



Predictiva.ai

Live #10

Os principais erros em Análise de Dados e como evitá-los

Parte 1

21/02/2023

Antes de mais nada, já se inscreveu em nosso canal para ter acesso aos materiais e avisos das lives?



<https://t.me/preditiva>

Erros acontecem...

A definição formal

Risco operacional

Probabilidade de ocorrência de perdas resultantes de eventos externos ou de falha, deficiência ou inadequação de processos internos, pessoas ou sistemas.

Erros acontecem...

A definição que eu gosto é...

Risco operacional

=

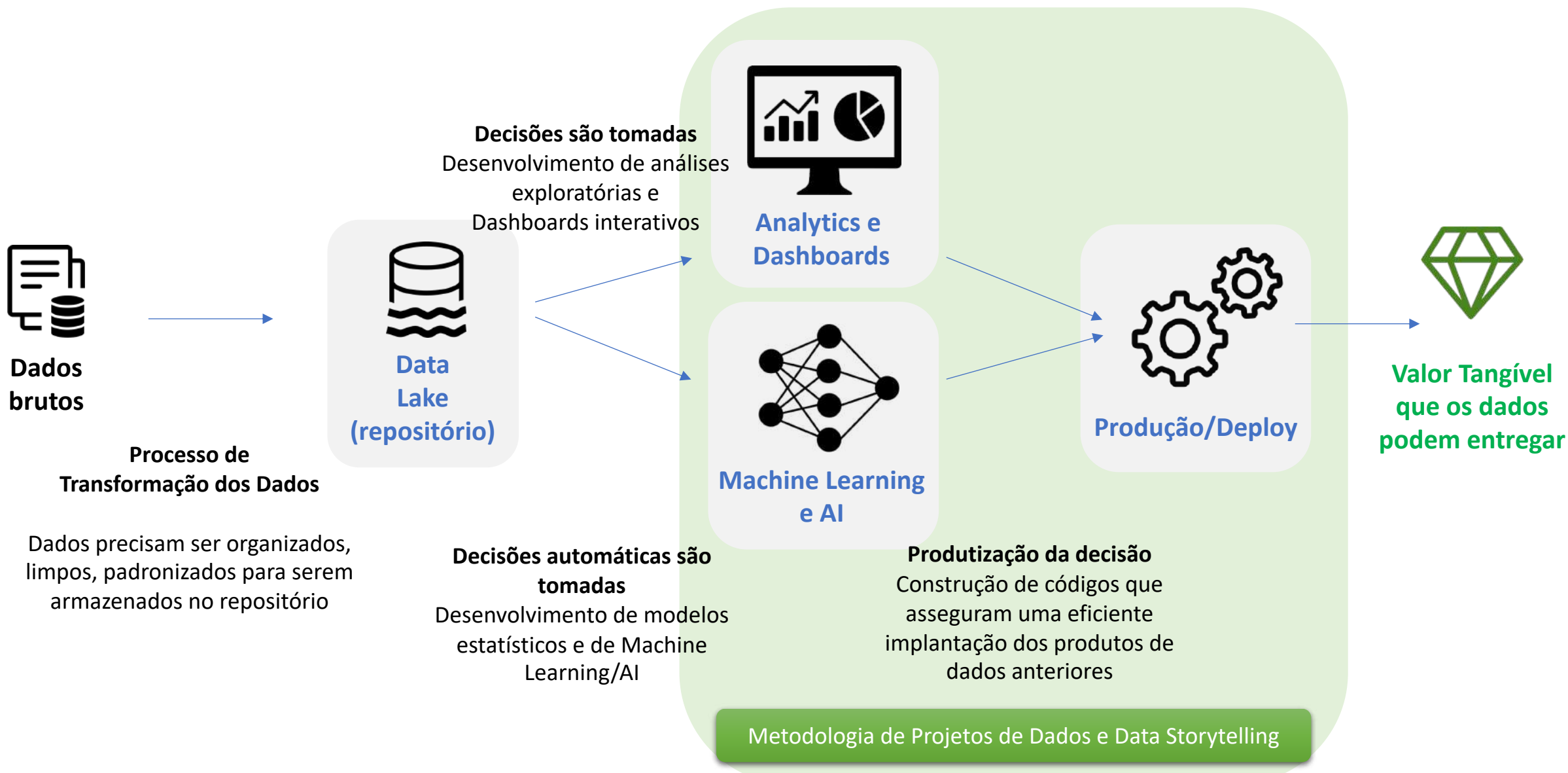
Probabilidade de você fazer uma CAGADA!



Vish

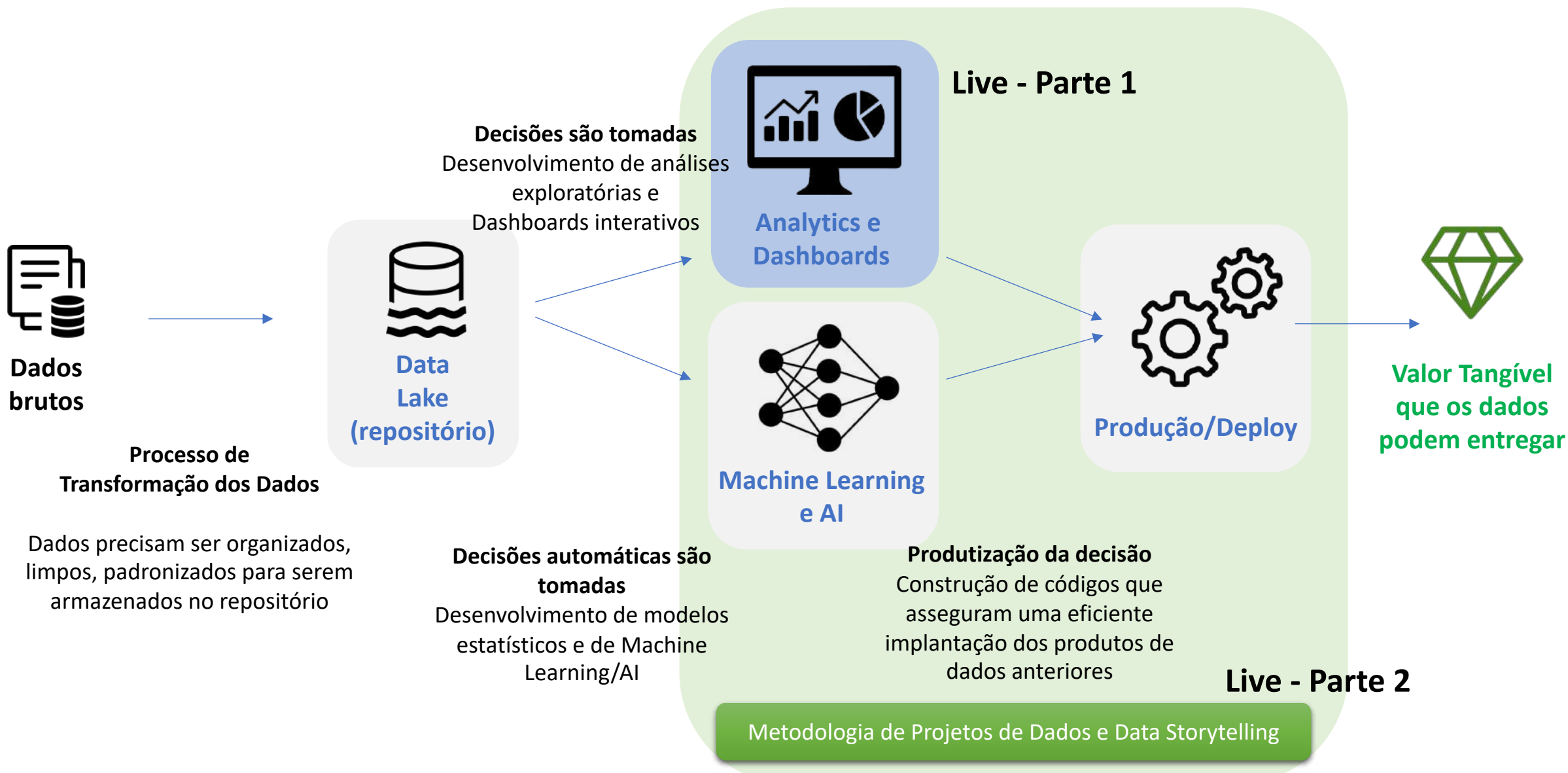
Jornada dos dados

Existem muitas chances de erros na jornada dos dados



Jornada dos dados

Existem muitas chances de erros na jornada dos dados



O mapa da CAGADA em Dados!

Parte 1



Medidas e Gráficos

Medidas estatísticas

Unidades da base

Medidas vs Gráficos

Correlações

Causalidade

Projeções

Frequências imprecisas

**Probabilidade e
Possibilidade**

Amostragem

Viés de seleção

Erros de Coleta

Tamanho de Amostra



O mapa da CAGADA em Dados!

