# Big Data Analytics: What is Big Data?
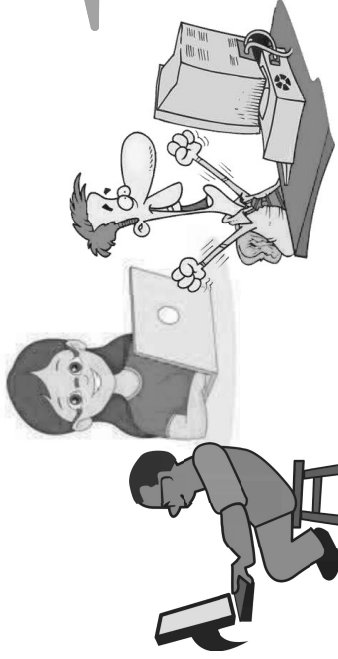
H. Andrew Schwartz

CSE545
Spring 2020
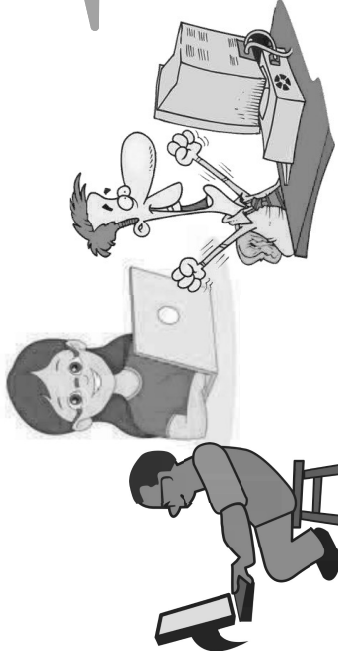
# Big Data, what is it?

data that will not fit in main memory.

traditional computer science

# Big Data, what is it?

data that will not fit in main memory.

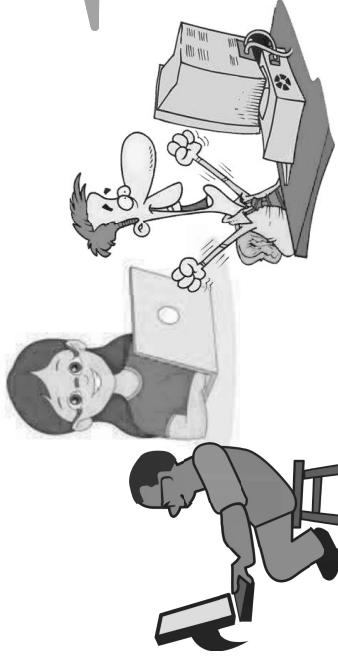traditional computer science

SSD Sequential Read: ~500 MB/s

For example...

busy web server access logs

graph of the entire Web

all of Wikipedia

daily satellite imagery over a year

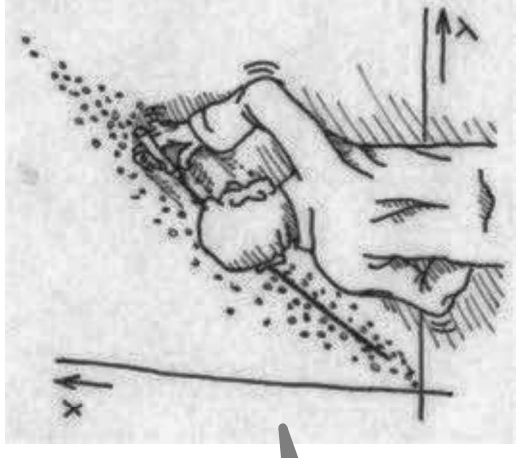# Big Data, what is it?

data that will not fit in main memory.

data with a *large* number of observations and/or features.

traditional computer science

statistics

# Big Data, what is it?

*Tall data:*

edge list of a large graph

rgb values per pixel location in large images

data with a large number of observations and/or features.

statistics

*Wide data:* mobile app usage statistics of 100 people

# Big Data, what is it?

data that will not fit in main memory.

data with a *large* number of observations and/or features.

traditional computer science

statistics

# Big Data, what is it?
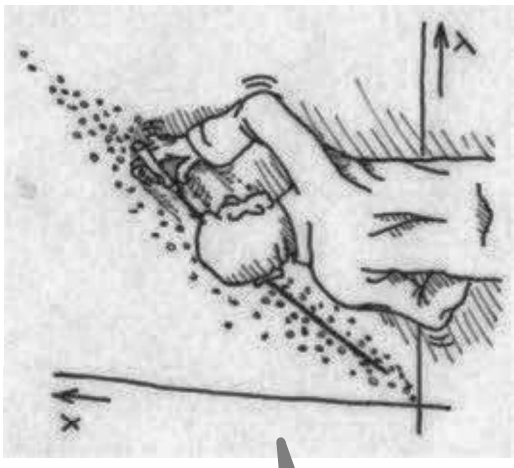
data that will not fit in main memory.

data with a *large* number of observations and/or features.

non-traditional sample size (i.e. > 100 subjects); can't analyze in stats tools (Excel).

