

Week 3

Essential Data Visualization



Agenda

- Data Visualisation - The good and the bad.
- The essential data visualizations
- Stylizing your plots



DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

From the October 2012 Issue

Summary Save Share Comment Print \$8.95 Buy Copies

When Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know anyone. So you just stand in the corner sipping your drink—and you probably leave early."

Goldman, PhD, a statistician from St. Louis, was invited to the LinkedIn headquarters in Menlo Park,

Sep 27, 2012, 10:49pm EDT

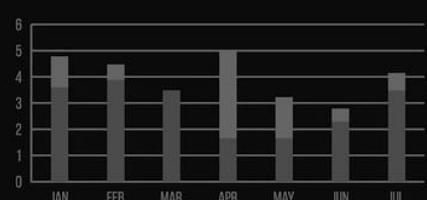
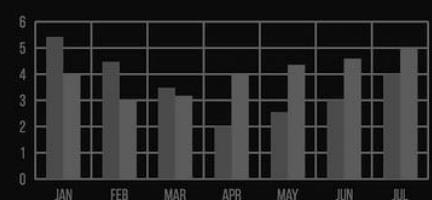
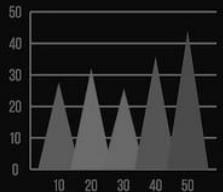
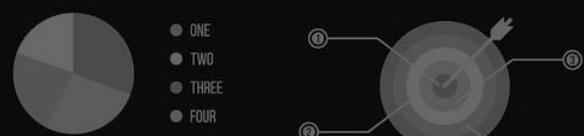
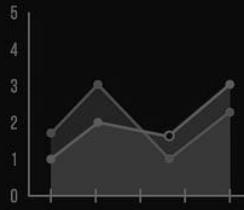
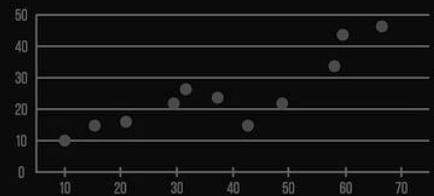
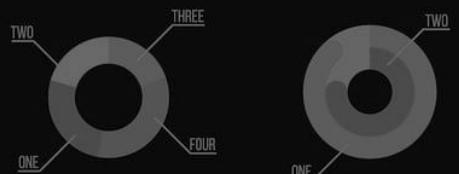
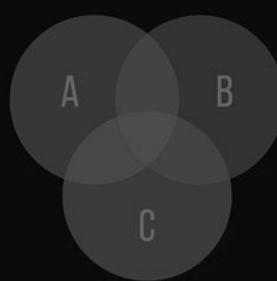
Data Scientists: The Definition of Sexy



Gil Press Senior Contributor

[Enterprise & Cloud](#)

I write about technology, entrepreneurs and innovation.

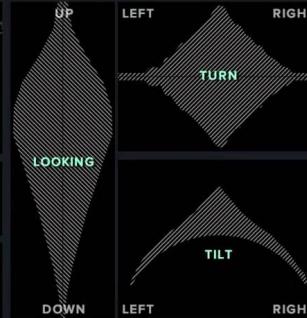


Good Data Visualisations

DEMOGRAPHICS



POSE



FEATURES



MOOD



3200 of 3200 selfies.

Normal

Crop

Crop & rotate



<http://selfiecity.net/#selfiexploratory>

Why do buses bunch?

INSTRUCTIONS EXPLANATION

Click and hold a bar below to delay its respective bus. Note how even a short delay causes the buses to bunch together after a while.

Hover over a stop to see its history. The area of the curve is cumulative wait time. Bunching makes the area grow.

PAUSE

RESET



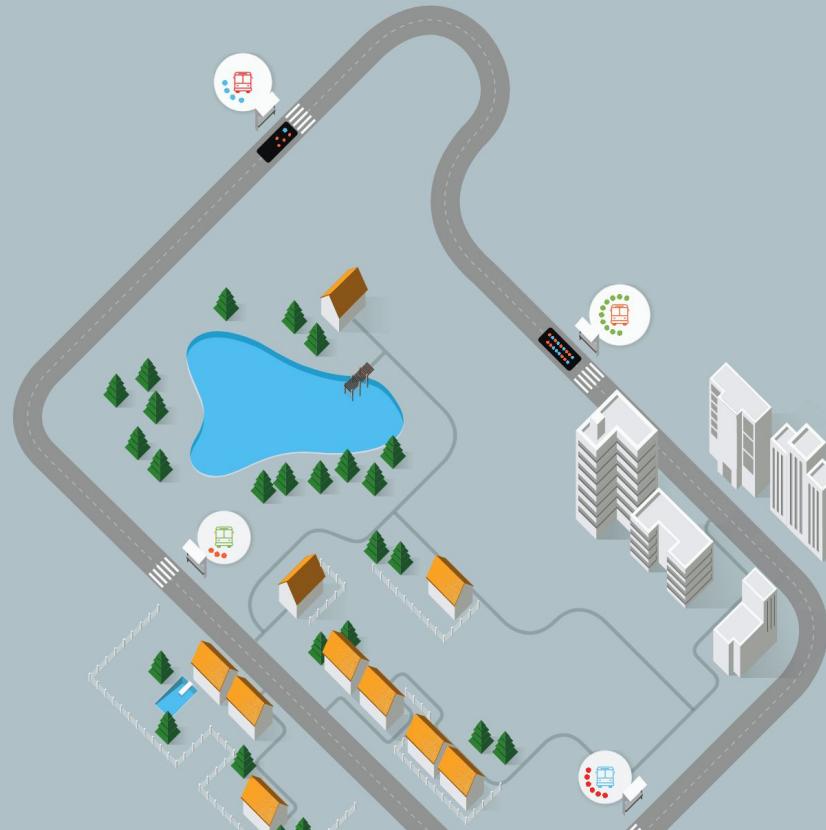
bus 1

17 passengers



bus 2

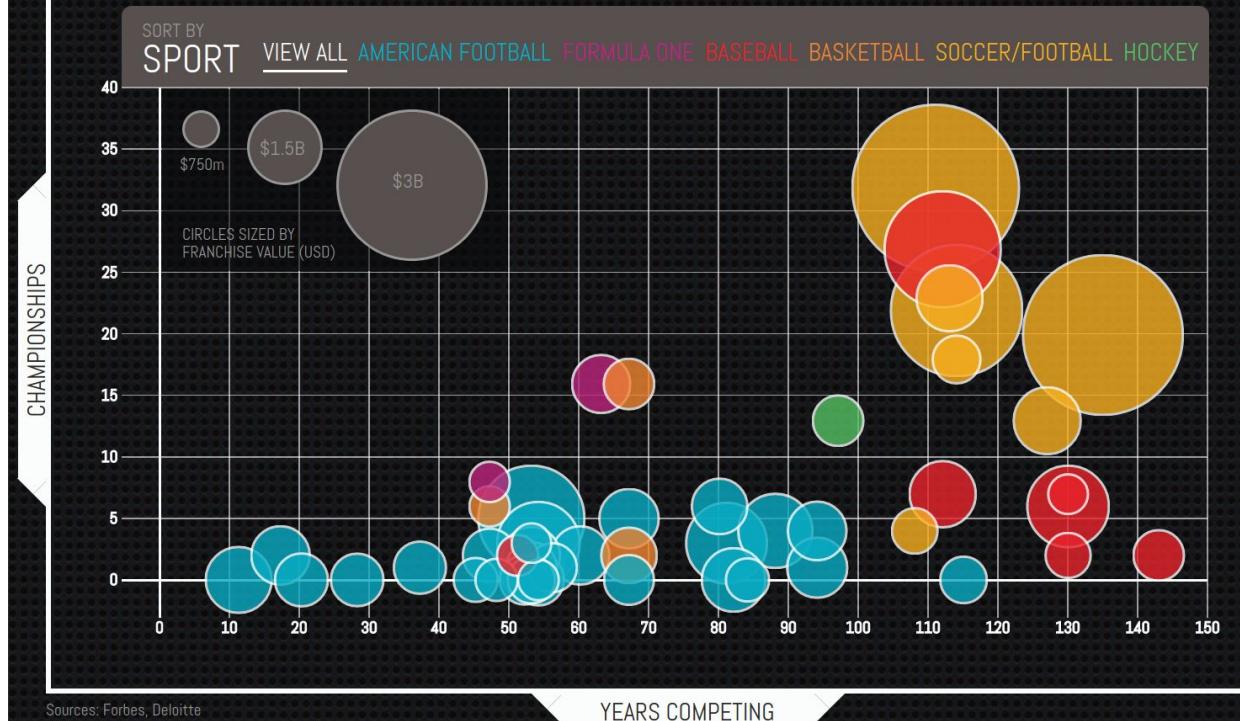
5 passengers



<https://setosa.io/bus/>

MOST VALUABLE PLAYERS

VISUALIZING FORBES' TOP 50 SPORTS FRANCHISES BY LONGEVITY AND SUCCESS



<http://mvp.columnfivemedia.com/>

routines.

SLEEP CREATIVE WORK DAY JOB/ADMIN FOOD/LEISURE EXERCISE OTHER

12 AM 1 2 3 4 5 6 7 8 9 10 11 12 PM 1 2 3 4 5 6 7 8 9 10 11 12



<https://podio.com/site/creative-routines>

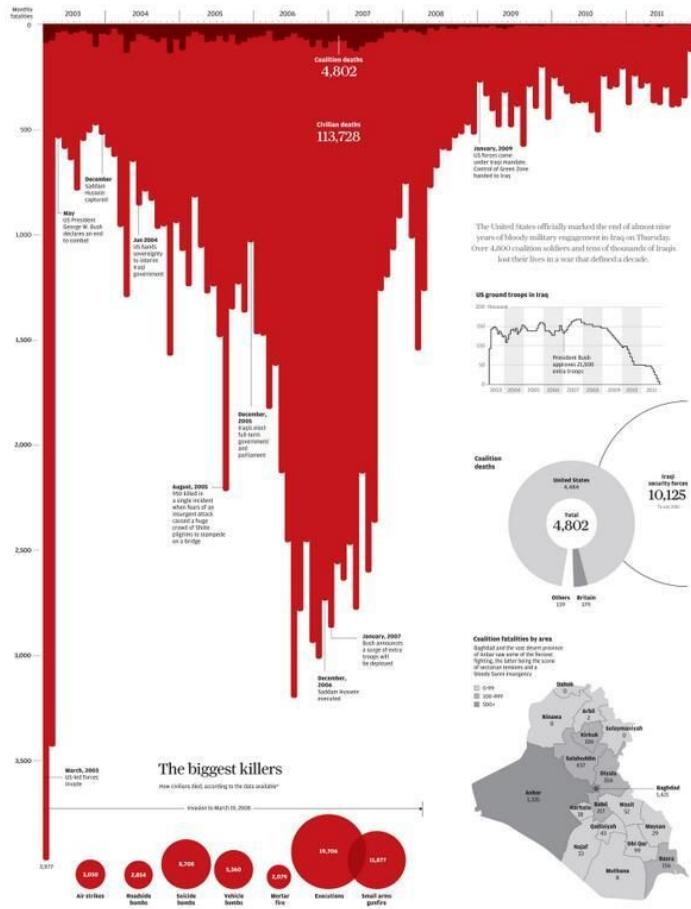
Austria Solar Energy Report



<https://www.behance.net/gallery/2986075/The-Solar-Annual-Report-powered-by-the-sun>

HOT OR NOT

Iraq's bloody toll



Last week, we asked whether
A levels are becoming harder
to pass. You said:



