

data

Análise Exploratória de  
Dados

Aimê @itsaime

# Conteúdo da aula

- Análise exploratória de dados
  - Análise intuitiva através da visualização de dados
  - Análise estatística
- Limpeza de dados
- Seleção de atributos



# Introdução

- É importante compreender a complexidade dos dados que estamos lidando
- Precisamos mensurar a:
  - propriedades dos dados;
  - qualidade dos dados;
  - e o poder preditivo dos dados.

```
1 "age";"job";"marital";"education";"default";"balance";"housing";"loan"  
2 58;"management";"married";"tertiary";"no";2143;"yes";"no";"unknown";5;  
3 44;"technician";"single";"secondary";"no";29;"yes";"no";"unknown";5;"n  
4 33;"entrepreneur";"married";"secondary";"no";2;"yes";"yes";"unknown";5;  
5 47;"blue-collar";"married";"unknown";"no";1506;"yes";"no";"unknown";5;  
6 33;"unknown";"single";"unknown";"no";1;"no";"no";"unknown";5;"may";198  
7 35;"management";"married";"tertiary";"no";231;"yes";"no";"unknown";5;"  
8 28;"management";"single";"tertiary";"no";447;"yes";"yes";"unknown";5;"  
9 42;"entrepreneur";"divorced";"tertiary";"yes";2;"yes";"no";"unknown";5;  
10 58;"retired";"married";"primary";"no";121;"yes";"no";"unknown";5;"may"  
11 43;"technician";"single";"secondary";"no";593;"yes";"no";"unknown";5;"  
12 41;"admin."; "divorced"; "secondary"; "no"; 270;"yes";"no";"unknown";5;"ma  
13 29;"admin."; "single"; "secondary"; "no"; 390;"yes";"no";"unknown";5;"may"  
14 53;"technician";"married";"secondary";"no";6;"yes";"no";"unknown";5;"n  
15 58;"technician";"married";"unknown";"no";71;"yes";"no";"unknown";5;"ma  
16 57;"services";"married";"secondary";"no";162;"yes";"no";"unknown";5;"n  
17 51;"retired";"married";"primary";"no";330;"yes";"no";"unknown";5;"may"
```



# Análise estatística de dados

Esse passo é usado para qualquer dataset, não só os usados para ML. Iremos avaliar a qualidade e o poder preditivo dos dados.

Conteúdo:

- Análise univariada (descritiva)
- Análise bivariada (correlação)
  - Qualitativa
  - Quantitativa
- Análise contextual



# Análise univariada (descritiva)

## Atributo comum

- **Tipo do dado** (importante!!!)
- Valores faltantes (printar tabela com porcentagem de cada coluna)

## Atributo categórico

- cardinalidade (baixa == muitos elementos duplicados)
- quantidade de valores únicos

## Atributo numérico

- Quantis estatísticos (Q1, Q2, Q3, min, max, range)
- Estatística descritiva (mean, mode, sd, median, kurtosis, skewness)
- Histograma de distribuição

## Atributo textual

- tokens
- document frequency (DF)/ term frequency (TF) sem palavras vazias (stop words)



# Passo crucial.

Verificar atributos comuns:  
valores faltantes e tipo de dado  
das colunas.

	Missing Values	Percentage	Data Types
Id	0	0.000000	int64
MSSubClass	0	0.000000	int64
MSZoning	0	0.000000	object
LotFrontage	259	17.739726	float64
LotArea	0	0.000000	int64
...	...	...	...
MoSold	0	0.000000	int64
YrSold	0	0.000000	int64
SaleType	0	0.000000	object
SaleCondition	0	0.000000	object
SalePrice	0	0.000000	int64
81 rows × 3 columns			

# Limpeza dos dados

- Remover duplicados
- Remover dados irrelevantes
- Padronizar capitalização
- Converter tipos de dado
- Remover erros e inconsistências
- Todos dados na mesma linguagem
- Lidar com dados faltantes
- Lidar com outliers



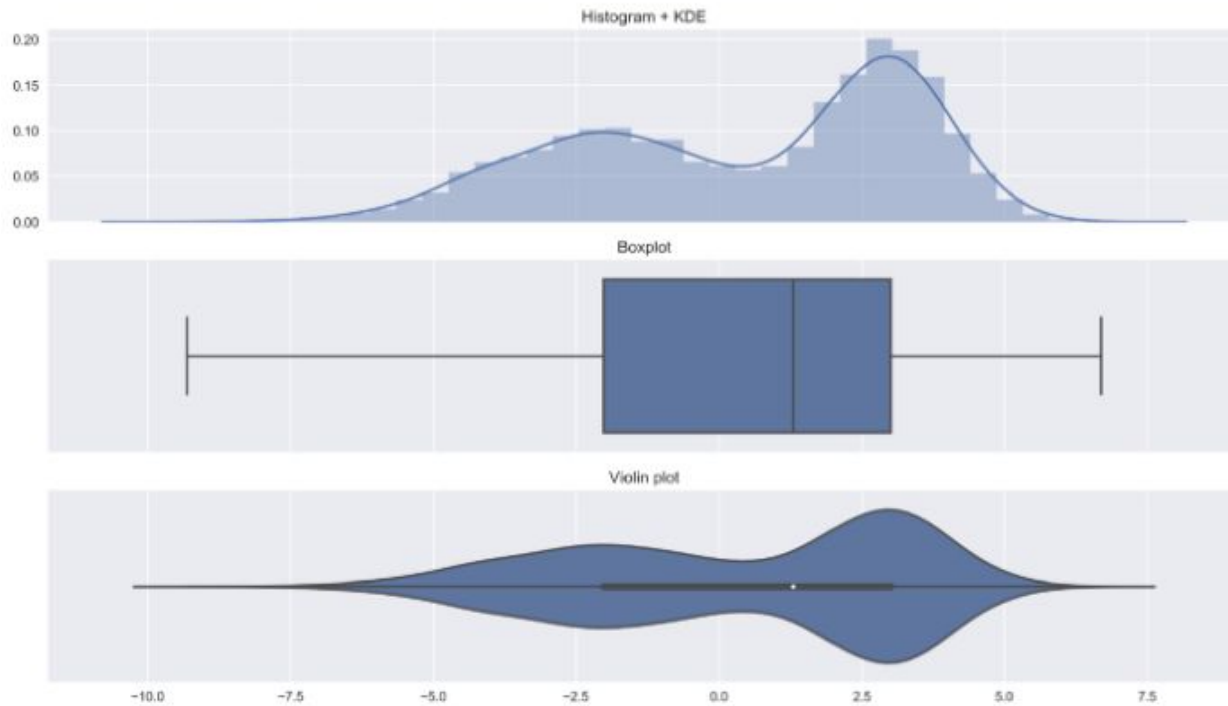
# Atributos numéricos

	max	range	IQR	mode	mad	kurtosis	skewness
age	90	73	20	36	11.189182	-0.166127	0.558743
functional_weight	1484705	1472420	119224	123011	77608.21854	6.218811	1.44698
education_num	16	15	3	9	1.903048	0.623444	-0.311676
capital_gain	99999	99999	0	0	1977.373437	154.799438	11.953848
capital_loss	4356	4356	0	0	166.462055	20.376802	4.594629
hours_per_week	99	98	5	40	7.583228	2.916687	0.227643

	mean	std	min	25%	50%	75%
age	38.581647	13.640433	17	28	37	48
functional_weight	189778.3665	105549.9777	12285	117827	178356	237051
education_num	10.080679	2.57272	1	9	10	12
capital_gain	1077.648844	7385.292085	0	0	0	0
capital_loss	87.30383	402.960219	0	0	0	0
hours_per_week	40.437456	12.347429	1	40	40	45



