

# Machine Learning Concepts

1.1

Some [slides](#) borrowed from the presentation of [Szilard Pafka](#).



Zoltán Papp  
Quant & Machine Learning  
Practitioner/Teacher  
<https://www.linkedin.com/in/pzoltan/>  
[@zozopp](https://twitter.com/zozopp)

The word "WHY?" is rendered in large, bold, red, three-dimensional letters. The letters have a metallic texture and are set against a plain white background. The perspective is from slightly below and to the side, giving a sense of depth.

WHY?

Machine Learning

# Why?

- Drive the car on a highway
- Park the car
- Autonomous driving expected in 2 years by Elon Musk



# Why?

May, 2016: Google's AlphaGo Defeats Lee Sedol, Chinese Go Master

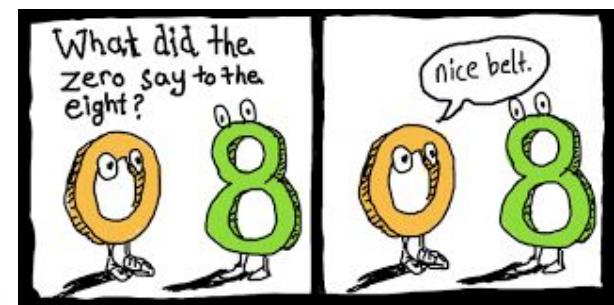


# Why?

In December, 2018: Google's DeepMind Defeats TLO and MANA top ranking professional Starcraft players



# Why?

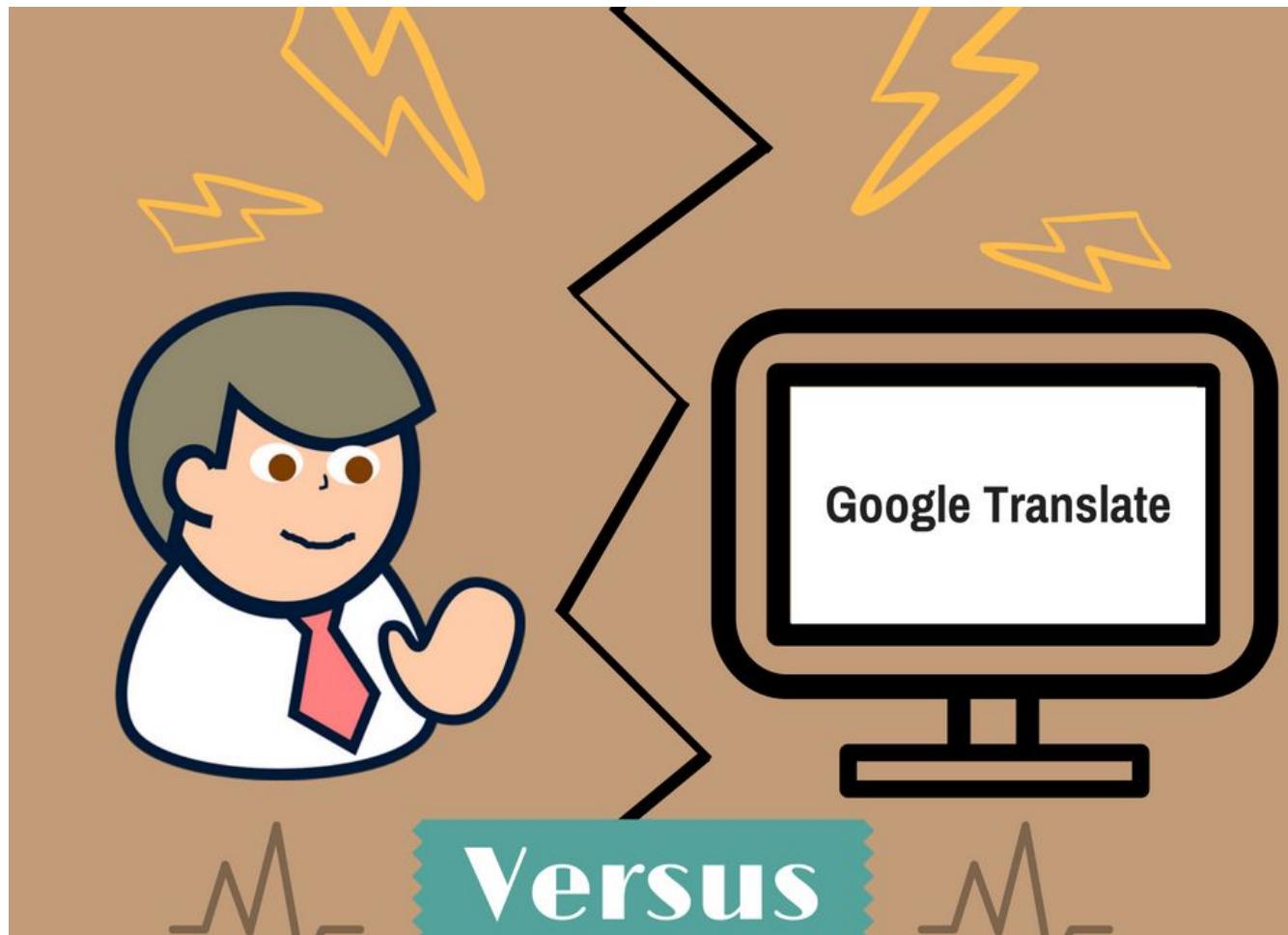


Amazon

# Why?

Google Translate now serves 200 million people daily

A billion translations a day



# What is needed for ML

*“The rate at which we’re generating data is rapidly outpacing our ability to analyze it”*

Professor Patrick Wolfe, Executive Director of the University College of London’s Big Data Institute

*“The trick here is to turn these massive data streams from a liability into a strength.”*



Data

Computing Capacity



So What's  
in it  
for me?

# What is your background?



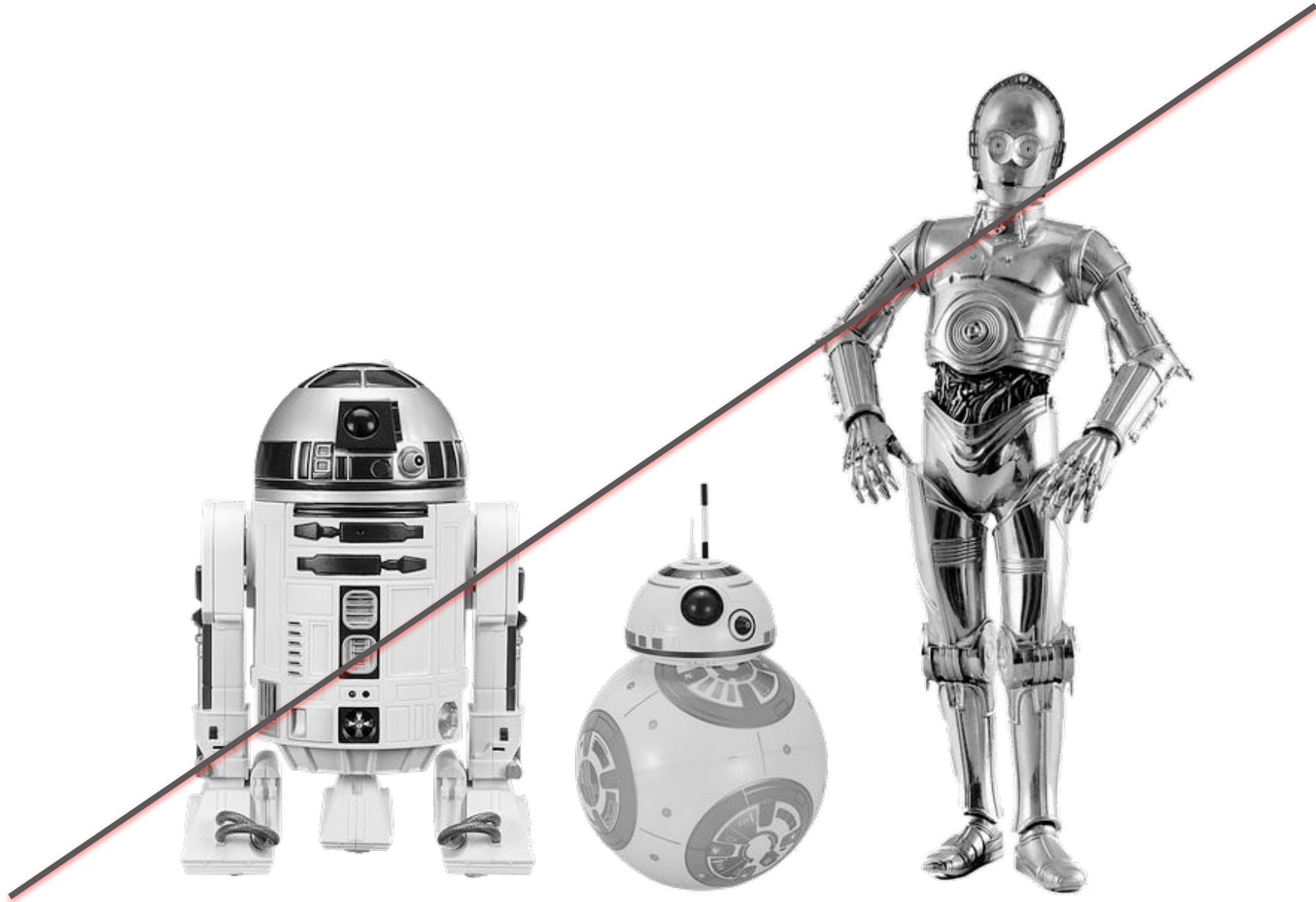
# What do you expect?



# Is it General Artificial Intelligence?



**It is not General Artificial Intelligence!**



# Course assignments and exam

Grading:

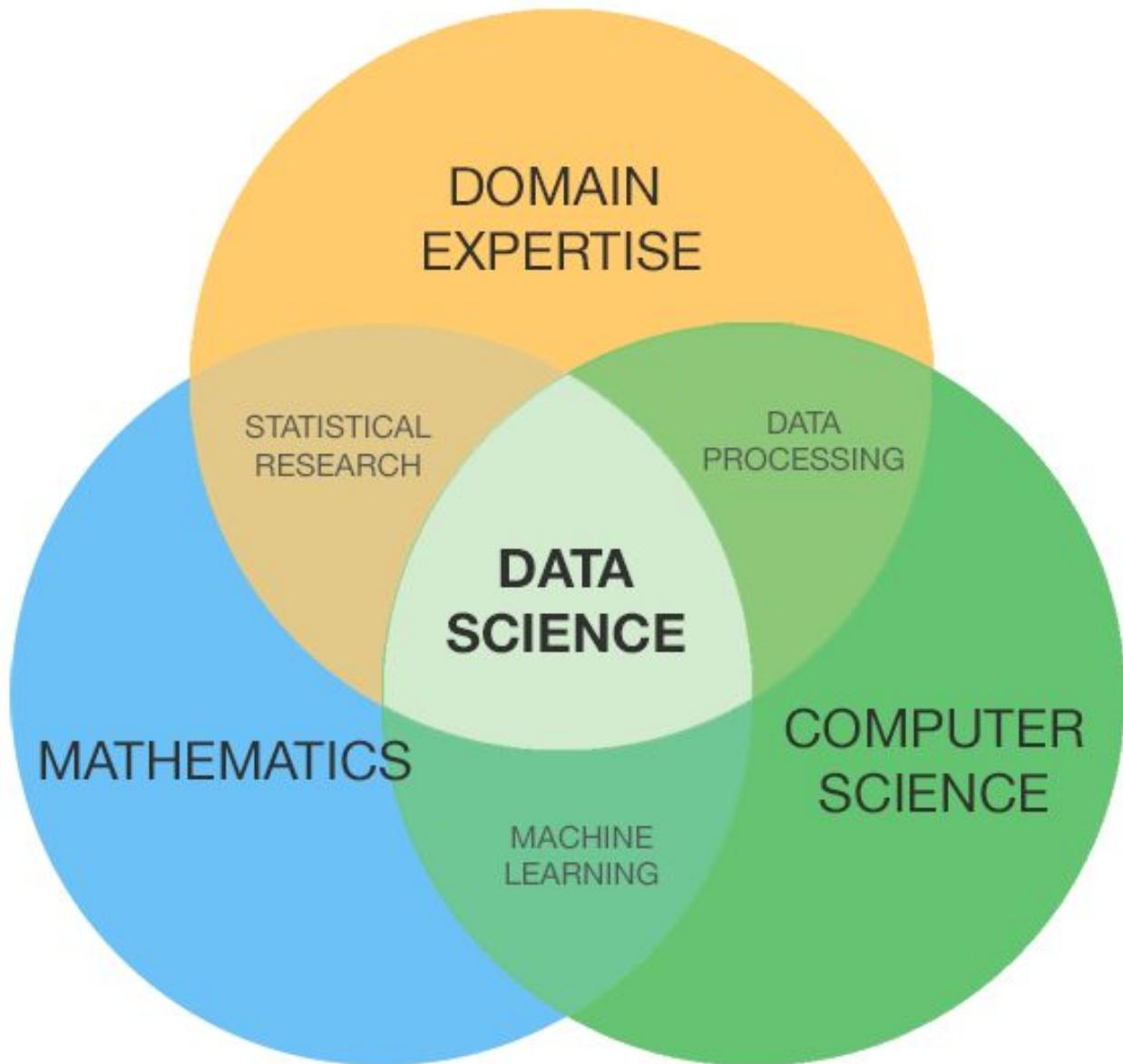
Exam and Assignment dates and deadlines:

Weekly assignment acceptance policy and achievable grades:

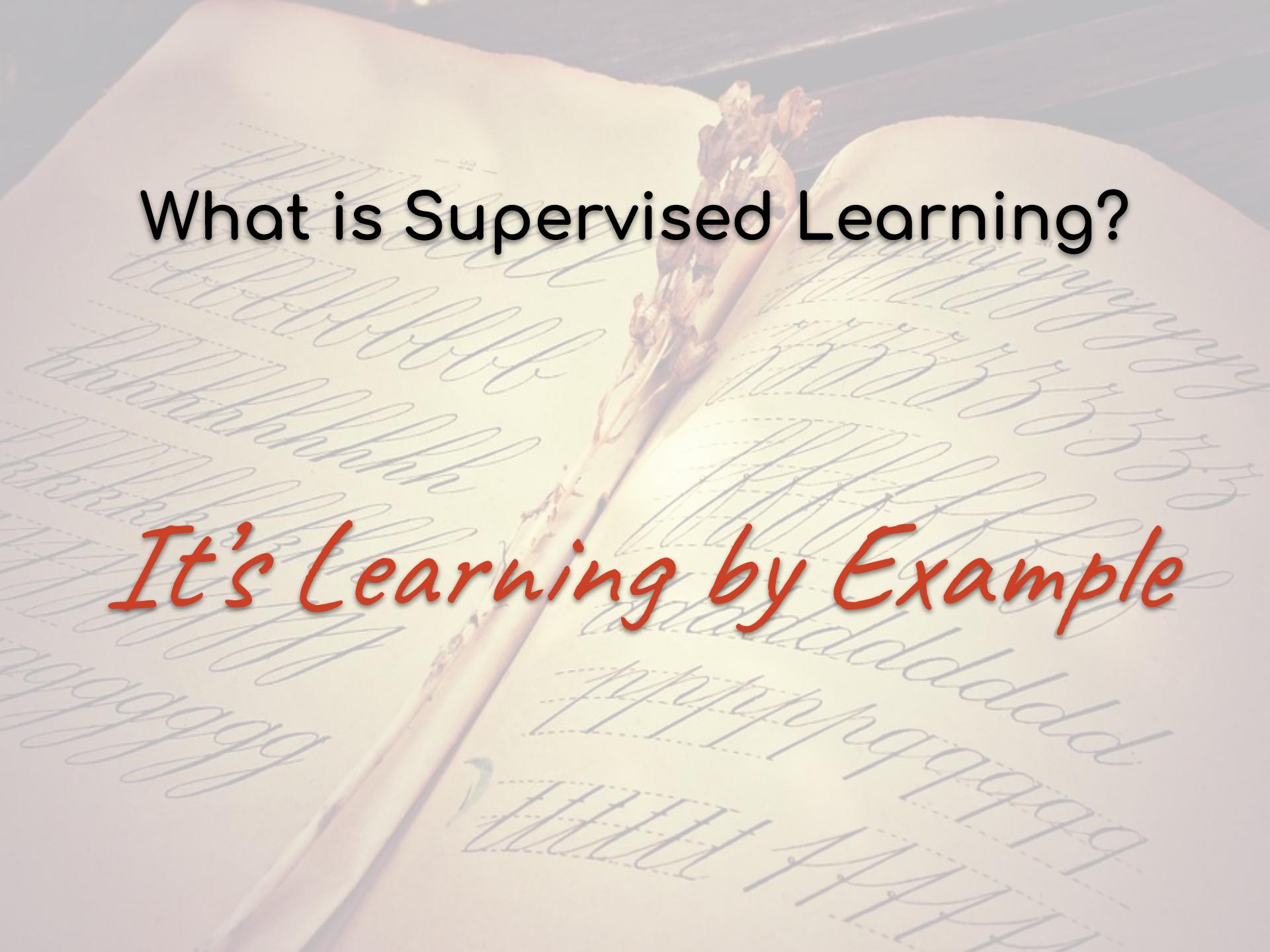
<https://github.com/pappzoltan>

The screenshot shows the GitHub repository page for 'pappzoltan / machine-learning-course'. The repository was forked from 'szilard/teach-ML-CEU-master-bizanalytics'. The main navigation bar includes links for Code, Pull requests, Actions, Projects, Wiki, Security, Insights, and Settings. Below the navigation bar, there are buttons for master, 4 branches, 2 tags, Go to file, Add file, and a green button for Code. The repository has 4 branches and 2 tags.