statistics

traditional computer science

other fields

# Big Data, what is it? *Government View*

**THE WORLD BANK** (2016)
IBRD • IDA | WORLD BANK GROUP

**Big**Data
UN Global Working Group

## 1. Survey of SDG-related Big Data projects

### Type of data source(s)



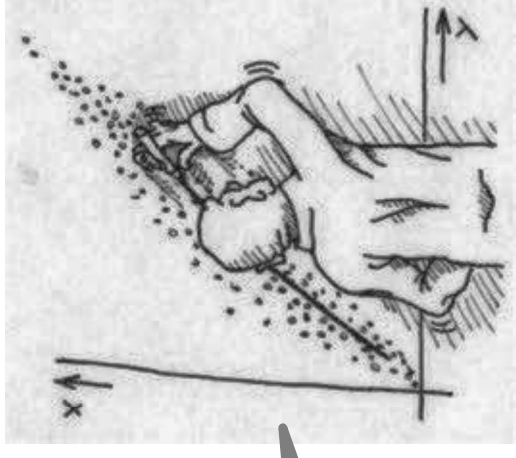| Data source | Count |
|---|---|
| Mobile phone data | 23 |
| Satellite imagery data and geodata | 20 |
| Web data | 17 |
| Twitter data | 12 |
| Other social networks | 12 |
| Financial transaction data | 11 |
| Scanner data | 11 |
| Facebook data | 8 |
| Sensor data | 6 |
| Smart meter data | 5 |
| Health records | 2 |
| Ships identification data | |
| Public transport usage data | |
| Credit card data | |

- Mobile (23), Satellite imagery (20) and social media (12+12+8) are the most prominent sources

# Big Data, what is it? *Industry View*

**Figure 2:** Sources of big data

Which of the following do you consider part of big data (regardless of whether your company uses each)?

| Source | Percentage |
|---|---|
| Large data files (20 terabytes or larger) | 65% |
| Advanced analytics or analysis | 60% |
| Data from visualization tools | 50% |
| Data from social networks | 48% |
| Unstructured data (e.g., video, open text, voice) | 43% |
| Geospatial/location information | 38% |
| Social media/monitoring/mapping | 37% |
| Telematics | 34% |
| Unstructured data/log files/free text | 25% |

Source: Accenture Big Success with Big Data Survey, April 2014

# Big Data, what is it? *Industry View*

**Figure 2:** Sources of big data

Which of the following do you consider part of big data (regardless of whether your company uses each)?

| | |
|---|---|
| Large data files (20 terabytes or larger) | 65% |
| Advanced analytics or analysis | 60% |
| Data from visualization tools | 50% |
| Data from social networks | 48% |
| Unstructured data (e.g., video, open text, voice) | 43% |
| Geospatial/location information | 38% |
| Social media/monitoring/mapping | 37% |
| Telematics | 34% |
| Unstructured data/log files/free text | 25% |

# Big Data, a type of analytics

**Figure 2:** Sources of big data

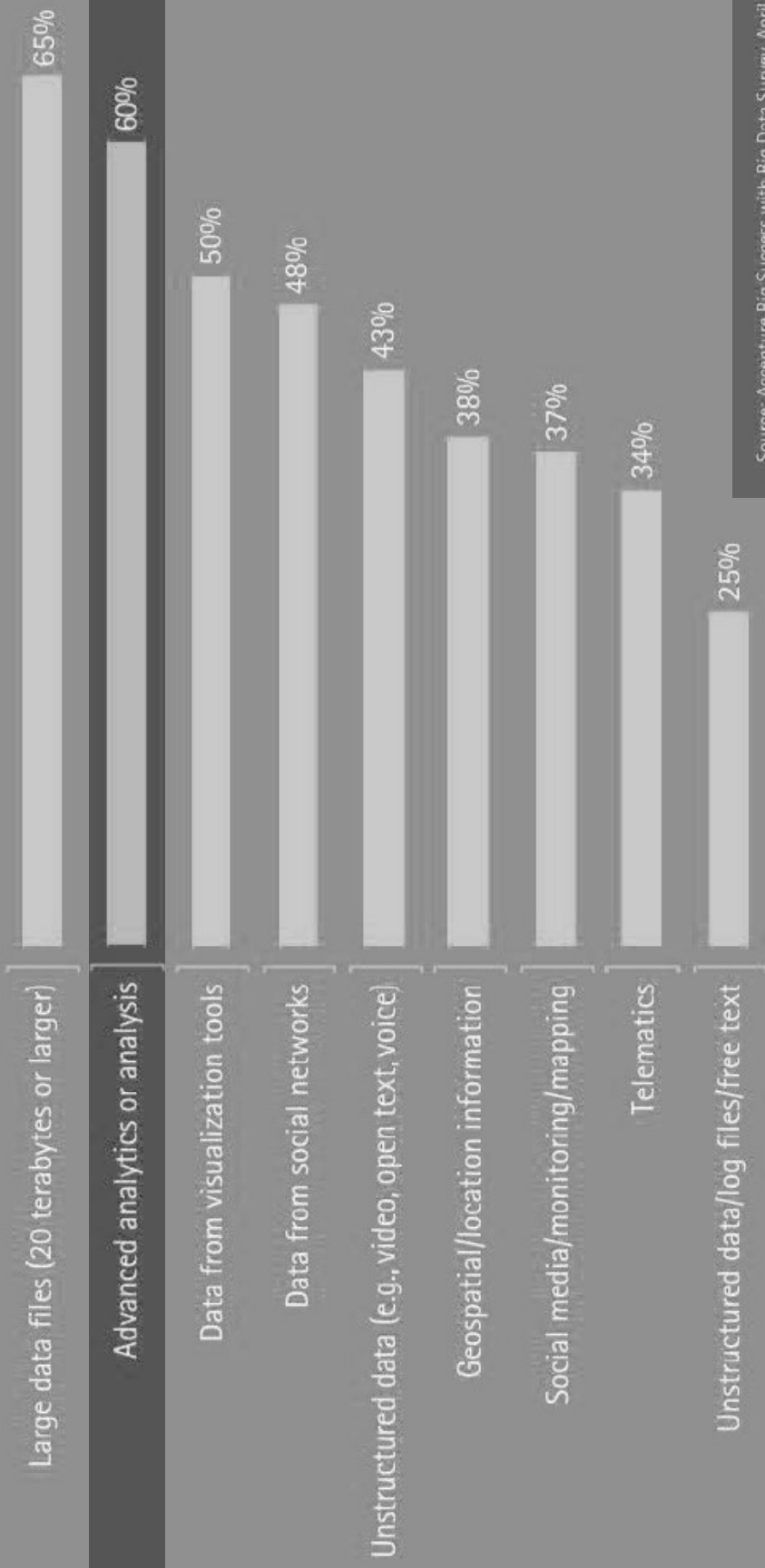Which of the following do you consider part of big data (regardless of whether your company uses each)?