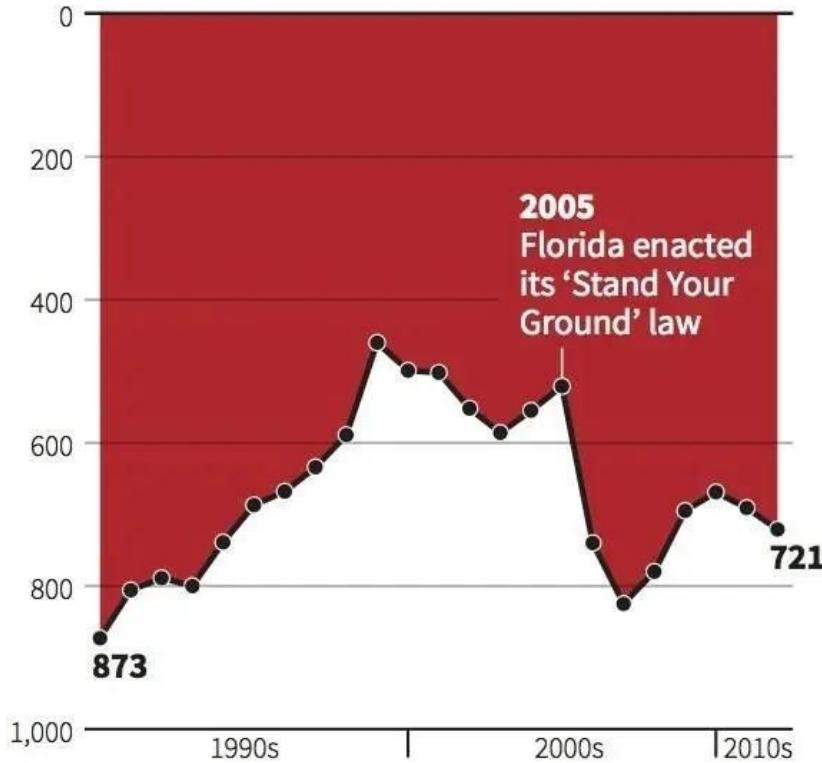
ch
rec
cli
C
S

Figures are rounded

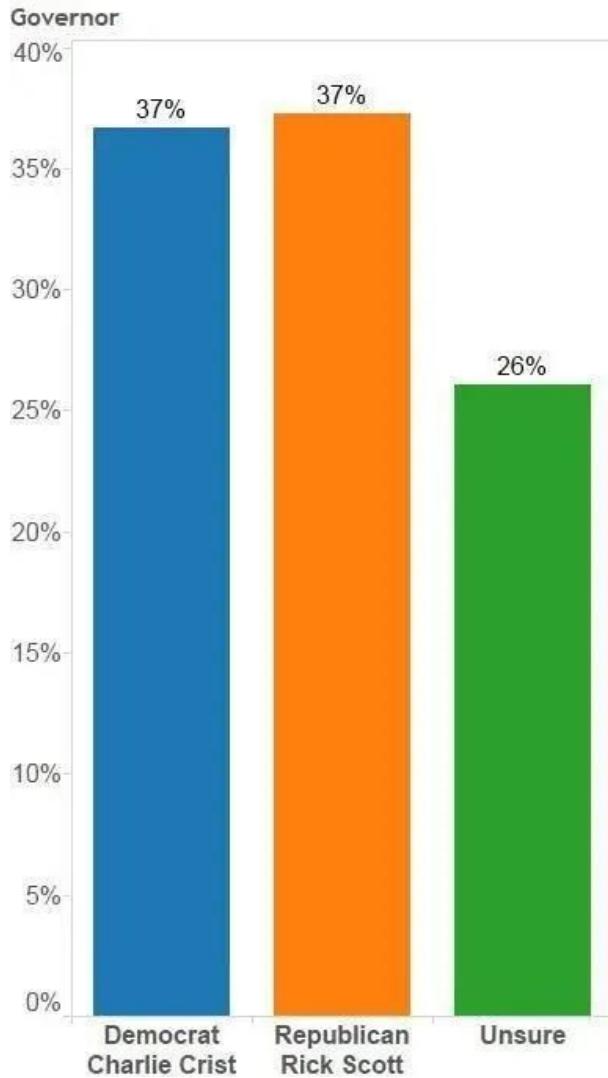
This week's question is:
Are GCSE entries a good

Gun deaths in Florida

Number of murders committed using firearms



Source: Florida Department of Law Enforcement



SCOTCEN POLL OF POLLS

SHOULD
SCOTLAND BE
INDEPENDENT?



NO

52%

YES

58%

LIVE

CNN

NAS ▲ 28.30

HAS KILLED MORE THAN 2,600 IN WEST AFRICA. WORLD HEALTH ORG. FEBROOKEBCNN

talks should happen soon.

BY
PRI

T
the
are
T
ry:
C
the
Am
Ho
by
the
You
the
Ho
day
C
I
Cla
St.
afte
No
I
tifi
asti
Gal
I

Question Of The Day

**What should cost less: a
gallon of gas or a gallon of milk?**

YES

43%



NO

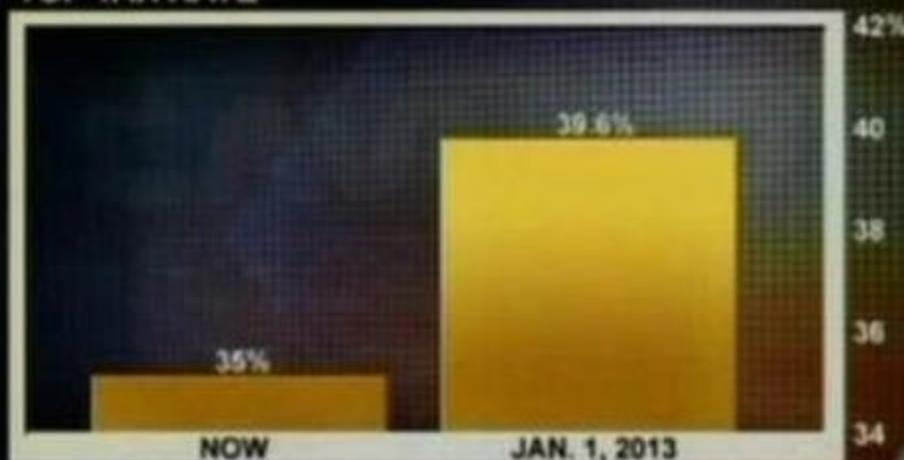
57%



To record your answer to the Question of the Day and other polls about today's news, visit www.post-journal.com.

IF BUSH TAX CUTS EXPIRE

TOP TAX RATE



8:01 p ET



TOP STORIES

TECHNOLOGY

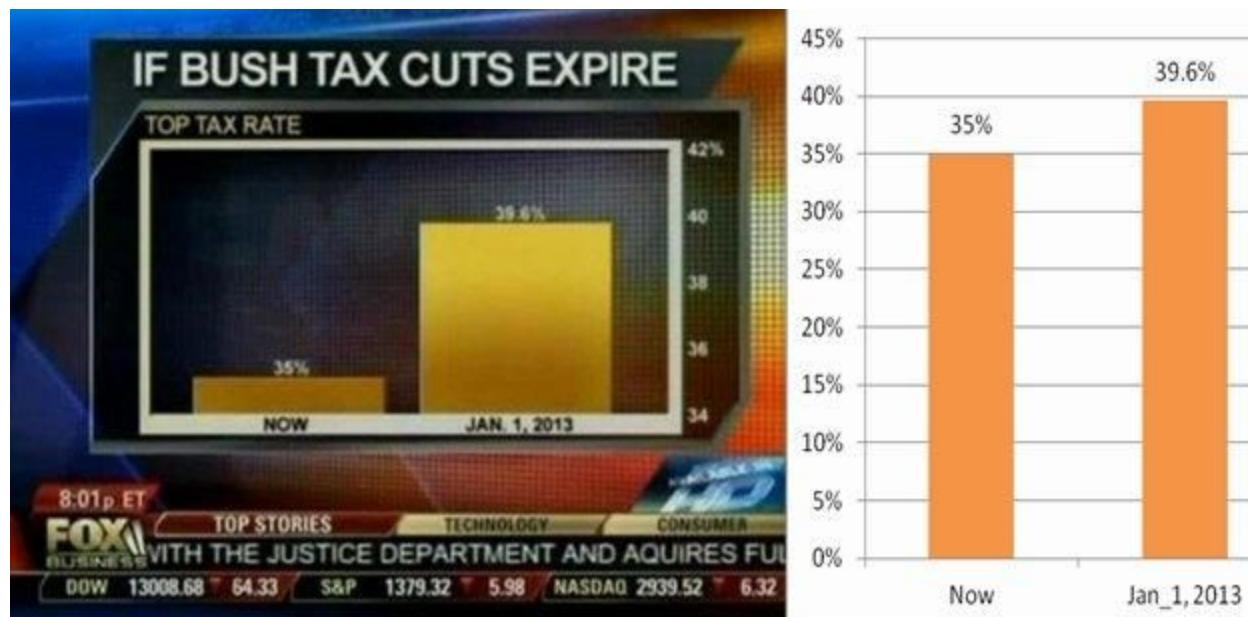
CONSUMER

WITH THE JUSTICE DEPARTMENT AND AQUIRES FULL T

DOW 13008.68 □ 64.33

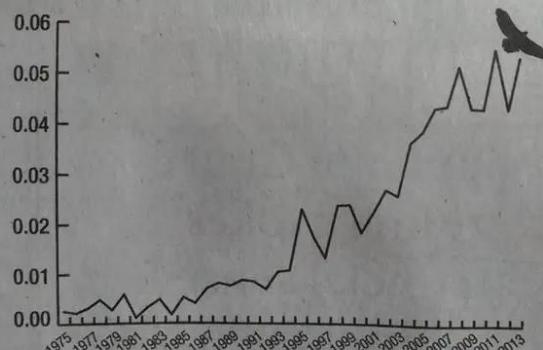
S&P 1379.32 □ 5.98

NASDAQ 2939.52 □ 6.32



Cooper's hawk population soars

A deadly combination of shootings and a pesticide, DDT, caused the Cooper's hawk population in Illinois to stay at low levels throughout the 20th century. However, over the past few years, the raptor has made a strong comeback.



SOURCE: NATURAL HISTORY SURVEY

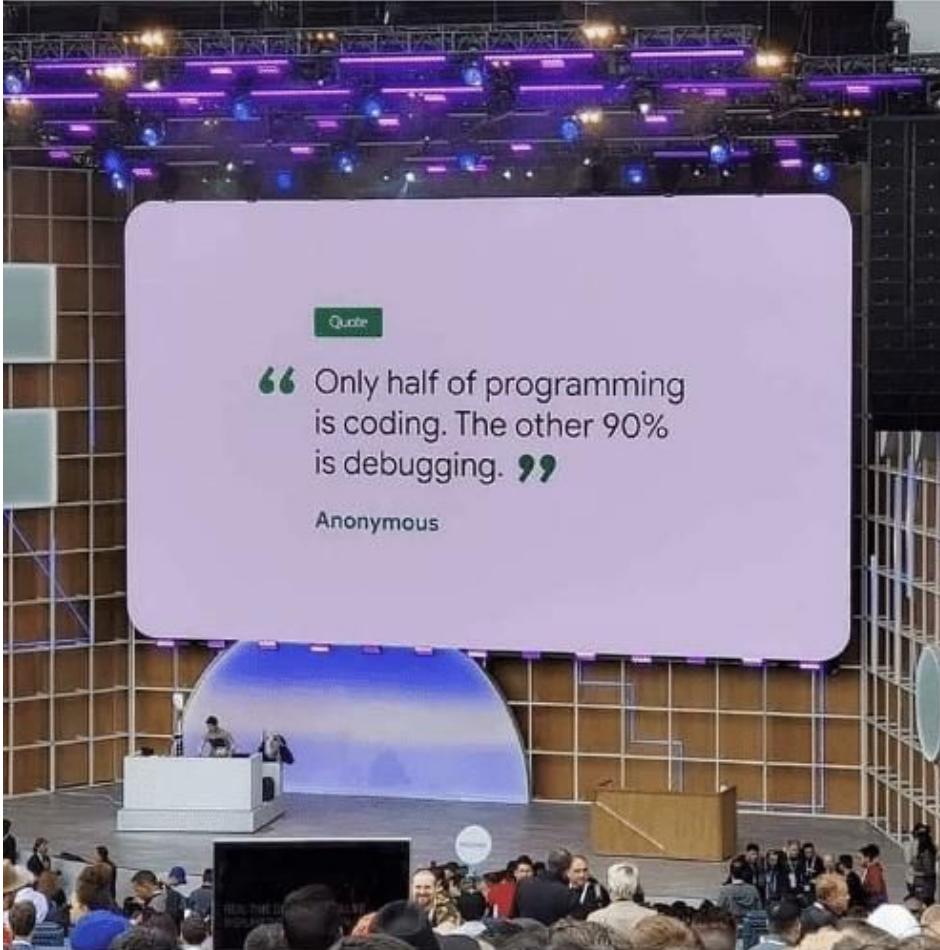
AUSTIN BAIRD THE DAILY ILLINI



Quote

“ Only half of programming
is coding. The other 90%
is debugging. ”