### Mixture of Gaussians - bimodal



Podem ser representados por gráficos

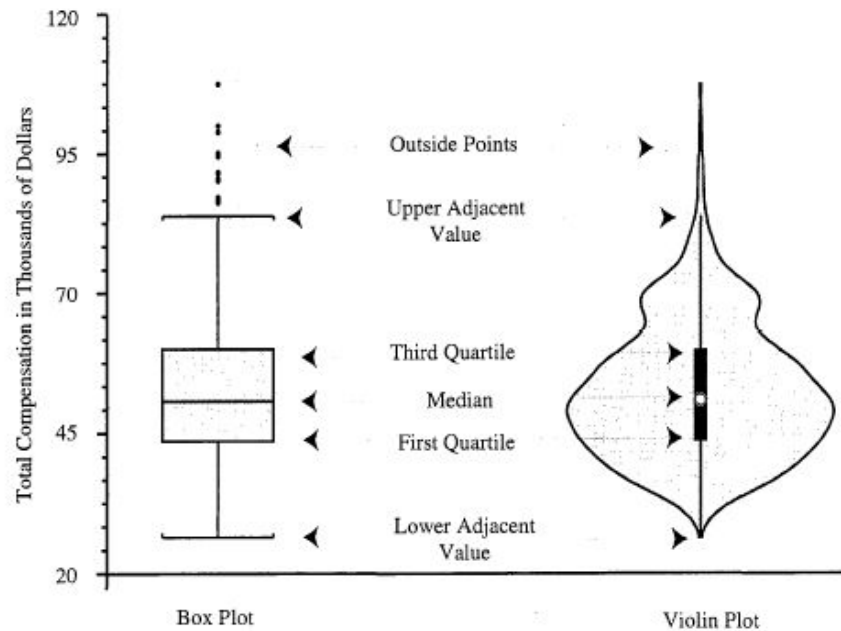
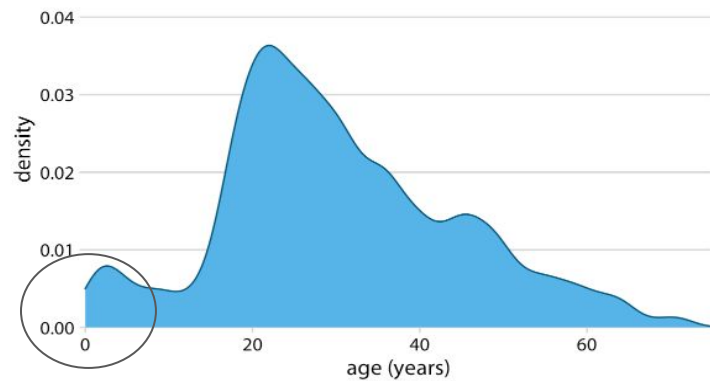
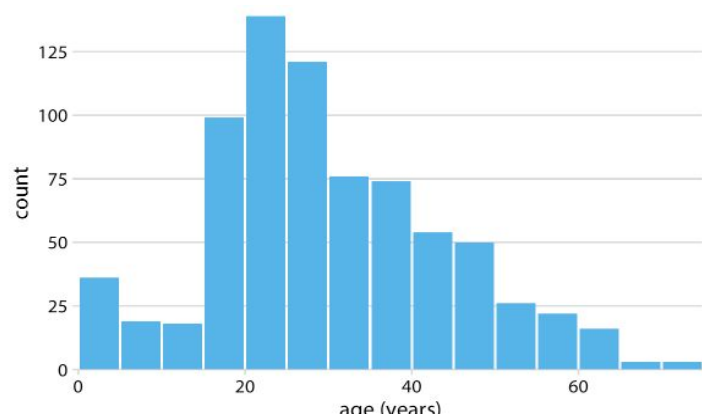


Figure 1. Common Components of Box Plot and Violin Plot. Total compensation for all academic ranks.

Plot de histograma, densidade, boxplot e violino.

# Análise bivariada (de correlação)

Examina correlação entre variáveis, dois a dois. Determina variáveis redundantes e/ou que não são preditivas em relação ao valor alvo.

Qualitativa:

- 2 categóricas: tabela de contingência
- categórica vs. numérica: estatística descritiva ou histograma

Problemas: enviesamento (verificado com métodos de amostragem e design de experimentos)

Quantitativa:

Teste de Hipóteses, onde a hipótese nula é de que não há relação entre as duas variáveis. Usamos p-valor de 0.05.



# Tabela de contingência

Item	Extra
Pizza	Cheese
Pizza	Sauce
Burger	Cheese
Burger	Cheese
Burger	Sauce



	Cheese	Sauce
Pizza	1	1
Burger	2	1

Handed- ness	Right-handed	Left-handed	Total
Sex			
Male	43	9	52
Female	44	4	48
Total	87	13	100



# Teste de hipóteses

X	Y	
	Categorical	Numerical
Categorical	<u>Chi-square test</u> <u>Information gain</u>	<u>Student T-test</u> <u>ANOVA</u> <u>Logistic regression</u> Discretize Y (left column)
Numerical	<u>Student T-test</u> <u>ANOVA</u> <u>Logistic regression</u> Discretize X (row above)	<u>Correlation</u> <u>Linear Regression</u> Discretize Y (left column) Discretize X (row above)



Handedness	Right-handed	Left-handed	Total
Sex			
Male	43	9	52
Female	44	4	48
Total	87	13	100

Handedness	Right-handed	Left-handed	Total
Sex			
Male	$100 \cdot 0.52 \cdot 0.87 = 45$	$100 \cdot 0.52 \cdot 0.13 = 7$	52
Female	$100 \cdot 0.48 \cdot 0.87 = 42$	$100 \cdot 0.48 \cdot 0.13 = 6$	48
Total	87	13	100

Test the statistics.

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = \frac{(43-45)^2}{45} + \frac{(44-42)^2}{42} + \frac{(9-7)^2}{7} + \frac{(4-6)^2}{6} = 1.42$$

$$p_{value} = Pr_{DF=1}(\chi^2 > 1.42) = 0.233$$

Exemplo de Teste Chi-Quadrado

# Análise Contextual

Entender o contexto no qual o dataset está inserido pode nos dar insights em relação aos dados.

Baseada em tempo:

- horas de rush
- quebra da bolsa de 2008
- datas de eleição
- black friday

Baseado em usuários.

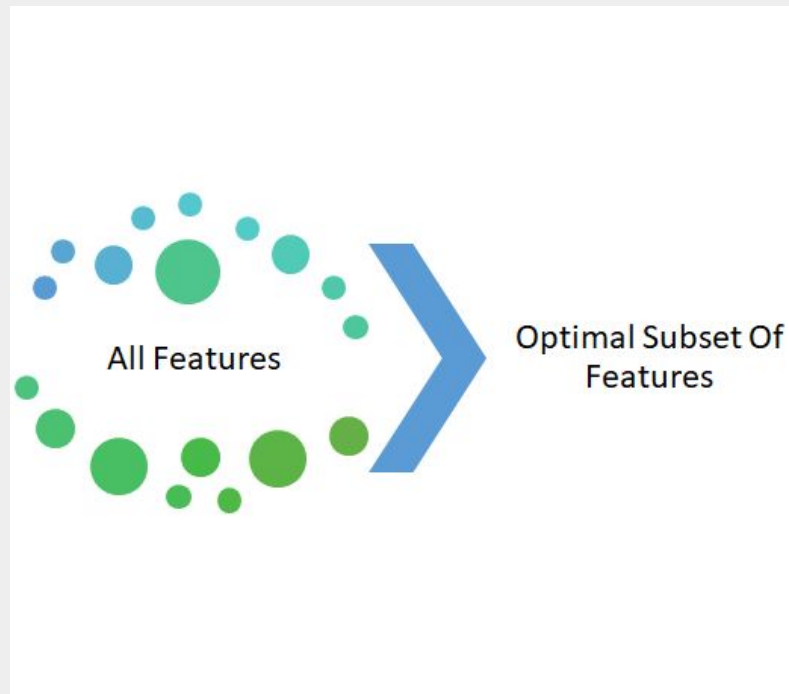
- verificar comportamento de cada usuário ou em conjunto
- achar anomalias



# Seleção de atributos

Baseado em:

- Análise descritiva
- Análise de correlação
- Análise contextual





# Seleção de atributos

## Baseado em análise descritiva

- Alta porcentagem de números faltantes não constituindo informação
- Pouca variância. Remover esses dados pode diminuir complexidade da ML
- Pouca entropia de categorias pode enviesar ML
  - Desbalanço de classe pode ser resolvida com amostragem
- Distribuição distorcida com calda longa. Pode ser resultado de erros que geram outliers.
- Alta cardinalidade, ou seja, muitos valores únicos

## Baseado em análise de correlação

- Correlação muito baixa pode significar que atributo não será útil para previsão desejada
- Alta correlação. Não acrescenta informação ao modelo e não é recomendada para modelos lineares.

## Baseado em análise contextual

Leva em conta conhecimento do contexto em que os dados estão inseridos.



# Visualização de dados

# Tabelas

Guiar o leitor para a  
informação que  
queremos passar.  
Procurar não usar  
bordas, pois poluem a  
visão.

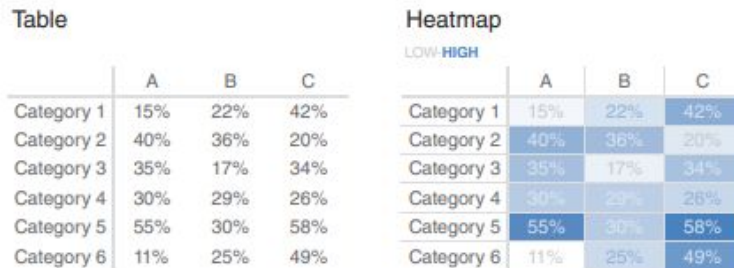


FIGURE 2.5 Two views of the same data



FIGURE 2.4 Table borders



# Escolhas sobre o design dos gráficos

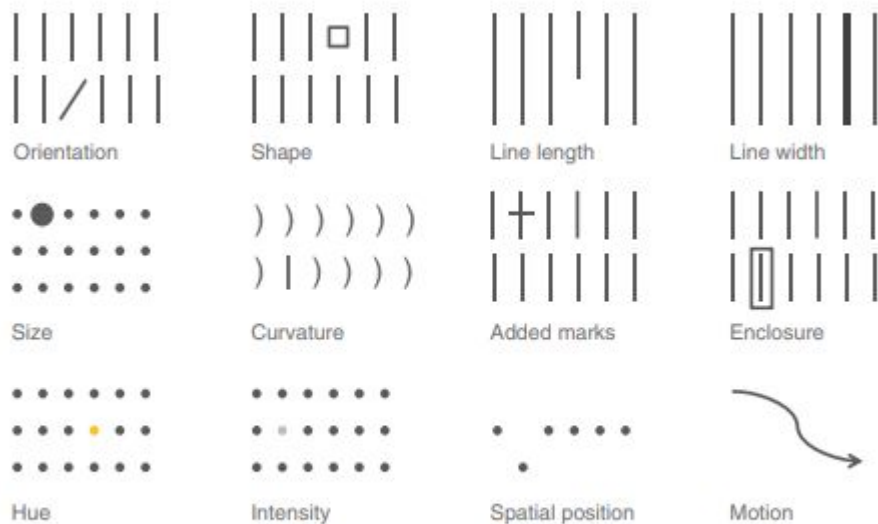
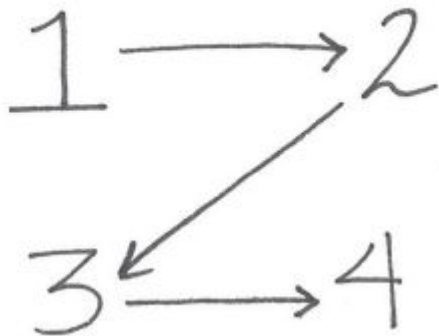


FIGURE 4.4 Preattentive attributes

Source: Adapted from Stephen Few's *Show Me the Numbers*, 2004.