- Lectures
  - a. Slides
  - b. Code samples
  - c. Data
- Labs
  - a. Code samples
  - b. Data



*Source: Palmer, Shelly. *Data Science for the C-Suite*. New York: Digital Living Press, 2015. Print.*



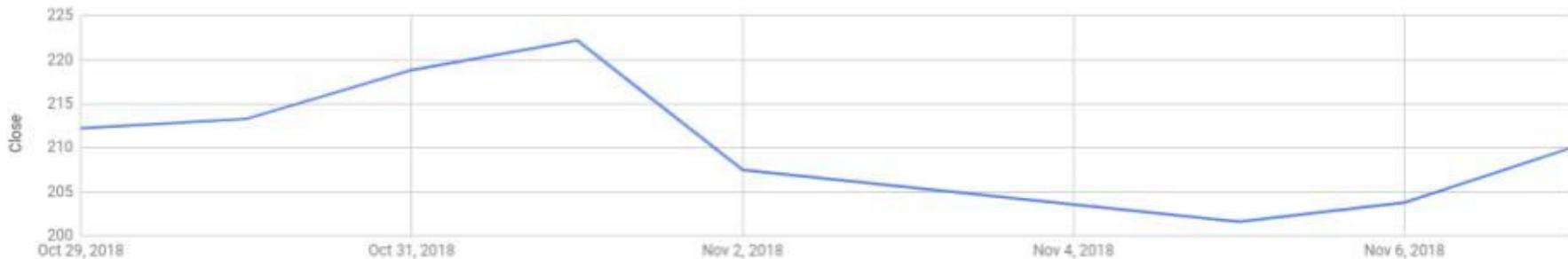
# What is Supervised Learning?

**It's Learning by Example**



# Financial panel data from Apple

IBM Close Price



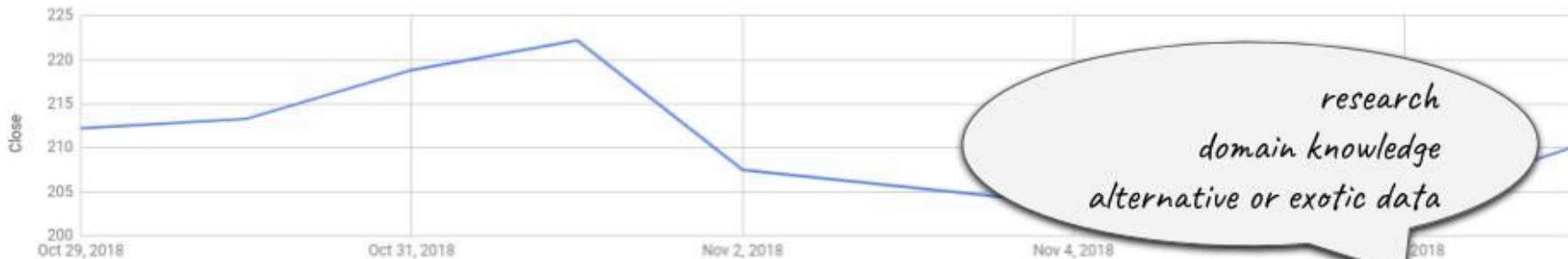
| Date         | Open   | High   | Low    | Close  | Volume     |
|--------------|--------|--------|--------|--------|------------|
| Oct 29, 2018 | 219.19 | 219.69 | 206.09 | 212.24 | 45,935,500 |
| Oct 30, 2018 | 211.15 | 215.18 | 209.27 | 213.3  | 36,660,000 |
| Oct 31, 2018 | 216.88 | 220.45 | 216.62 | 218.86 | 38,358,900 |
| Nov 01, 2018 | 219.05 | 222.36 | 216.81 | 222.22 | 58,323,200 |
| Nov 02, 2018 | 209.55 | 213.65 | 205.43 | 207.48 | 91,328,700 |
| Nov 05, 2018 | 204.3  | 204.39 | 198.17 | 201.59 | 66,163,700 |
| Nov 06, 2018 | 201.92 | 204.72 | 201.69 | 203.77 | 31,882,900 |
| Nov 07, 2018 | 205.97 | 210.06 | 204.13 | 209.95 | 33,424,400 |

| Date         | Close  | Previous Close | Previous Volume | Features |
|--------------|--------|----------------|-----------------|----------|
| Oct 29, 2018 | 212.24 | NaN            | NaN             |          |
| Oct 30, 2018 | 213.91 | 212.24         | 45,935,500      |          |
| Oct 31, 2018 | 218.86 | 213.3          | 36,660,000      |          |
| Nov 01, 2018 | 222.22 | 218.86         | 38,358,900      |          |
| Nov 02, 2018 | 207.48 | 222.22         | 58,323,200      |          |
| Nov 05, 2018 | 201.59 | 207.48         | 91,328,700      |          |
| Nov 06, 2018 | 203.77 | 201.59         | 66,163,700      |          |
| Nov 07, 2018 | 209.95 | 203.77         | 31,882,900      |          |



# Financial panel data from Apple

IBM Close Price



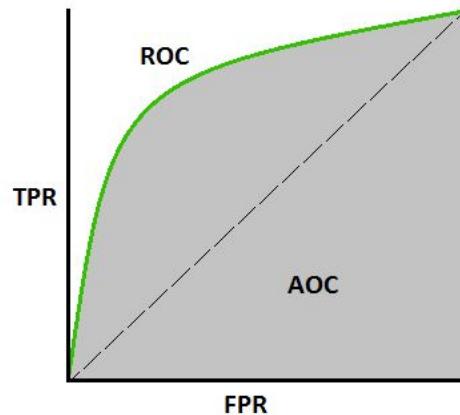
| Date         | Open   | High   | Low    | Close  | Volume     |
|--------------|--------|--------|--------|--------|------------|
| Oct 29, 2018 | 219.19 | 219.69 | 206.09 | 212.24 | 45,935,500 |
| Oct 30, 2018 | 211.15 | 215.18 | 209.27 | 213.3  | 36,660,000 |
| Oct 31, 2018 | 216.88 | 220.45 | 216.62 | 218.86 | 38,358,900 |
| Nov 01, 2018 | 219.05 | 222.36 | 216.81 | 222.22 | 58,323,200 |
| Nov 02, 2018 | 209.55 | 213.65 | 205.43 | 207.48 | 91,328,700 |
| Nov 05, 2018 | 204.3  | 204.39 | 198.17 | 201.59 | 66,163,700 |
| Nov 06, 2018 | 201.92 | 204.72 | 201.69 | 203.77 | 31,882,900 |
| Nov 07, 2018 | 205.97 | 210.06 | 204.13 | 209.95 | 33,424,400 |

| Date         | Close  | Previous Close | Previous Volume | Features |
|--------------|--------|----------------|-----------------|----------|
| Oct 29, 2018 | 212.24 | NaN            | 45,935,500      | ...      |
| Oct 30, 2018 | 213.3  | 212.24         | 36,660,000      | ...      |
| Oct 31, 2018 | 218.86 | 213.3          | 38,358,900      | ...      |
| Nov 01, 2018 | 222.22 | 218.86         | 58,323,200      | ...      |
| Nov 02, 2018 | 207.48 | 222.22         | 91,328,700      | ...      |
| Nov 05, 2018 | 201.59 | 207.48         | 66,163,700      | ...      |
| Nov 06, 2018 | 203.77 | 201.59         | 31,882,900      | ...      |
| Nov 07, 2018 | 209.95 | 203.77         | 33,424,400      | ...      |

# Classification? Regression? Performance

Classification

Is it going up or down?



Regression

How much up or down?

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

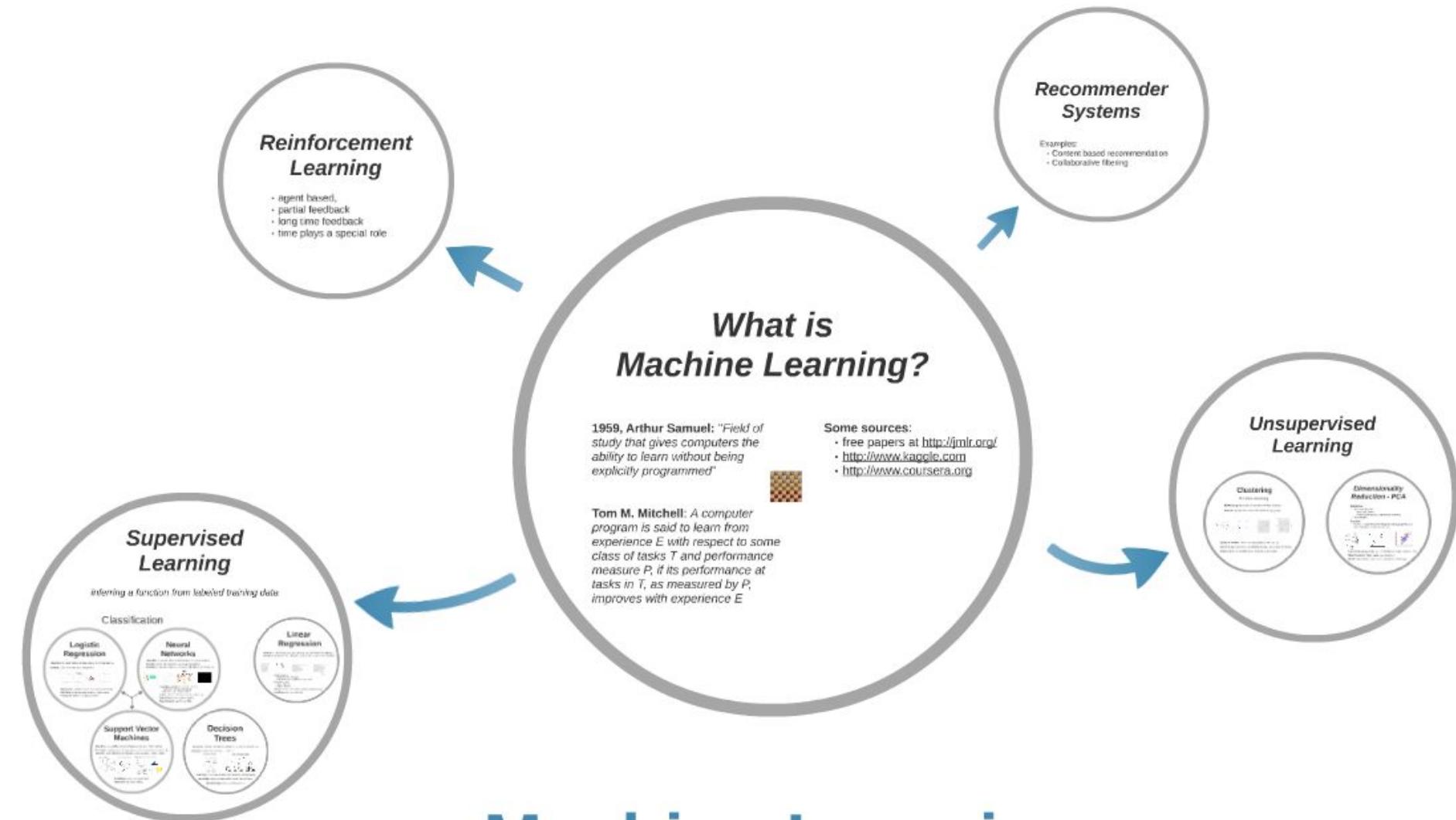
Classification Examples:

- Is mail a spam or not?
- Is it a cat on the image?
- Is it the same person on two pictures?

Regression Examples:

- Next day return for Apple?
- Temperature tomorrow?
- Expected bike sharing demand at a station tomorrow?

Let's see some **supervised** and  
**unsupervised** learning algorithms



# Machine Learning

# Keywords to summarize topics

Supervised Learning

Data: X (n observations, p features), y (labels)

Regression, classification

Train/learn/fit f from data (model)

Score: for new x, get f(x)

Algos: LR, k-NN, DT, RF, GBM, NN/DL, SVM, NB...

Goal: max acc/min err new data

Metrics: MSE, AUC (ROC)

Bad: measure on train set. Need: test/CV

Hyperparameters, model capacity, overfitting

Regularization

Model selection

Hyperparameter search (grid, random)

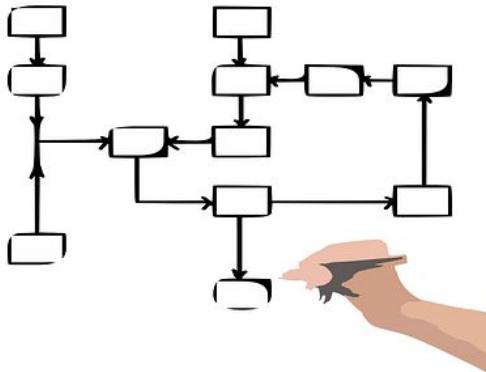
Algos: LR, k-NN, DT, RF, GBM, NN/DL, SVM, NB...

ML Tools: R packages, glmnet, h2o, xgboost, lightgbm, keras, Tensorflow...