Anonymous



по данным ЦИК

Ростовская область

1	Единая Россия	58,99%
2	КПРФ	32,96%
3	ЛДПР	23,74%
4	Справедливая Россия	19,41%
5	Яблоко	9,32%
6	Патриоты России	1,46%
7	Правое дело	0,59%

В Госдуму пройдут партии, набравшие более 7% голосов

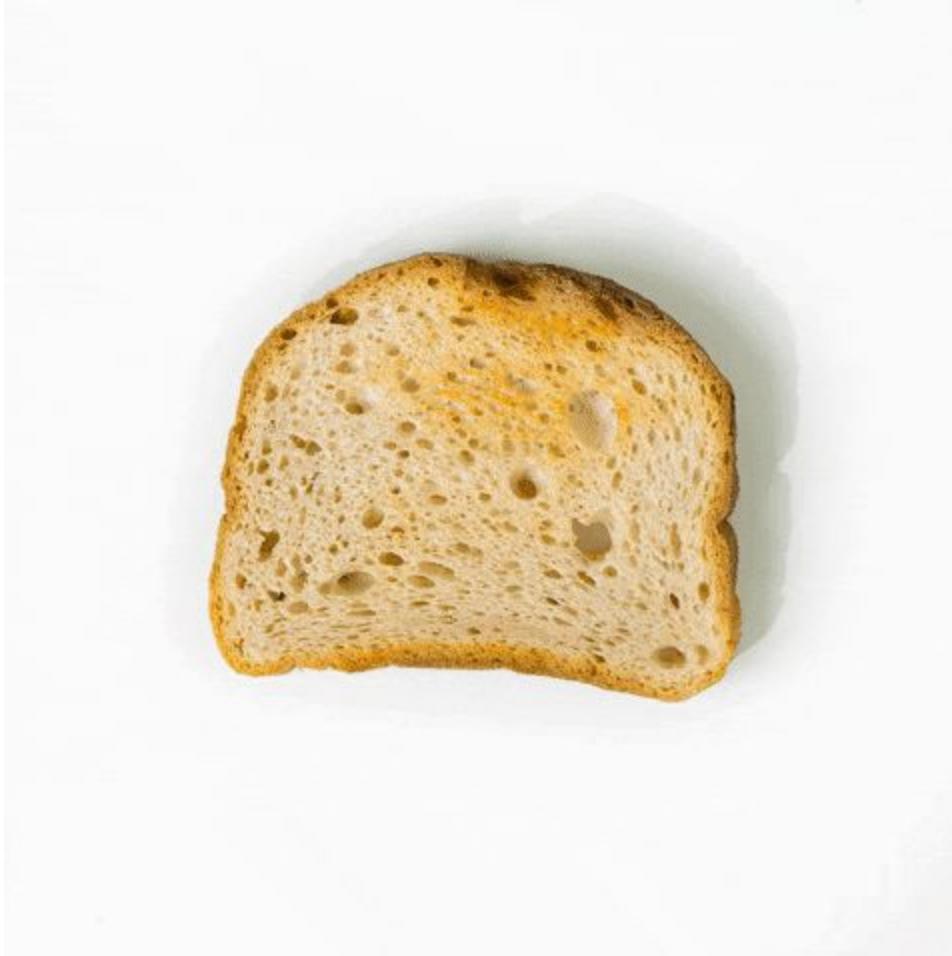
22:58 | ЛУКОЙЛ ММВБ

1704,00 | -0,06%

Setup Google Collab



plotly







Dataset



1248

Avocado Prices

Historical data on avocado prices and sales volume in multiple US markets



Justin Kiggins • updated 2 years ago (Version 1)

Data

Tasks

Kernels (167)

Discussion (10)

Activity

Metadata

Download (2 MB)

New Notebook



Relevant columns: Date

- **Date** - The date of the observation
- | | Date |
|---------------------|-------------|
| ○ First: 2015-01-04 | 2015-12-27 |
| ○ Last: 2018-03-25 | 2015-12-20 |
| | 2015-12-13 |
| | 2015-12-06 |
| | 2015-11-29 |
| | 2015-11-22 |
| | 2015-11-15 |

Relevant columns: AveragePrice

- AveragePrice - the average price of a single avocado

- Median: 1.37
- Mean: 1.40
- std: 0.40
- min: 0.44
- max: 3.2

AveragePrice
1.33
1.35
0.93
1.08
1.28
1.26
0.99

Relevant columns: Type

- **Type** - conventional or organic
 - Equal distribution of both types

type
organic
conventional
organic
conventional
organic
organic
conventional

QUIZ: What type of data is column 'type'

Relevant columns: region

- **Region** - the city or region of the observation
 - 54 unique values

region
Chicago
PhoenixTucson
RaleighGreensboro
Southeast
Midsouth
WestTexNewMexico

Relevant columns: Total Volume

- **Total Volume** - Total number of avocados sold
 - mean: 850644.0
 - median: 107376.76
 - std: 3453545.0
 - min: 84.56
 - max : 2505650.0

Relevant columns: 4225

- **4225** - Total number of avocados with PLU 4225 sold



Relevant columns: 4770

- 4770 - Total number of avocados with PLU 4770 sold



Relevant columns: 4046

- 4046 - Total number of avocados with PLU 4046 sold





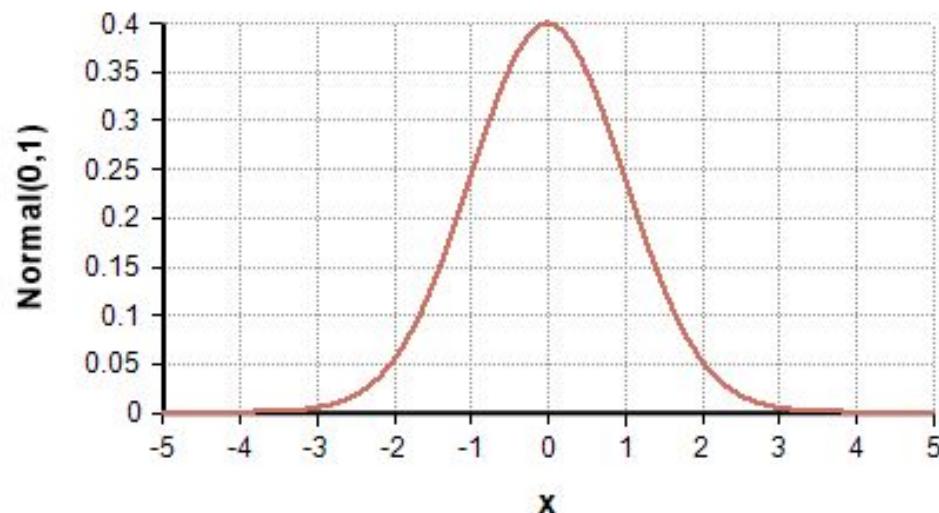
plotly

Plotly Express

Plotly Express is the easy-to-use, high-level interface to Plotly, which operates on "tidy" data and produces easy-to-style figures.

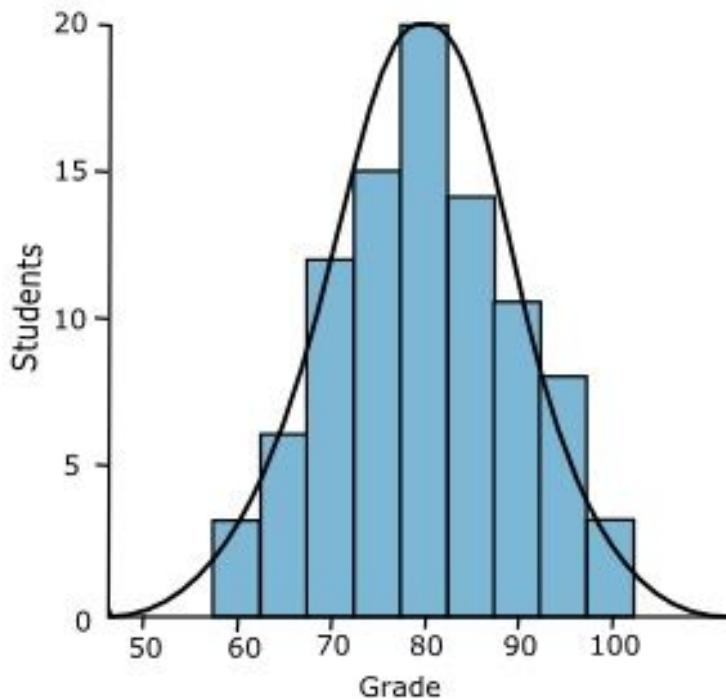


Remember this?



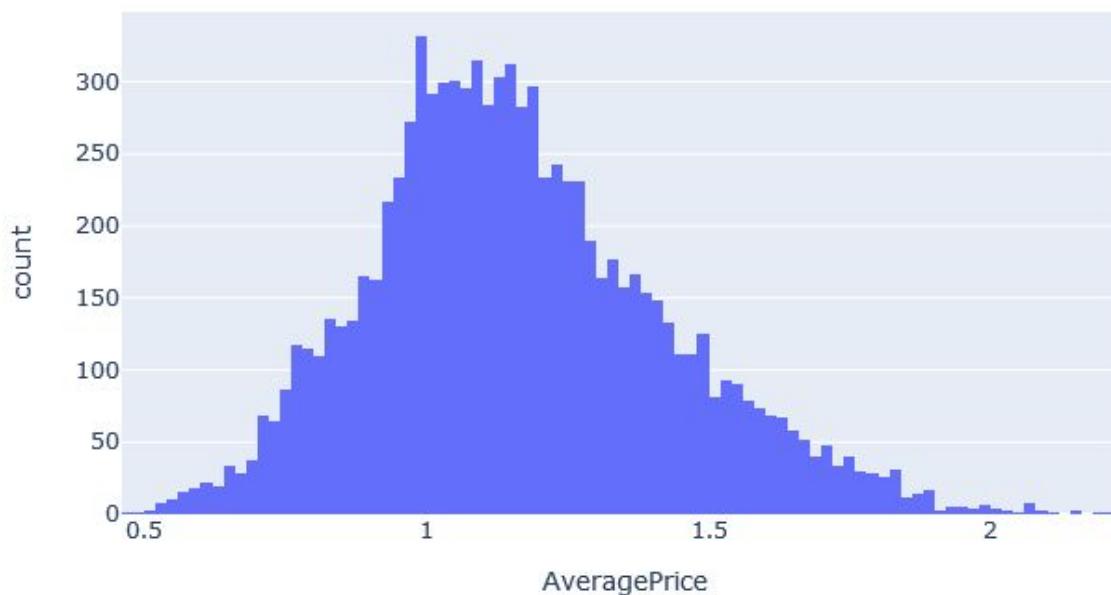
Histogram

A histogram is a graphical display of data using bars of different heights. In a histogram, each bar groups numbers into ranges. Taller bars show that more data falls in that range. A histogram displays the shape and spread of continuous sample data.