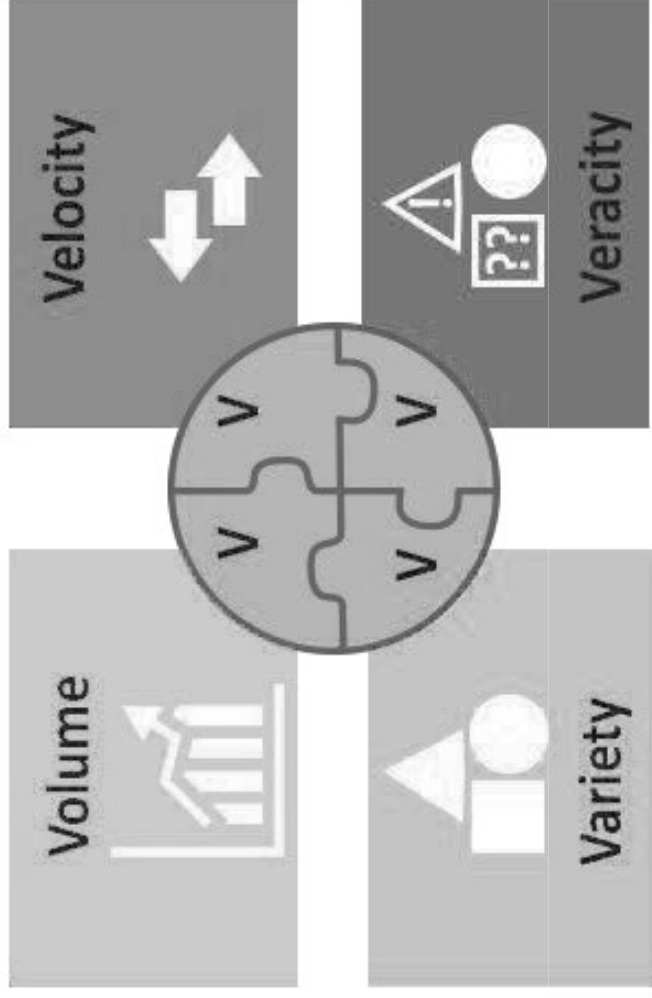| Source | Percentage |
|---|---|
| Large data files (20 terabytes or larger) | 65% |
| Advanced analytics or analysis | 60% |
| Data from visualization tools | 50% |
| Data from social networks | 48% |
| Unstructured data (e.g., video, open text, voice) | 43% |
| Geospatial/location information | 38% |
| Social media/monitoring/mapping | 37% |
| Telematics | 34% |
| Unstructured data/log files/free text | 25% |

# Big Data, a type of analytics

Analyses which can handle the "3 Vs":

1. Volume – large quantity

2. Velocity – arriving quickly

3. Variety – [un]structured, multi-modal

# Big Data, a type of analytics

*Analyses which can handle the "34 Vs":*

# Big Data, a type of analytics

# Big Data, a type of analytics

Data

Insights!