Unsupervised Learning

Best practices

# Topics to cover



Data Mining Process



Exploratory Data Analysis



Data Visualization

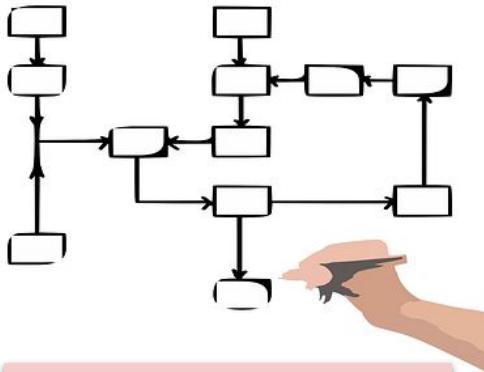


Data Science Tools



Reproducible Research

# Topics to cover



Data Mining Process



Exploratory Data Analysis



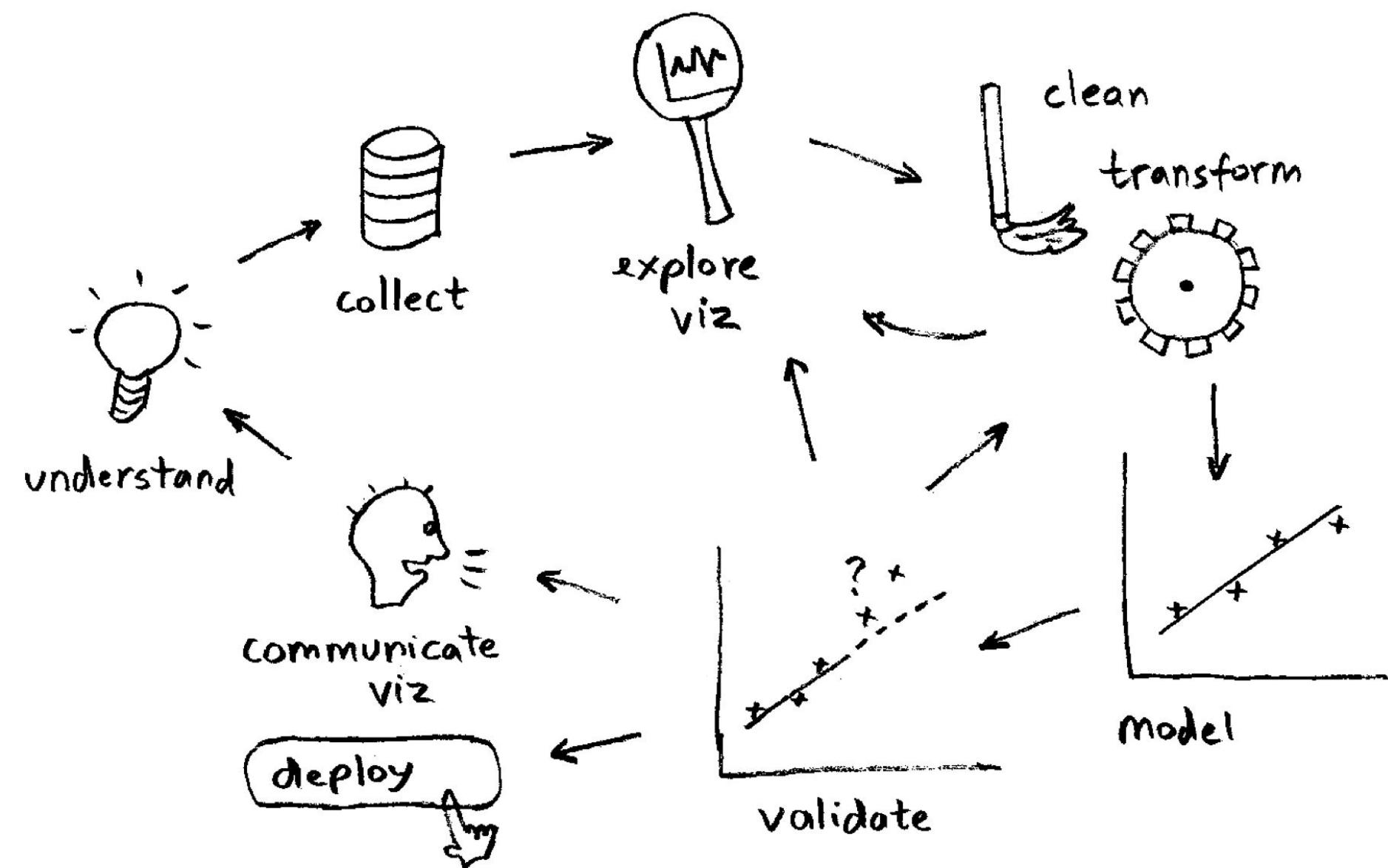
Data Visualization



Data Science Tools

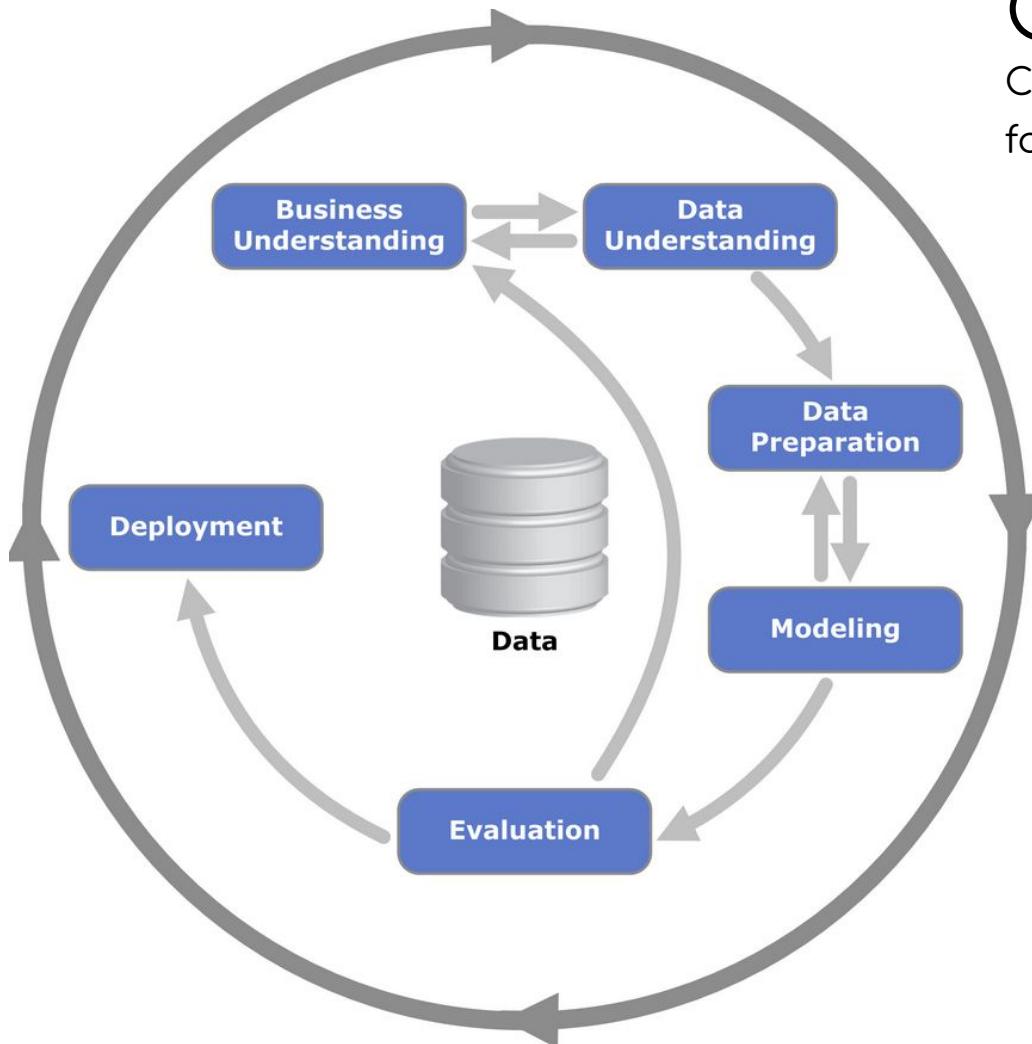


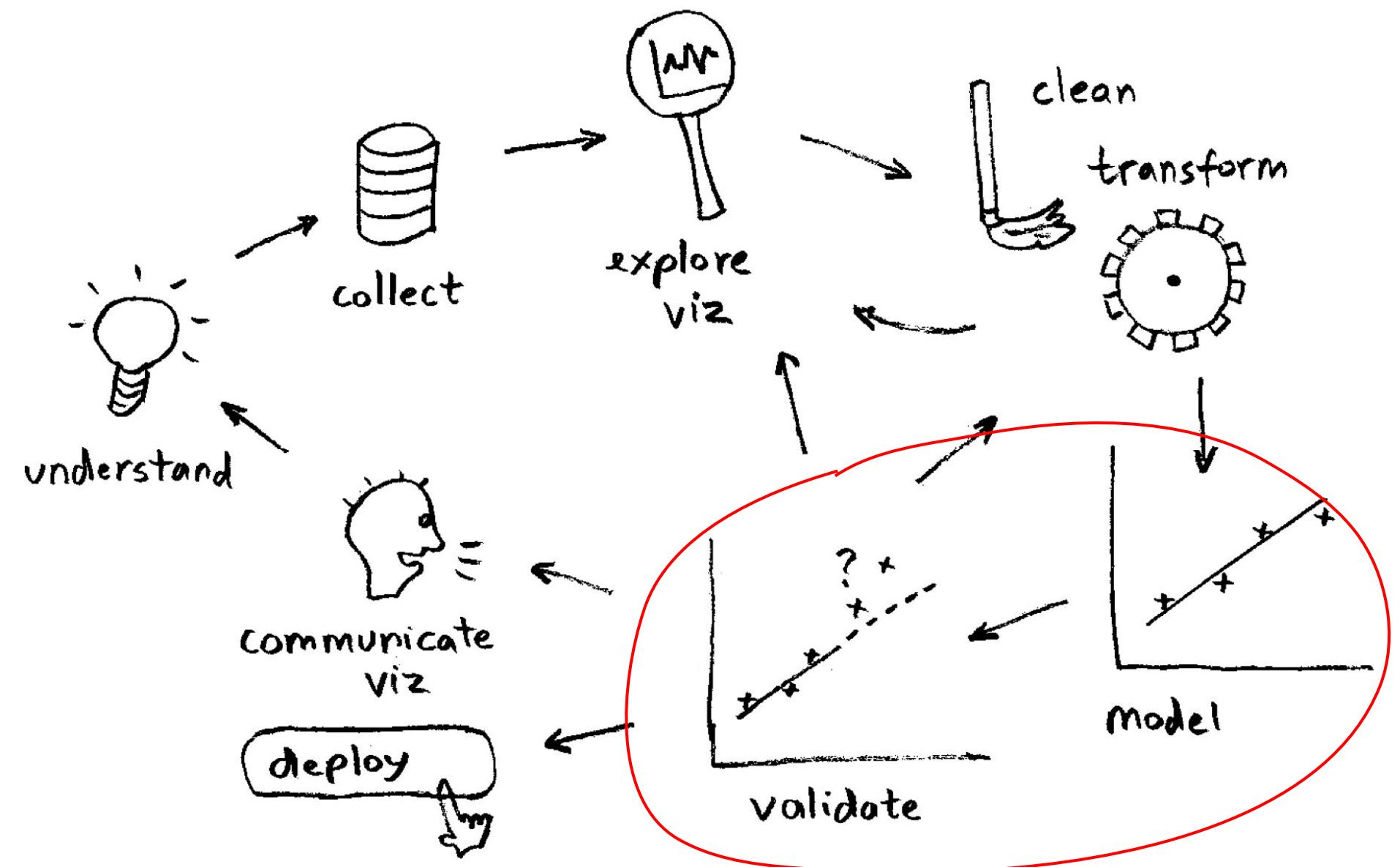
Reproducible Research

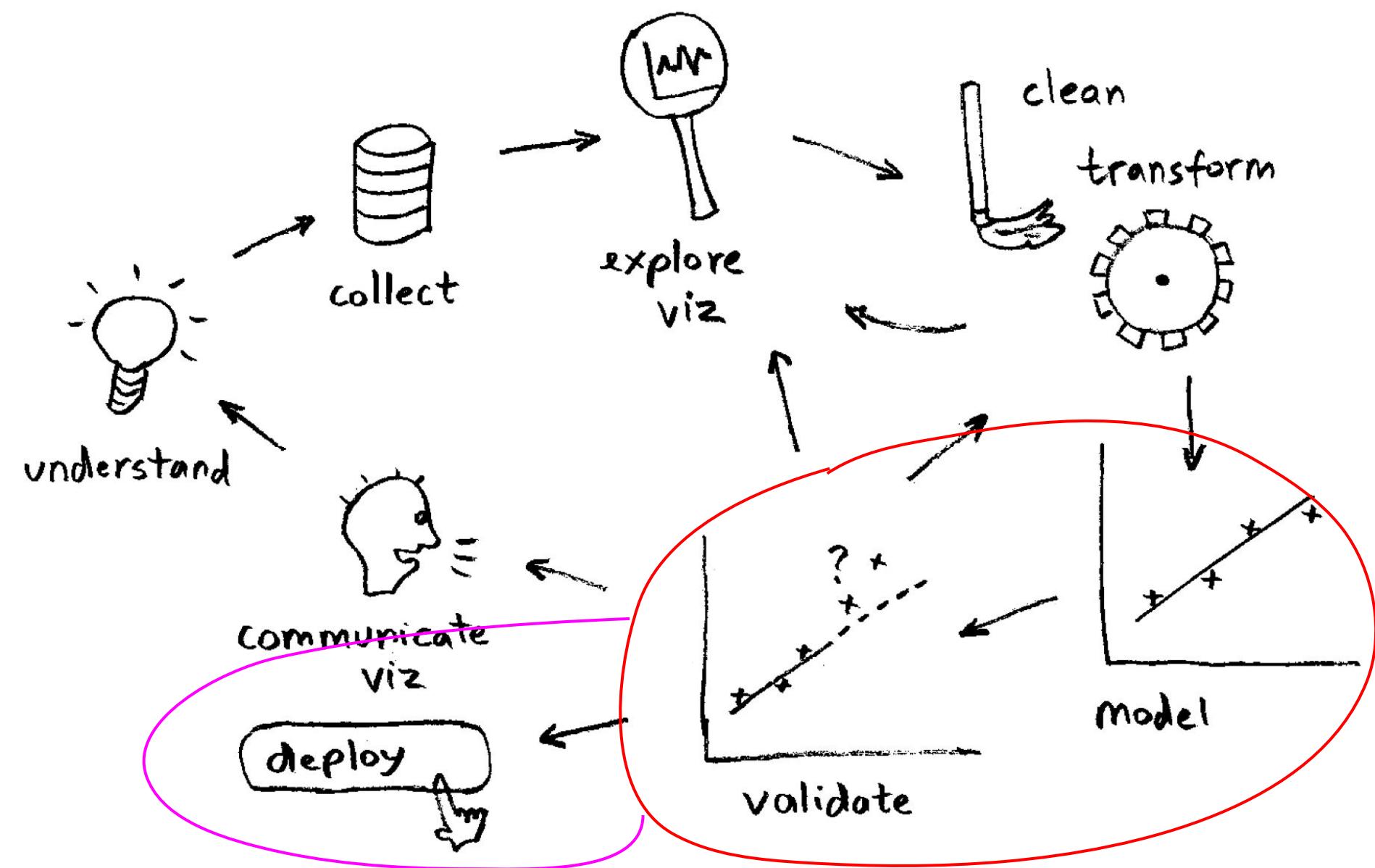


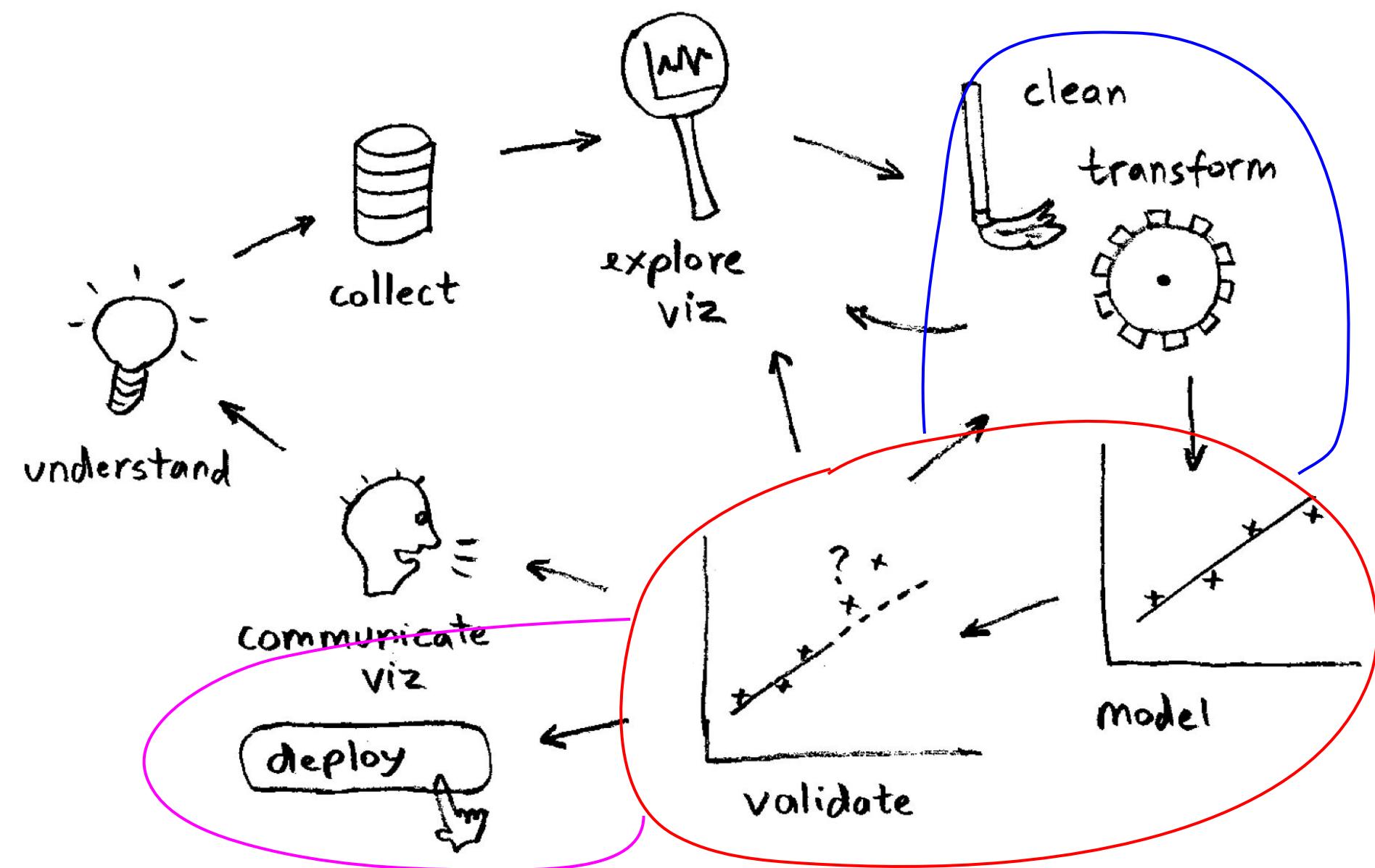
# CRISP-DM, 1999

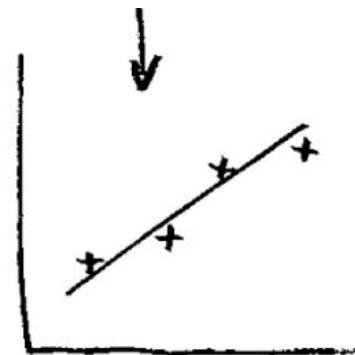
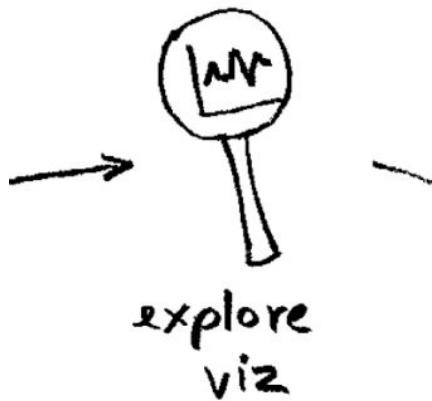
Cross-industry standard process  
for data mining











# David Donoho: 50 years of Data Science

GDS1: Exploratory data analysis and data manipulation/preparation and cleaning

GDS2: Systems/databases

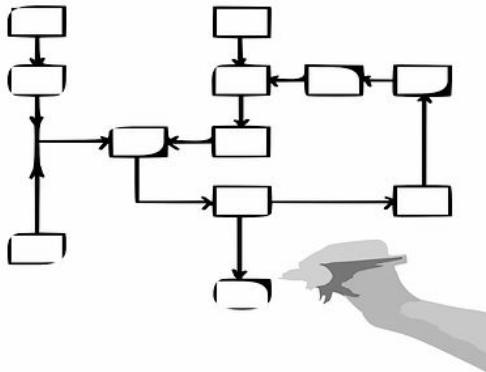
GDS3: Programming: R

GDS4: Data visualization

GDS5: Modeling (statistics/machine learning)

GDS6: Science of Data Science

# Topics to cover



Data Mining Process



Exploratory Data Analysis



Data Visualization



Data Science Tools



Reproducible Research

John W. Tukey

# EXPLORATORY DATA ANALYSIS





**Big Data Borat**

@BigDataBorat

Follow

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.