Histogram of Conventional Avocado Prices

Distribution of prices on Conventional Avocado

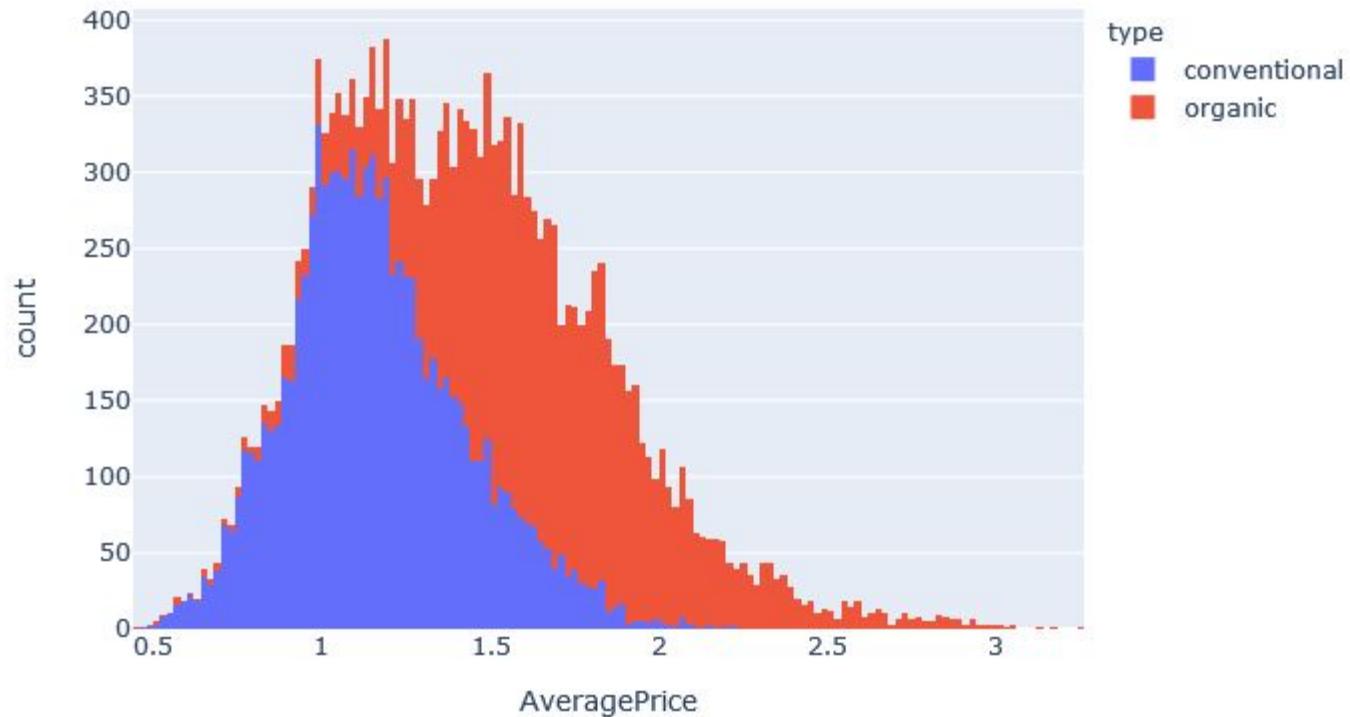


DEMO

**What is price range of the highest count?
(organic)**

1.58-1.59

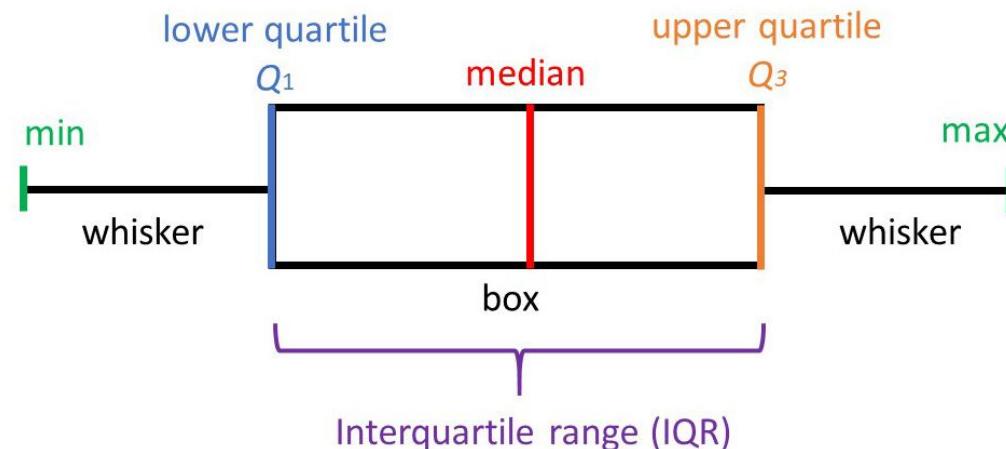
Comparing two histograms

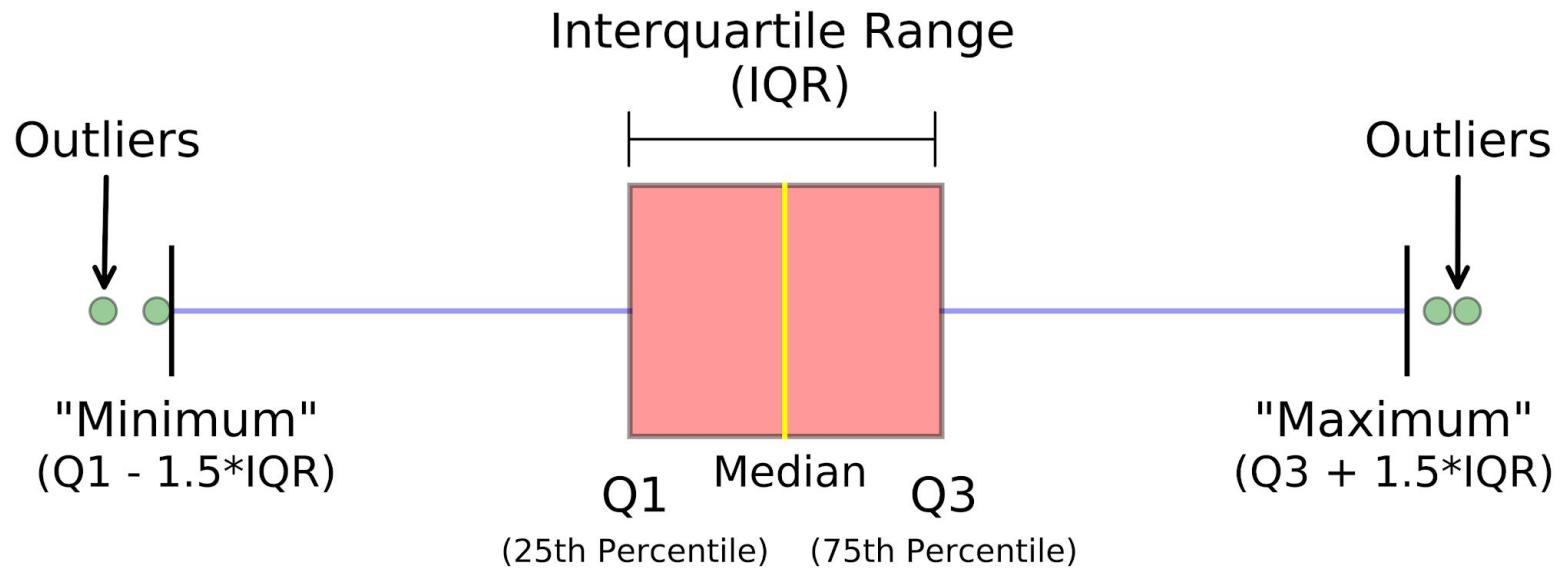


DEMO

Box Plot

A box plot displays the five-number summary of data set. The five-number summary is the minimum, first quartile, median, third quartile, and maximum.





-4

-3

-2

-1

0

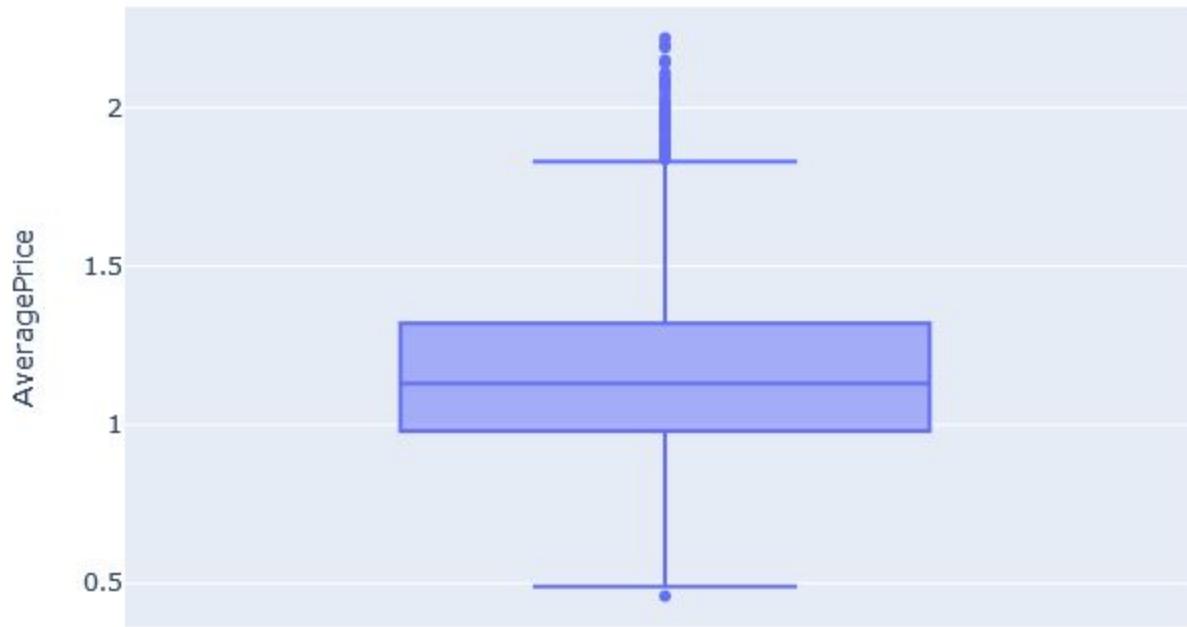
1

2

3

4

Box Plot of Conventional Avocado Prices



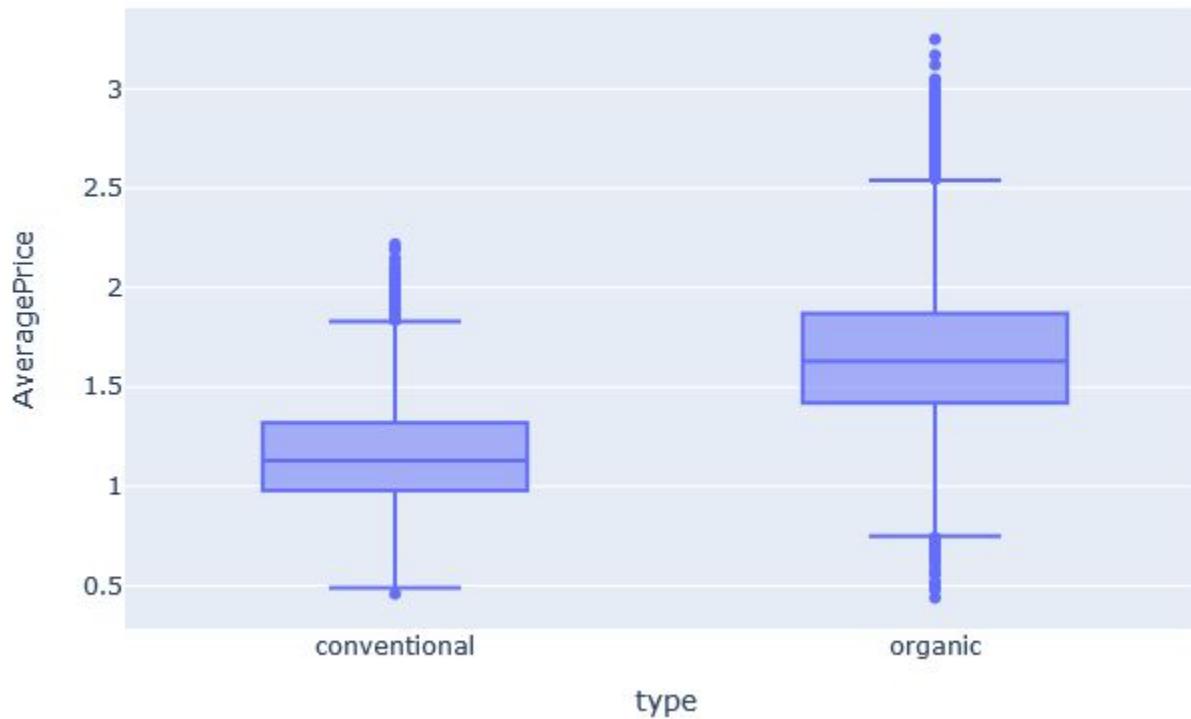
DEMO

Quiz Plot box plot for organic avocado

What is value of upper fence?