# Big Data, a type of analytics

## The Big Data Challenge

# Big Data, a buzz word?

VISIBILITY

Peak of Inflated Expectations

Plateau of Productivity

Slope of Enlightenment

Trough of Disillusionment

Technology Trigger

TIME

(Gartner Hype Cycle)

# Big Data, a buzz word?



2018

2012

2011

2011

2010

2008

# Big Data, a buzz word?

POPULAR
DATA IS POWER

VISIBILITY

Peak of Inflated Expectations

Flu Trends Criticized (2014)

Plateau of Productivity

Slope of Enlightenment

Trough of Disillusionment

Google Flu Trends (2008)
Technology Trigger

TIME

(Gartner Hype Cycle)

# Big Data, a buzz word?



Peak of Inflated Expectations

Flu Trends Criticized (2014)

Plateau of Productivity

Big Data Analytics Adoption Soared In The Enterprise In 2018

Louis Columbus
Enterprise & Cloud

TWEET THIS

Big data adoption in enterprises soared from 17% in 2015 to 59% in 2018, reaching a Compound Annual Growth Rate (CAGR) of 36%.

Slope of Enlightenment

Trough of Disillusionment

Google Flu Trends (2008)
Technology Trigger

TIME

(Gartner Hype Cycle)

# Big Data, a buzz word?

- main-stream study being established
  - Realization of what subfields are really doing "big data" (i.e. data mining, ML, Statistics, computational social sciences).
  - Best practices being established.

Peak of Inflated Expectations

Flu Trends Criticized (2014)

Plateau of Productivity

Slope of Enlightenment

Big Data Analytics Adoption Soared In The Enterprise In 2018

Trough of Disillusionment

Google Flu Trends (2008)
Technology Trigger

TIME

(Gartner Hype Cycle)

# Big Data, in demand?

**Figure 3:** Main challenges with big data projects

What are the main challenges to implementing big data in your company?



| Challenge | % |
|---|---|
| Security | 51% |
| Budget | 47% |
| Lack of talent to implement big data | 41% |
| Lack of talent to run big data and analytics on an ongoing basis | 37% |
| Integration with existing systems | 35% |
| Procurement limitations on big data vendors | 33% |
| Enterprise not ready for big data | 27% |

Source: Accenture Big Success with Big Data Survey, April 2014

# Big Data, in demand?

**Figure 6:** Big data's competitive significance



Big data will revolutionize the way we do business to a degree similar to the advent of the Internet in the 1990s
- Strongly Agree: 51%
- Agree: 38%
- Neither Agree nor Disagree: 10%
- Disagree: 1%

Big data will dramatically change the way we do business in the future
- Strongly Agree: 39%
- Agree: 46%
- Neither Agree nor Disagree: 13%
- Disagree: 2%

Companies that do not embrace big data will lose their competitive position and may even face extinction
- Strongly Agree: 37%
- Agree: 42%
- Neither Agree nor Disagree: 19%
- Disagree: 2%

We feel we are ahead of our peers in using big data and this creates a competitive advantage for us
- Strongly Agree: 37%
- Agree: 46%
- Neither Agree nor Disagree: 12%
- Disagree: 4%