Reply



Retweet



Favorite



More

# Data Scientist: The **Sexiest** Job of the 21st Century

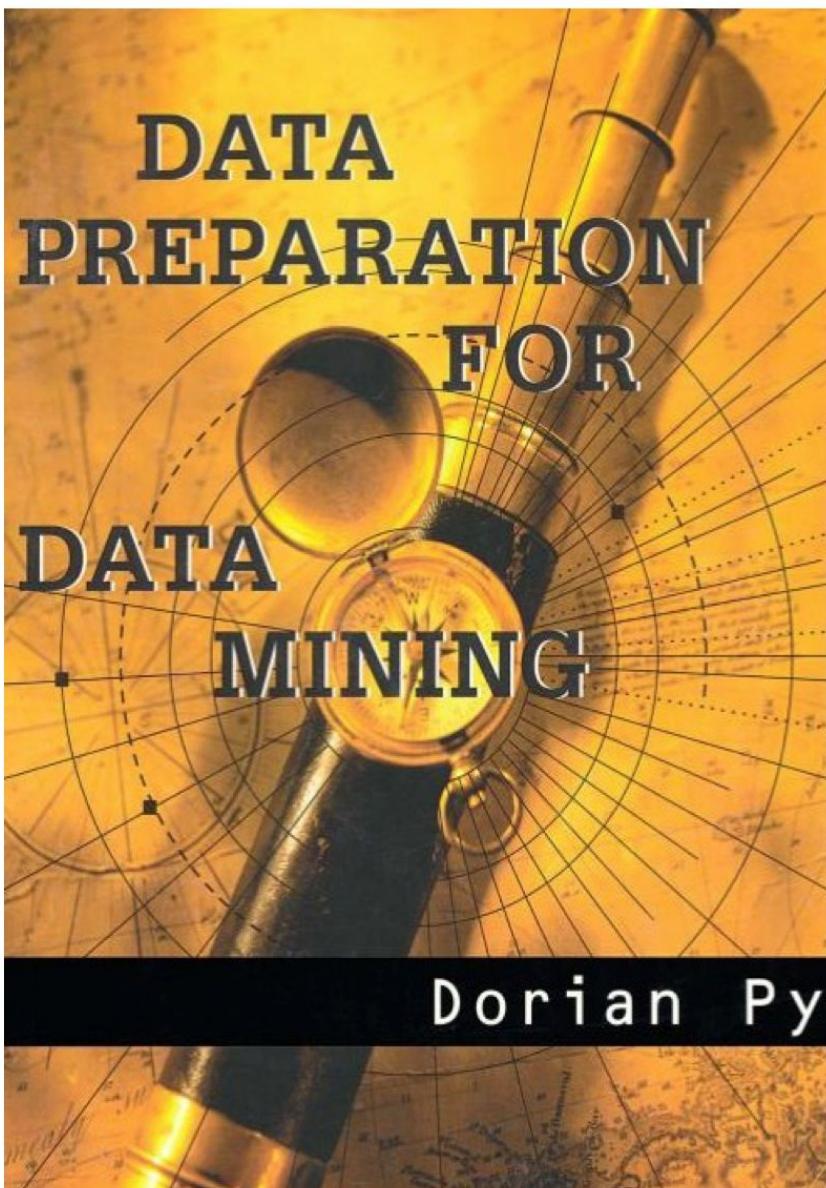
Harvard  
Business  
Review

# Data Scientist: The **Sexiest** Job of the 21st Century

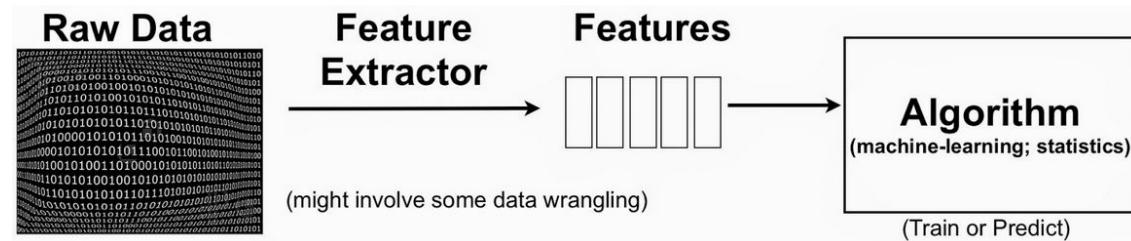
Harvard  
Business  
Review

*For Big-Data Scientists, ‘Janitor’ Work Is Key Hurdle to Insights*

The New York Times



“the data preparation process prepares  
both the data and the modeler”



# Bike Sharing Demand

## Data Fields

**datetime** - hourly date + timestamp

**season** - 1 = Q1, 2 = Q2, 3 = Q3, 4 = Q4

**holiday** - whether the day is considered a holiday

**workingday** - whether the day is neither a weekend nor holiday

**weather**

1. Clear, Few clouds, Partly cloudy, Partly cloudy
2. Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
3. Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds
4. Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog

**temp** - temperature in Celsius

**atemp** - "feels like" temperature in Celsius

**humidity** - relative humidity

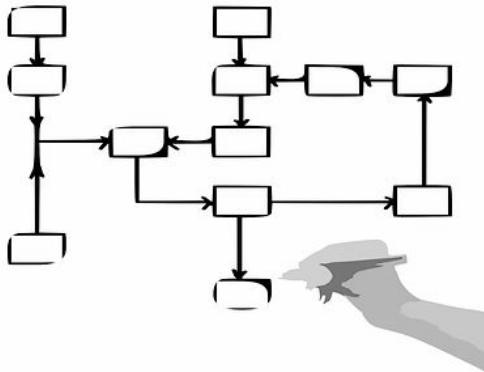
**windspeed** - wind speed

**casual** - number of non-registered user rentals initiated

**registered** - number of registered user rentals initiated

**count** - number of total rentals

# Topics to cover



Data Mining Process



Exploratory Data Analysis



Data Visualization



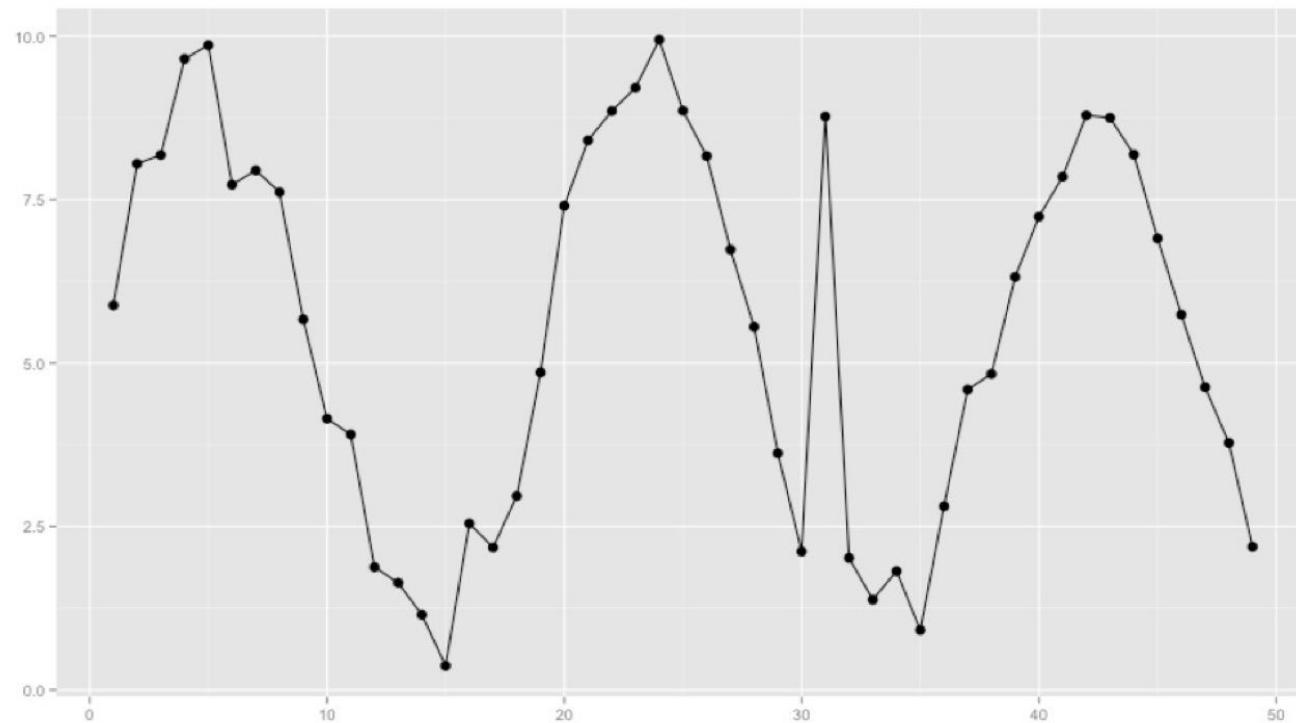
Data Science Tools



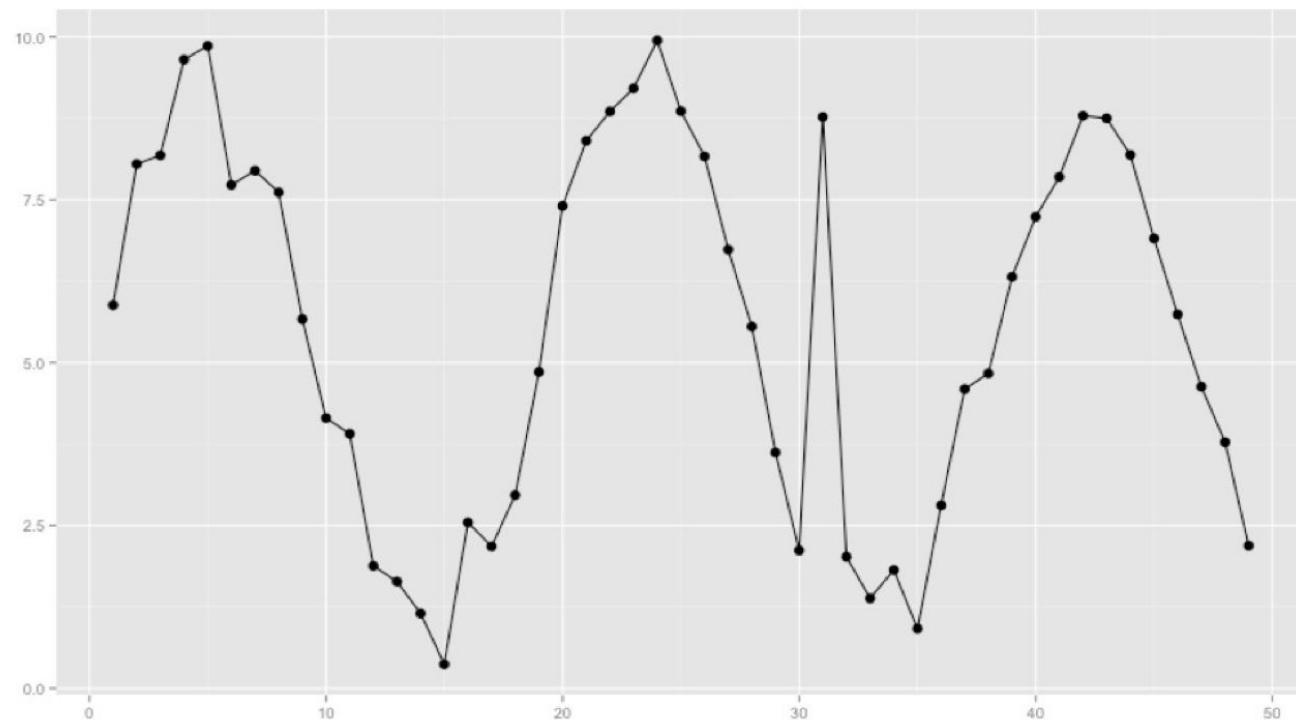
Reproducible Research

|           |           |           |           |           |           |           |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 5.8839338 | 8.0500895 | 8.1838378 | 9.6537864 | 9.8625664 | 7.7283027 | 7.9485545 |
| 7.6139286 | 5.6673501 | 4.1509576 | 3.9085585 | 1.8794583 | 1.6390833 | 1.1494471 |
| 0.3701523 | 2.5463324 | 2.1793825 | 2.9664571 | 4.8563495 | 7.4056122 | 8.4070250 |
| 8.8555943 | 9.2110448 | 9.9459725 | 8.8605880 | 8.1672658 | 6.7366060 | 5.5535530 |
| 3.6201724 | 2.1181429 | 8.7715833 | 2.0190151 | 1.3814497 | 1.8169363 | 0.9166560 |
| 2.8093003 | 4.5931840 | 4.8333278 | 6.3170302 | 7.2390577 | 7.8509665 | 8.7915221 |
| 8.7507326 | 8.1899304 | 6.9060409 | 5.7413247 | 4.6312077 | 3.7803116 | 2.1903559 |

|           |           |           |           |           |           |           |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 5.8839338 | 8.0500895 | 8.1838378 | 9.6537864 | 9.8625664 | 7.7283027 | 7.9485545 |
| 7.6139286 | 5.6673501 | 4.1509576 | 3.9085585 | 1.8794583 | 1.6390833 | 1.1494471 |
| 0.3701523 | 2.5463324 | 2.1793825 | 2.9664571 | 4.8563495 | 7.4056122 | 8.4070250 |
| 8.8555943 | 9.2110448 | 9.9459725 | 8.8605880 | 8.1672658 | 6.7366060 | 5.5535530 |
| 3.6201724 | 2.1181429 | 8.7715833 | 2.0190151 | 1.3814497 | 1.8169363 | 0.9166560 |
| 2.8093003 | 4.5931840 | 4.8333278 | 6.3170302 | 7.2390577 | 7.8509665 | 8.7915221 |
| 8.7507326 | 8.1899304 | 6.9060409 | 5.7413247 | 4.6312077 | 3.7803116 | 2.1903559 |



|           |           |           |           |           |           |           |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 5.8839338 | 8.0500895 | 8.1838378 | 9.6537864 | 9.8625664 | 7.7283027 | 7.9485545 |
| 7.6139286 | 5.6673501 | 4.1509576 | 3.9085585 | 1.8794583 | 1.6390833 | 1.1494471 |
| 0.3701523 | 2.5463324 | 2.1793825 | 2.9664571 | 4.8563495 | 7.4056122 | 8.4070250 |
| 8.8555943 | 9.2110448 | 9.9459725 | 8.8605880 | 8.1672658 | 6.7366060 | 5.5535530 |
| 3.6201724 | 2.1181429 | 8.7715833 | 2.0190151 | 1.3814497 | 1.8169363 | 0.9166560 |
| 2.8093003 | 4.5931840 | 4.8333278 | 6.3170302 | 7.2390577 | 7.8509665 | 8.7915221 |
| 8.7507326 | 8.1899304 | 6.9060409 | 5.7413247 | 4.6312077 | 3.7803116 | 2.1903559 |