Answer: 2.54

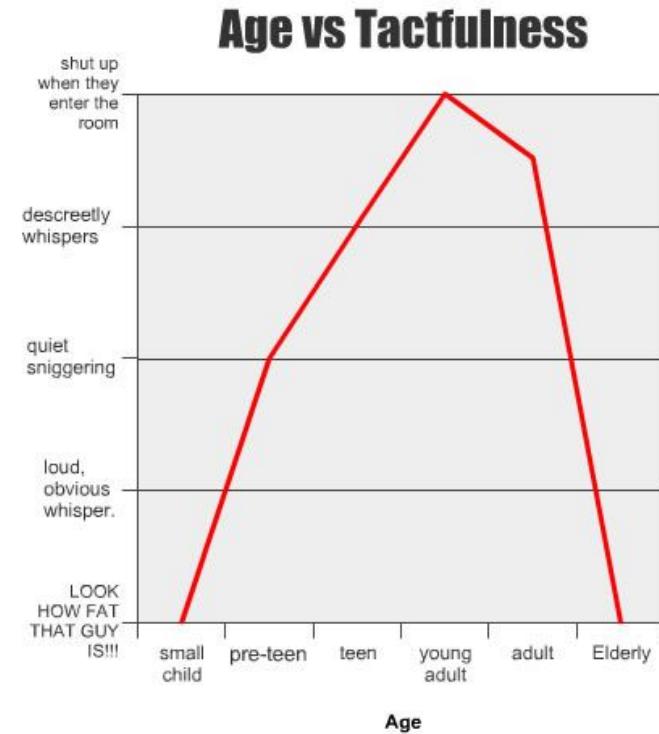
Box Plot for each type of Avocado



DEMO

Line Plot

A line chart or line plot or line graph or curve chart is a type of chart which displays information as a series of data points called 'markers' connected by straight line segments.



Average Price of Conventional Avocado Over Time



DEMO

Quiz: make organic plot. What is the date of the highest peak.

27. Aug 2017

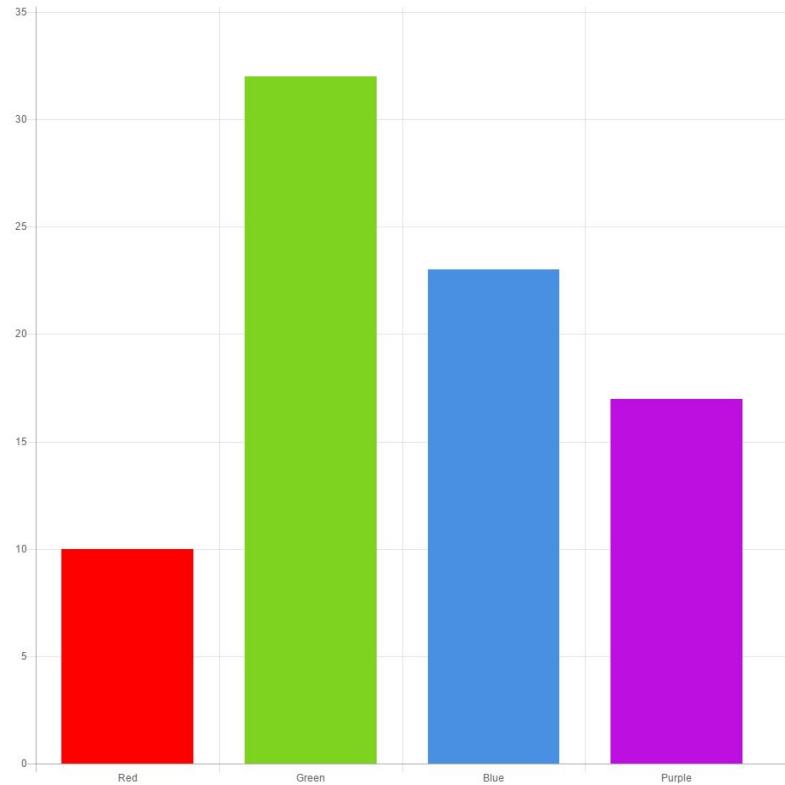
Average Price of Avocado Over Time



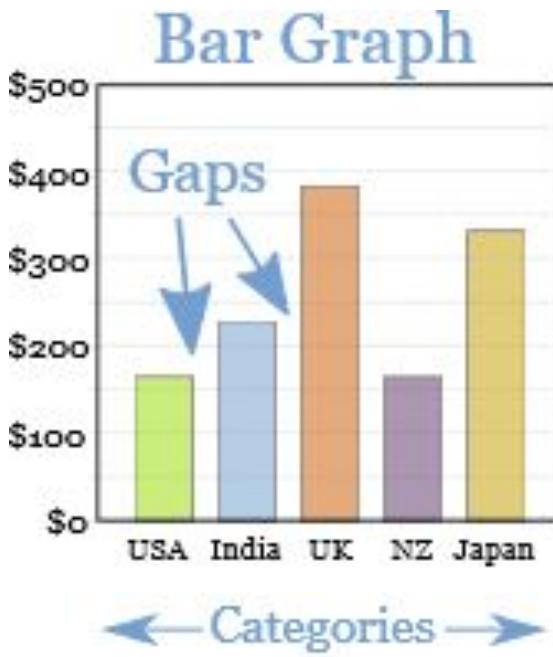
DEMO

Bar plot

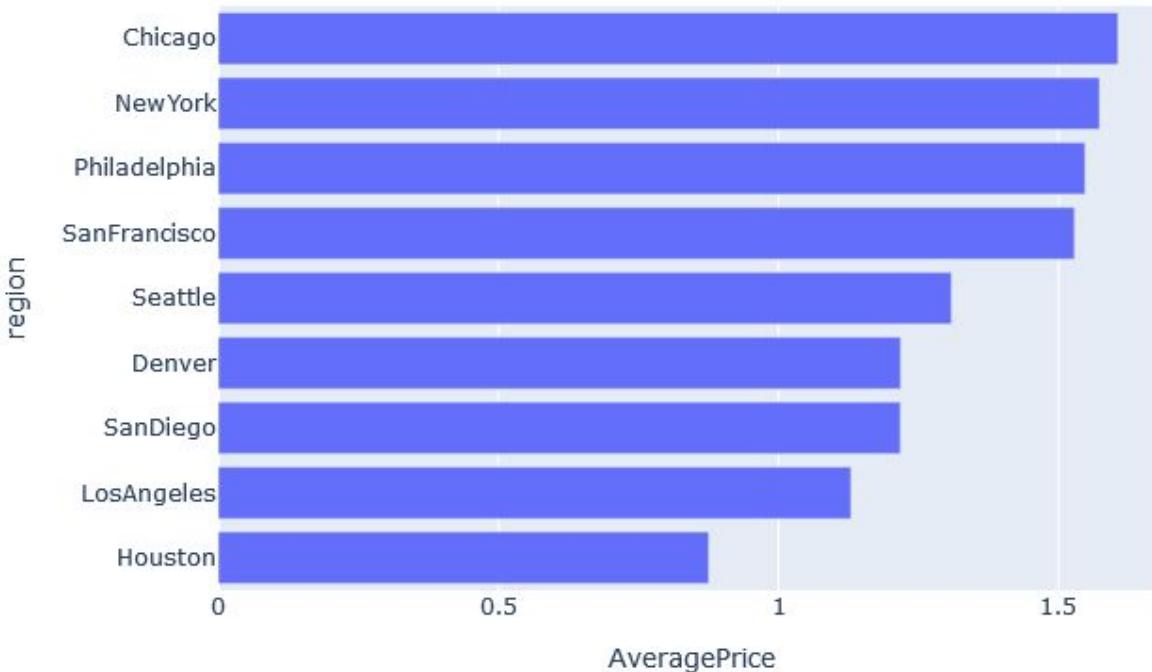
A bar chart or bar graph is a chart or graph that presents categorical data with rectangular bars with heights or lengths proportional to the values that they represent. The bars can be plotted vertically or horizontally.



Bar Plot vs Histogram



Average Price of Conventional Avocado in 2017



DEMO

Quiz df_2017. How many rows?

Answer: 5722

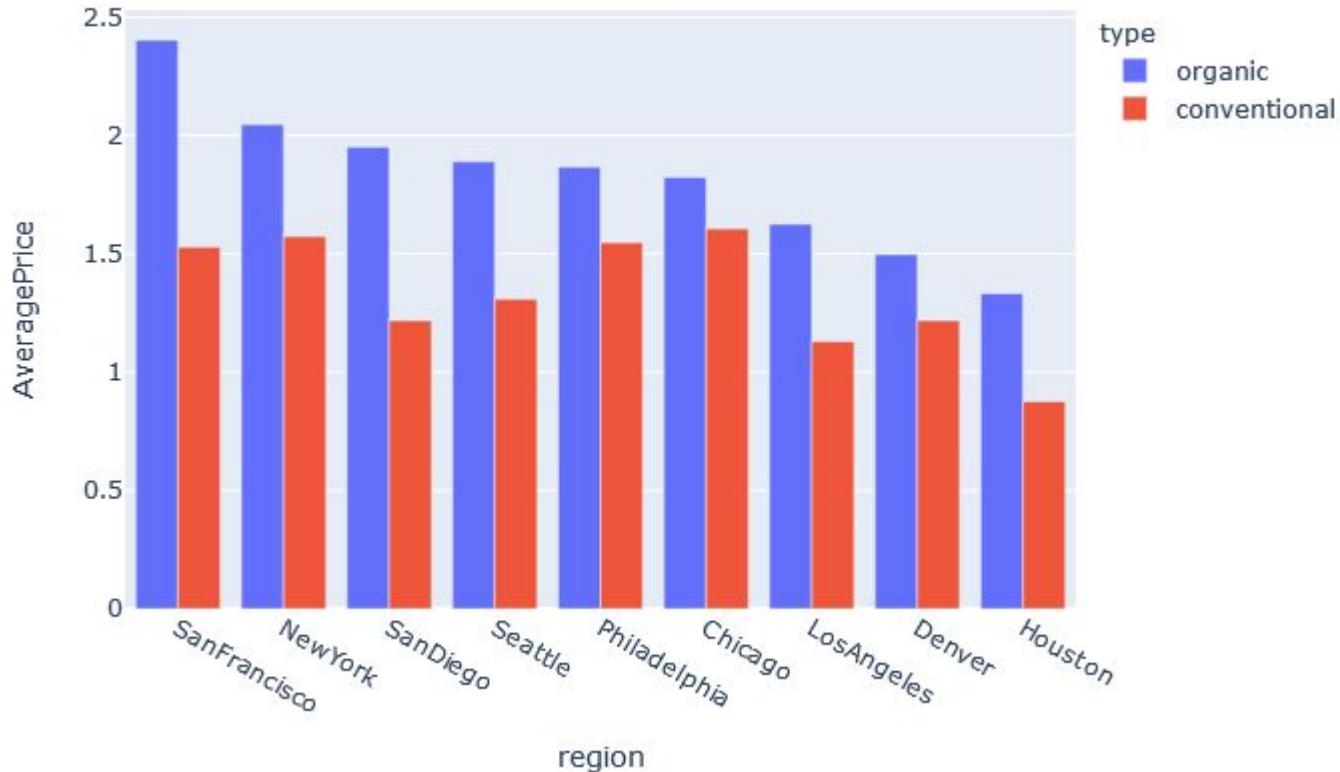
Groupby Exercice.