Caminho do foco do leitor

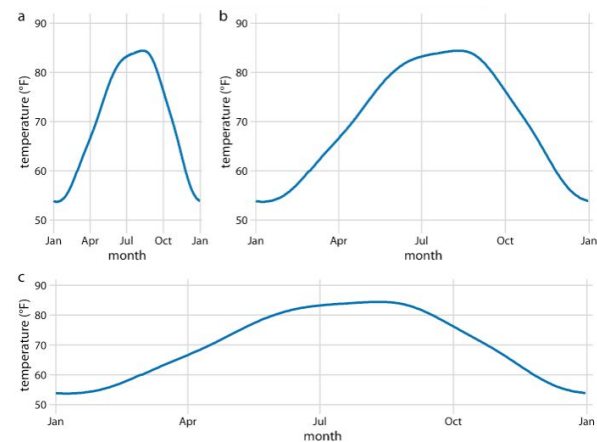
Recursos para destacar atributos.



# Escolhas sobre o design dos gráficos



Barras do gráfico



Escala do gráfico



Okabe Ito



ColorBrewer Dark2



ggplot2 hue



ColorBrewer Blues



Heat



Viridis

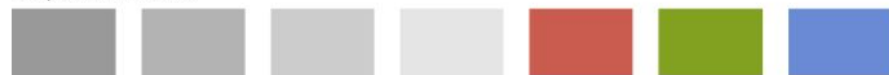


## Categórico

Okabe Ito Accent



Grays with accents



ColorBrewer Accent



## Destaque

## Quantitativo



## Amigável a daltônicos



Círculo cromático

# Que informação devemos apresentar?

**Quem:** a quem iremos falar

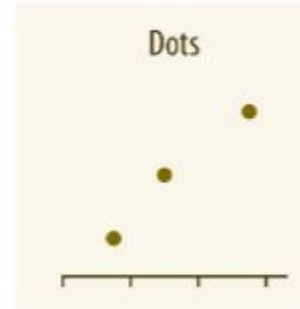
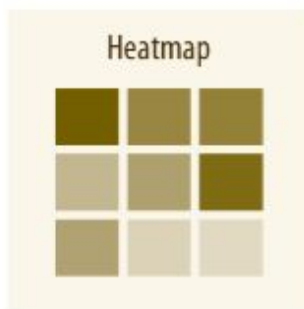
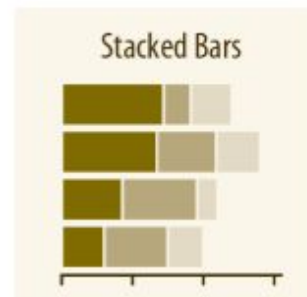
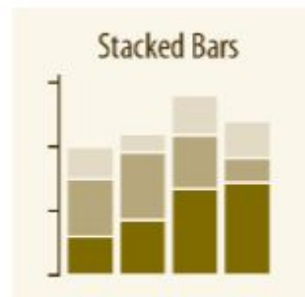
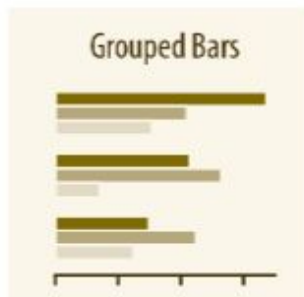
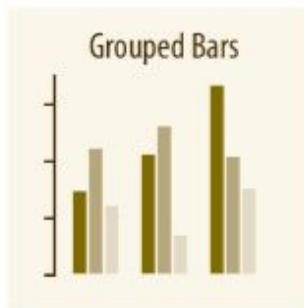
**O que:** ideia que queremos passar

**Como:** escolha da informação apresentada

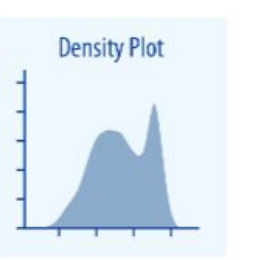
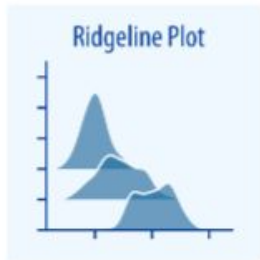
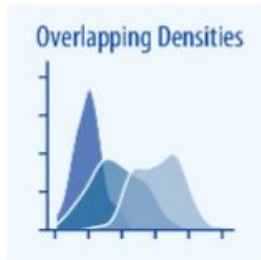
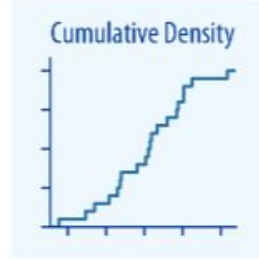
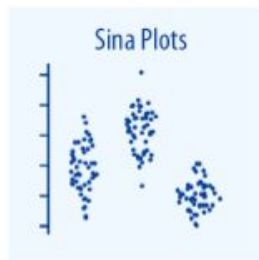
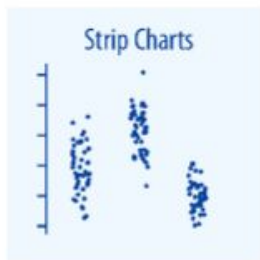
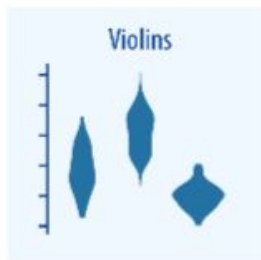
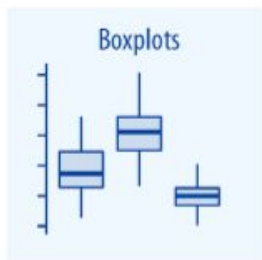




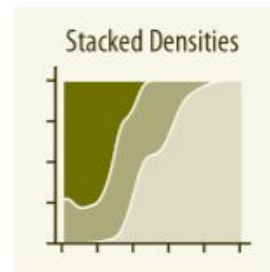
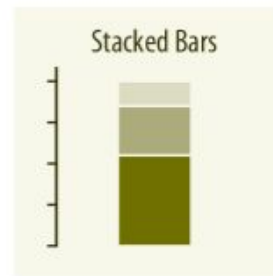
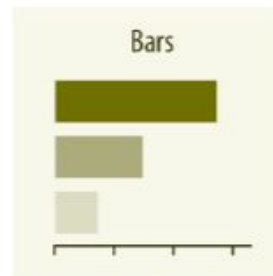
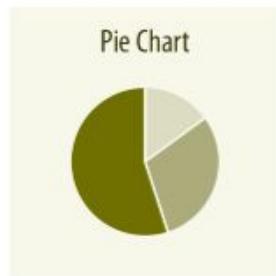
# Visualizando quantidades



# Visualizando distribuições



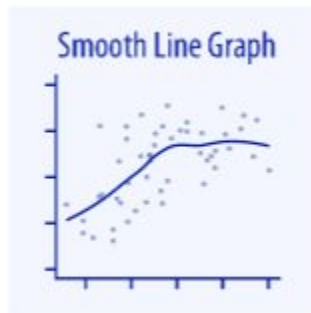
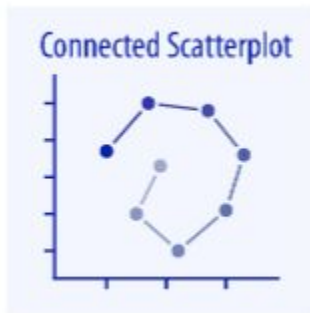
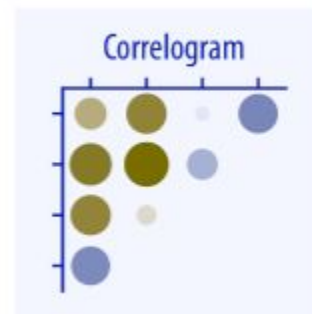
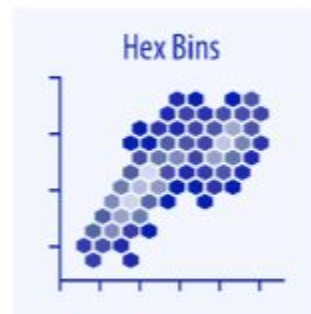
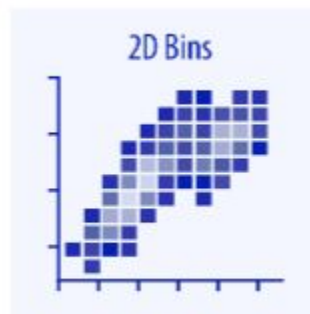
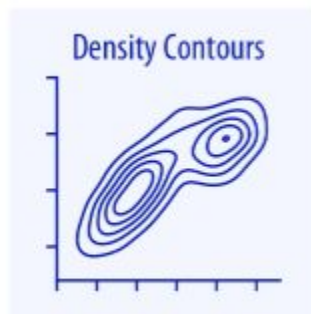
# Visualizando proporções