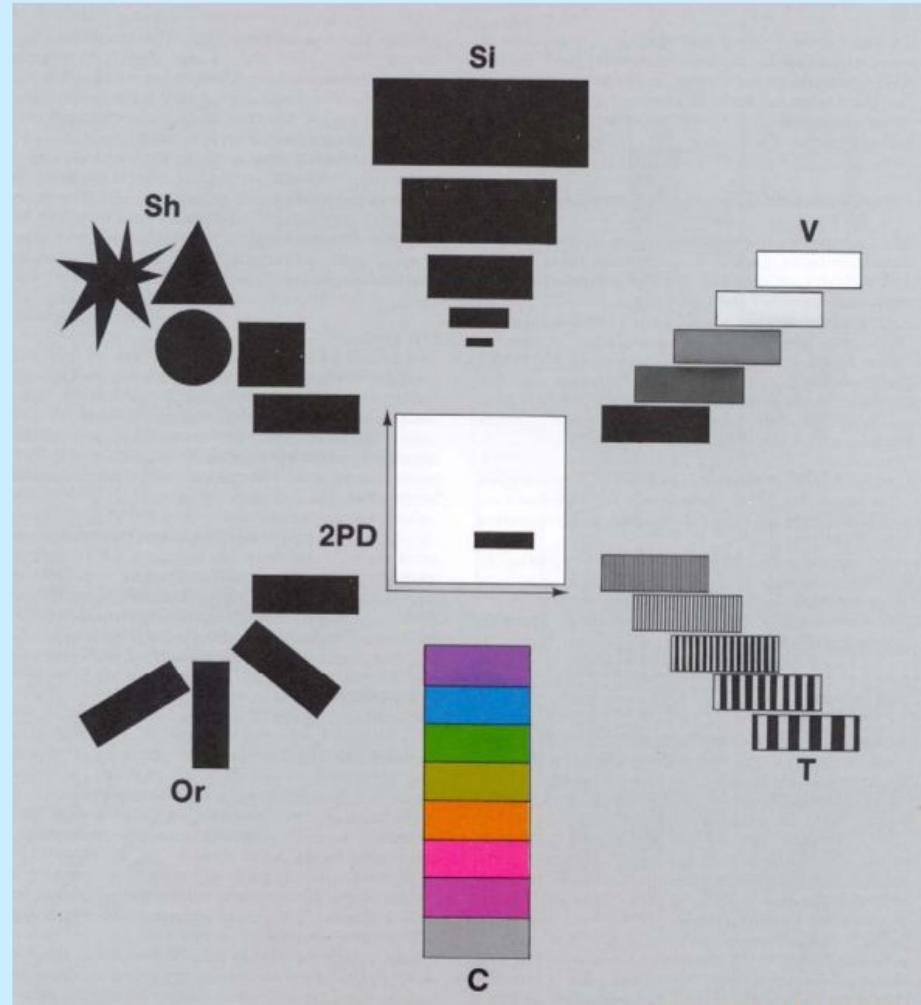
# Visual Representation from Semiology of Graphics by J. Bertin

## The Plane

- Points
  - “A point represents a location on the plane that has no theoretical length or area. This signification is independent of the size and character of the mark which renders it visible.”
  - a location
  - marks that indicate points can vary in all visual variables
- Lines
  - “A line signifies a phenomenon on the plane which has measurable length but no area. This signification is independent of the width and characteristics of the mark which renders it visible.”
  - a boundary, a route, a connection
- Areas
  - “An area signifies something on the plane that has measurable size. This signification applies to the entire area covered by the visible mark.”
  - an area can change in position but not in size, shape or orientation without making the area itself have a different meaning

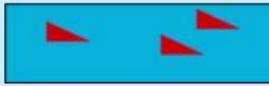
# Visual Representation from Semiology of Graphics by J. Bertin

## Visual Variables



# Visual Representation from Semiology of Graphics by J. Bertin

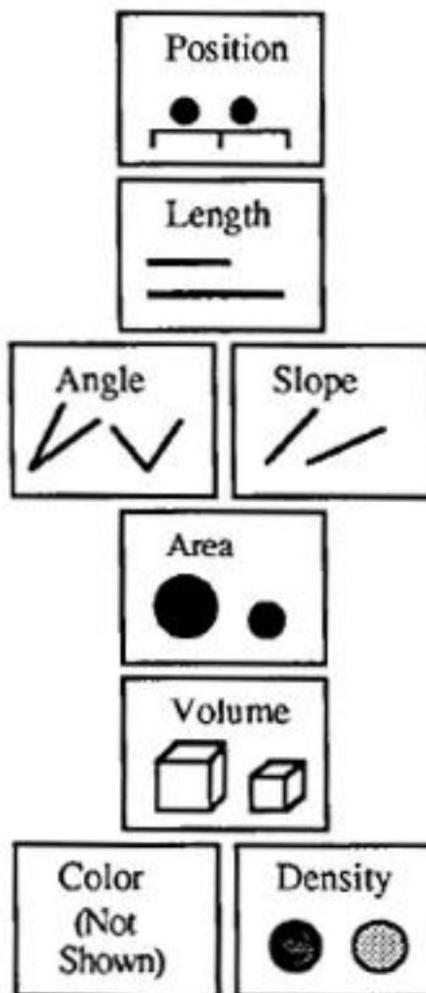
## Visual Variables

- **position**
  - changes in the x, y, (z) location
- **size**
  - change in length, area or repetition
- **shape**
  - infinite number of shapes
- **value**
  - changes from light to dark
- **orientation**
  - changes in alignment
- **colour**
  - changes in hue at a given value
- **texture**
  - variation in pattern
- **motion**

More accurate

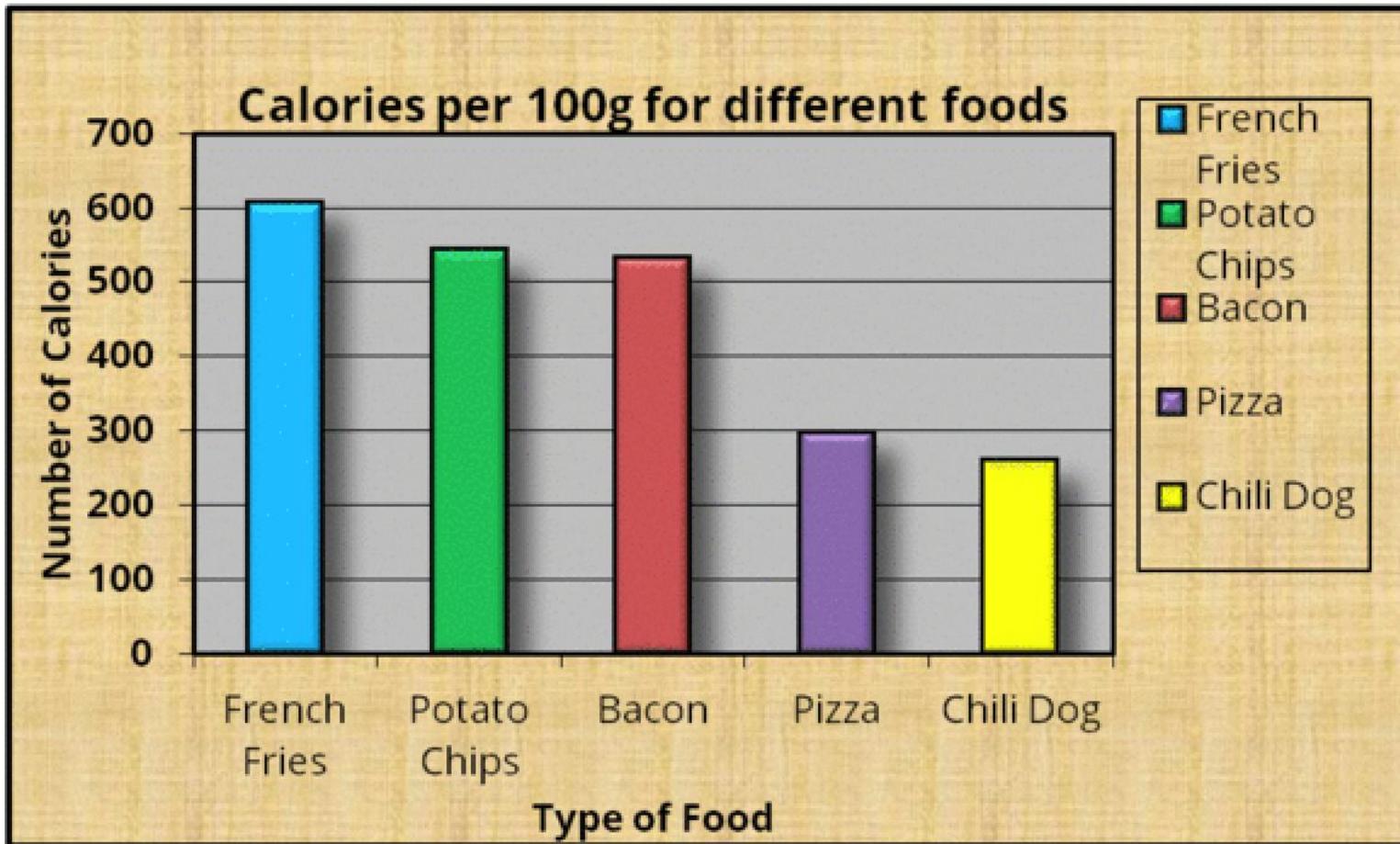


Less accurate

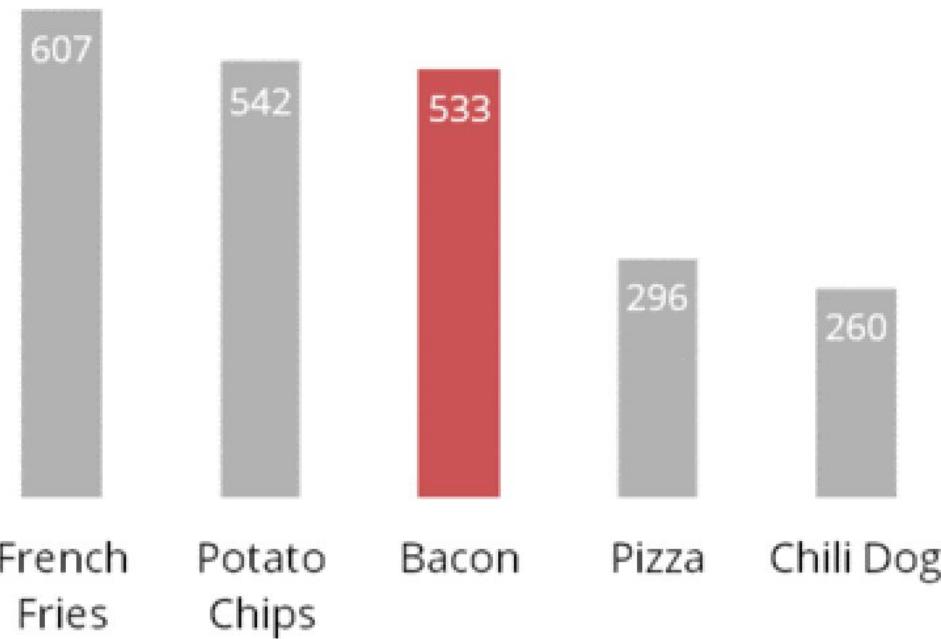


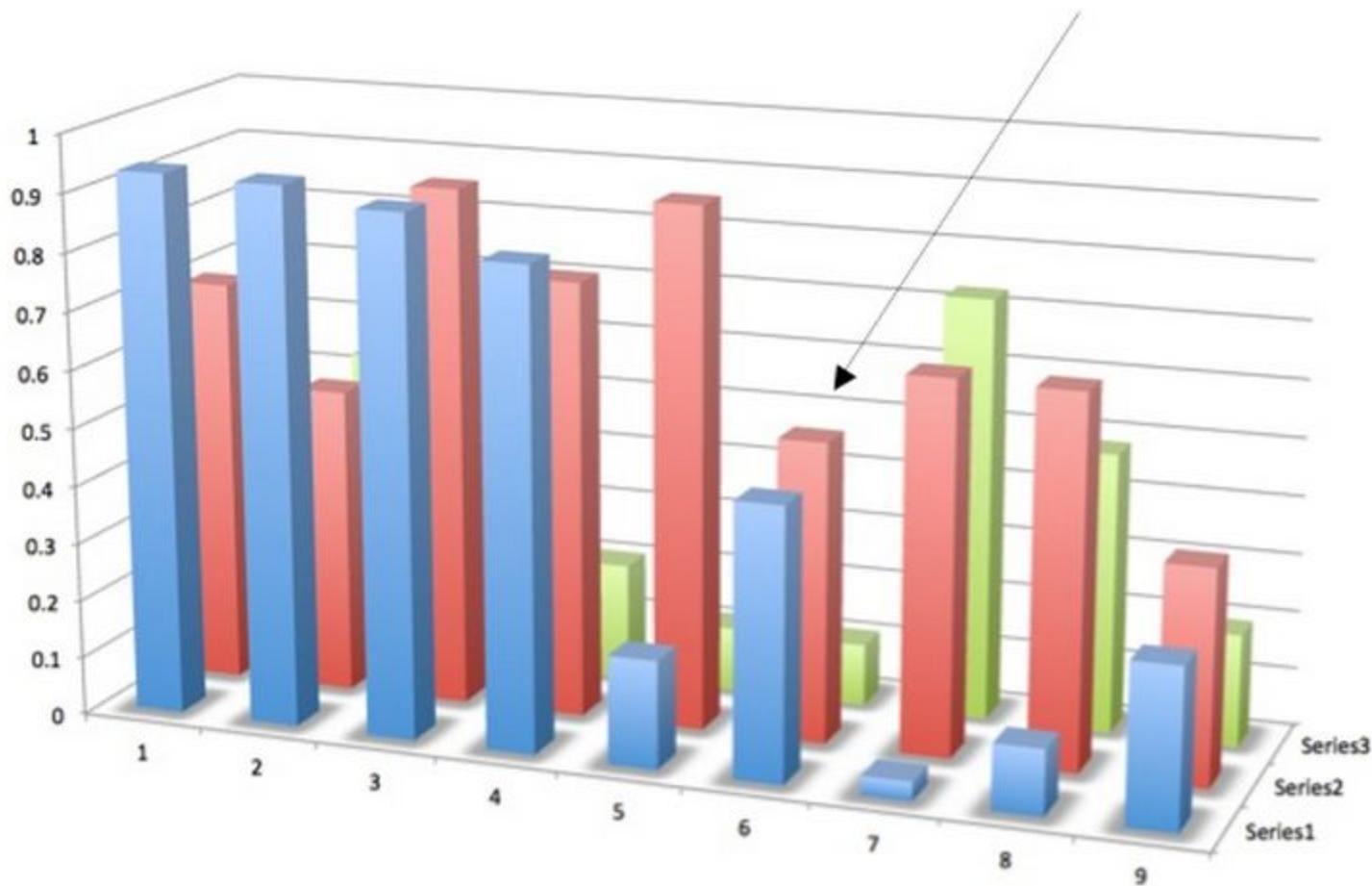
**Data-ink**

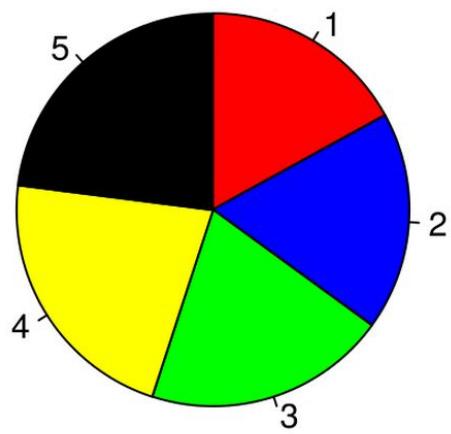
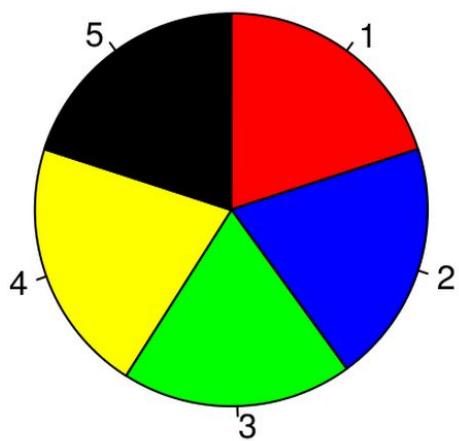
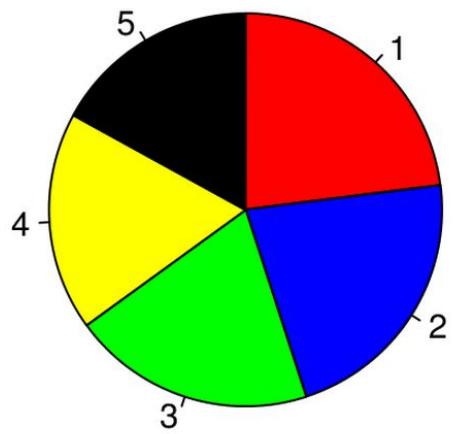
**Data-ink ratio =**  $\frac{\text{Data-ink}}{\text{Total ink used to print the graphic}}$