Legend: Strongly Agree | Agree | Neither Agree nor Disagree | Disagree

# Big Data, in demand?



**Adoption of Big Data 2015-2018**
(Copyright 2018 – Dresner Advisory Services)

Legend: 2015, 2016, 2017, 2018

Categories: Yes. We use big data today | We may use big data in the future | No. We have no plans to use big data at all
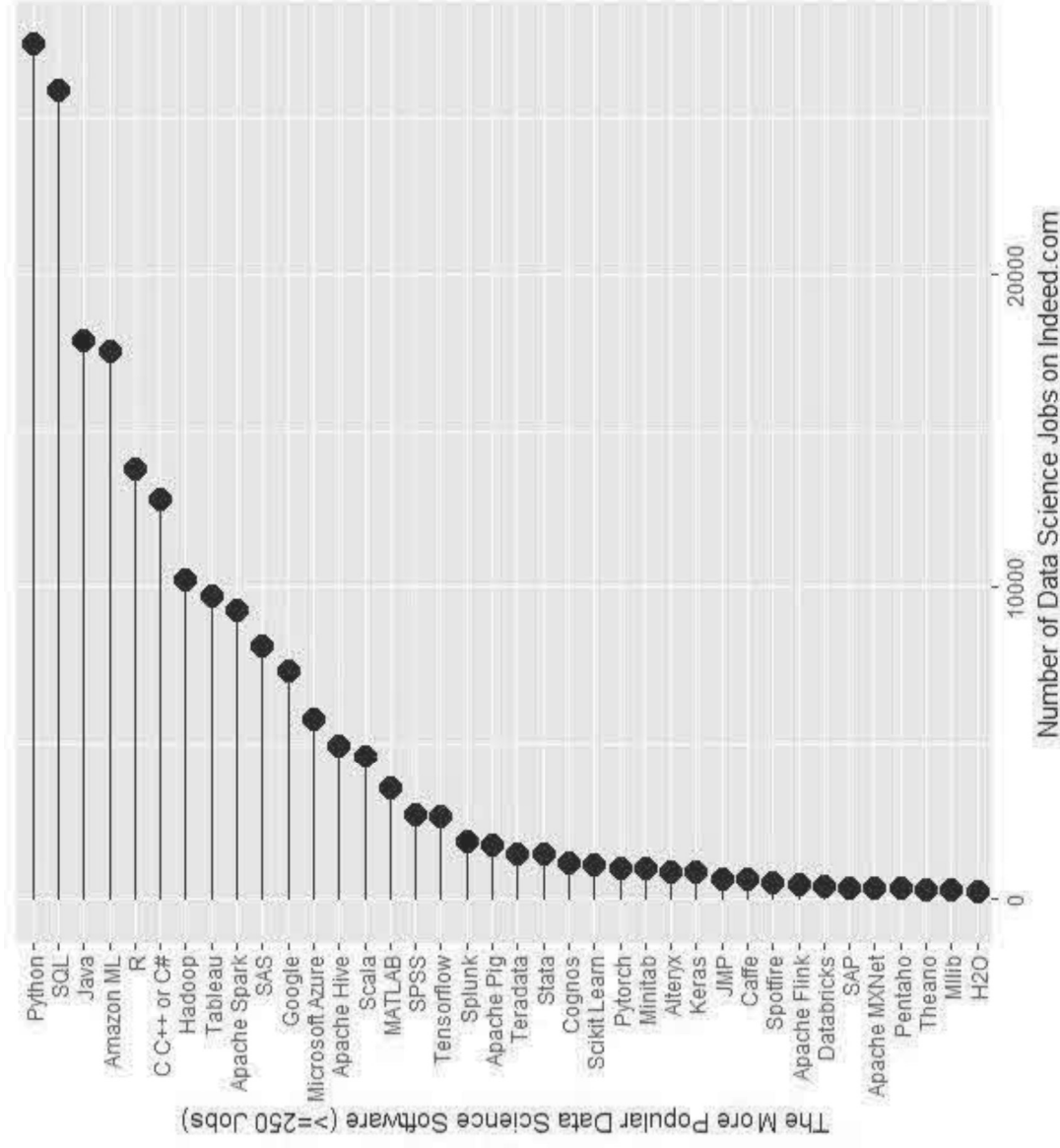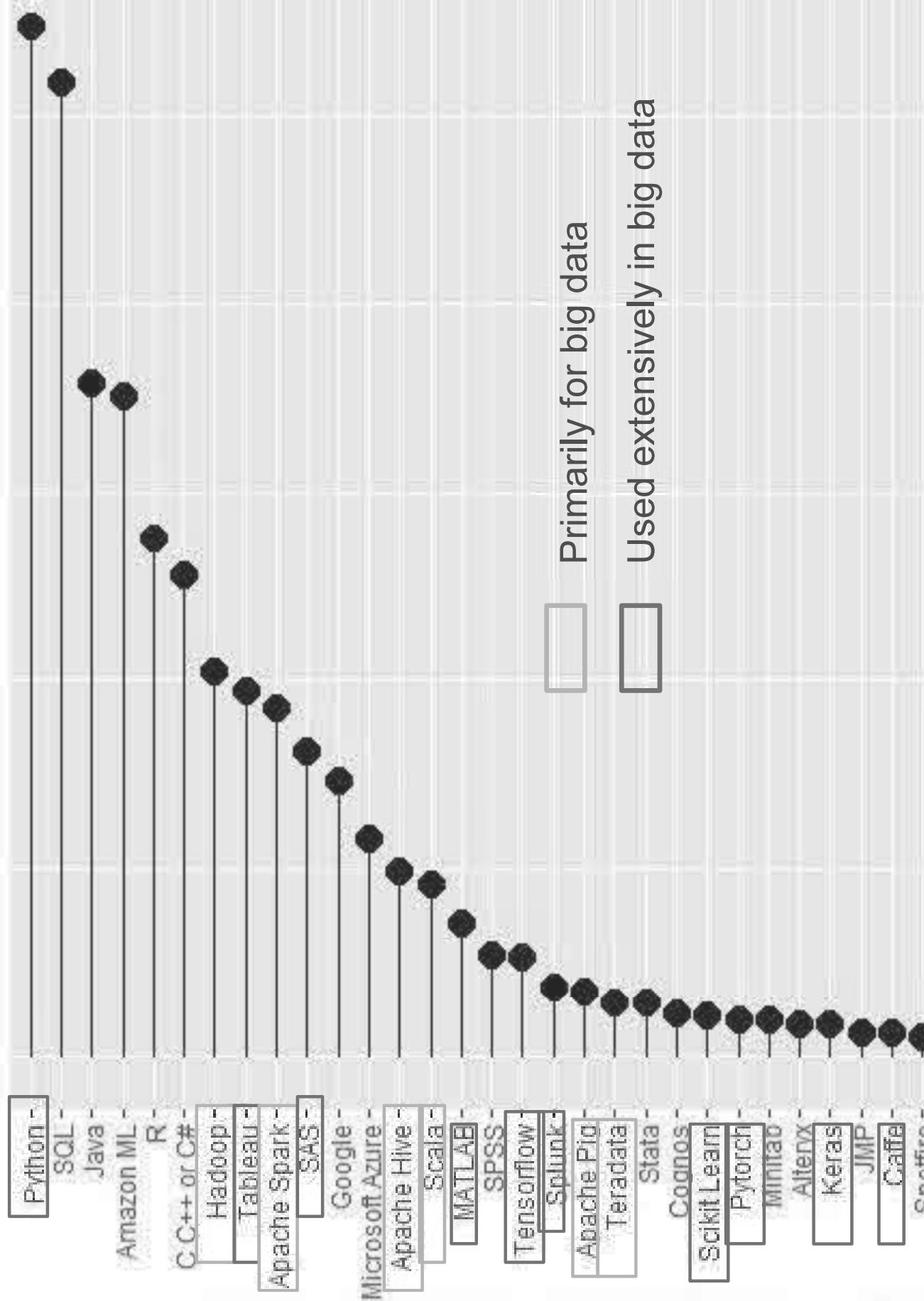
Axis: 0%, 10%, 20%, 30%, 40%, 50%, 60%, 70%