How many rows and columns does "df_2017_region_price" have?

108,3

DEMO

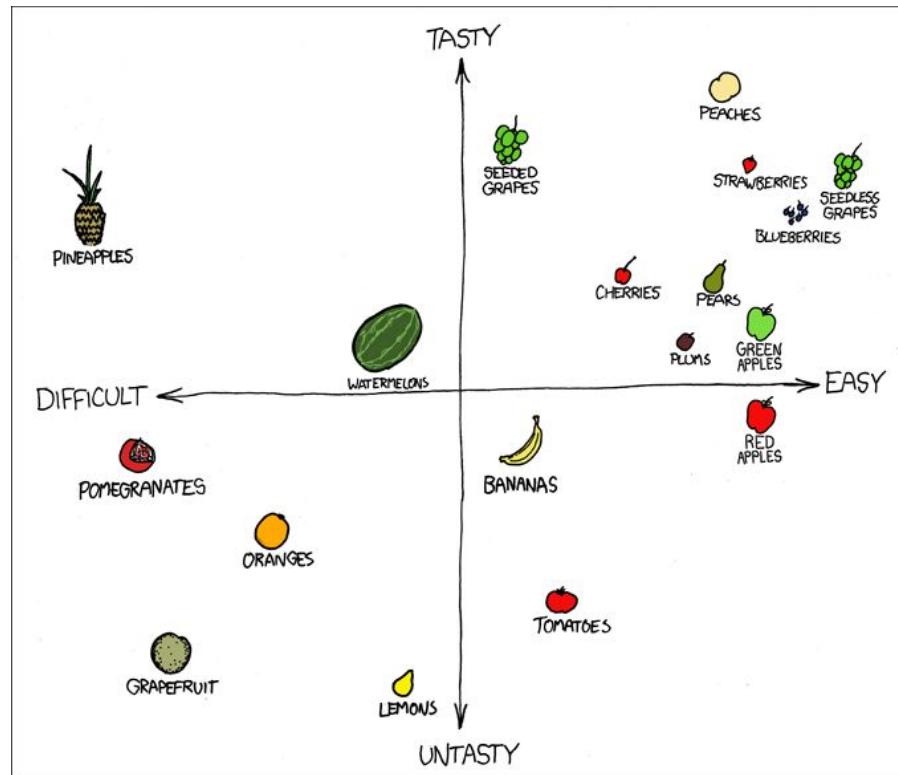
Average Price of Avocado in 2017



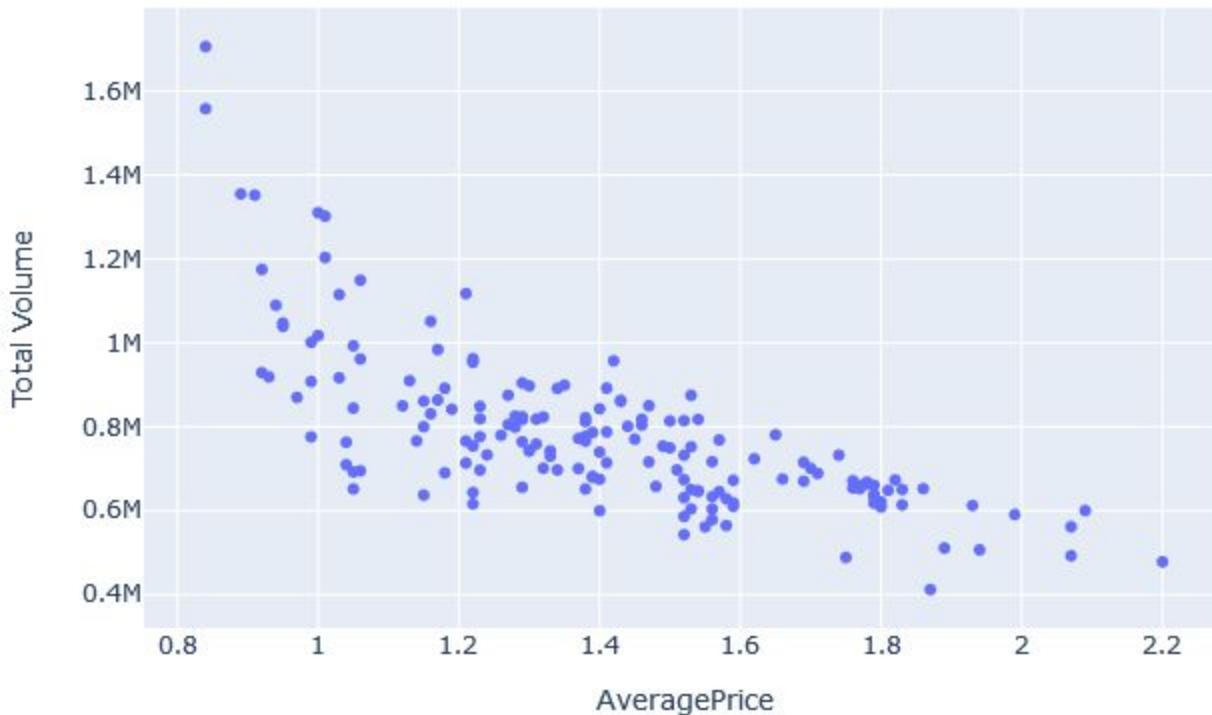
DEMO

Scatter Plot

A scatter plot uses coordinates to display values for typically two variables for a set of data. If the points are coded (color/shape/size), one additional variable can be displayed. The data are displayed as a collection of points, each having the value of one variable determining the position on the horizontal axis and the value of the other variable determining the position on the vertical axis.



Correlation between Average Price and Total Volume

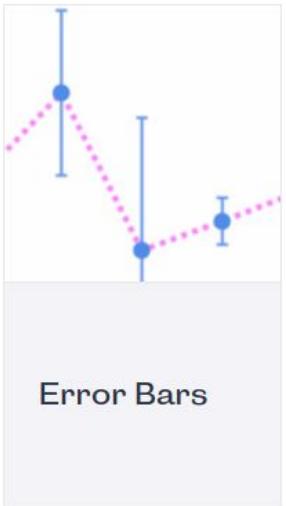


DEMO

Correlation in Houston and San Francisco



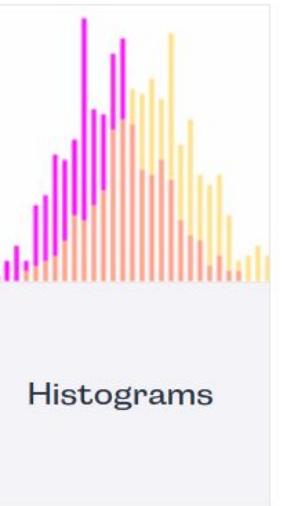
DEMO



Error Bars



Box Plots



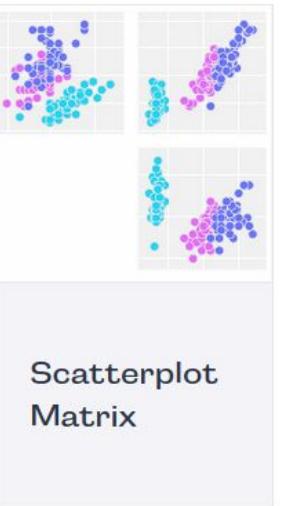
Histograms



Distplots



2D Histograms



Scatterplot Matrix



Facet and Trellis Plots



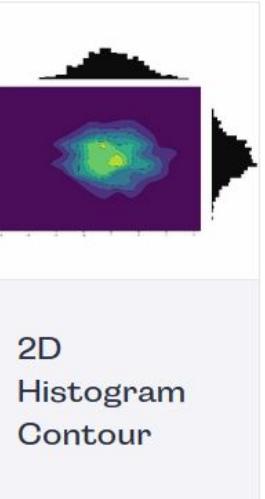
Parallel Categories Diagram



Tree-plots



Violin Plots



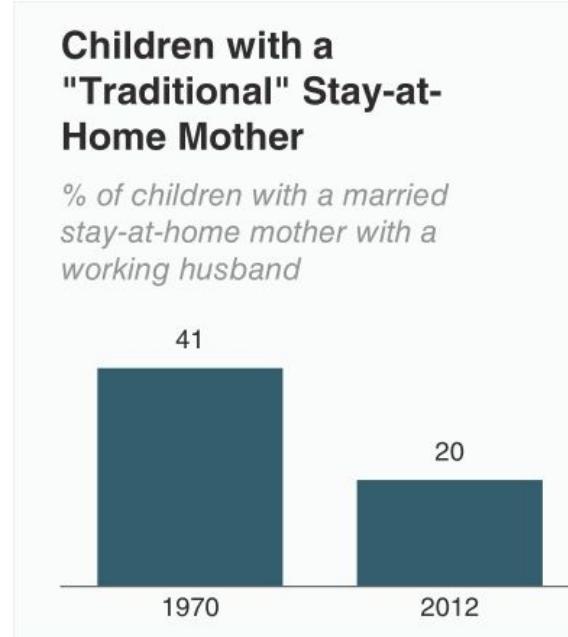
2D Histogram Contour



Linear and Non-Linear Trendlines

Simple text

simple text can be a great way to communicate



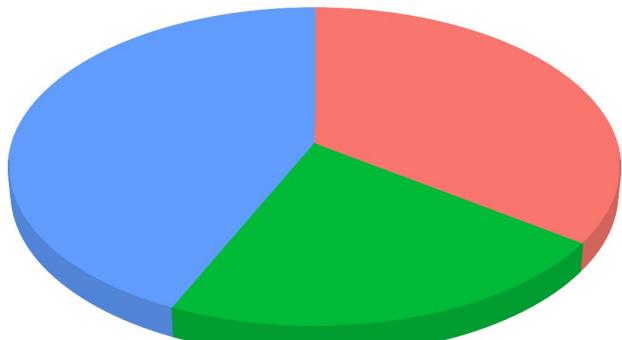
Simple text

simple text can be a great way to communicate

20%

of children had a
traditional stay-at-home mom
in 2012, compared to 41% in 1970

Never use 3D



as.factor(cyl) 4 6 8

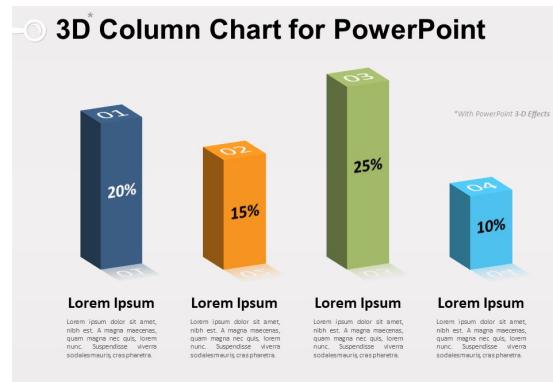


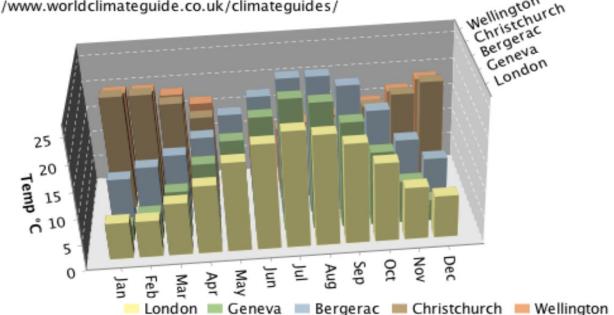
Chart rotation demo
Test options by dragging the sliders below