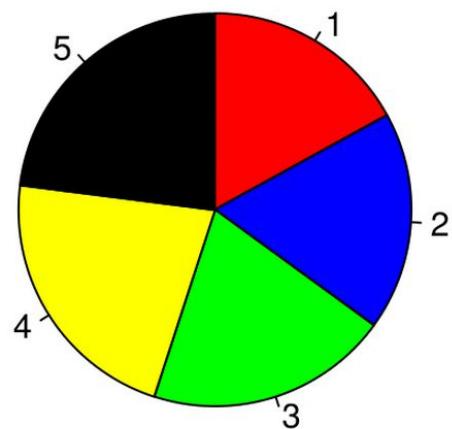
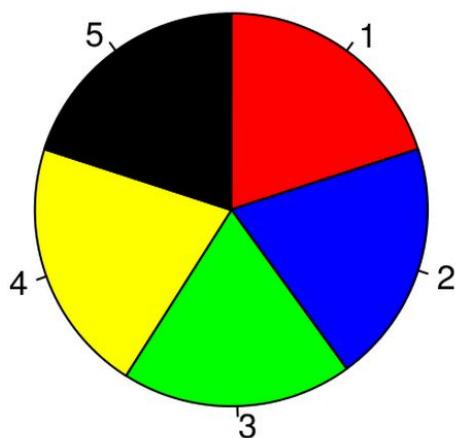
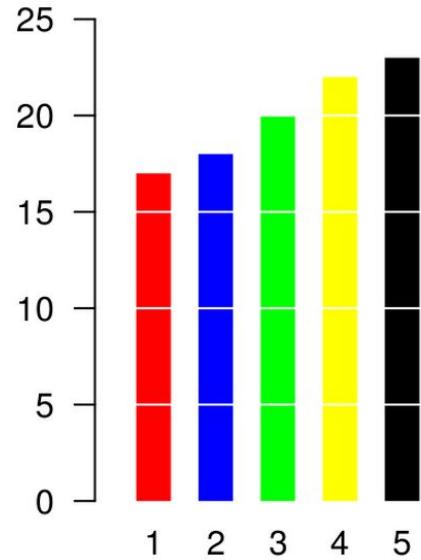
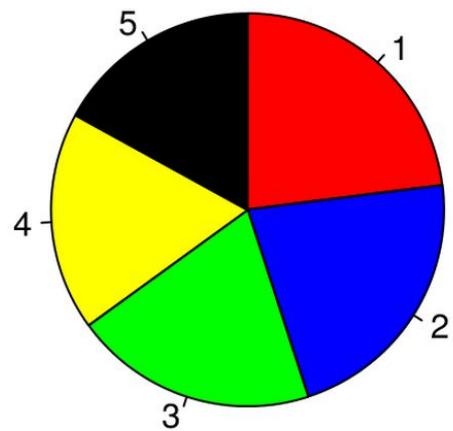


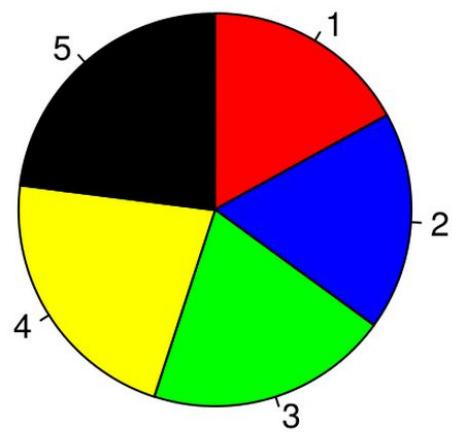
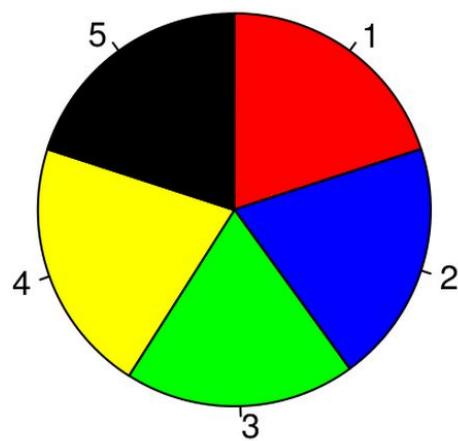
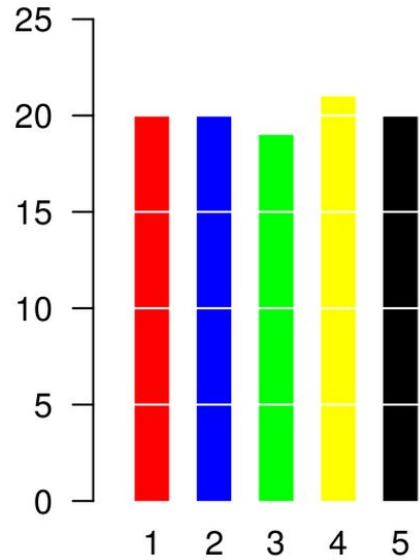
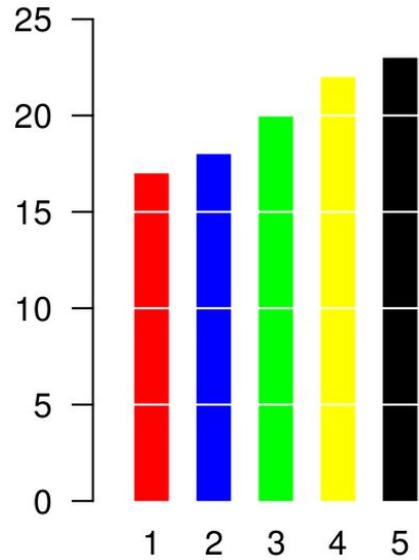
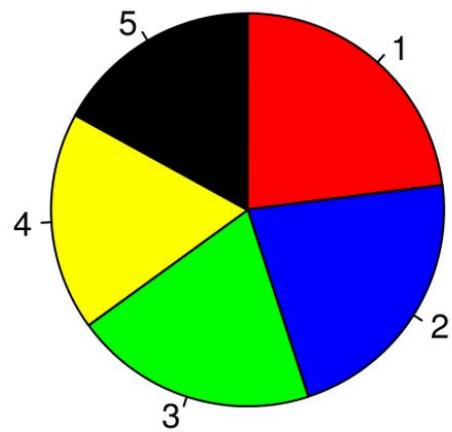
Calories per 100g

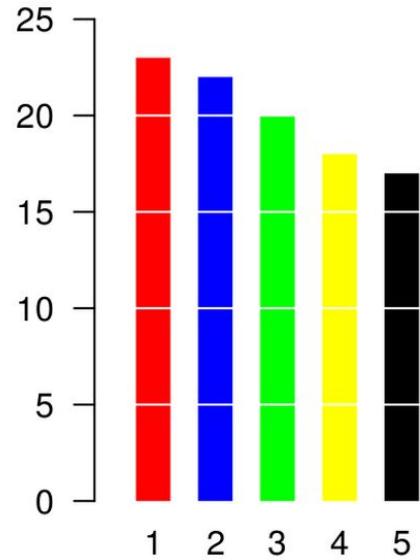
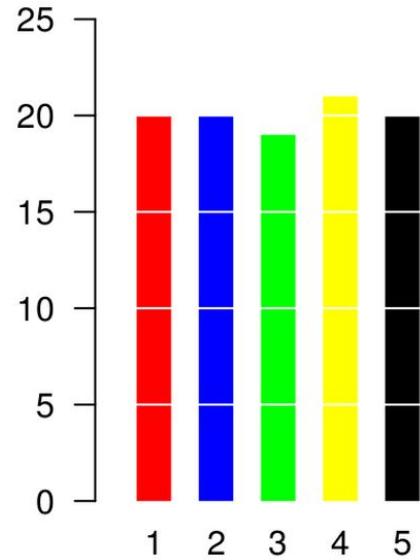
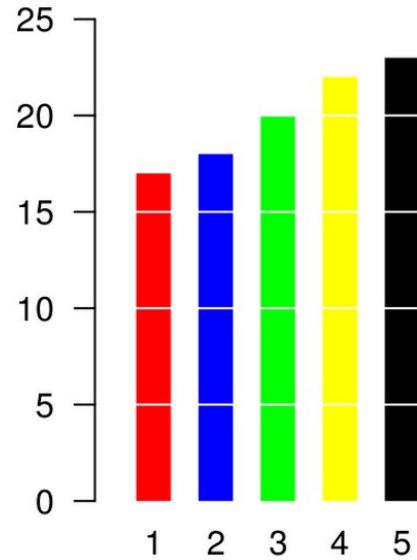
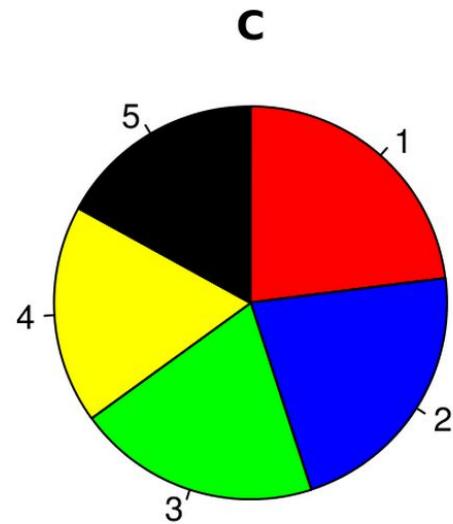
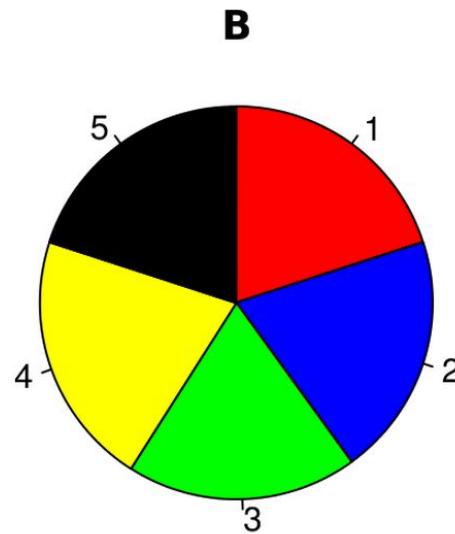
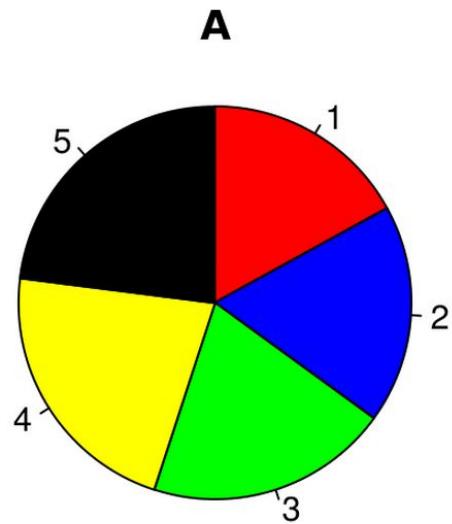




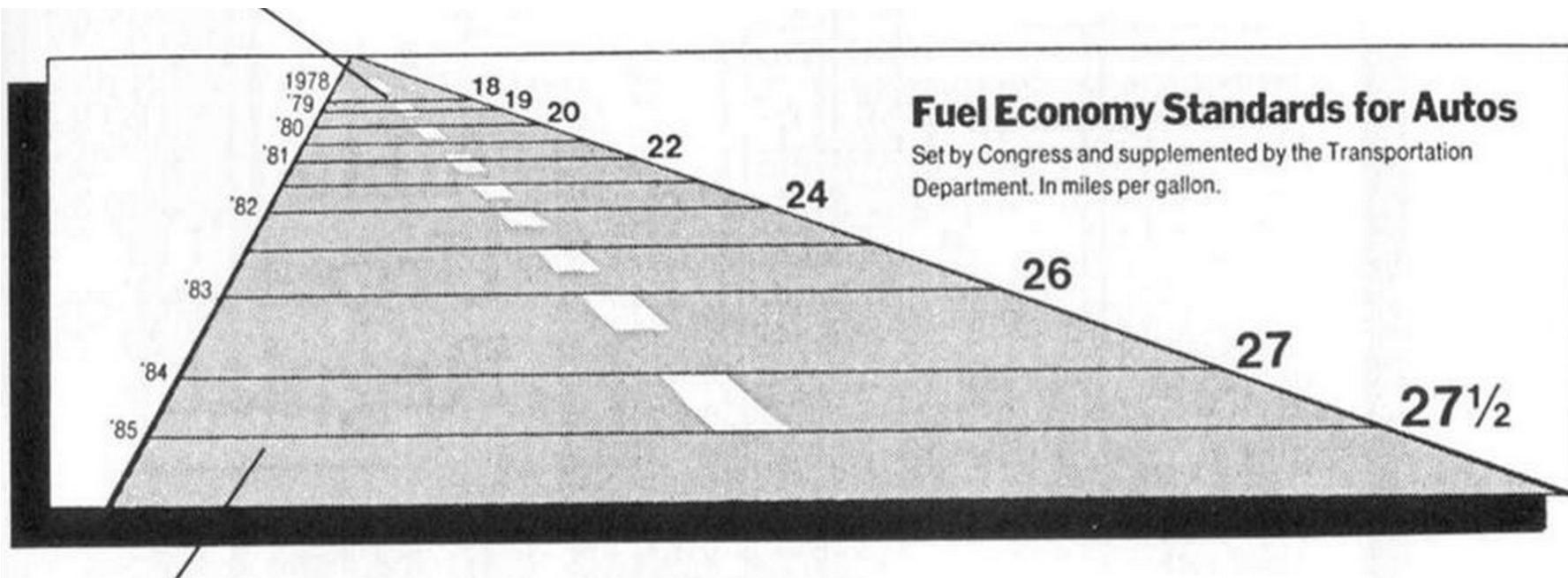
**A****B****C**

**A****B****C**

**A****B****C**



# Chart junk

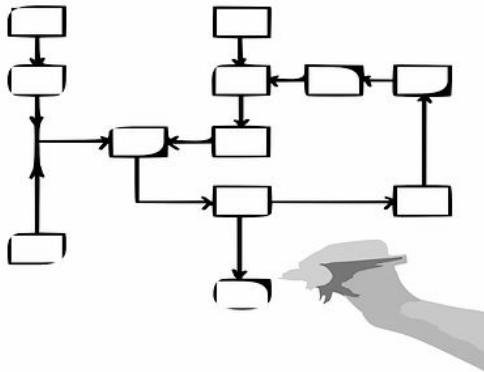


# MONSTROUS COSTS

Total House and Senate campaign expenditures,  
in millions



# Topics to cover



Data Mining Process



Exploratory Data Analysis



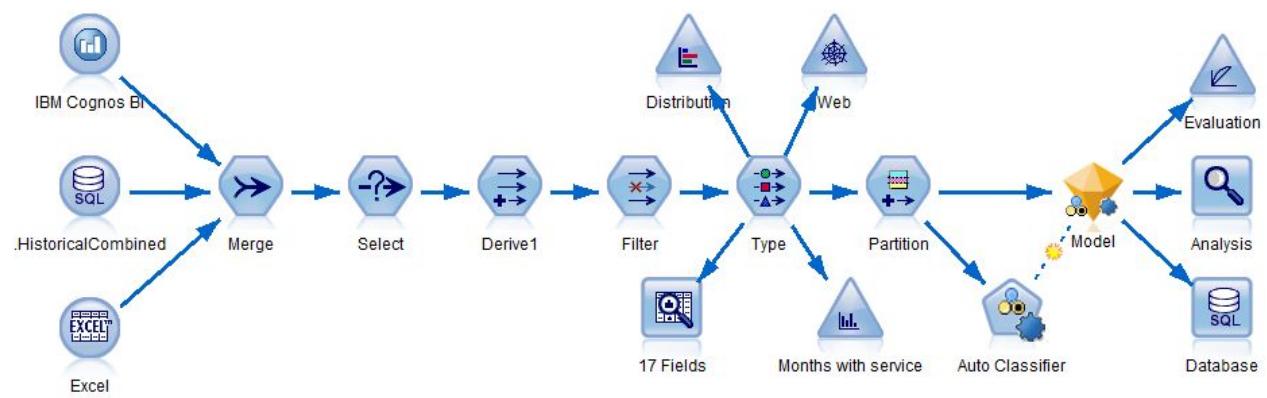
Data Visualization



Data Science Tools



Reproducible Research



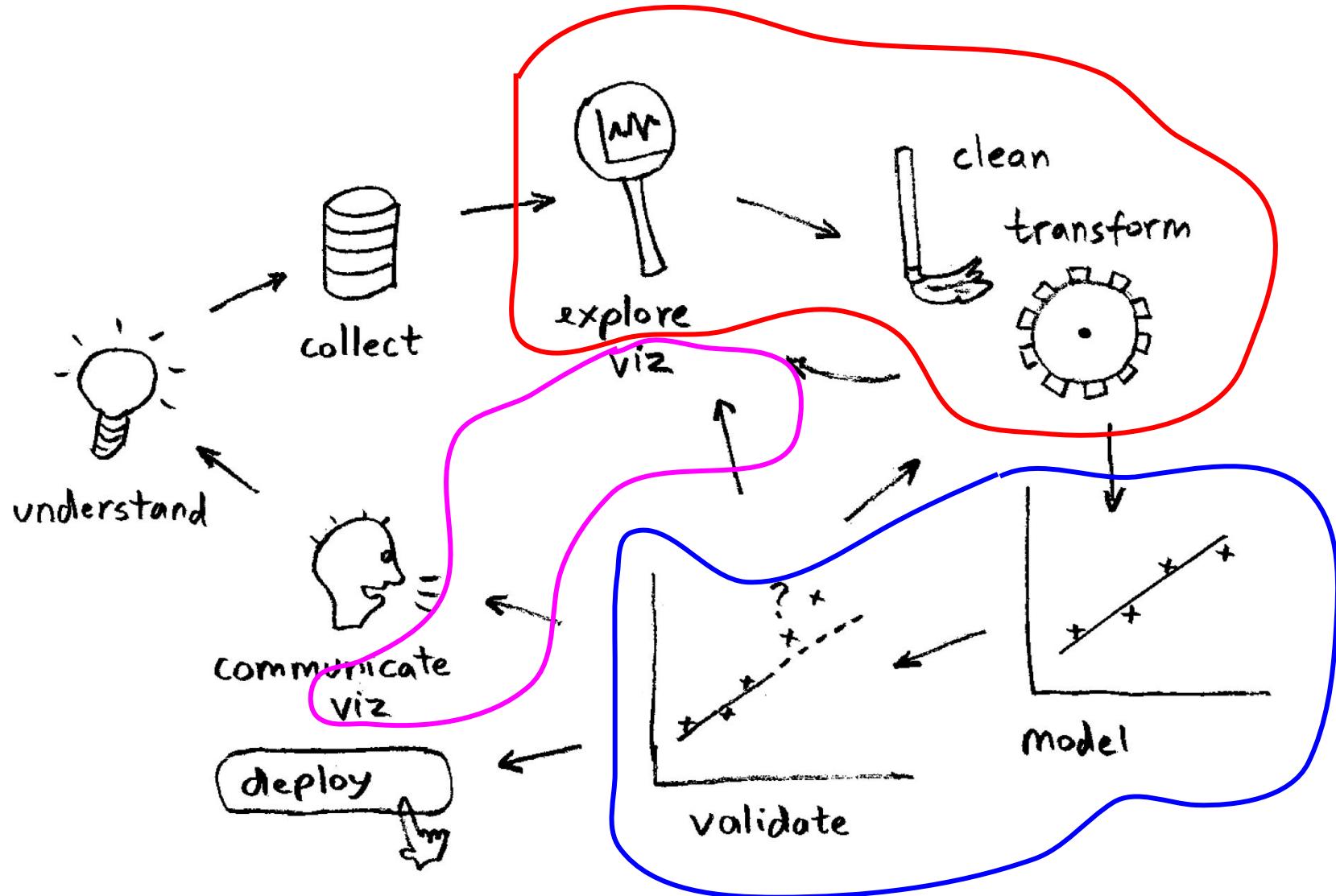


Alice Zhao  
@adashofdata

+ Follow

The Great Debate: If You Can't Code, You  
Can't Be a Data Scientist **#TeamCode** wins!  
**#Strataconf**