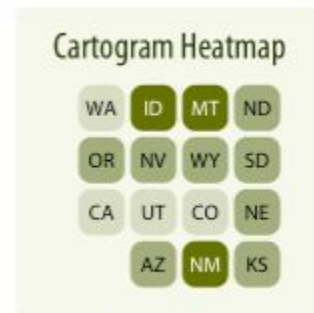
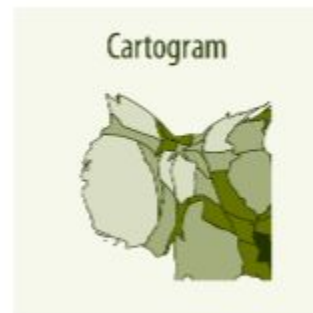
# Visualizando proporções



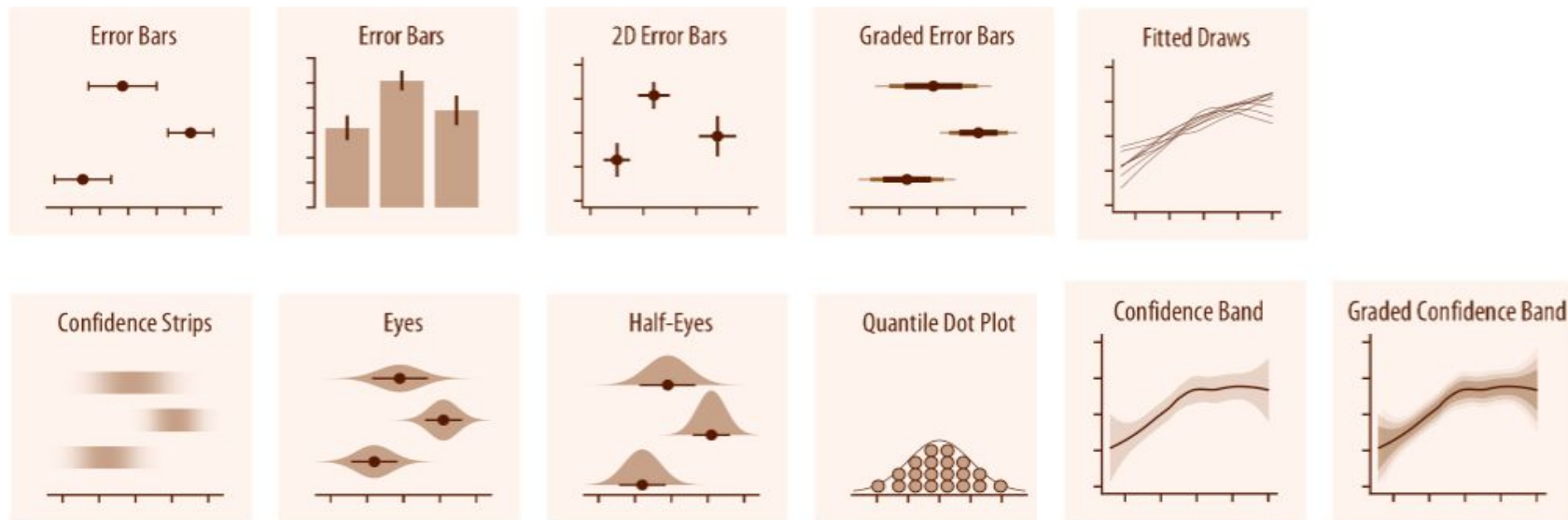
# Relações entre duas variáveis



# Dado geoespacial



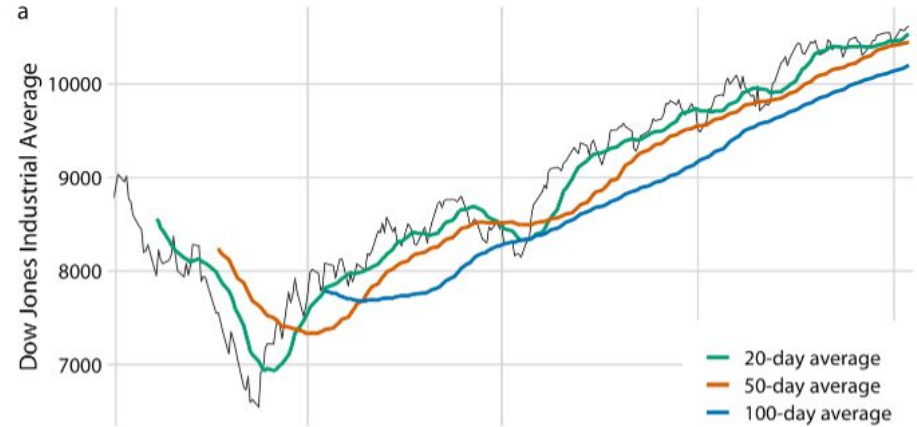
# Visualizando incerteza



# Séries temporels

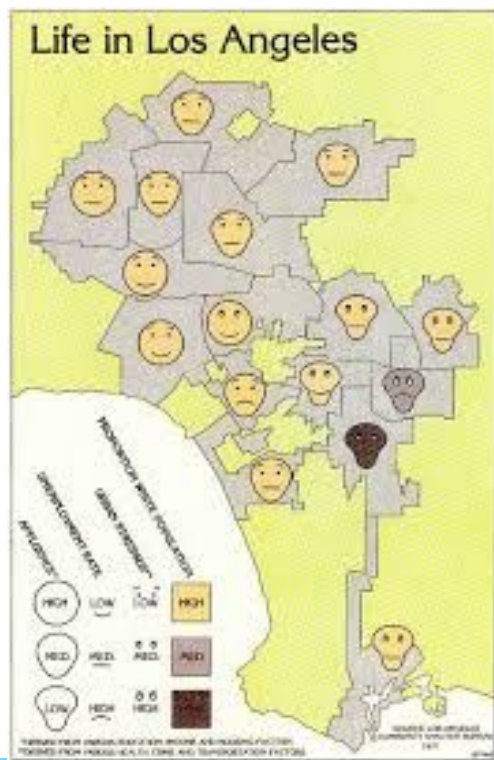


Figure 14.1: Daily closing values of the Dow Jones Industrial Average for the year 2009. Data source: Yahoo! Finance





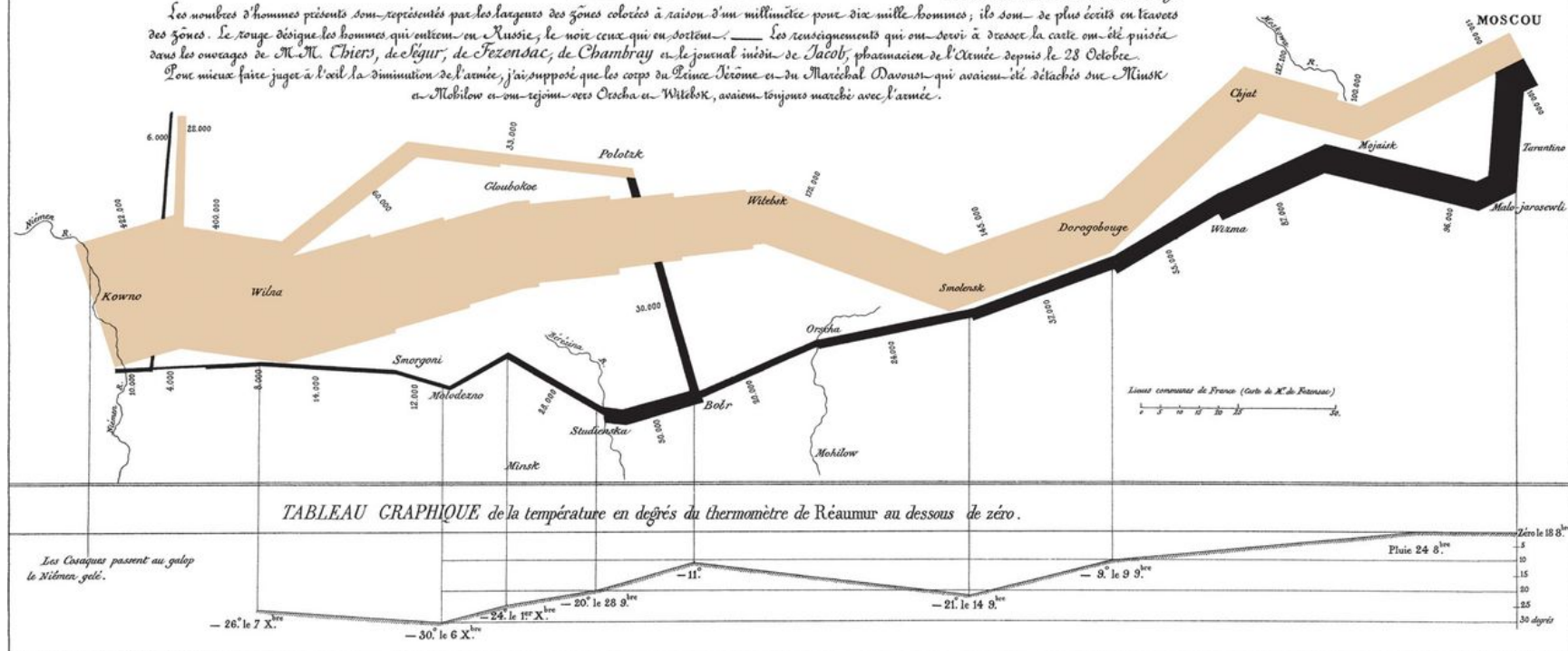
# Faces de Chernoff



# Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.

Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite. Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en caractères des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. Thiers, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre. Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout, qui avaient été détachés sur Minsk et Mohilew et s'en retournèrent vers Orscha et Witebsk, avaient toujours marché avec l'armée.



Asses. par Regnier, 1. Par. 5<sup>me</sup> Maria 5<sup>me</sup> 0<sup>me</sup> à Paris.

Imp. Lith. Regnier et Doucet.

Diagrama de Sankey

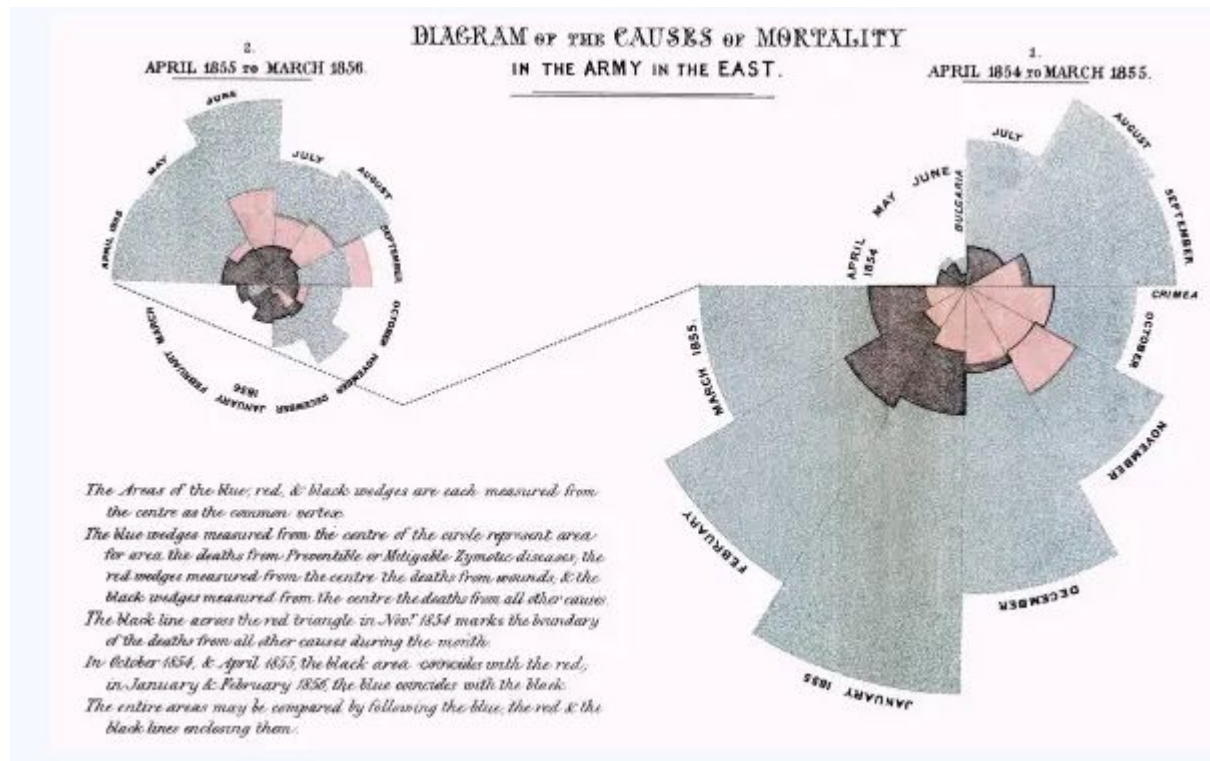
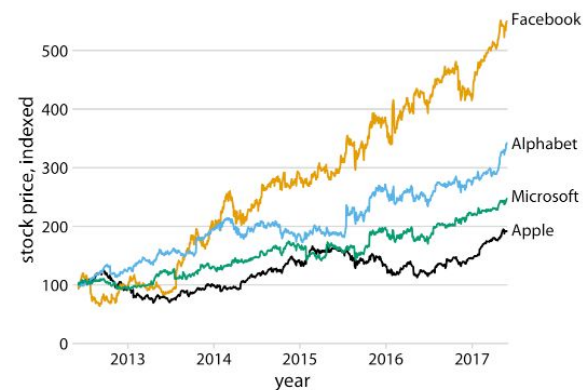
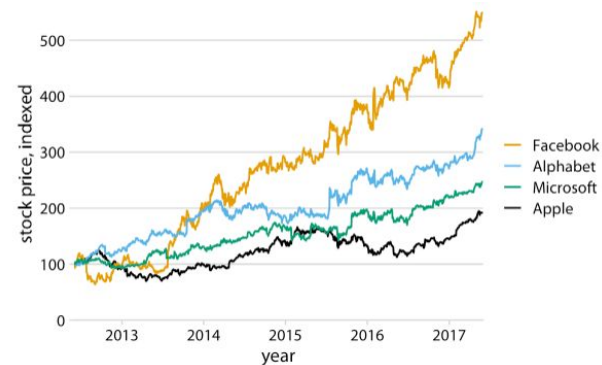
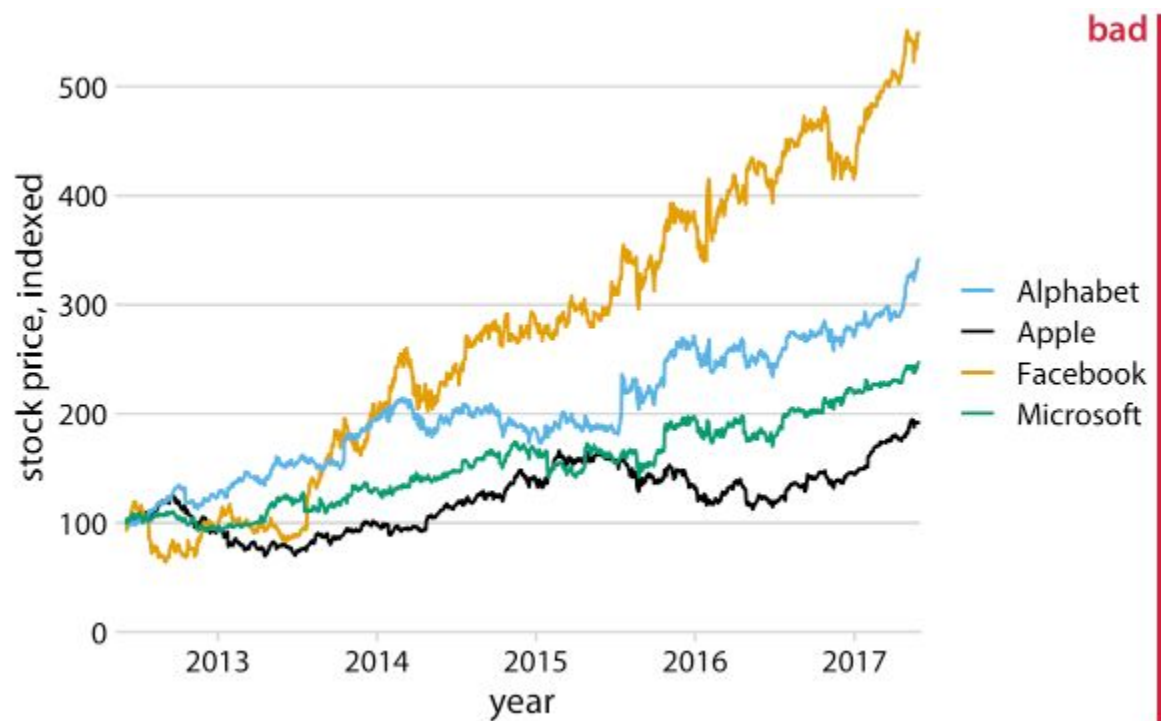


Gráfico de Florence Nightingale (Coxcomb Plot)

# Atenção às legendas



# Importância da escolha do gráfico

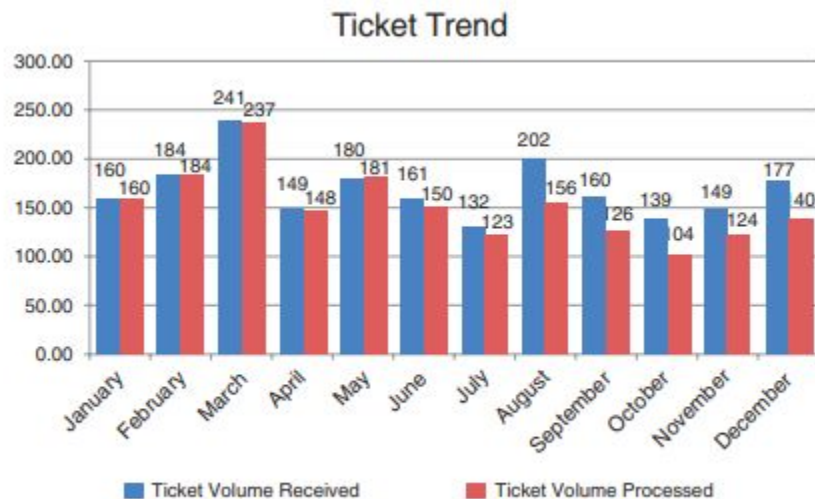
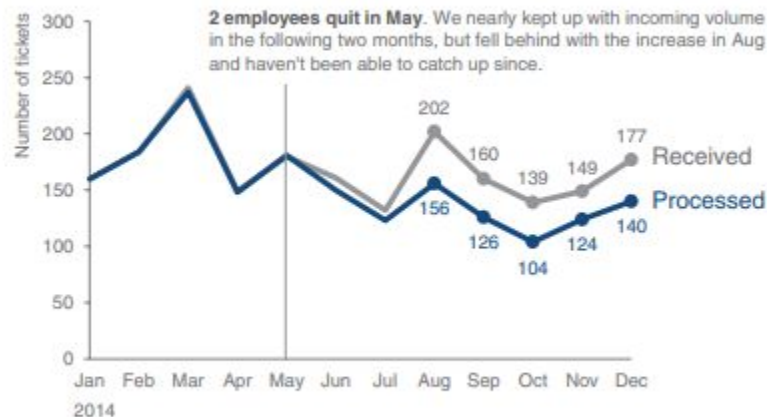


FIGURE 0.2 Example 1 (before): showing data

Ticket volume over time



Data source: XYZ Dashboard, as of 12/31/2014 | A detailed analysis on tickets processed per person and time to resolve issues was undertaken to inform this request and can be provided if needed.