# Big Data, in demand?

By the requirements
in job ads.
(Muenchen,2019)



The More Popular Data Science Software (>=250 Jobs)

Number of Data Science Jobs on Indeed.com

More Popular Data Science Software (>=250 Jobs)

Python, SQL, Java, Amazon ML, R, C C++ or C#, Hadoop, Tableau, Apache Spark, SAS, Google, Microsoft Azure, Apache Hive, Scala, MATLAB, SPSS, Tensorflow, Splunk, Apache Pig, Teradata, Stata, Cognos, Scikit Learn, Pytorch, Minitab, Alteryx, Keras, JMP, Caffe, Spotfire

Primarily for big data

Used extensively in big data

# Big Data, What is it?

*Short Answer:*

*Big Data ≈ Data Mining ≈ Predictive Analytics ≈ Data Science*

(Leskovec et al., 2014)

# Big Data, What is it?

*Short Answer:*

*Big Data ≈ Data Mining ≈ Predictive Analytics ≈ Data Science*
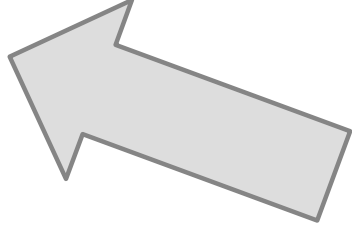
(Leskovec et al., 2014)

*CSE545 focuses on:*

How to analyze data that is mostly too large for main memory.

Analyses only possible with a *large* number of observations or features.

# Big Data, What is it?

**Goal:** Generalizations
A *model* or *summarization* of the data.

Analyses only possible with a *large* number of observations or features.

How to analyze data that is mostly too large for main memory.

# Big Data, What is it?

> **Goal:** Generalizations
>
> A *model* or *summarization* of the data.

E.g.

- Google's PageRank: *summarizes* web pages by a single number.
- Twitter financial market predictions: *Models* the stock market according to shifts in sentiment in Twitter.
- Distinguish tissue type in medical images: *Summarizes millions of pixels into clusters.*
- Mental health diagnosis in social media: *Models* presence of diagnosis as a distribution (a summary) of linguistic patterns.
- Frequent co-occurring purchases: *Summarize billions of purchases as items that frequently are bought together.*

# Big Data, What is it?

**Goal: Generalizations**
*A model or summarization of the data.*

1. Descriptive analytics
   Describe (*generalizes*) the data itself

2. Predictive analytics
   Create something *generalizeable* to new data

# Big Data Analytics, The Class

**Core Data Science Courses**

CSE 519: Data Science Fundamentals

CSE 544: Prob/Stat for Data Scientists

**CSE 545: Big Data Analytics**

CSE 512: Machine Learning

CSE 537: Artificial Intelligence

CSE 548: Analysis of Algorithms

CSE 564: Visualization

**Applications of Data Science**

CSE 527:
   Computer Vision

CSE 538:
   Natural Language Processing

CSE 549:
   Computational Biology

...

# Big Data Analytics, The Class

## Core Data Science Courses

CSE 519: Data Science Fundamentals

CSE 544: Prob/Stat for Data Scientists

**CSE 545: Big Data Analytics**

CSE 512: Machine Learning

CSE 537: Artificial Intelligence

CSE 548: Analysis of Algorithms

CSE 564: Visualization

## Applications of Data Science

CSE 527:
Computer Vision

CSE 538:
Natural Language Processing

CSE 549:
Computational Biology

...

**Key Distinction:**
Focus on scalability and algorithms / analyses not possible without large data.

# Big Data Analytics, The Class

**Goal:** Generalizations
A model or summarization of the data.

*Data Frameworks*

*Algorithms and Analyses*

# Big Data Analytics, The Class

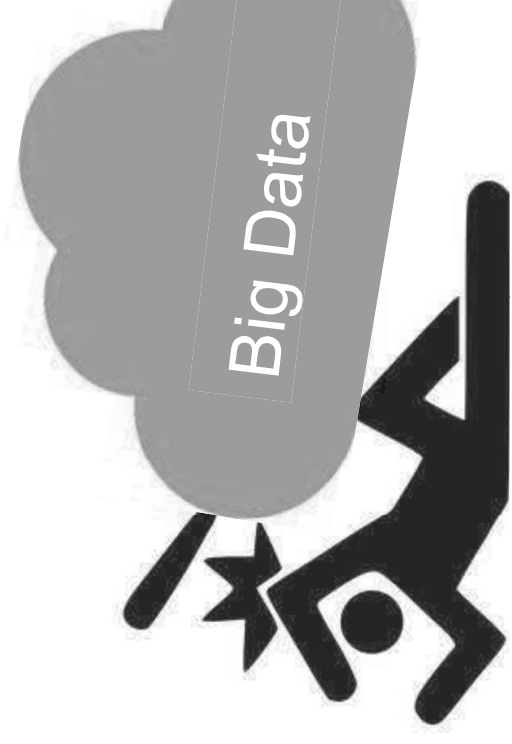**Goal:** Generalizations
A model or summarization of the data.