Parte 1



Medidas e Gráficos

Medidas estatísticas

Unidades da base

Medidas vs Gráficos

Correlações

Causalidade

Projeções

Frequências imprecisas

**Probabilidade e
Possibilidade**

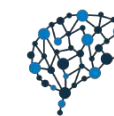
Amostragem

Viés de seleção

Erros de Coleta

Tamanho de Amostra





O que você faria na seguinte situação?

Você é convidado para trabalhar em uma startup com 15 funcionários e, segundo o

RH o salário **médio** dos funcionários é

R\$4.200,00.

Você atualmente ganha R\$1.000,00.

Funcionário	Salário
1	1.000,00
2	1.000,00
3	1.000,00
4	1.000,00
5	1.000,00
6	1.000,00
7	1.000,00
8	1.000,00
9	1.000,00
10	1.000,00
11	1.000,00
12	1.000,00
13	1.000,00
14	10.000,00
15	40.000,00
Média	4.200,00

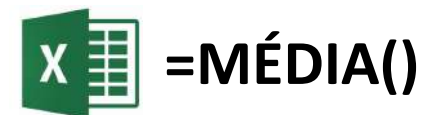
Medidas e Gráficos

Medidas Resumo: Medidas de Posição - Média

ID	Salário
1	5.130,00
2	4.193,00
3	3.468,00
4	3.068,00
5	2.670,00
6	2.693,00
7	9.526,00
8	3.068,00
9	5.237,00
10	9.980,00
11	2.426,00
12	2.911,00

A **Média** é uma **medida de tendência central**. Isto é, ela tenta nos dar uma noção de onde o valor central dos dados deve aparecer. Veja:

$$\text{Média} = \frac{54.370,00}{12} = 4.530,83$$



Porém, muitas vezes **ela não é a medida resumo mais indicada**. Veremos isso mais adiante.

Medidas e Gráficos

Medidas Resumo: Medidas de Posição - Mediana



A **Mediana**, assim como a média, também é uma medida de tendência central. Porém, no seu caso, ela realmente representa o centro do seu conjunto de dados.

Ou seja, 50% dos valores são inferiores à **Mediana** e 50% dos valores são superiores à **Mediana**.

Obs: Quando o número de observações é par (no exemplo, temos 12 salários), a **Mediana** é a média entre os 2 valores centrais. Se fosse ímpar, seria o próprio valor central.



ID	Salário
11	2.426,00
5	2.670,00
6	2.693,00
12	2.911,00
4	3.068,00
8	3.068,00
3	3.468,00
2	4.193,00
1	5.130,00
9	5.237,00
7	9.526,00
10	9.980,00

Mediana
3.268,00

O que você faria na seguinte situação?

Você está doente e só existe um remédio que pode te ajudar. Segundo a bula, o tempo de sobrevida **mediano** é de 8 semanas. Todos que tomam o remédio tem severos efeitos colaterais. Vale a pena tomar o remédio?

Tempo de Sobrevida (em semanas)	% Relativa
2	15%
4	15%
8	20%
16	5%
32	5%
64	5%
128	5%
256	10%
512	20%

O erro:

Usar apenas **uma** medida para resumir os dados.

Como evitá-lo?

Quanto mais medidas resumo usar, melhor conseguirá entender os dados analisados. Exemplos de medidas:

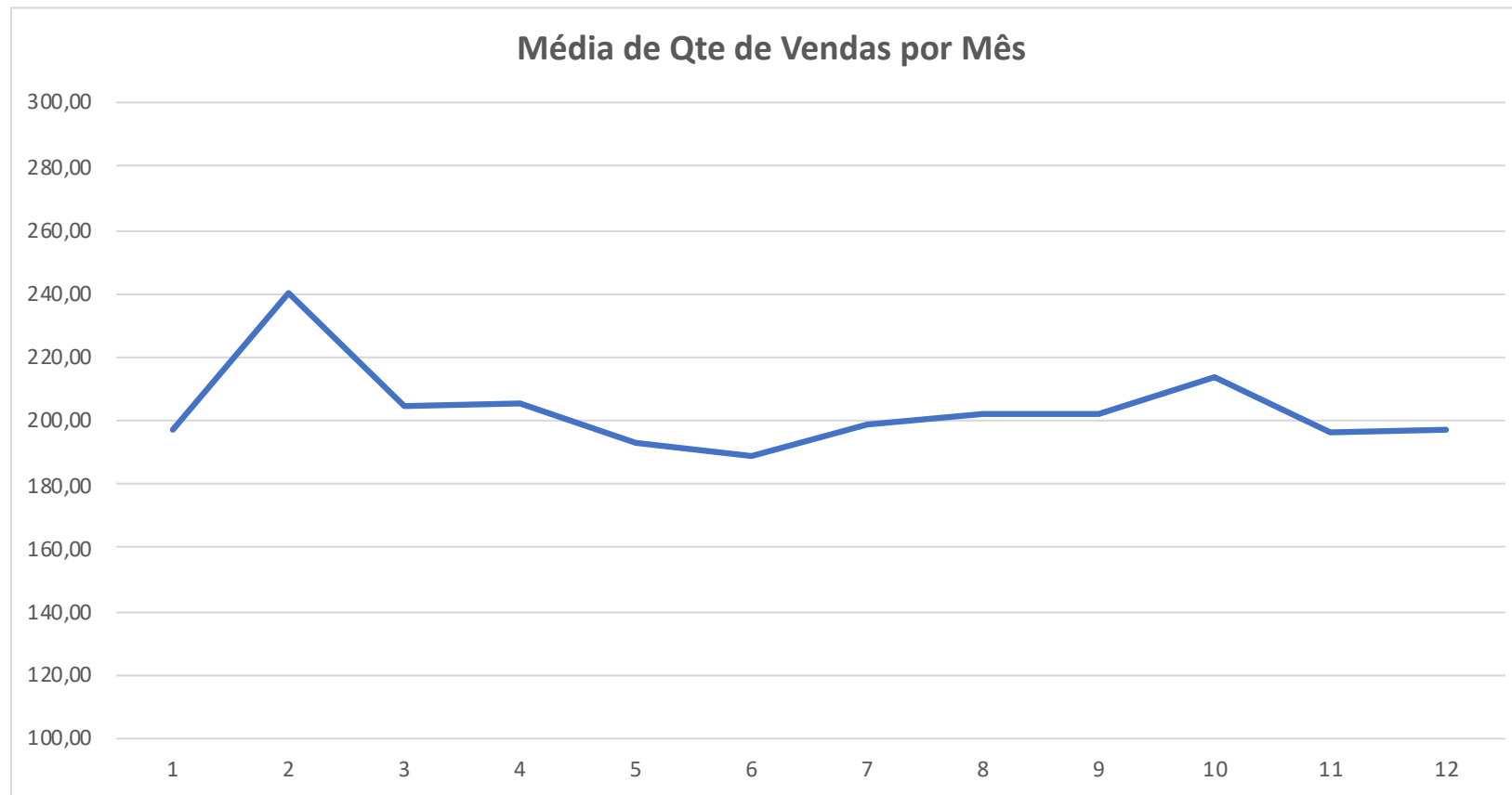
- Média
- Mediana
- Quartis
- Mínimo e Máximo
- Desvio Padrão

Medidas e Gráficos

Exemplo 2



Você pede para alguém um relatório de vendas anual. O analista entrega isso...



Conclusões mais comuns:

1. Fevereiro foi um mês atípico de vendas.
2. De março a setembro ficamos estáveis em cerca de 200 vendas por mês.

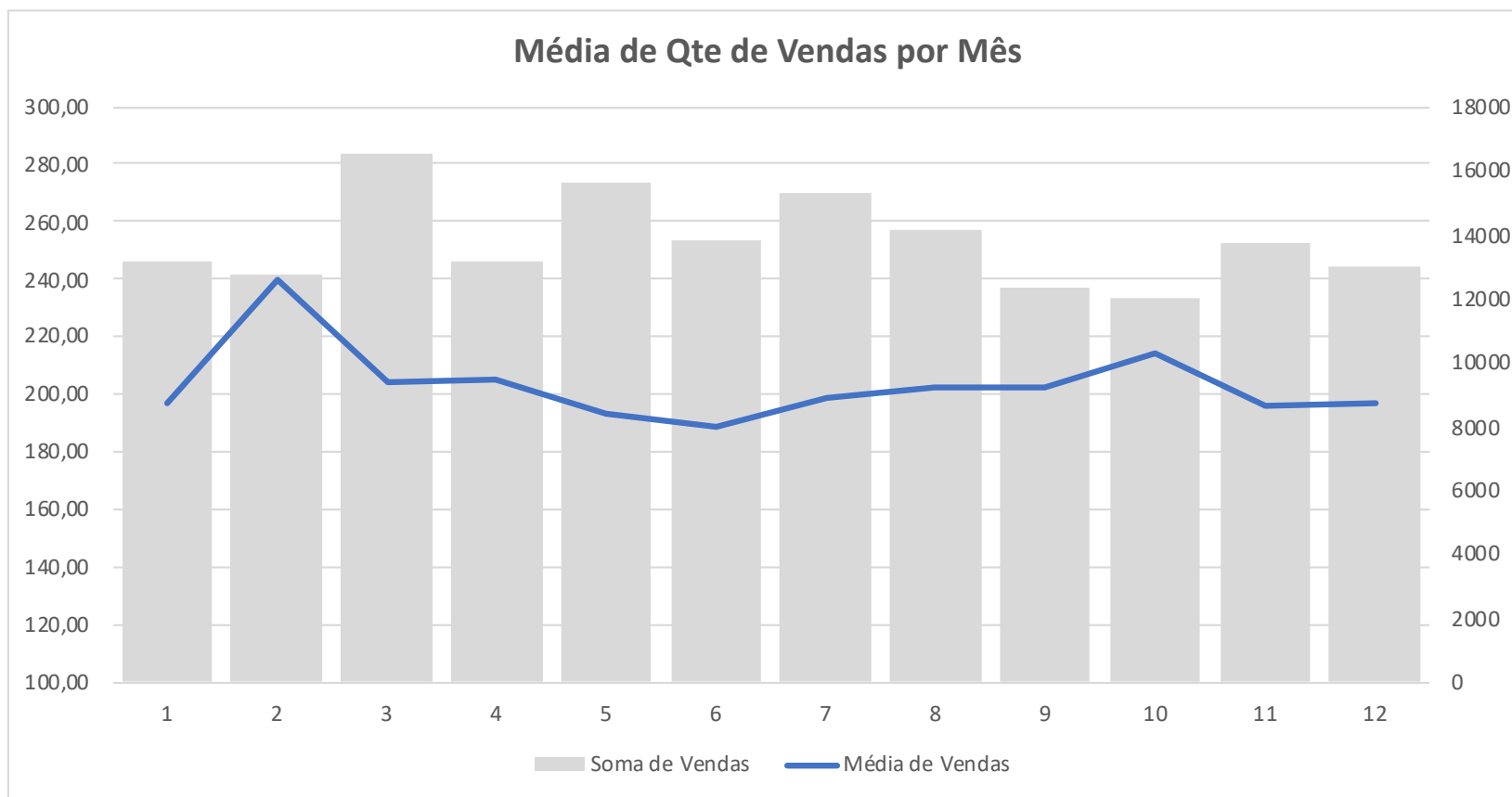
Medidas e Gráficos

Exemplo 2



Preditiva.ai

Para tentar entender melhor, você pede a soma de vendas por mês. O analista te entrega isso...



Sua conclusão mudou de alguma forma?

Será que conseguimos realmente entender nossas vendas com esse resumo?

E se tivéssemos acesso aos dados brutos para um *Double Check*?

O erro:

Assumir a unidade da base utilizada **sem** perguntar para o analista gerador da análise.

Como evitá-lo?

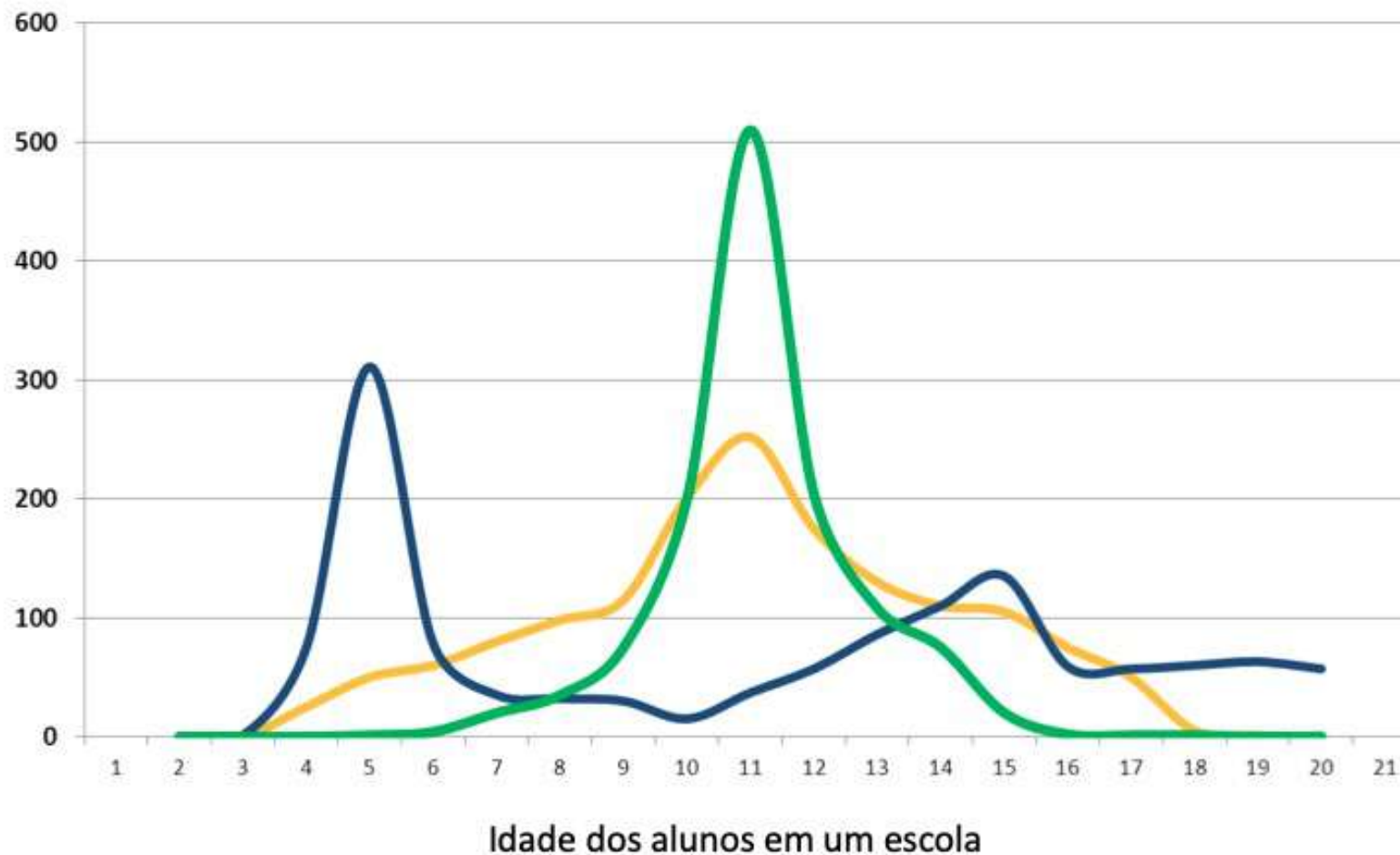
Deixar claro qual unidade da base deve ser trabalhada (tanto se você for o cliente do trabalho quanto você for o criador da análise).

Pergunte também:

- Filtro de tempo usado (Diário, Semanal, Mensal ?)
- Se existem outros tipos de filtros ocultos (ex: produto, segmentos, regiões etc)

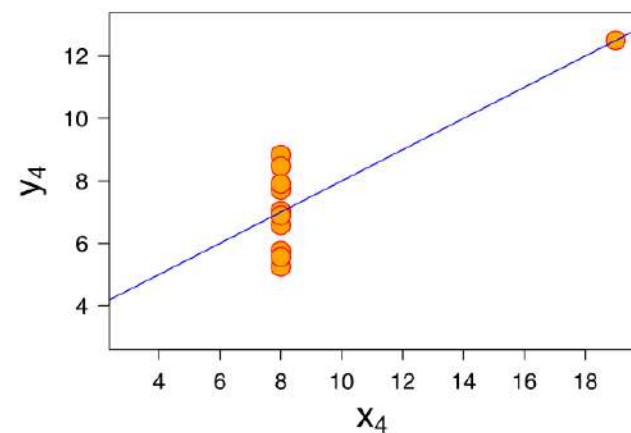
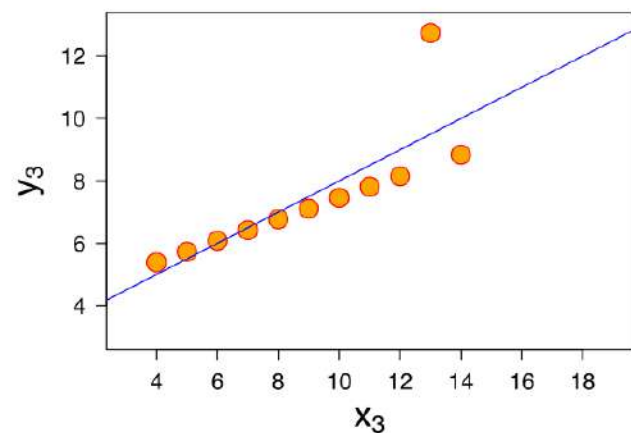
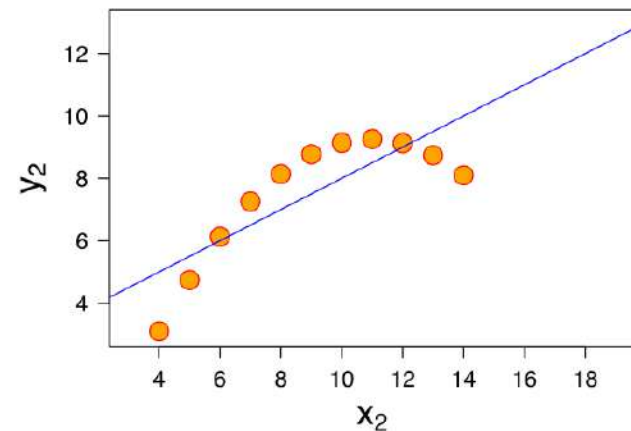
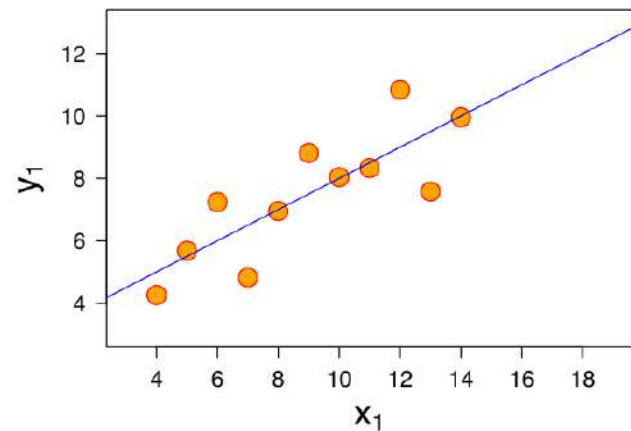