# Importância da escolha do gráfico

Average Retail Product Price per Year

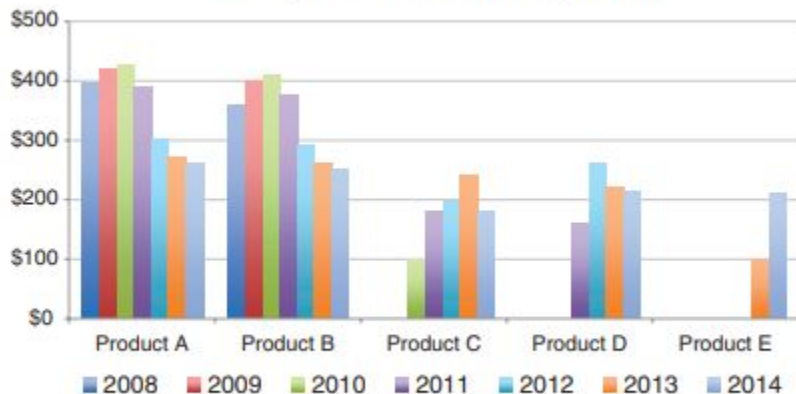


FIGURE 0.6 Example 3 (before): showing data

To be competitive, we recommend introducing our product *below* the \$223 average price point in the \$150–\$200 range

Retail price over time by product

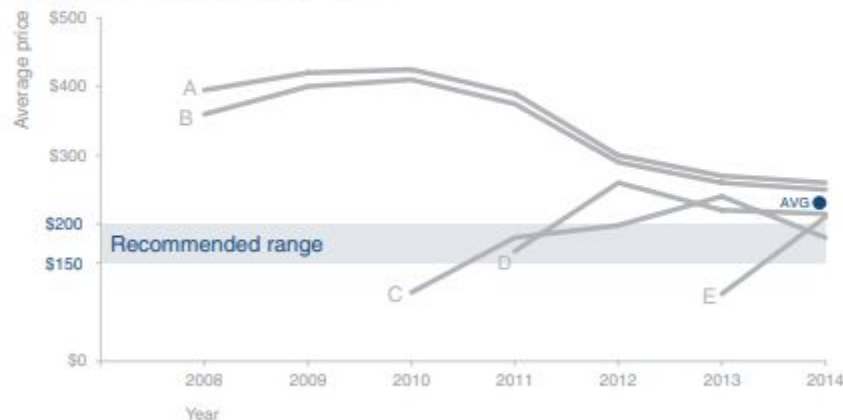


FIGURE 0.7 Example 3 (after): storytelling with data





# Importância da escolha do gráfico

## Survey Results

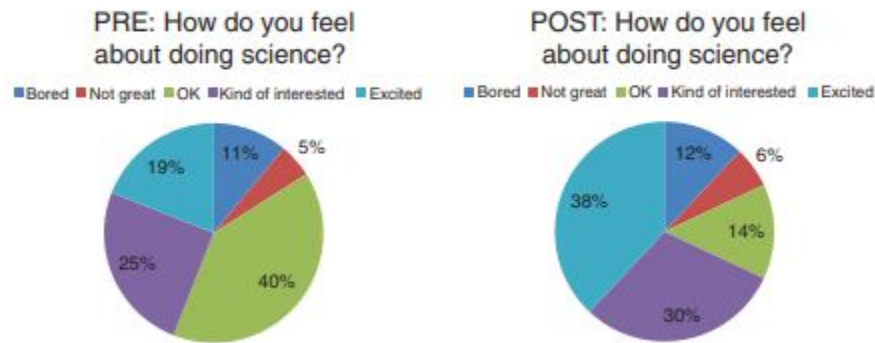
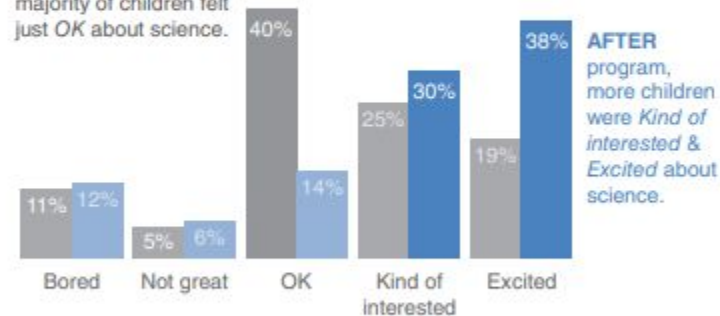


FIGURE 0.4 Example 2 (before): showing data

## Pilot program was a success

How do you feel about science?

**BEFORE** program, the majority of children felt just OK about science.



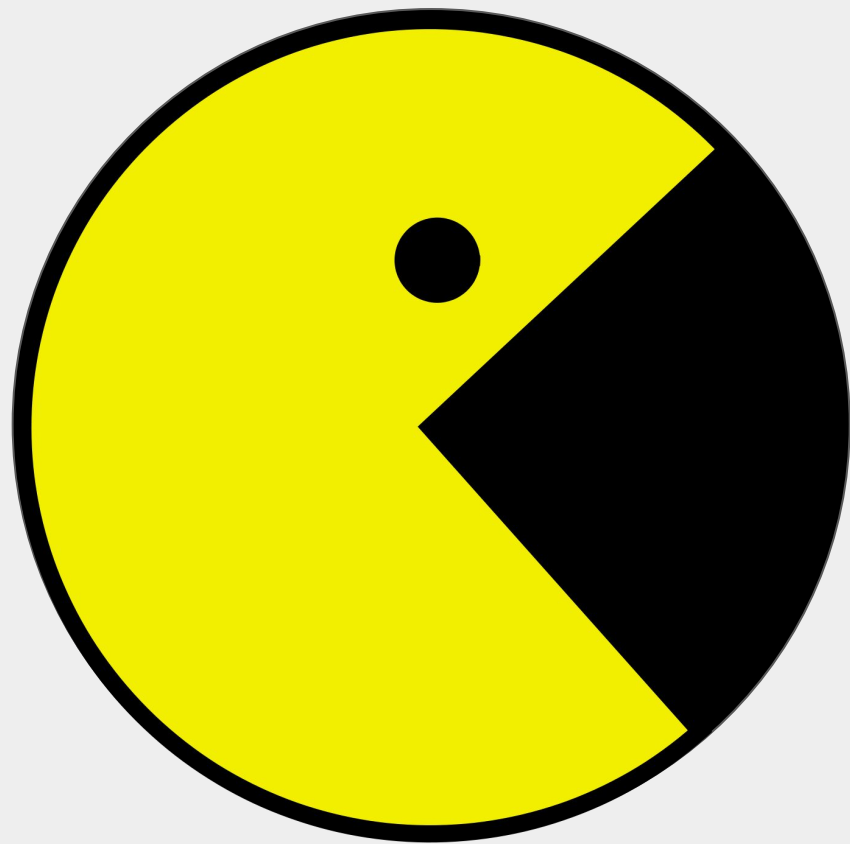
Based on survey of 100 students conducted before and after pilot program (100% response rate on both surveys).

FIGURE 0.5 Example 2 (after): storytelling with data



# Gráfico de pizza

Evite o gráfico de pizza. Se fizer questão, saiba que não passa muita informação. Use então a regra do Pacman.





# 75%

Muito mais impactante que um gráfico de pizza.