Algorithms and Analyses

Data Frameworks

Hadoop File System    Spark

Streaming    TensorFlow

MapReduce

# Big Data Analytics, The Class

**Goal:** Generalizations
A model or summarization of the data.

Data Frameworks

Algorithms and Analyses

Hadoop File System   Spark
Streaming
MapReduce   Tensorflow

Similarity Search   Linear Modeling
Recommendation Systems
Graph Analysis   Deep Learning

# Big Data Analytics, The Class

Big Data

# Preliminaries

Ideas and methods that will repeatedly appear:

- Bonferroni's Principle
- Normalization (TF.IDF)
- Power Laws
- Hash functions
- IO Bounded (Secondary Storage)
- Unstructured Data

- *Parallelism*
- *Functional Programming*
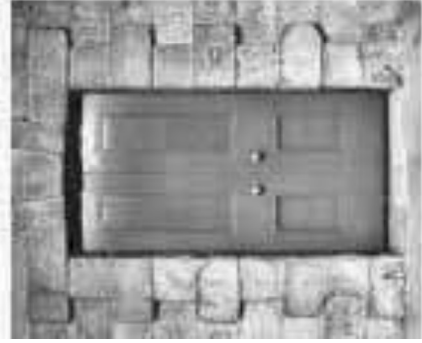
# Statistical Limits.      Goal:   **Generalization**

Bonferroni's Principle

A to consider goal of generalization:

Find events that didn't just happen *by chance*.

# Statistical Limits.

## Bonferroni's Principle; an example:

# Statistical Limits.

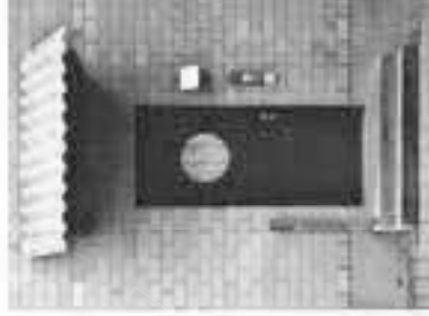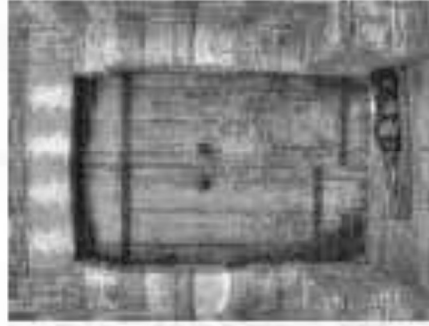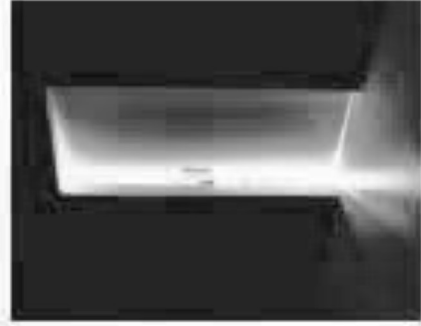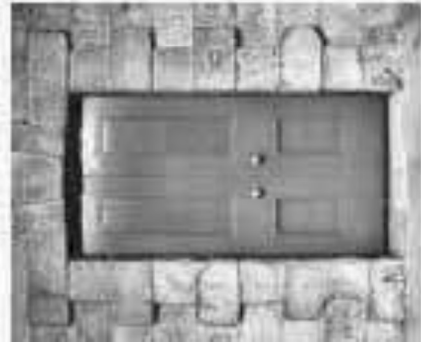## Bonferroni's Principle; an example:

Statistical Limits.

Bonferroni's Principle

Goal:

**Generalization
(i.e. not by chance)**

# Bonforroni's Principle; Task Example

snazzyphones.com wants to know which case to eliminate.

6 total cases:



Red

Green

Blue

Teal

Purple

Yellow
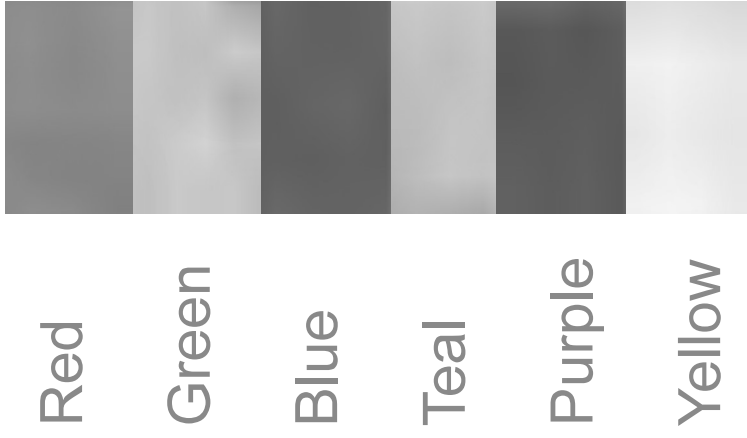
# Bonforroni's Principle; Task Example

snazzyphones.com wants to know which case to eliminate.