Qual linha tem a maior média?



Medidas e Gráficos

Exemplo 4 - O que esses quatro conjuntos têm em comum?



Propriedade	Valor
Média em x	9
Variância em x	11
Média em y	7.50
Variância em y	4.125
Correlação entre x e y	0.816
Regressão linear	$y = 3.00 + 0.500x$
R^2	0.67

Quarteto de Ascombe

O erro:

As vezes **nem muitas medidas resumo são suficientes** para entender seus dados.

Como evitá-lo?

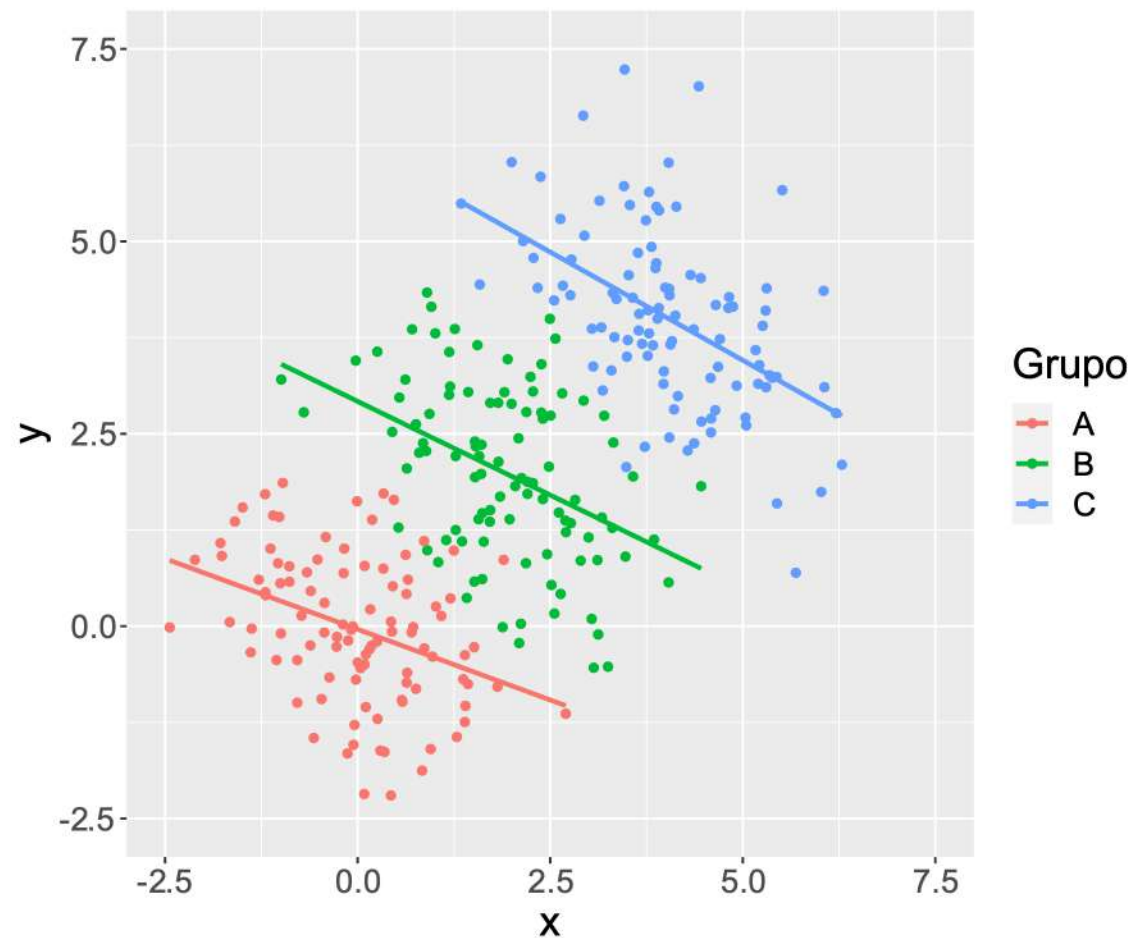
Sempre que possível plotar os dados em um gráfico analítico e comparar com as medidas resumo.

Gráficos mais indicados:

- Histogramas
- Boxplots
- Gráficos de Dispersão (Scatter Plots)



Paradoxo de Simpson



O erro:

A **correlação entre seus dados pode mudar** ao agrupá-los de formas diferentes.

Como evitá-lo?

Sempre plote o gráfico para cada grupo mais importante de sua base de dados.

Agrupamentos possíveis:

- Por tempo
- Por produto
- Por região
- Por dados demográficos

Medidas e Gráficos

Exemplo 6



Você pede para um analista um relatório de vendas de um determinado produto. A ideia é saber se o gênero explica a compra desse produto. O analista te entrega isso:

Gênero	Unidades Vendidas	Frequência de Compra (%)
Masculino	400	20%
Feminino	1600	80%

2000

Conclusões mais comuns:

Produto claramente voltado para o público feminino.

Medidas e Gráficos

Exemplo 6

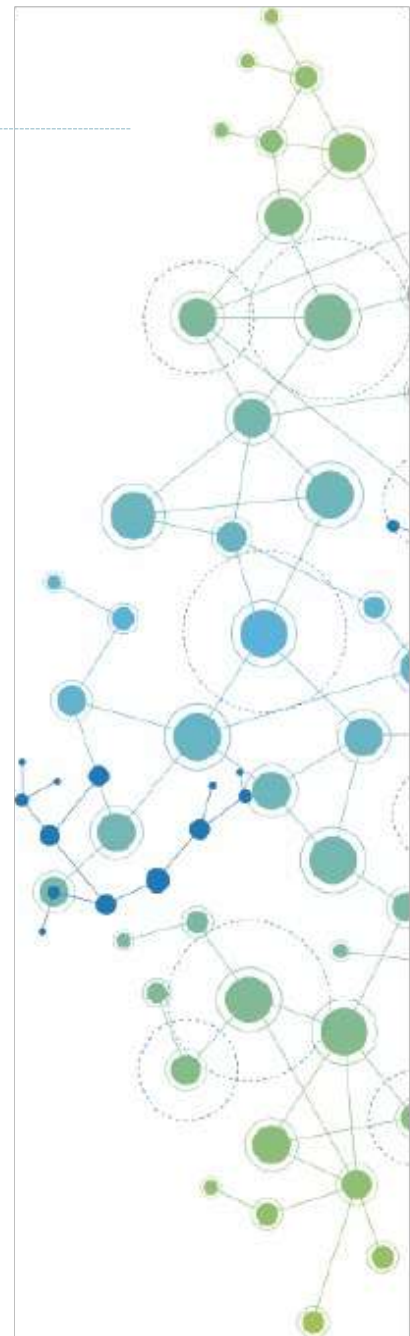
Essa análise simples não considera o fator **contrafactual**, ou seja, não considera as compras **não** realizadas por cada um dos gêneros. Teria a mesma proporção? Vejamos:

Gênero	Unidades enviadas para o carrinho de compras	Frequência de Compra (%)	Frequência de Não Compra (%)	Diferença
Masculino	3600	20%	19%	1%
Feminino	7200	80%	81%	-1%

6000

E agora? A conclusão muda?