# Scatterplot

Cost per mile by miles driven

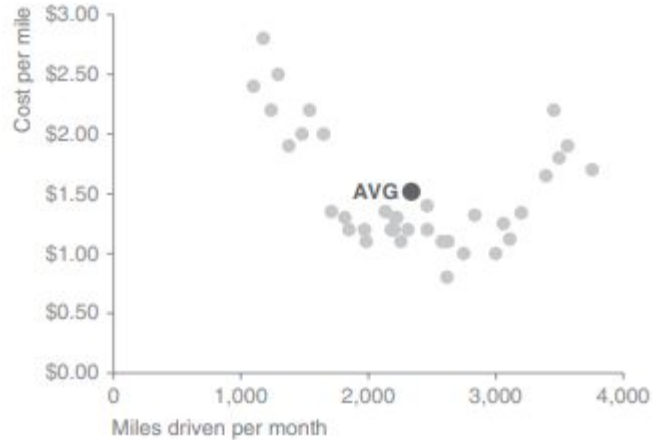


FIGURE 2.6 Scatterplot

Cost per mile by miles driven

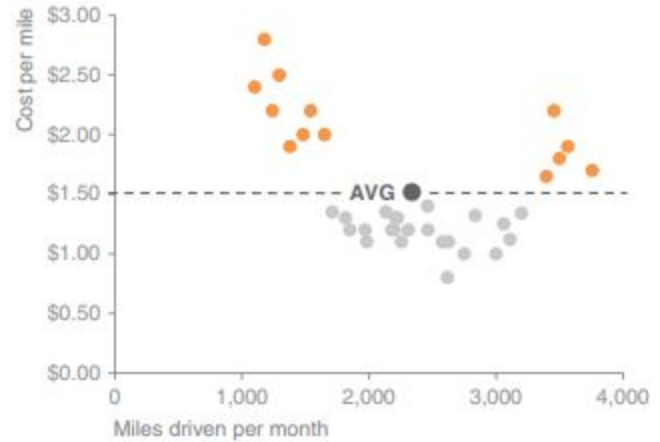


FIGURE 2.7 Modified scatterplot



# Slopegraph

Employee feedback over time



Employee feedback over time

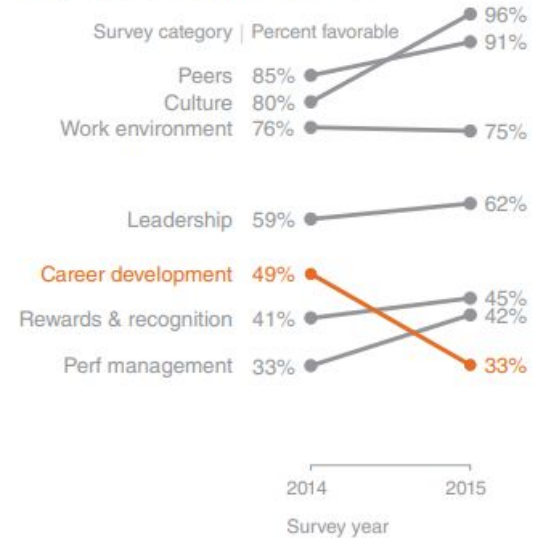


FIGURE 2.11 Modified slopegraph



# Elimine a poluição visual

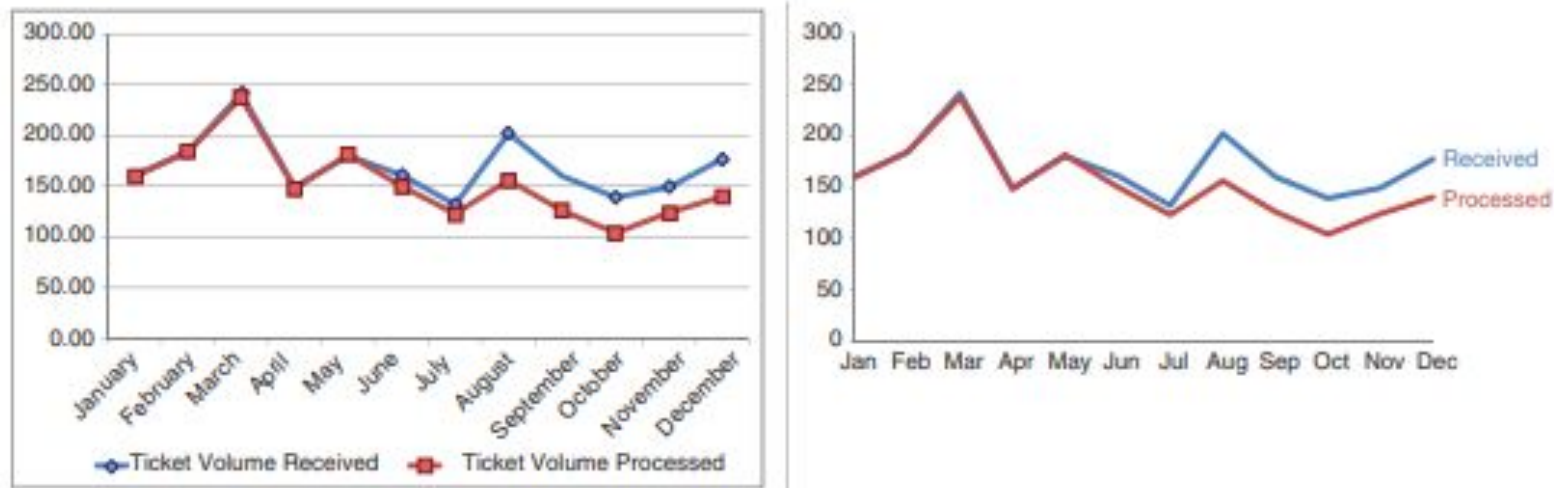


FIGURE 3.24 Before-and-after



# Foco direccionado

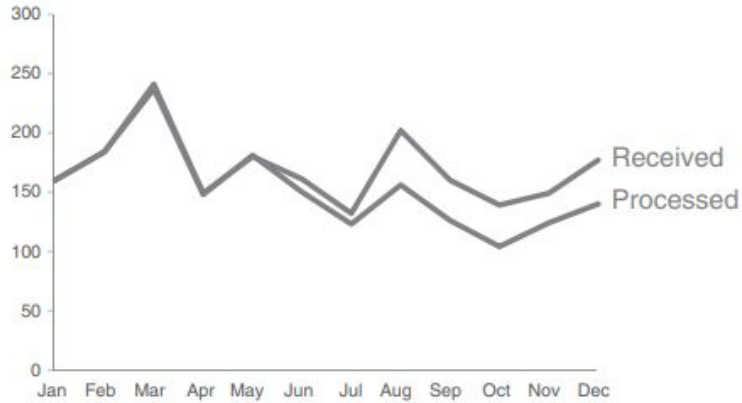


FIGURE 4.11 First, push everything to the background



FIGURE 4.14 Data labels used sparingly help draw attention



# Acréscimo de informações no gráfico

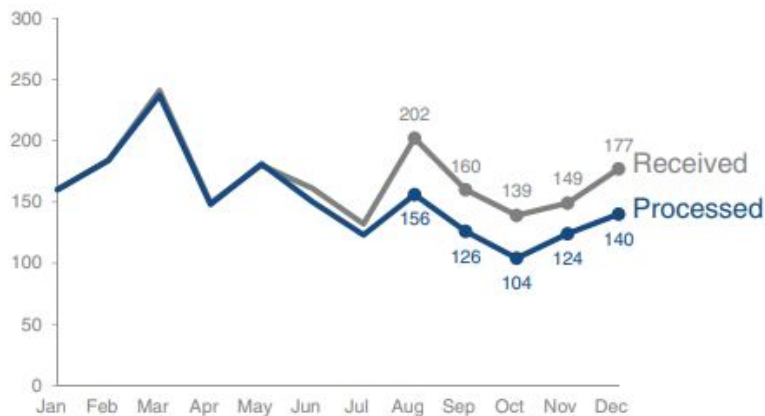
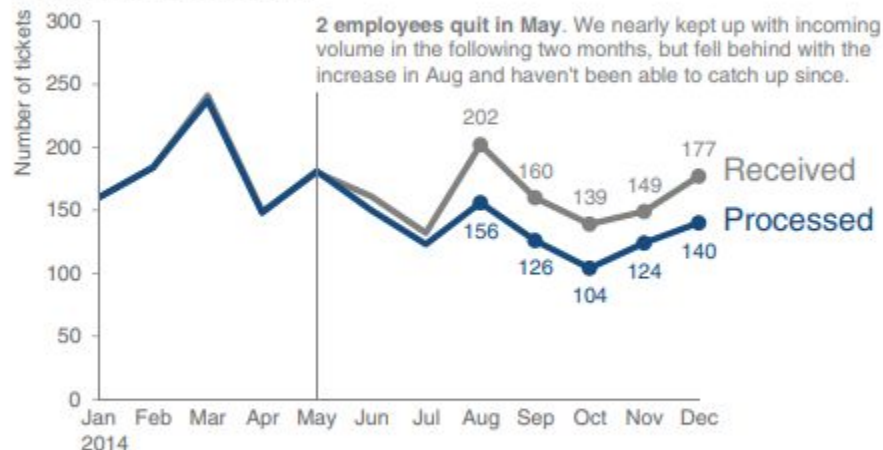


FIGURE 4.14 Data labels used sparingly help draw attention

## Please approve the hire of 2 FTEs

to backfill those who quit in the past year

### Ticket volume over time



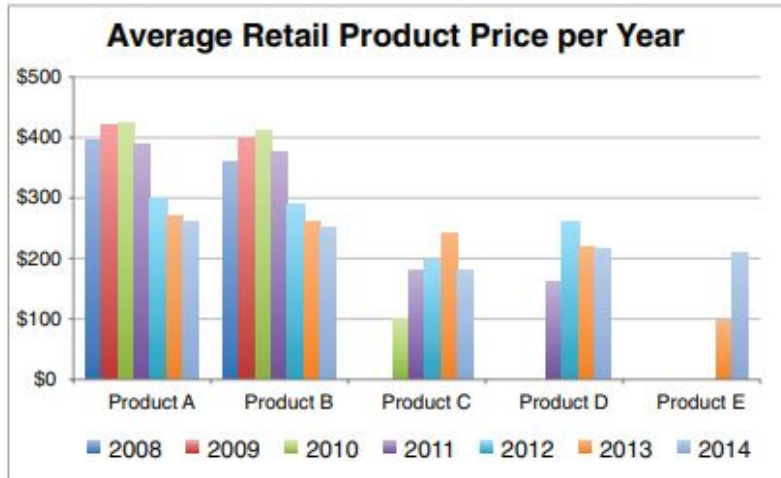
Data source: XYZ Dashboard, as of 12/31/2014 | A detailed analysis on tickets processed per person and time to resolve issues was undertaken to inform this request and can be provided if needed.

