortimento e previsão diferenciadas.

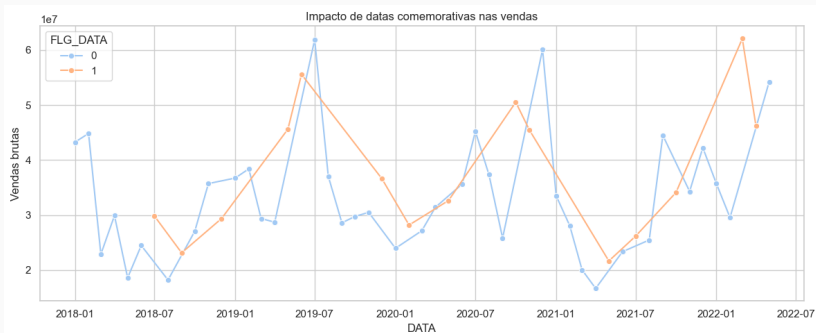
# Evolução temporal por canal e categoria



- Vendas variam muito de um período para outro, sem seguir um padrão simples  $\Rightarrow$  exige modelos flexíveis, ou seja, que consigam adaptar-se rapidamente a mudanças. Em alguns casos até apresenta certa sazonalidade.
- SARIMAX com variáveis exógenas (como eventos) pode capturar essas variações;
- Segmentar os dados por canal/categoria, também pode capturar os padrões específicos de cada grupo.



# Impacto de datas comemorativas



- Eventos sazonais (Natal, dia das mães, etc.) têm forte impacto e podem ser usados como variáveis exógenas no SARIMAX para melhorar previsibilidade nos picos.

## Observação

Variáveis exógenas (ou variáveis explicativas externas) são variáveis que influenciam a série temporal que queremos prever, mas não são influenciadas por ela  $\Rightarrow$  são fornecidas ao modelo como *inputs* adicionais, pois carregam informação relevante para a previsão.

# Objetivos da modelagem

- Prever vendas com boa acurácia;
- Fornecer *input* confiável para o processo de S&OP;
- Avaliar robustez estatística por canal e categoria;
- Detectar limitações e propor melhorias.

- Modelo: **SARIMAX** com variável exógena (datas comemorativas);
- Segmentação por canal e por categoria;
- Métricas: MAE, MASE, wAPE, Bias;
- Testes estatísticos: ADF, KPSS, Shapiro-Wilk, Ljung-Box.

## Importância das métricas para o negócio:

- **MAE** impacto direto em volume de produção e custo logístico;
- **MASE**: razão entre erro do modelo e erro do naïve (referência mínima de desempenho)  $\Rightarrow$  mostra o quanto o modelo melhora a previsão;
- **wAPE**: indica percentual de erro em relação ao volume real  $\Rightarrow$  críticos para orçamento e faturamento  $\Rightarrow$  soma total dos valores reais, o que faz com que as observações com vendas maiores “pesem” mais no erro final;
- **Bias**: identifica tendência sistemática de subestimação ou superestimação.

## Importância dos testes estatísticos:

- **ADF / KPSS:** testam estacionariedade (propriedades estatísticas não mudam ao longo do tempo  $\Rightarrow$  requisito para SARIMAX;
- **Shapiro-Wilk:** avalia normalidade dos resíduos  $\Rightarrow$  fundamental que o modelo capture bem os padrões estruturais da série;
- **Ljung-Box:** verifica se resíduos são ruído branco  $\Rightarrow$  garante que o modelo não deixou padrões não explicados.

## Observação

Ao ajustar um modelo de séries temporais (como SARIMAX), espera-se que os resíduos (erros) atendam a dois critérios fundamentais:

- Não tenham padrão temporal, isto é, que se comportem como um ruído branco;
- Não sejam autocorrelacionados  $\Rightarrow$  o erro de hoje não deve depender do erro de períodos anteriores.

## ■ Shapiro-Wilk:

- anon\_S7, S12:  $p < 0.05 \Rightarrow$  não normal;
- demais:  $p > 0.05 \Rightarrow$  normal

## ■ Ljung-Box:

- Todos os segmentos com  $p > 0.25 \Rightarrow$  sem autocorrelação significativa.



### Segmentos com bom/médio desempenho:

- **anon\_S30**:  $MASE = 0.686$ ,  $wAPE = 22.49\%$ ,  $Bias = 37636.43$ , resíduos normais;
- **anon\_S18**:  $MASE = 0.751$ ,  $wAPE = 28.49\%$ ,  $Bias = -305771.62$ , resíduos normais;
- $Bias(\text{anon\_S30}) \Rightarrow$  o modelo está ligeiramente superestimando as vendas;
- $Bias(\text{anon\_S18}) \Rightarrow$  o modelo está subestimando as vendas em cerca de R\$ 300 mil/mês.

## **Fatores que explicam este desempenho:**

- Volume de vendas e variância mais estáveis;
- Padrões sazonais regulares;
- Resíduos sem autocorrelação;
- $MASE < 1$ .

## Segmentos com erro elevado:

- **anon\_S0, S7, S5, S12, S2, S11**: wAPE acima de 50%;
- Bias negativo (na casa dos milhões)  $\Rightarrow$  subestimação elevadíssima;
- $MASE > 2.0$  em vários casos  $\Rightarrow$  modelo perde para baseline naïve.

## Possíveis causas:

- Cada segmento reage de forma diferente à variável exógena utilizada.

- Erro maior em períodos de pico extremo, mesmo com variável exógena;
- Subestimação sistemática das vendas;
- Erros elevados em segmentos com alta volatilidade;

## Próximos passos recomendados

1. Variável binária FLG\_DATA mais granular, com diferenciações entre os diversos tipos de datas comemorativas, i.e., trabalhar com mais variáveis exógenas pode melhorar os resultados bons, médios e até mesmo aqueles que foram insatisfatórios;
2. Avaliar modelos alternativos, como o Prophet, por exemplo.