Demonstração desse erro



Medidas e Gráficos

Exemplo 7 – Correlação de Pearson não é a única forma de medir



Turn Over?	Idade	Frequência de Viagens	Distância do trabalho
Sim	41	Viaja raramente	1
Não	49	Viaja frequentemente	8
Sim	37	Viaja raramente	2
Não	33	Viaja frequentemente	3
Não	27	Viaja raramente	2
Não	32	Viaja frequentemente	2
Não	59	Viaja raramente	3
Não	30	Viaja raramente	24
Não	38	Viaja frequentemente	23
Não	36	Viaja raramente	27
Não	35	Viaja raramente	16
Não	29	Viaja raramente	15
Não	31	Viaja raramente	26
Sim	34	Viaja raramente	19
Não	28	Viaja raramente	24
Não	29	Viaja raramente	21
Não	32	Viaja raramente	5
Não	22	Não viaja	16
Não	53	Viaja raramente	2
Não	38	Viaja raramente	2
Sim	24	Não viaja	11
Não	36	Viaja raramente	9
Não	34	Viaja raramente	7
Não	21	Viaja raramente	15
Sim	34	Viaja raramente	6
Não	53	Viaja raramente	5
Sim	32	Viaja frequentemente	16
Não	42	Viaja raramente	16

Information Value:

Uma das melhores técnicas de análise de dados

Muito trabalho, só fazer uma matriz de correlação

13w 8 likes Reply See translation

preditiva.analytics a matriz só funciona quando a variável de interesse for numérica. No exemplo que demos a variável é categórica e binária.

13w Reply See translation

@preditiva.analytics transformá o que dá em 1 e 0

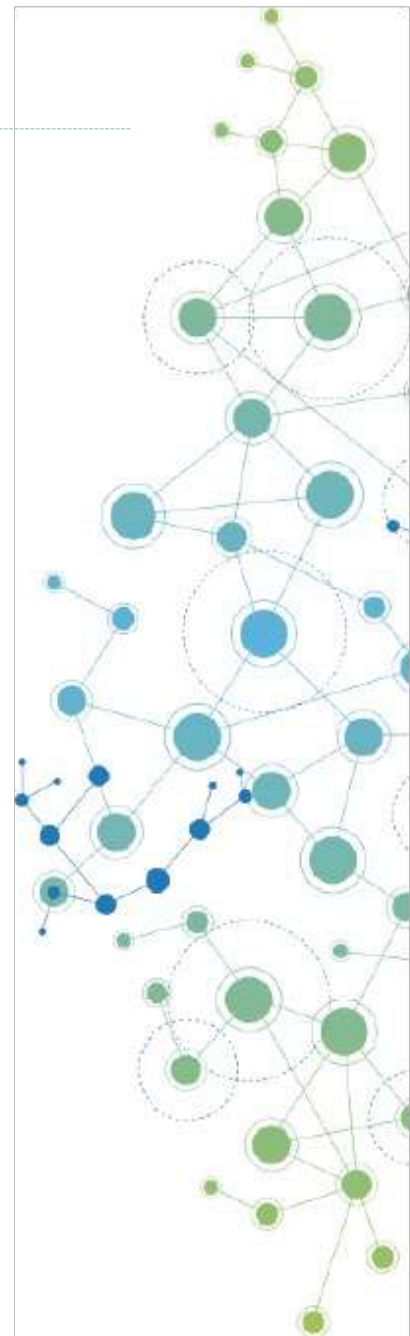
13w Reply See translation

Faz um heatmap da Matriz de Correlação das variáveis po

6w Reply See translation

Muito trabalho, só fazer uma matriz de correlação

Demonstração desse erro



É fundamental dominarmos a **diferença entre esses 2 conceitos** para não cairmos em algumas **armadilhas** de Analytics.

Vejamos a definição destes 2 termos:

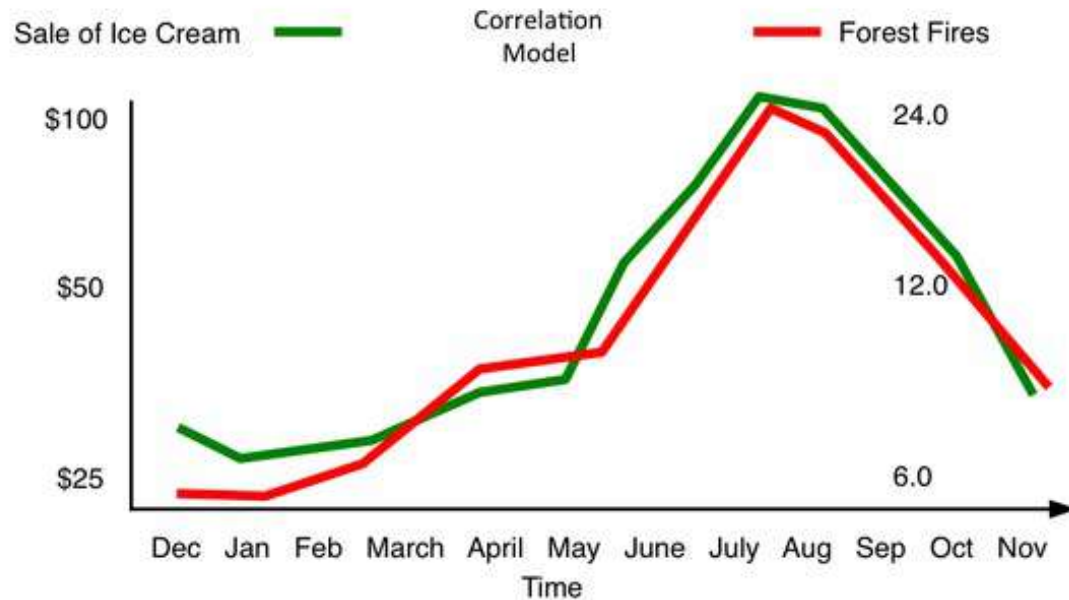
- **Correlação**: relação de **dependência** ou **associação** entre duas variáveis.
- **Causalidade**: relação entre um evento A e um evento B, sendo que o evento B é **consequência** do evento A.

Ou seja, **Correlação** está relacionada com a **dependência ou associação** e a **Causalidade** relacionada a **consequência**.

Medidas e Gráficos

Exemplo 8 - Causalidade

Vamos avaliar a **Correlação** entre **Venda de Sorvetes** e **Incêndio nas Florestas**:



Fonte: <https://www.decisionskills.com/blog/how-ice-cream-kills-understanding-cause-and-effect>

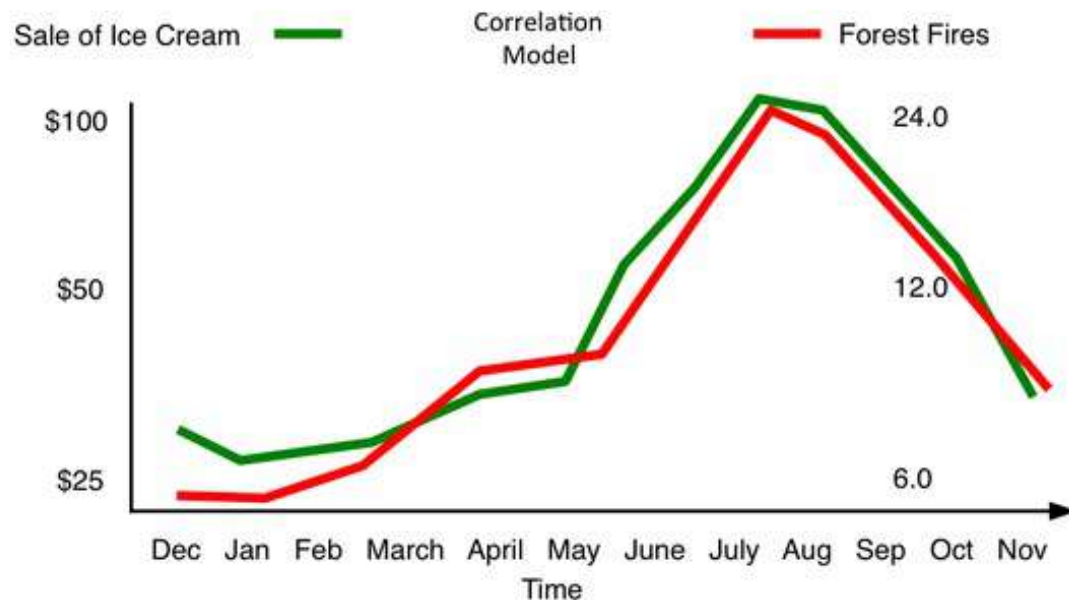
Você acha que a **venda de sorvetes** pode **causar** **incêndios nas florestas**?