FIGURE 5.10 Add action title and annotation



# Exemplo

Price has declined for all products on the market since the launch of Product C in 2010



To be competitive, we recommend introducing our product *below* the \$223 average price point in the **\$150–\$200 range**

Retail price over time

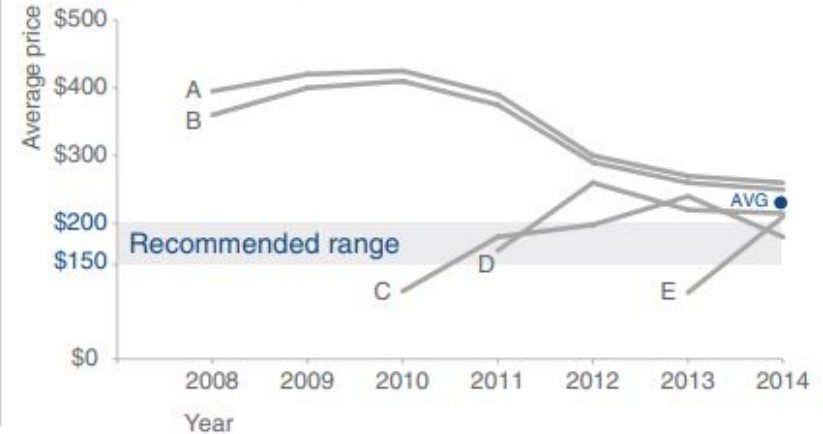
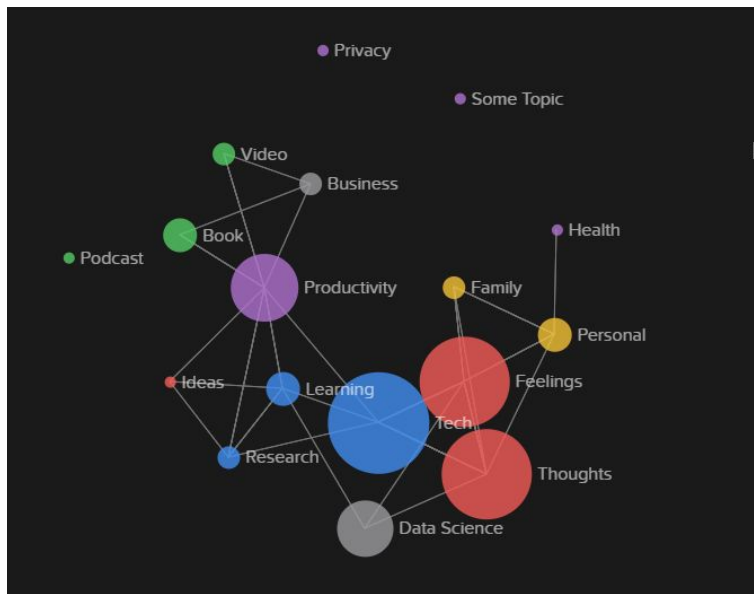


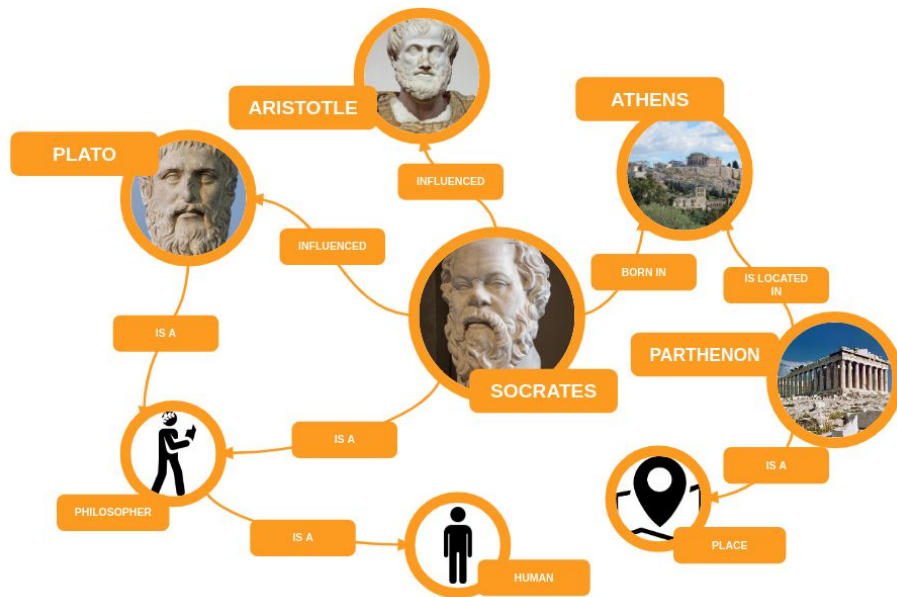
FIGURE 8.20 Before-and-after



# Extra: Visualização informação através de grafos



Grafo do aplicativo Obsidian

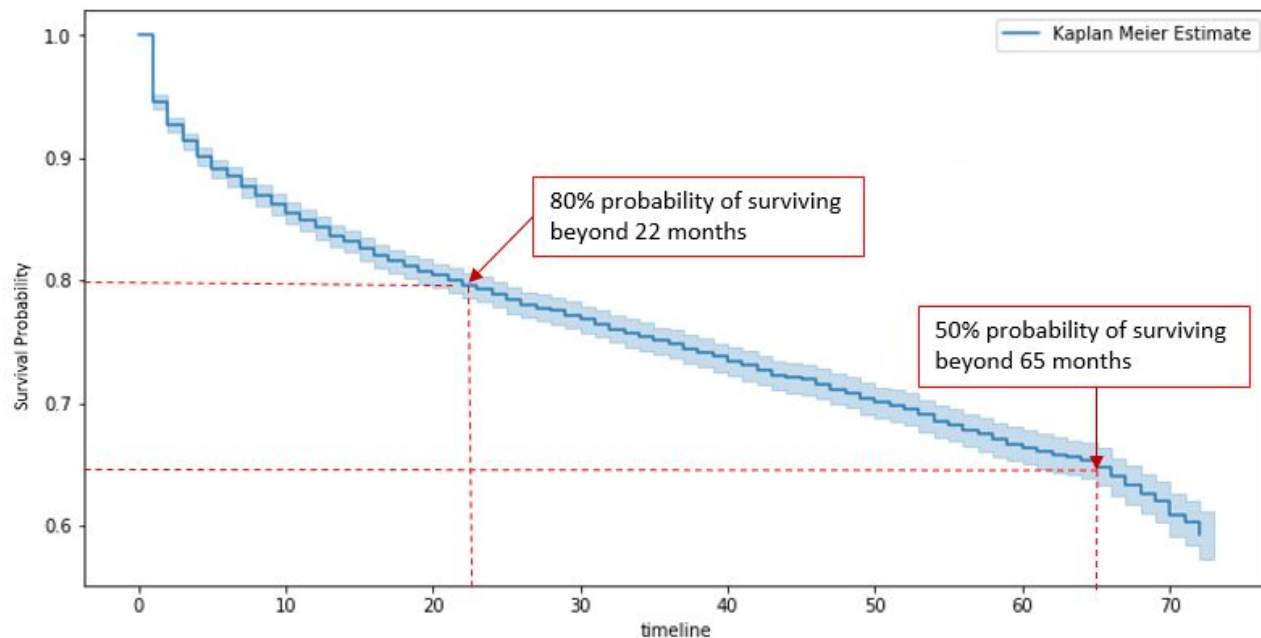


Grafo de Conhecimento



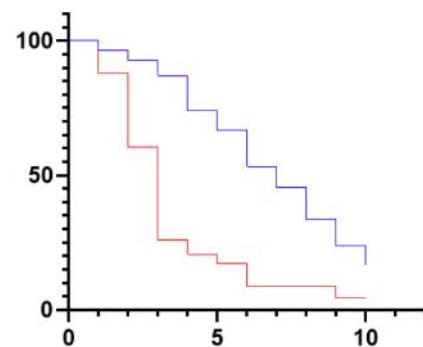


# Extra: Análise de sobrevivência



Biblioteca python:  
lifelines

Mais a fundo: usar R



# Ferramentas

- Python
  - Seaborn
  - Matplotlib
  - Plotly
  - geoplotlib
  - Folium
- R
  - Gráficos mais específicos de estatística, como QQ plot
  - Gráficos interativos

Obs: olhar sempre a documentação.

## Ferramentas de DashBoard

- Python
- Data Studio (Google)
- PowerBI
- Tableau
- InDesign
- Slides (Google)

## Inspiração

- Pinterest
- Behance (Adobe)
- Data Studio Report Gallery



# Bibliografia

Fundamentals of Data Visualization - Claus O. Wilke

Exploratory Data Analysis for Feature Selection in Machine Learning - Google Cloud

<https://datavizcatalogue.com/>

Storytelling with data - Cole Nussbaumer Knafl





dúvidas?