# Tools I Use



knitr

R Studio

`#!/bin/bash`

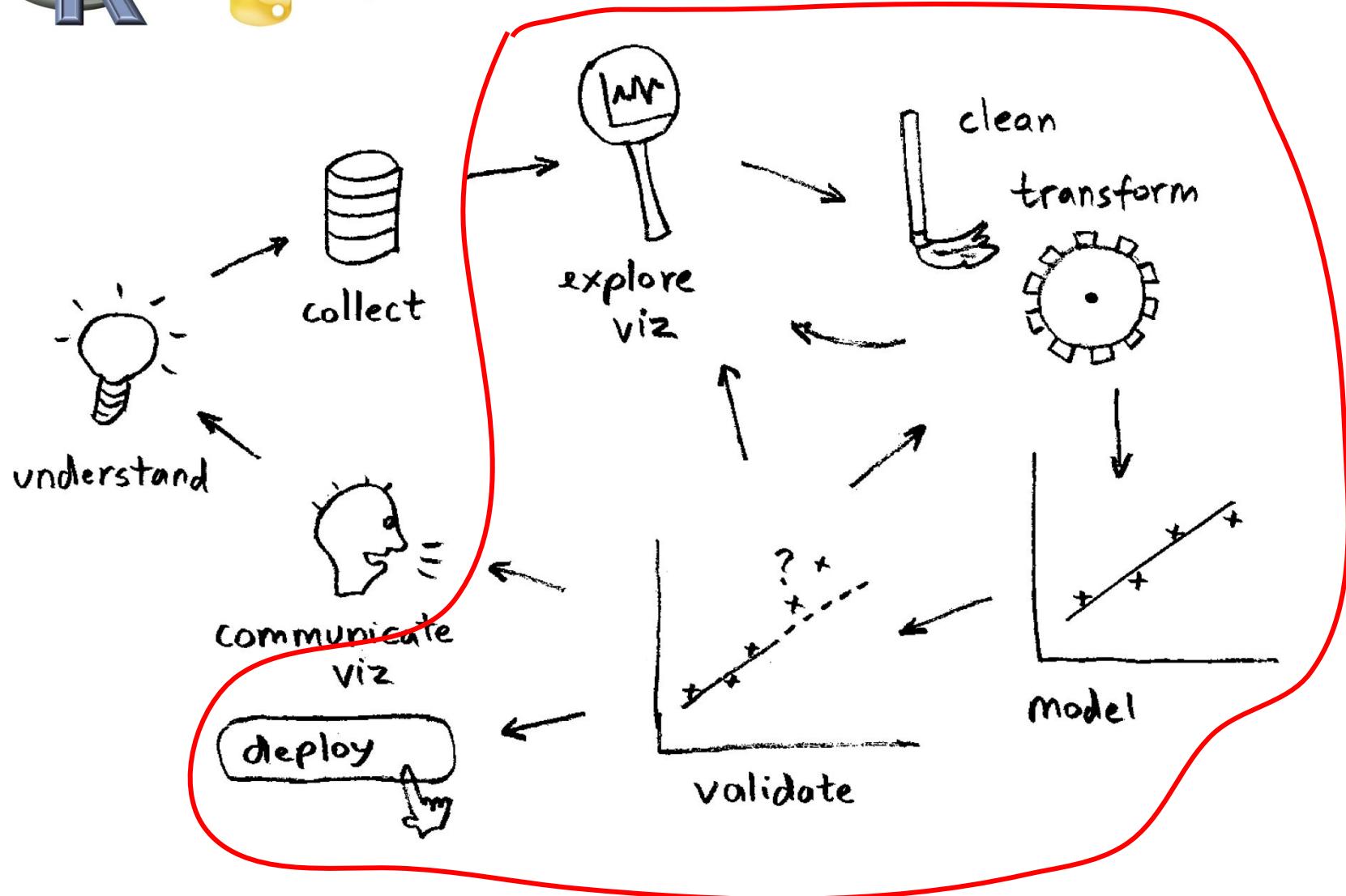


kx

# Tools Others Use (Survey)

LA Data Science/ML meetup, Apr 2014, 200 people

- **Data munging:** R 60%, Python 50%, SQL 40%, Hadoop (mostly Hive) 30%, Unix shell 20%, Excel 10% + Perl, Matlab, SAS, Impala, Pig, Shark...
- **Visualization:** R 40%, Python 30%, Tableau 10%, Javascript 10% + Matlab, Excel...
- **Machine learning/modeling:** R 30%, Python 30% + Vowpal Wabbit, Matlab, Mahout, SAS, SPSS...



# API

```
Facet(gender, height).hist()
```

**mplfacet**

```
fig, axes = plt.subplots(nrows=1, ncols=2, sharex=True, sharey=True)

axes[0].hist(height[gender == 'Male'])
axes[0].set_title('Male')
axes[1].hist(height[gender == 'Female'])
axes[1].set_title('Female')
```

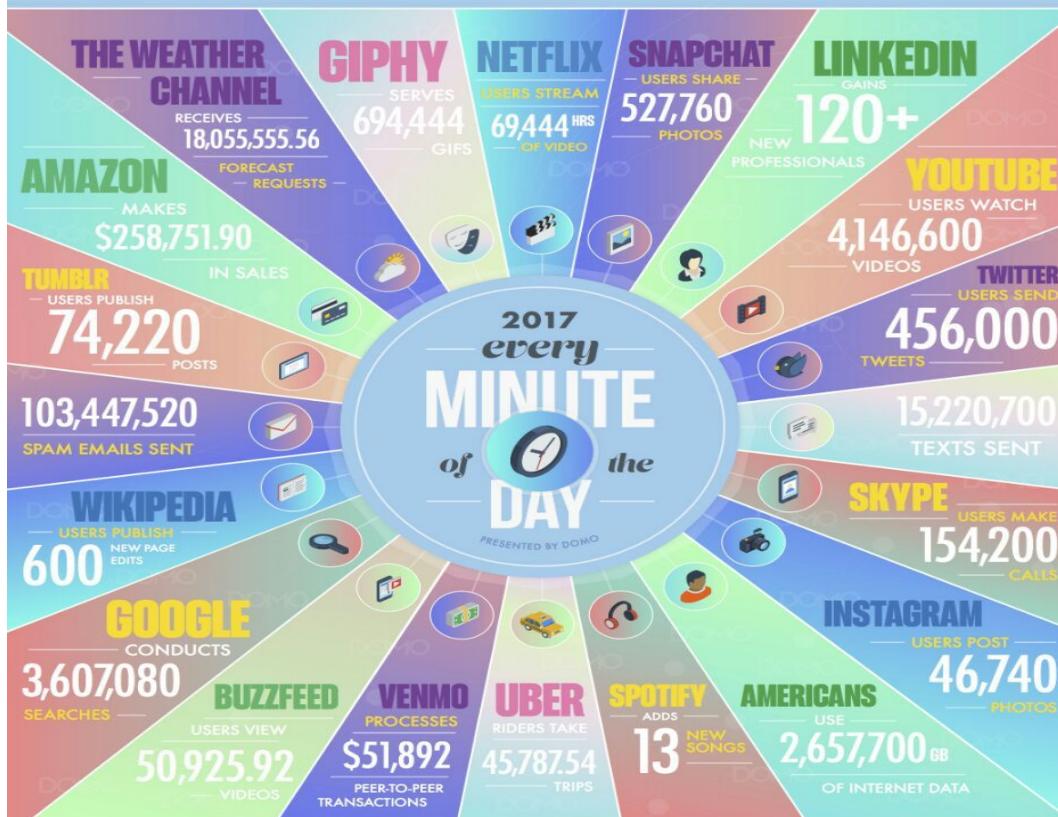
Systems,  
Databases  
Big data

DOMO

# DATA NEVER SLEEPS 5.0

How much data is generated *every minute*?

90% of all data today was created in the last two years—that's 2.5 quintillion bytes of data per day. In our 5th edition of Data Never Sleeps, we bring you the latest stats on just how much data is being created in the digital sphere—and the numbers are staggering.



The world internet population has grown 7.5% from 2016 and now represents 3.7 billion people.



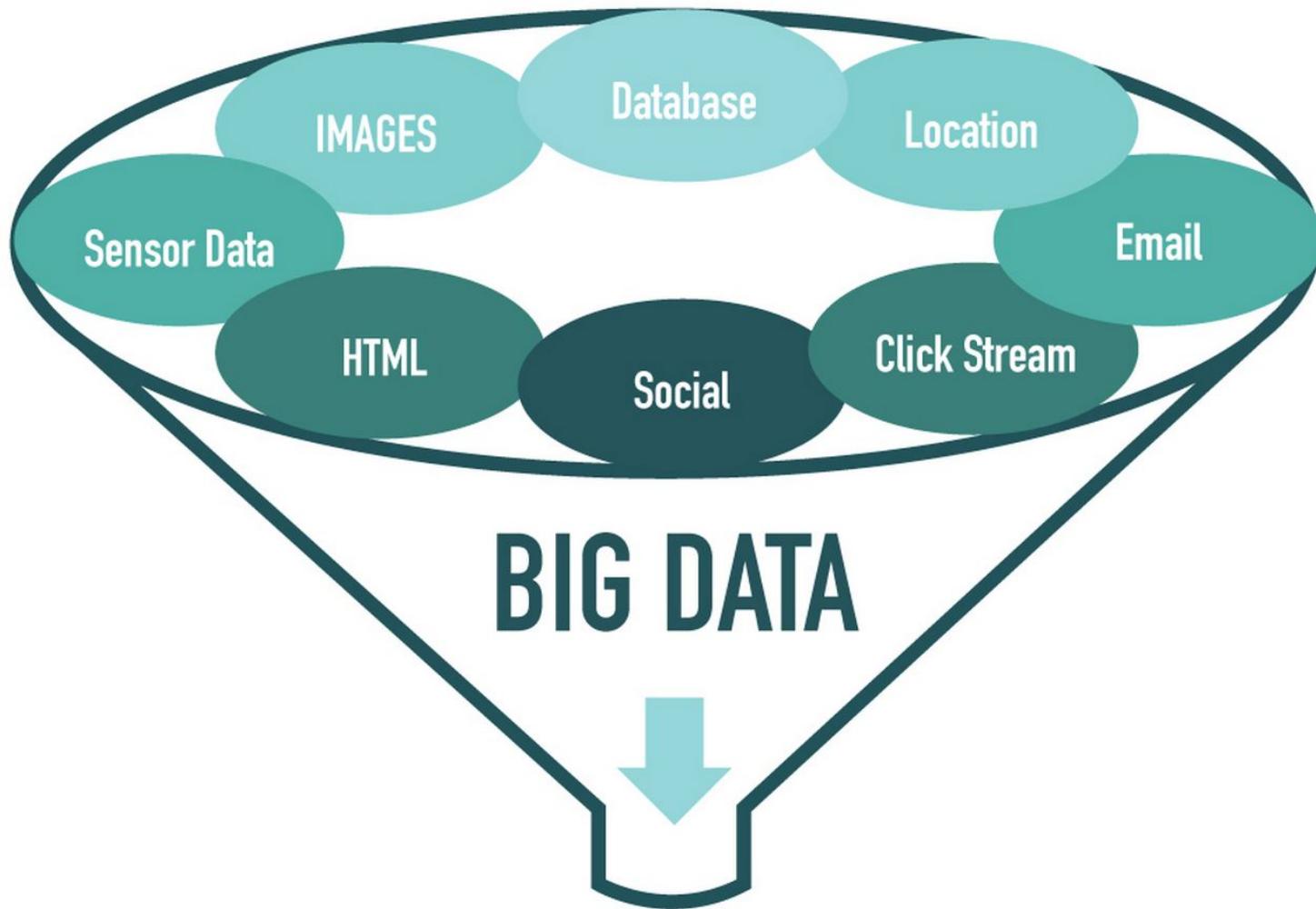
GLOBAL INTERNET POPULATION GROWTH 2012–2017  
(IN BILLIONS)

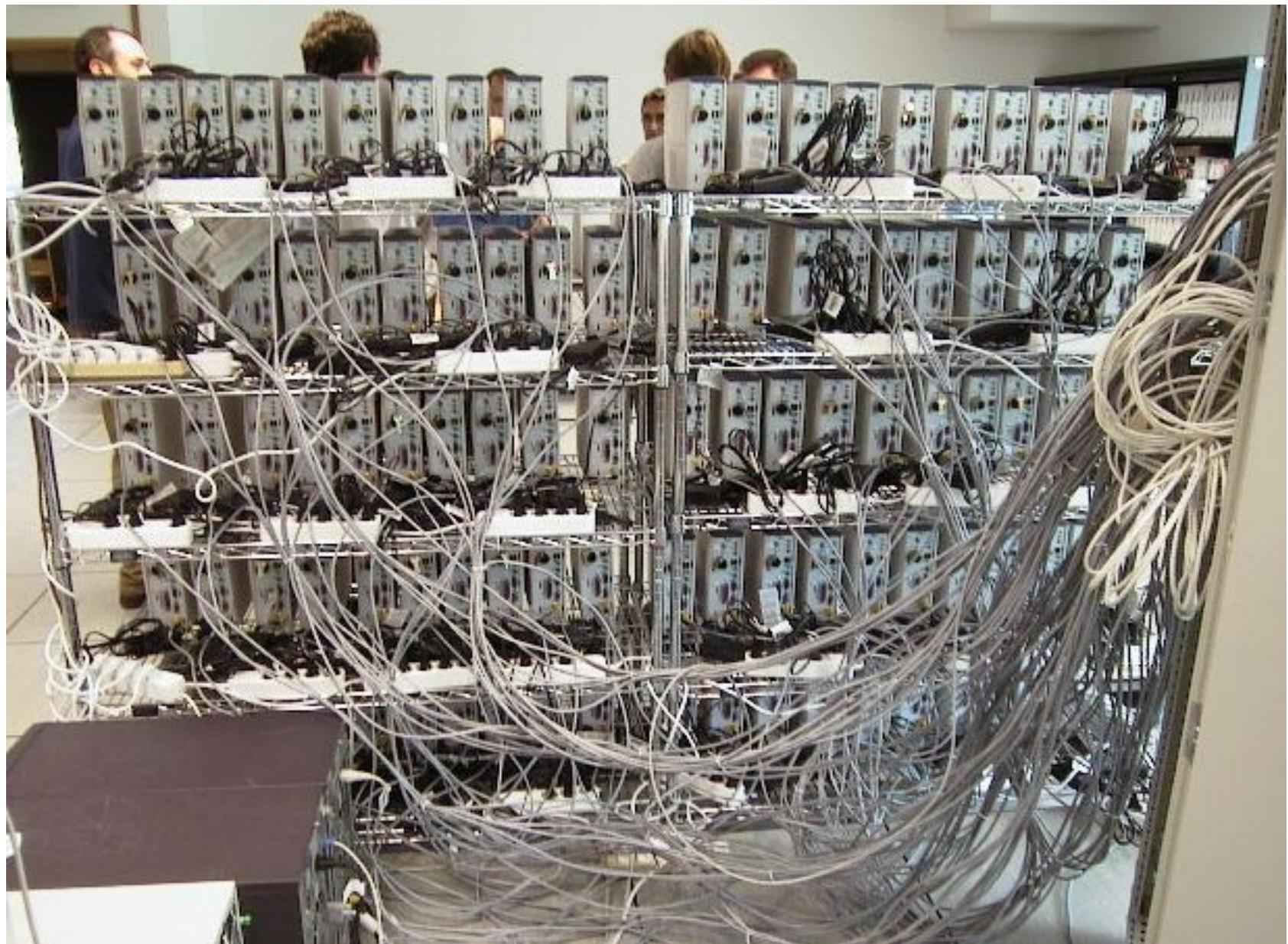
With each click, swipe, share, and like, businesses are using data to make decisions about the future. Domo gives everyone in your business real-time access to data from virtually any data source in a single platform for smarter decision-making at any moment.