Neste caso há uma **3ª variável não avaliada** e que faz mais sentido ser a **causadora** do aumento no **consumo de sorvete** e dos **incêndios nas florestas**: **o clima quente!**

Medidas e Gráficos

Exemplo 8 - Causalidade

Vamos avaliar a **Correlação** entre **Venda de Sorvetes** e **Incêndio nas Florestas**:



Fonte: <https://www.decisionskills.com/blog/how-ice-cream-kills-understanding-cause-and-effect>

A **Correlação** entre **Venda de Sorvetes** e **Incêndio nas Florestas** é conhecida como **Correlação Espúria**.

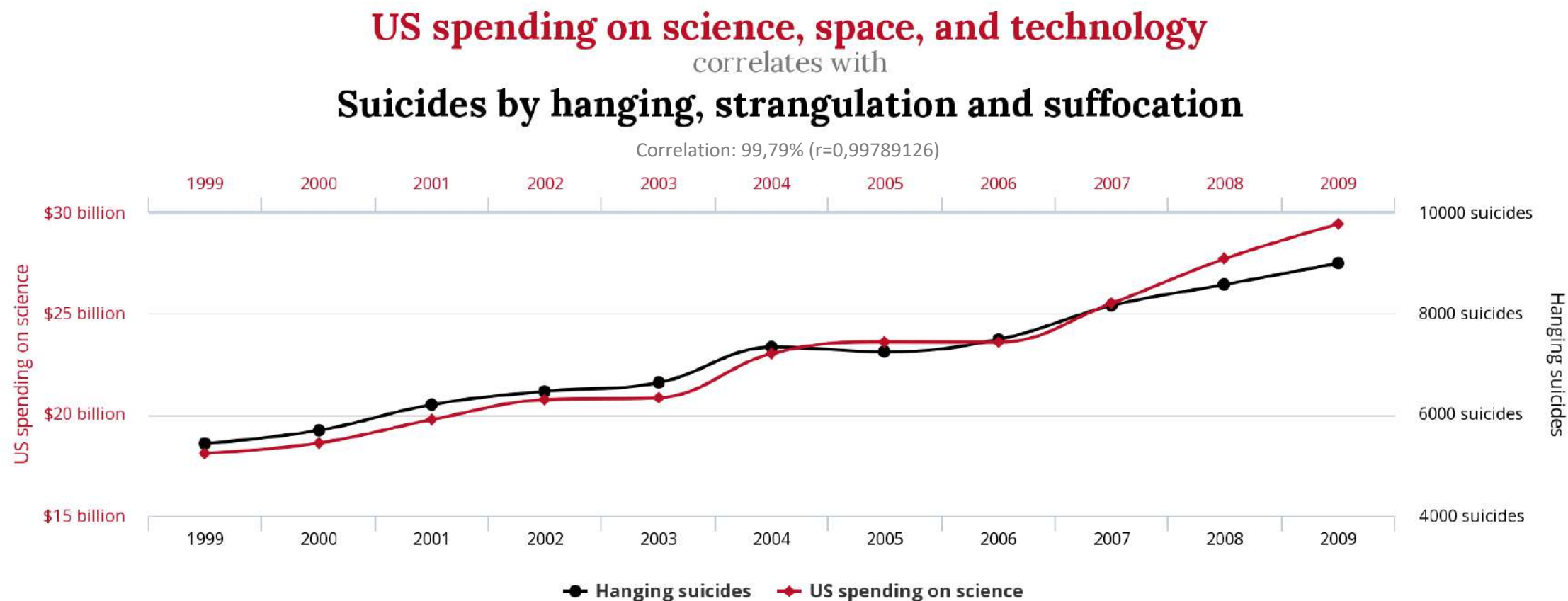
As **Correlações Espúrias** podem ser uma armadilha para **falsas conclusões**.

Vejamos alguns outros exemplos.

Medidas e Gráficos

Exemplo 8 - Causalidade

Exemplo 1: Gasto em Pesquisa no EUA vs. Suicídios

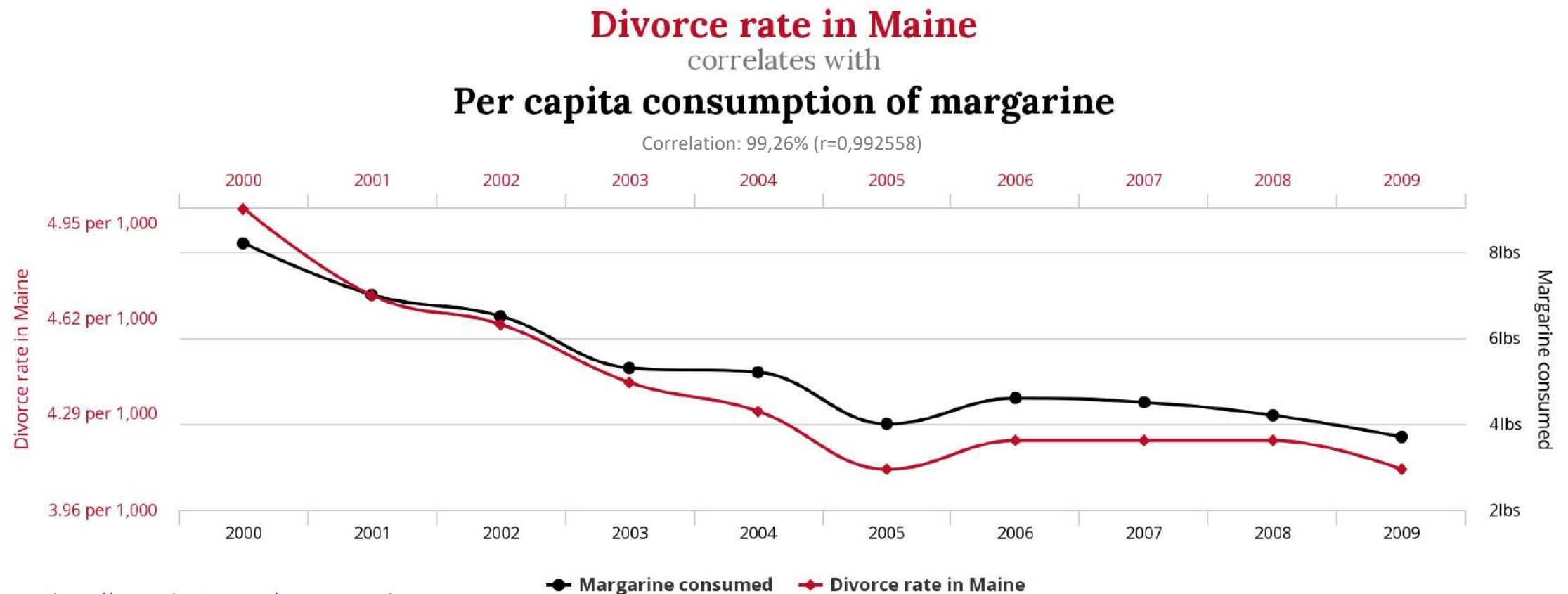


Fonte: <https://www.tylervigen.com/spurious-correlations>

Medidas e Gráficos

Exemplo 8 - Causalidade

Exemplo 2: Divórcios em Maine vs. Consumo de margarina



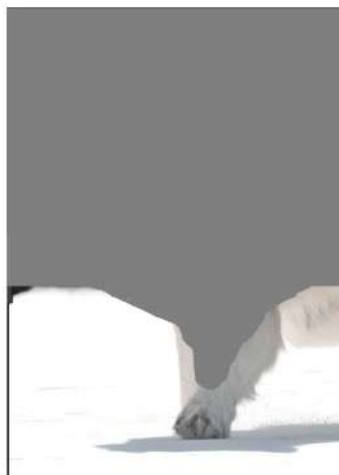
Fonte: <https://www.tylervigen.com/spurious-correlations>