Highcharts.com

Average Maximum Temperature

<http://www.worldclimateguide.co.uk/climateguides/>

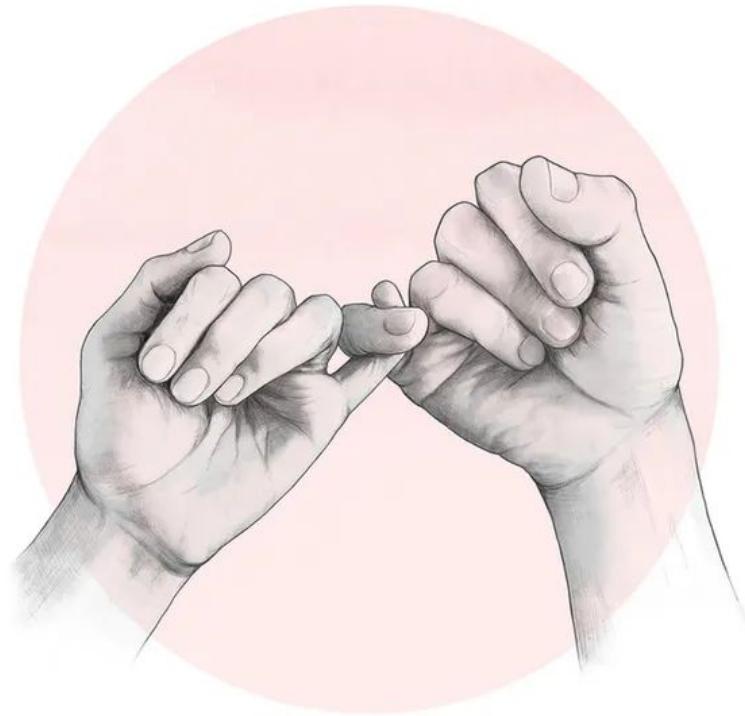


Never use 3D

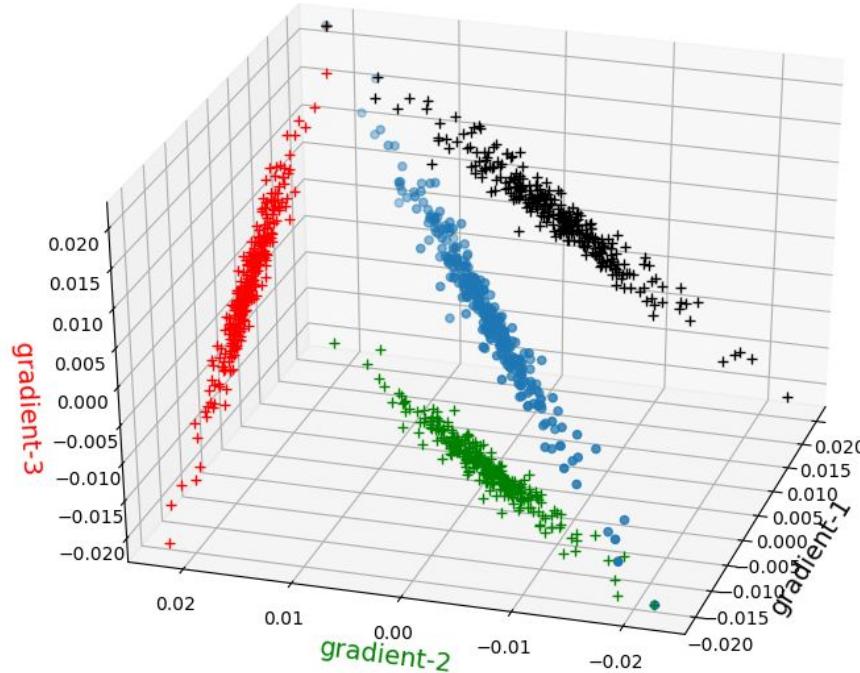
Repeat after me

Repeat after me

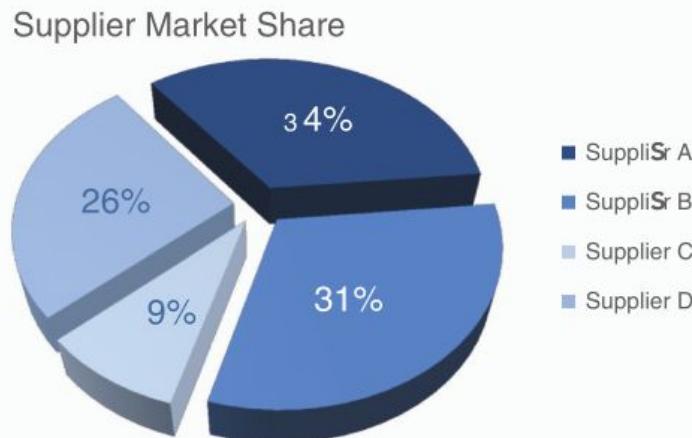
I will **never** use **3D Plots** to impress
my managers.



Exception: 3d Scatter Plots



Bar plot as alternative to pie charts.



The secret of good plots:

The secret of good plots:
Remove the clutter.

**If you don't remember anything
Remember just this.**





“Discard
everything that
does not spark
joy.”