Learn more at [domo.com](http://domo.com)

SOURCES: EXPANDERAMBLINGS.COM, WEARESOCIAL.COM, WIKIPEDIA, FORBES, ADWEEK.COM, FORTUNE.COM, BLOOMBERG.COM, ONEREACH.COM, IBM, BUZZFEED, INTERNET LIVE STATS, INTERNET WORLD STATS, BBC

DOMO







Google Cloud Platform

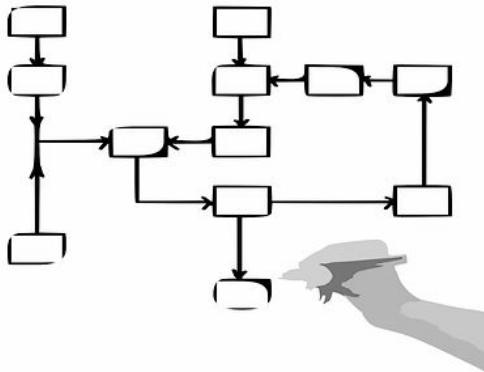


**SAY BIG DATA**



**ONE MORE TIME**

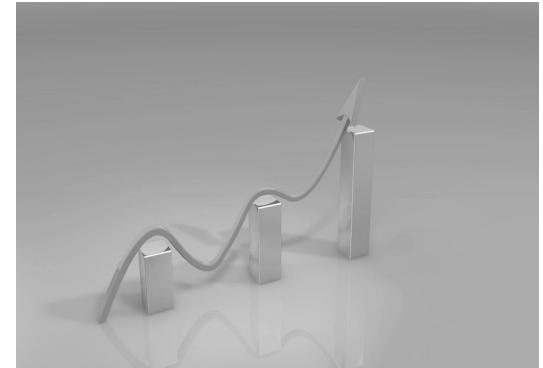
# Topics to cover



Data Mining Process



Exploratory Data Analysis



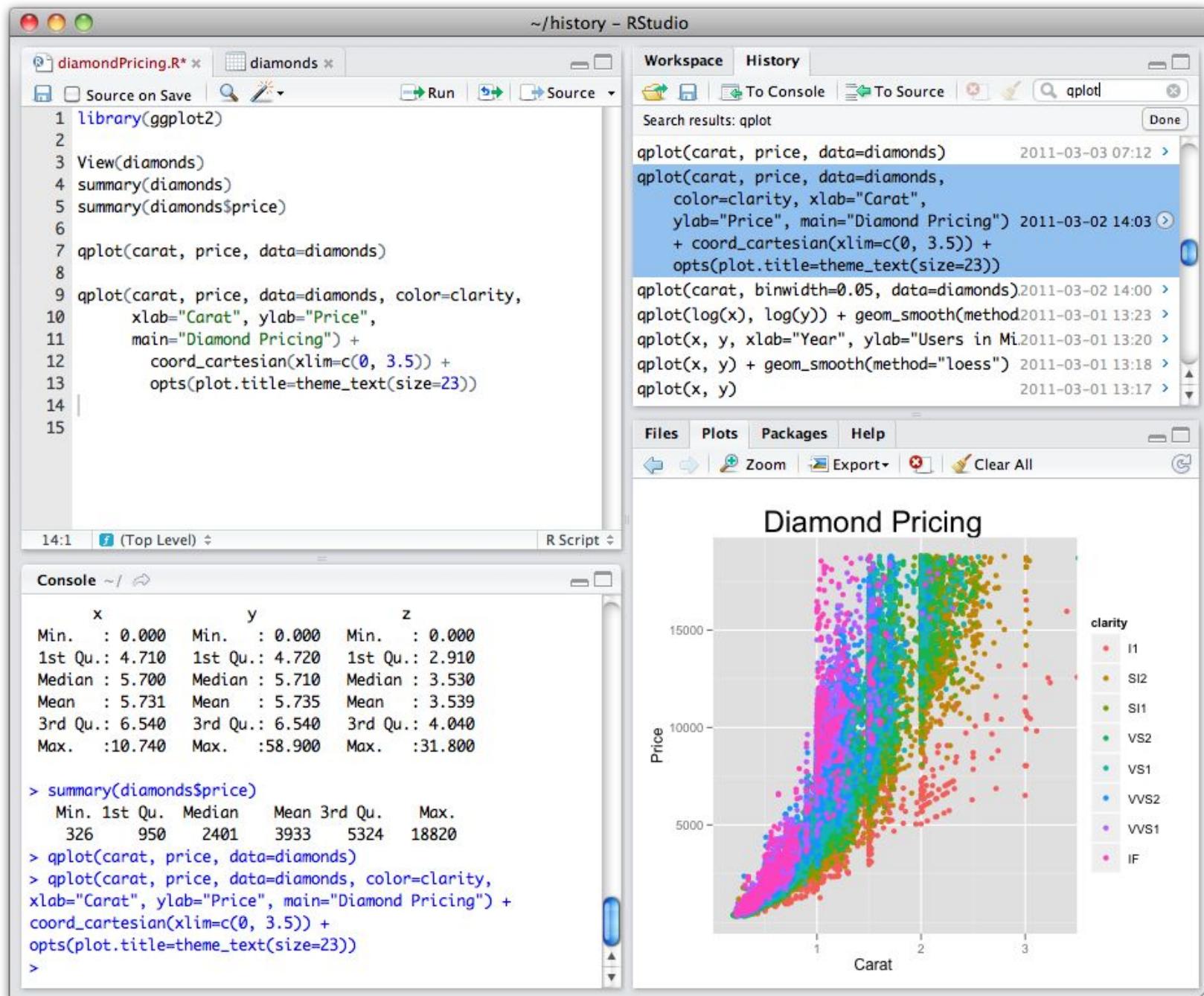
Data Visualization



Data Science Tools



Reproducible Research



IPy IPython Dashboard × IPy spectrogram × +

127.0.0.1:8888/a5222740-848b-4ac1-b212-d732c9f8f78b

IP[y]: Notebook spectrogram Last saved

File Edit View Insert Cell Kernel Help

Markdown

## Simple spectral analysis

An illustration of the [Discrete Fourier Transform](#)

$$X_k = \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi i}{N} kn} \quad k = 0, \dots, N-1$$

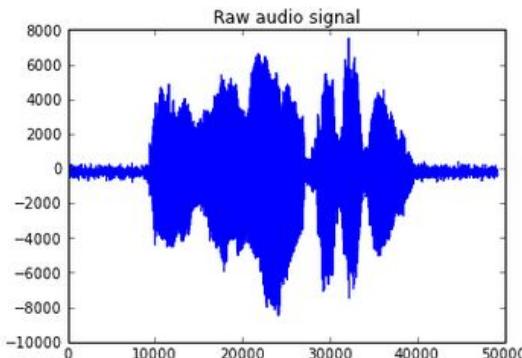
using windowing, to reveal the frequency content of a sound signal.

We begin by loading a datafile using SciPy's audio file support:

```
In [1]: from scipy.io import wavfile
rate, x = wavfile.read('test_mono.wav')
```

And we can easily view its spectral structure using matplotlib's builtin `specgram`:

```
In [2]: fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12,
ax1.plot(x); ax1.set_title('Raw audio signal')
ax2.specgram(x); ax2.set_title('Spectrogram'))
```



R Studio

chunks.Rmd × Knit HTML Chunks

1 R Code Chunks

2 =====

3

4 With R Markdown, you can insert R code

5 chunks including plots:

6 

```
```{r qplot, fig.width=4, fig.height=3,
| message=FALSE}
7 # quick summary and plot
8 library(ggplot2)
9 summary(cars)
10 qplot(speed, dist, data=cars) +
11   geom_smooth()
```

12

13

With R Markdown, you can insert R code chunks including plots:

# quick summary and plot

library(ggplot2)

summary(cars)

## speed dist

## Min. : 4.0 Min. : 2

## 1st Qu.:12.0 1st Qu.: 26

## Median :15.0 Median : 36

## Mean :15.4 Mean : 43

## 3rd Qu.:19.0 3rd Qu.: 56

## Max. :25.0 Max. :120

qplot(speed, dist, data = cars) + geom\_smooth()

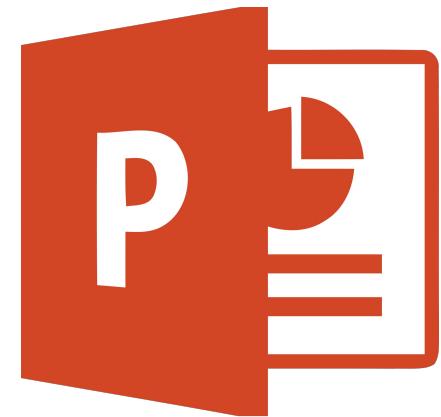
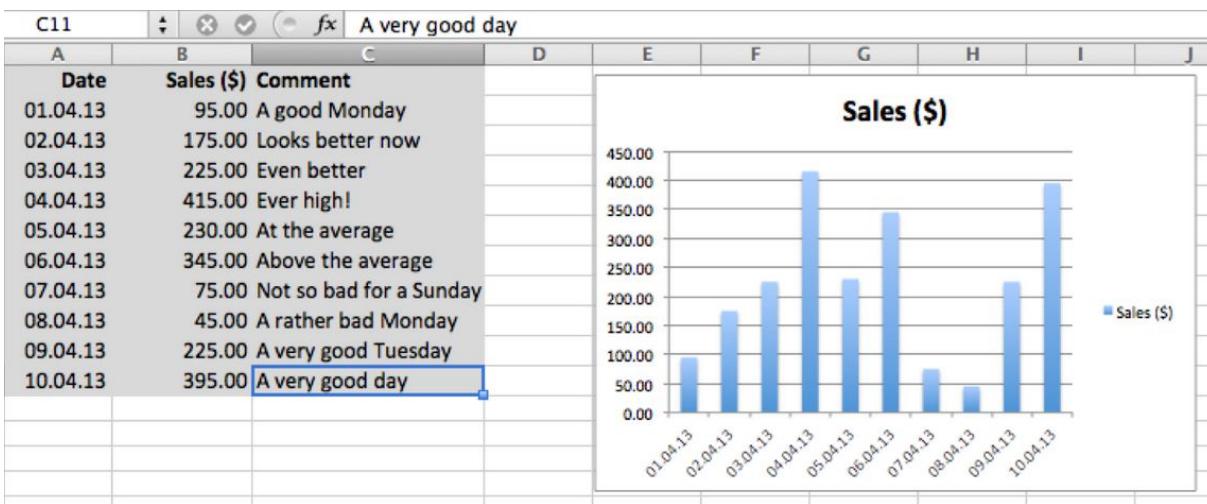
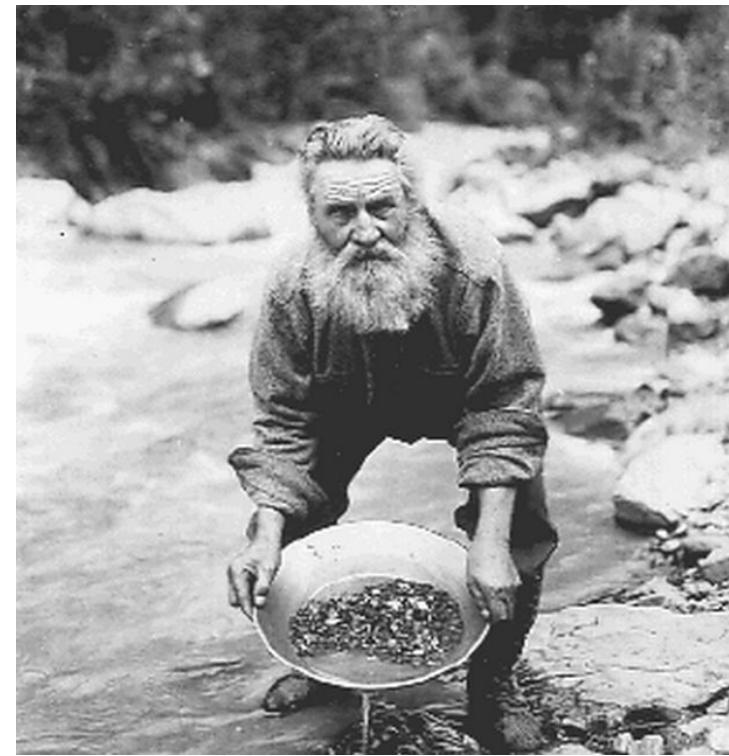
Query 1

```

1 •  SELECT `actor`.`actor_id`,
2     `actor`.`first_name`,
3     `actor`.`last_name`,
4     `actor`.`last_update`
5  FROM `sakila`.`actor`;
6
7 •  SELECT `film`.`film_id`,
8     `film`.`title`,
9     `film`.`description`,
10    `film`.`release_year`,
11    `film`.`language_id`,
12    `film`.`original_language_id`,
13    `film`.`rental_duration`,
14    `film`.`rental_rate`,
15    `film`.`length`,
16    `film`.`replacement_cost`,
17    `film`.`rating`,
18    `film`.`special_features`,
19    `film`.`last_update`
20

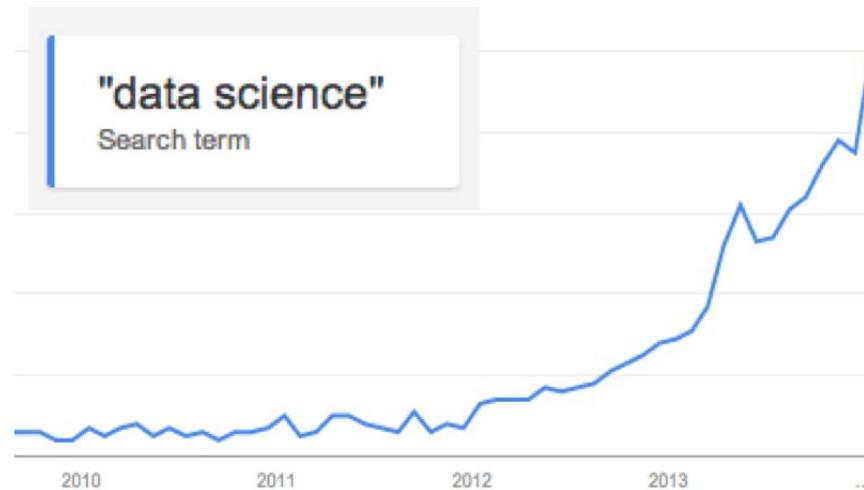
```

Result Set Filter:



# Data Science New?

# Is Data Science New?



Data Science is a newly emerging field dedicated to analyzing and manipulating data to derive insights and build data products.

<https://www.kaggle.com/wiki/WhatIsDataScience>



## **Data Scientist: The Sexiest Job of the 21st Century**

The United States alone faces a shortage of 140,000 to 190,000 people with deep analytical skills. [2011]

[http://www.mckinsey.com/insights/business\\_technology/big\\_data\\_the\\_next\\_frontier\\_for\\_innovation](http://www.mckinsey.com/insights/business_technology/big_data_the_next_frontier_for_innovation)

Four major influences act on **data analysis today**:

1. The formal theory of **statistics**
2. Revolutionary developments in **computers** and display devices
3. The challenge, in many fields, of more and ever **larger** bodies of **data**
4. The accelerating emphasis on quantification in an ever **wider variety of disciplines**

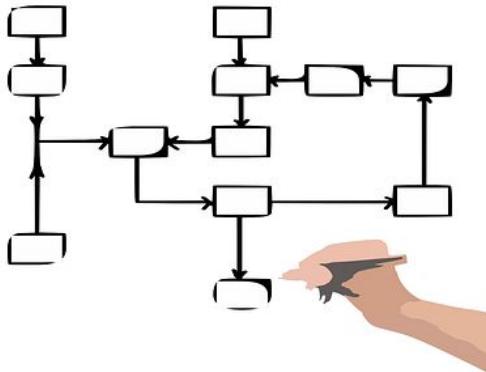
Four major influences act on **data analysis today**:

1. The formal theory of **statistics**
2. Revolutionary developments in **computers** and display devices
3. The challenge, in many fields, of more and ever **larger** bodies of **data**
4. The accelerating emphasis on quantification in an ever **wider variety of disciplines**

## Tukey & Wilk, 1965

Tukey, J.W., & Wilk, M.B. (1965). Data analysis and statistics: techniques and approaches    Reprinted in  
The Collected Works of John W. Tukey, Vol. V, Graphics 1965- 1985, 1-22 (1988)

# Topics covered



Data Mining Process



Exploratory Data Analysis



Data Visualization



Data Science Tools



Reproducible Research