Medidas e Gráficos

Exemplo 8 - Causalidade



Lobo



Lobo



Cachorro

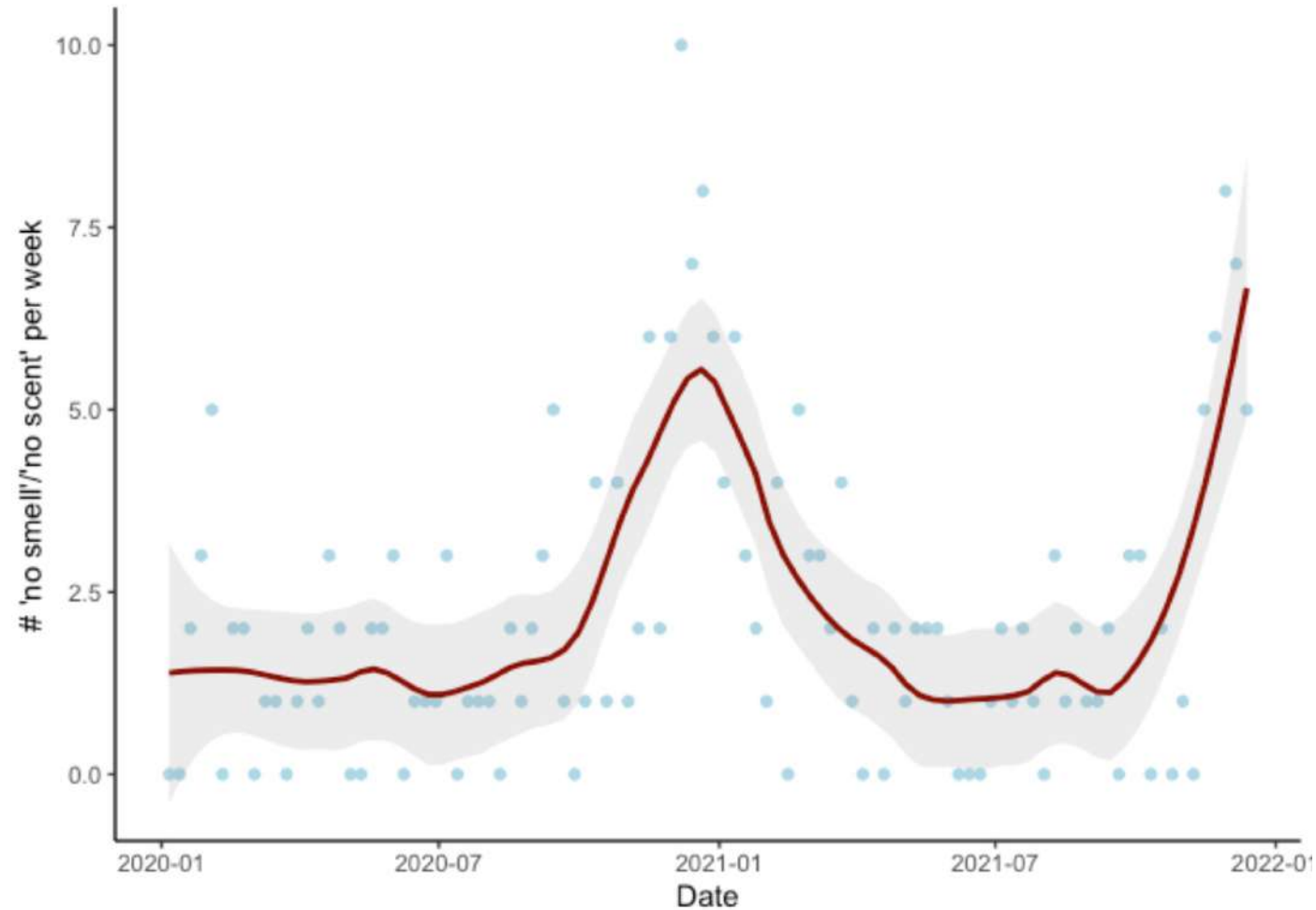


Cachorro

Medidas e Gráficos

Exemplo 8 - Causalidade

Número de reclamações na Amazon para uma vela aromática (sem cheiro)



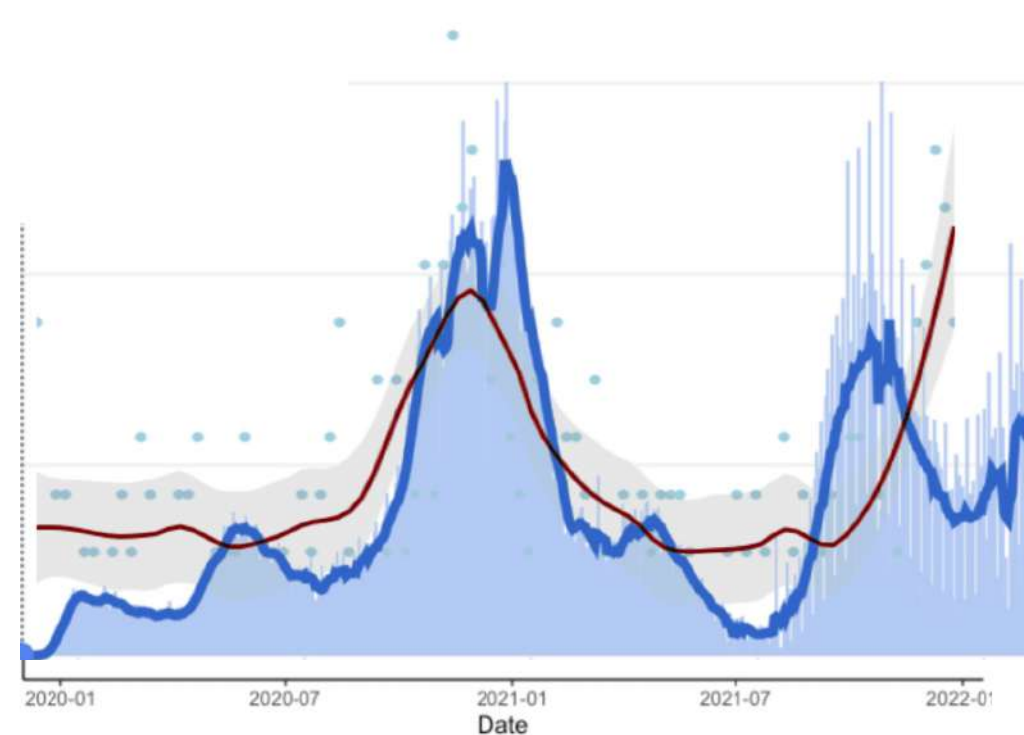
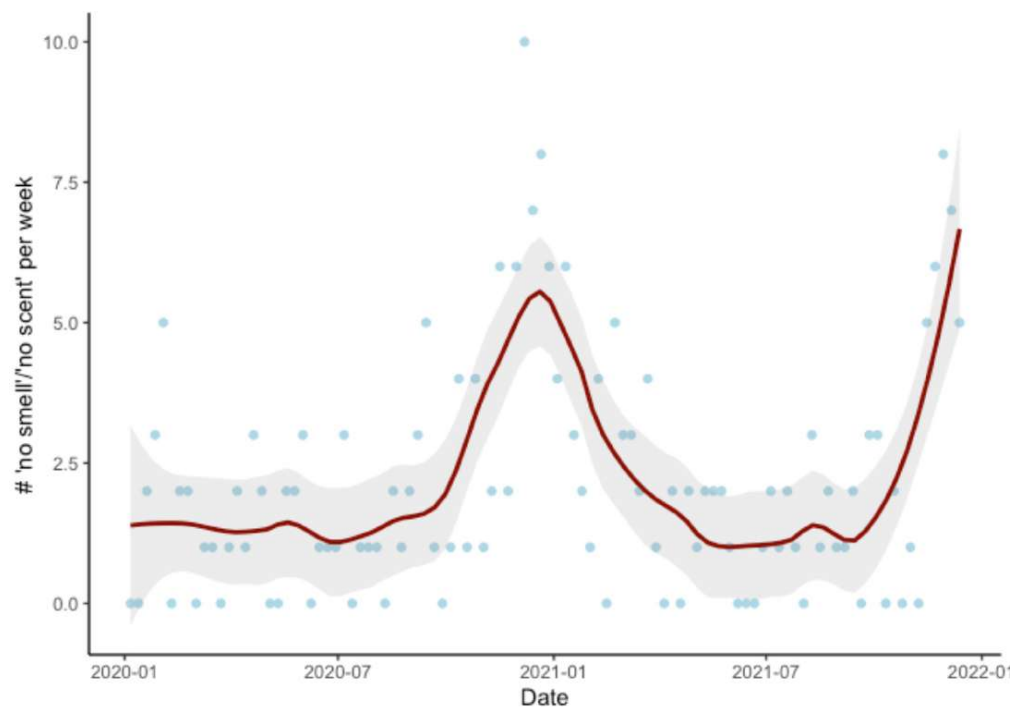
Medidas e Gráficos

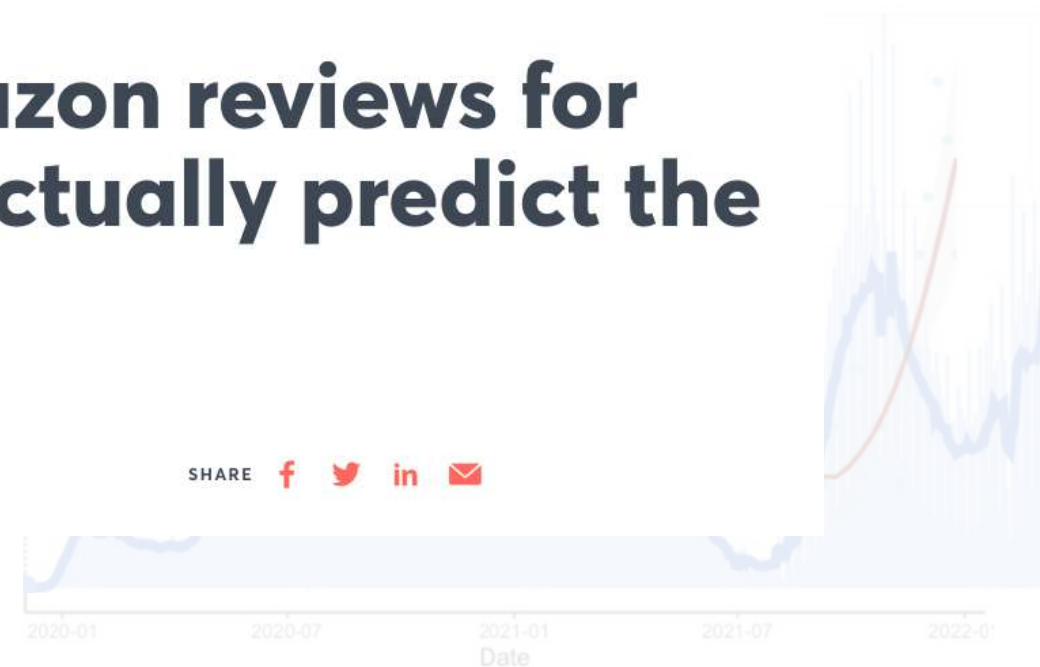
Exemplo 8 - Causalidade



Preditiva.ai

Número de reclamações na Amazon para uma vela aromática (sem cheiro)