MARIE KONDO

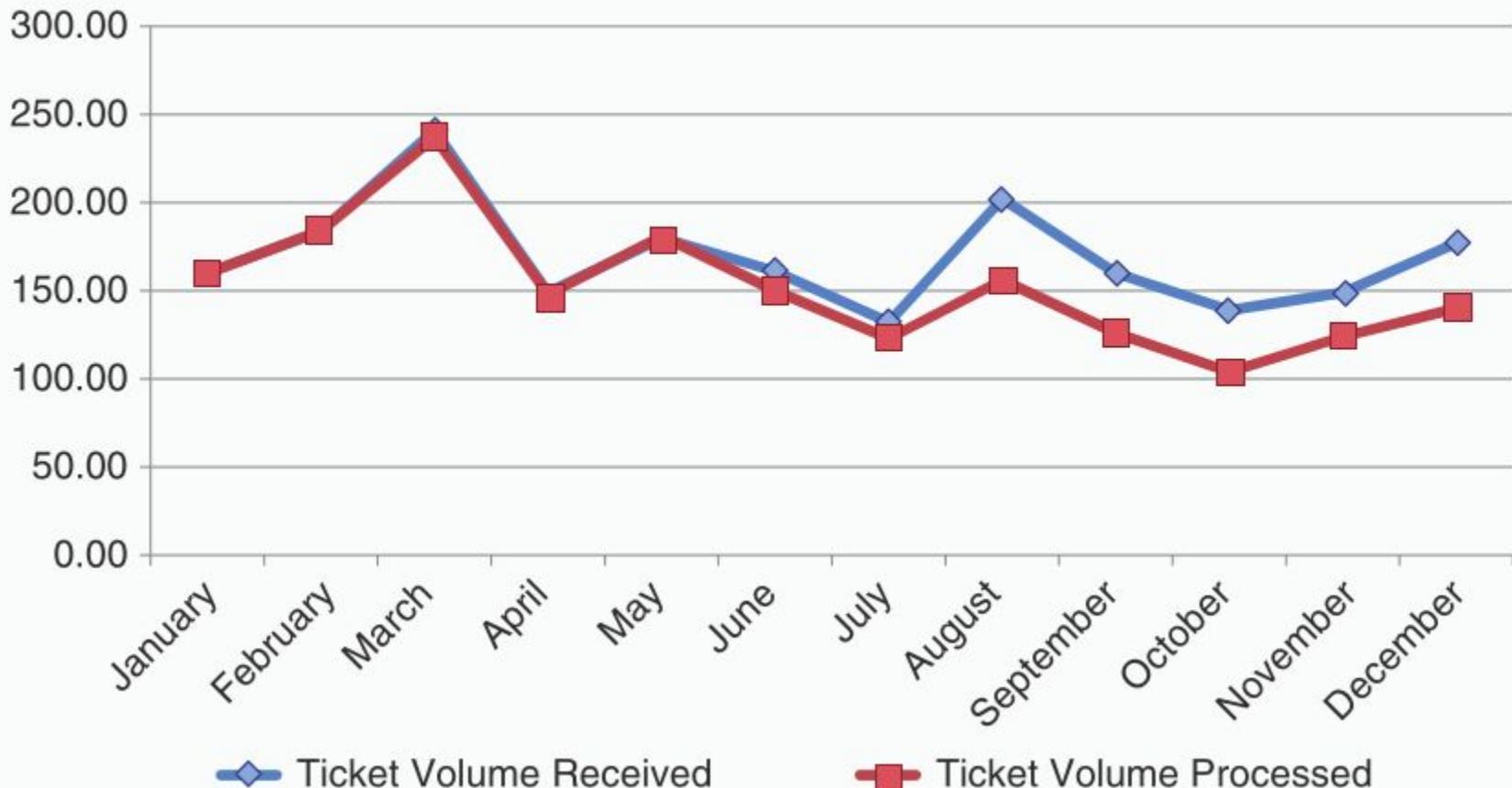
yourtango

Step 1

Remove what has no purpose

Step 1

Remove what has no purpose

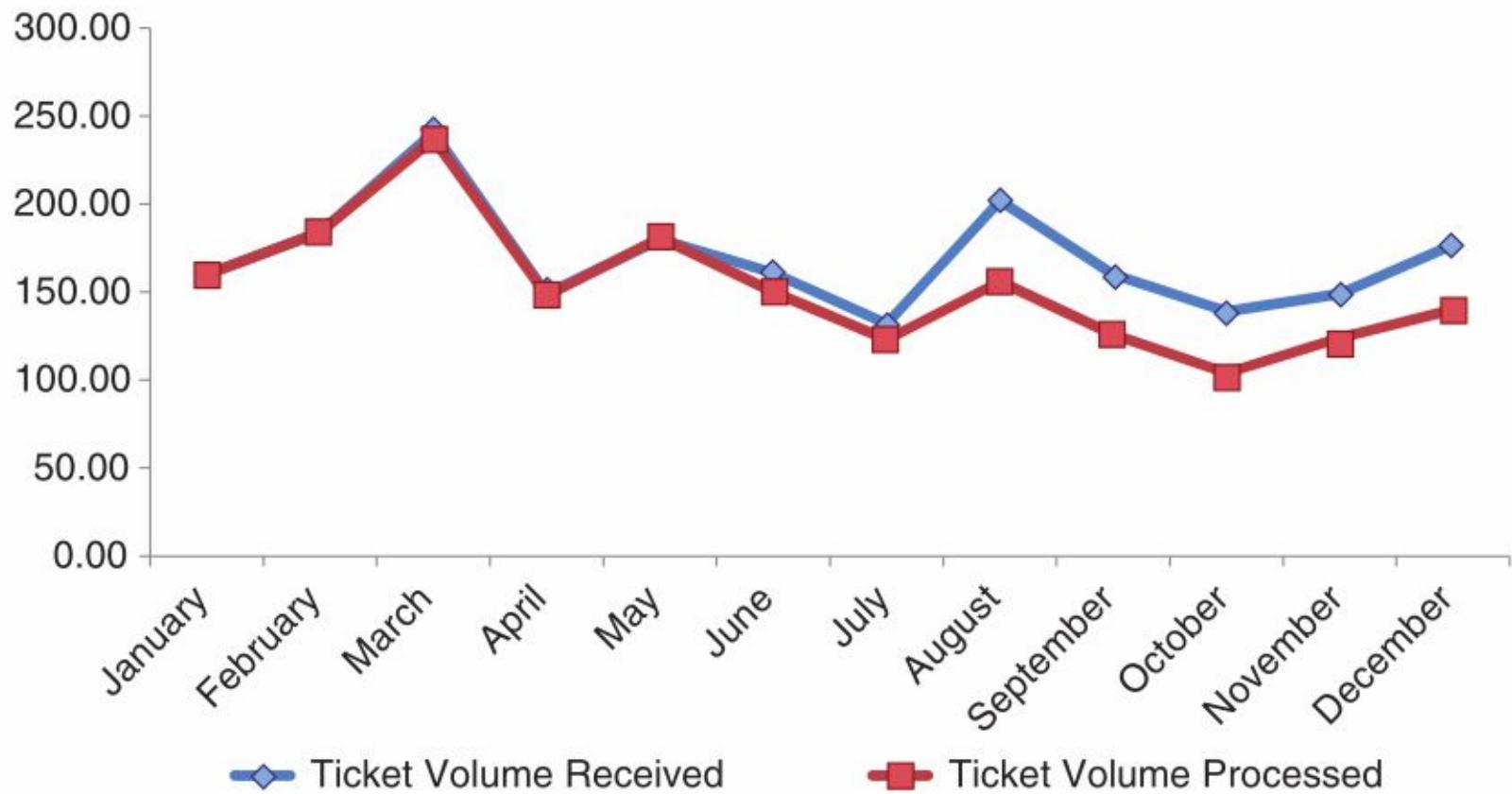


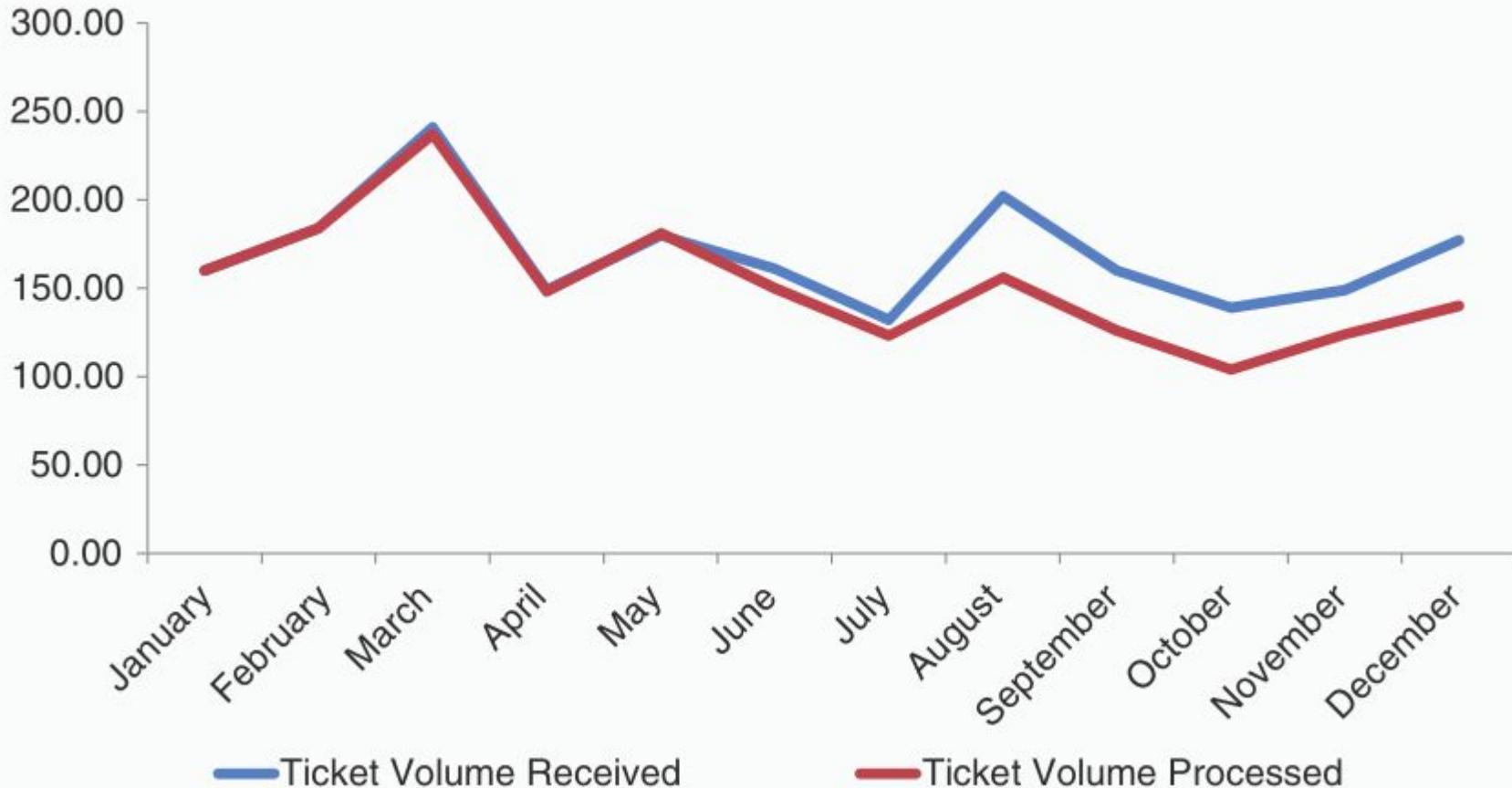
Knafllic, C. N. (2015). *Storytelling with data: a data visualization guide for business professionals*. Hoboken, New Jersey: John Wiley & Sons, Inc.

Gridlines?

Do they help the viewer?

No Gridlines = Better Contrast





Step 2
Clean your axis.



Step 3

The more Obvious The Better.

Step 3
The more **Obvious**
The Better.

3483075861872364872136
4786328756321876407321
6087301647261358768721
5982374872134986213075
7329847983215032814798
3217490257982173498732
1498278749832174903217

How many 0s?

3483**0**75861872364872136

47863287563218764**0**7321

6087301647261358768721

5982374872134986213**0**75

7329847983215**0**32814798

321749**0**257982173498732

14982787498321749**0**3217



Orientation



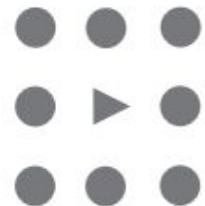
Length



Width



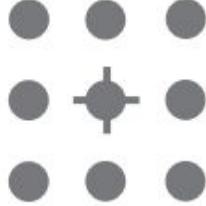
Size



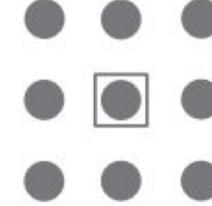
Shape



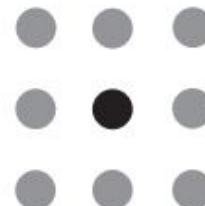
Curvature



Added Marks



Enclosure



Color Value



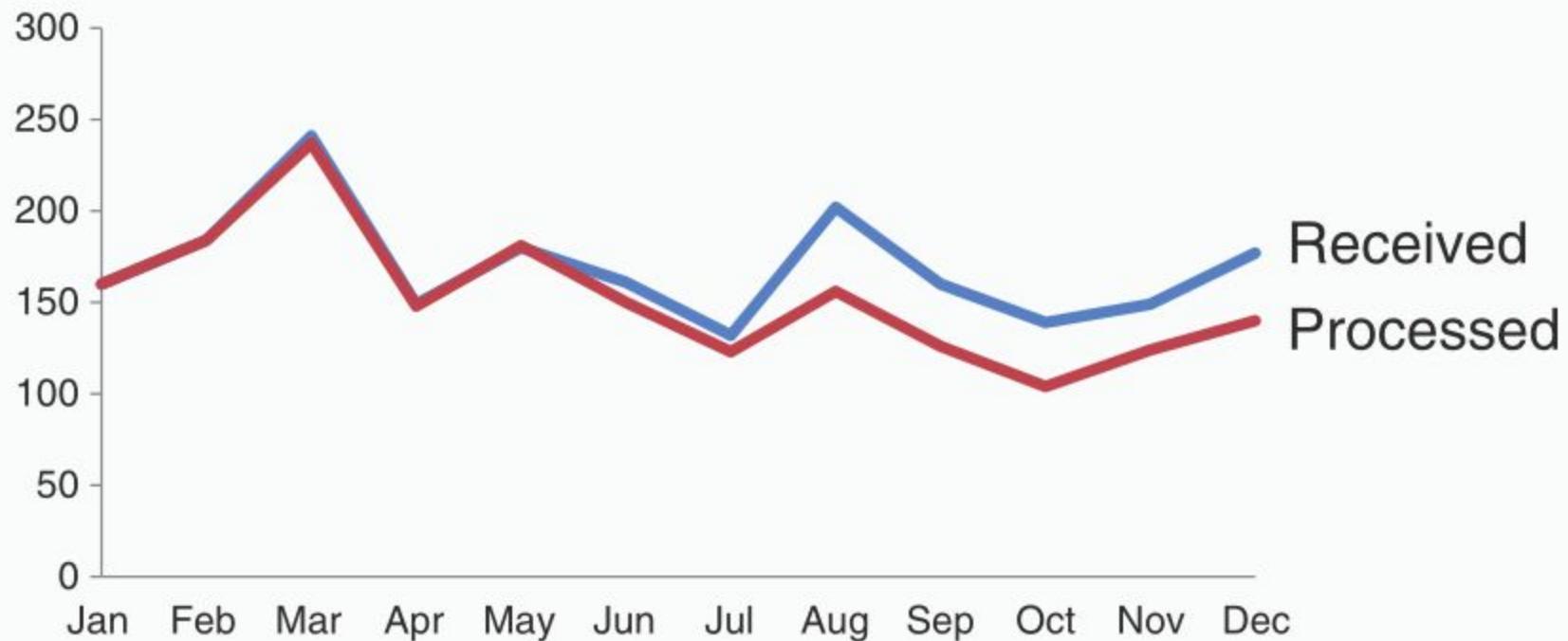
Color Hue



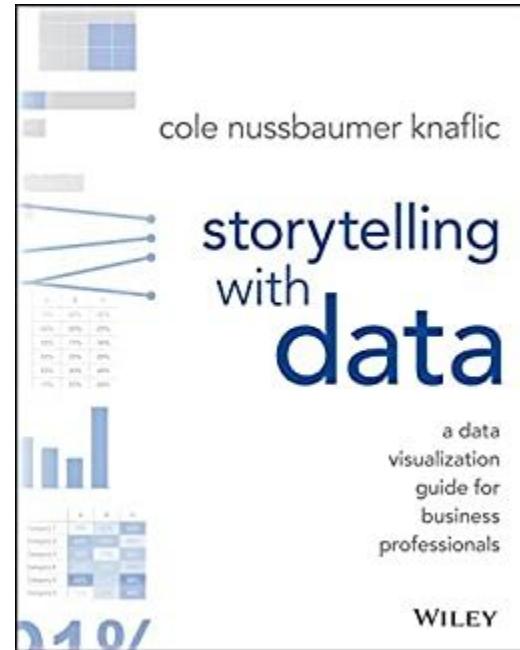
Position



Spatial Grouping



Book Recommendation



And no, dont listen the audiobook.