Conclusões:

- Se 2 variáveis estão **correlacionadas**, **pode ou não** haver **causalidade**
- Se houver **correlação** e não houver **causalidade** entre essas 2 variáveis, possivelmente há uma **3ª variável que não foi observada**
- **Mantenha-se cético**: busque **fortes evidências** para assumir a **causalidade**
- Antes de assumir a **causalidade** responda as seguintes perguntas:
 - **Por que** a variável A **causa** a variável B?
 - **Como** a variável A **causa** a variável B?

O mapa da CAGADA em Dados!

Parte 1



Medidas e Gráficos

Medidas estatísticas

Unidades da base

Medidas vs Gráficos

Correlações

Causalidade

Projeções

Frequências imprecisas

Probabilidade e
Possibilidade

Amostragem

Viés de seleção

Erros de Coleta

Tamanho de Amostra

Essas cagadas ficam para a parte 2 ... rs



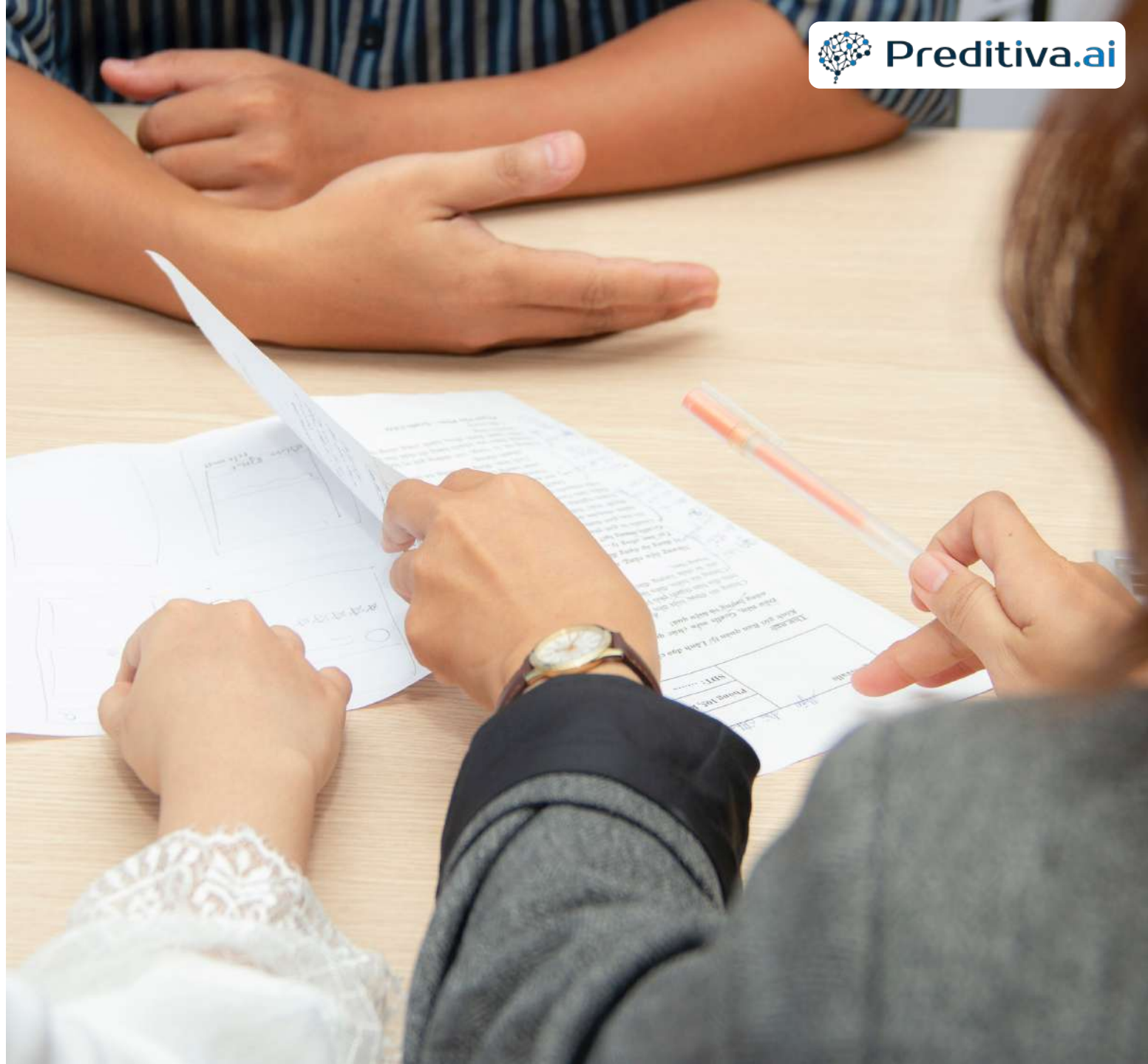
Resumindo os erros da live (Parte 1)

Guarda esse check-list



Cagada 🍌		Como evitar?
1	Usar apenas uma medida para resumir os dados.	Quanto mais medidas resumo usar, melhor conseguirá entender os dados analisados.
2	Assumir a unidade da base utilizada sem perguntar para o analista gerador da análise.	Deixar claro qual unidade da base deve ser trabalhada (tanto se você for o cliente do trabalho quanto você for o criador da análise).
3	As vezes nem muitas medidas resumo são suficientes para entender seus dados.	Sempre que possível plotar os dados em um gráfico analítico e comparar com as medidas resumo.
4	A correlação entre seus dados pode mudar ao agrupá-los de formas diferentes. (Paradoxo de Simpson)	Sempre plote o gráfico para cada grupo mais importante de sua base de dados.
5	Assumir que uma variável causa a outra só porque as variáveis estão correlacionadas. Correlação não significa causalidade.	Mantenha-se cético: busque fortes evidências para assumir a causalidade. Se não estiver confiante da causalidade, evite usar essa correlação.

Você quer
acelerar seu
desenvolvimento
e se sentir **ainda
mais confiante**
para resolver
problemas como
mostrei aqui ?



Como **complemento** à nossa formação principal
“**Gerando Valor com Dados**”, criamos o programa:



ACELERA

Evolua 6 meses de experiência em Dados em apenas 1 mês

Como funciona?

ACELERA

Evolua 6 meses de experiência em dados em apenas 1 mês



○ Acelera é um **programa ao vivo** onde você vai:

- 1** **Resolver problemas reais** utilizando as **técnicas, ferramentas e metodologias** aprendidas na Plataforma Preditiva com o **apoio de mentores experientes**.
- 2** Ter a experiência de trabalhar em um **projeto de dados colaborativo**, com colegas que têm o mesmo objetivo que você e o **apoio de mentores experientes**.
- 3** Treinar seu **Data Storytelling** **apresentando suas análises** para os gestores de diversas áreas e obtendo feedback em tempo real.
- 4** Aprender em **apenas 1 mês**, de forma **100% prática** a extrair **insights dos dados**, e estar pronto para buscar sua transição de carreira.

ACELERA

Evolua 6 meses de experiência em dados em apenas 1 mês



Carga Horária: 15h

Serão 5 encontros ao vivo das 19h30 às 22h30 (3h cada encontro)



Certificado

Participando de **pelo menos 4 encontros**, você conquistará **seu certificado em Projetos de Dados**



Calendário

- **Dia 1:** 03-Abr-23: Onboarding, Team Building e Início do projeto
- **Dia 2:** 10-Abr-23: Desenvolvimento das Análises – Parte 1
- **Dia 3:** 17-Abr-23: Desenvolvimento das Análises – Parte 2
- **Dia 4:** 24-Abr-23: Avaliação dos Resultados e Preparação da Apresentação
- **Dia 5:** 08-Mai-23: Apresentação Final

Inscreva-se agora

Inscrições até 28/02



Preditiva.ai