6 total cases:

first day, 17 sales:

*What is the data telling you?*

Red

Green

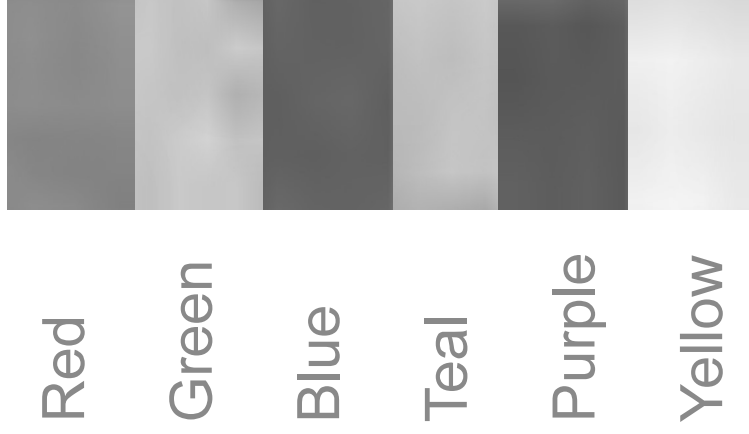Blue

Teal

Purple

Yellow

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

# Bonforroni's Principle; Task Example

snazzyphones.com wants to know which case to eliminate.

6 total cases:

first day, 17 sales:

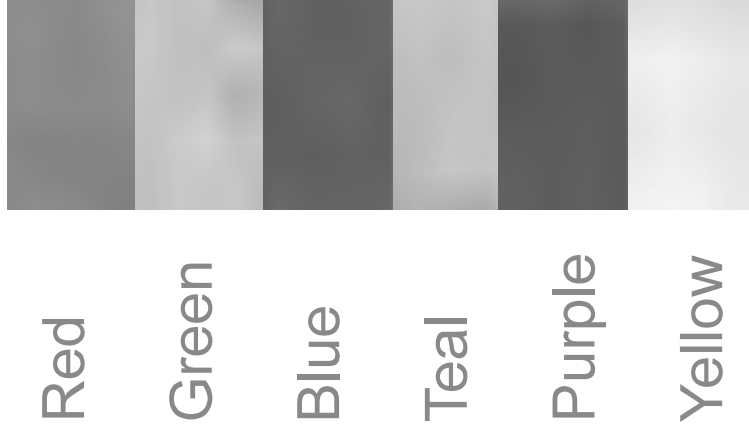*What is the data telling you?*

*The blue isn't selling?*

Red

Green

Blue

Teal

Purple

Yellow

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

# Statistical Limits.        Goal: **Generalization**

Bonferroni's Principle

Roughly, calculating the probability of any of n *findings* being true requires n times the probability as testing for 1 finding.

https://xkcd.com/882/

In brief, one can only look for so many patterns (i.e. features) in the data before one finds something just by chance (i.e. finding something that does **not** generalize).

"Data mining" is a bad word in some communities!

# Statistical Limits.    Goal: **Generalization**

Note: *Bonferroni's principle* is simply an abstract idea inspired by a precisely defined method of hypothesis testing called "Bonferroni correction".

We will go over this <u>correction method</u> later. The *<u>principle</u>* is the more important idea to understand as a big data practitioner.

In brief, one can only look for so many patterns (i.e. features) in the data before one finds something just by chance (i.e. finding something that does **not** generalize).

"Data mining" is a bad word in some communities!

# Normalizing

Count data often need *normalizing* -- putting the numbers on the same "scale".

Prototypical example: TF.IDF

# Normalizing

Count data often need *normalizing* -- putting the numbers on the same "scale".

Prototypical example: TF.IDF of word *i* in document *j*:

Term Frequency:

$$tf_{ij} = \frac{count_{ij}}{\max_k count_{kj}}$$

$$tf.idf_{ij} = tf_{ij} \times idf_i$$

Inverse Document Frequency:

$$idf_i = log_2\left(\frac{docs_*}{docs_i}\right) \propto \frac{1}{\frac{docs_i}{docs_*}}$$

where docs is the number of documents containing word *i*.

# Normalizing

Count data often need *normalizing* -- putting the numbers on the same "scale".

Prototypical example: TF.IDF of word *i* in document *j:*

Term Frequency:

$$tf_{ij} = \frac{count_{ij}}{\boxed{\max_k count_{kj}}}$$

$$tf.idf_{ij} = tf_{ij} \times idf_i$$

Inverse Document Frequency:

$$idf_i = \boxed{log_2} \left( \frac{docs_*}{docs_i} \right) \propto \frac{1}{\boxed{\frac{docs_i}{docs_*}}}$$

where docs is the number of documents containing word *i*.

# Normalizing

**Standardize**: puts different sets of data (typically vectors or random variables) on the same scale with the same center.

- Subtract the mean (i.e. "mean center")

- Divide by standard deviation

$$z_i = \frac{x_i - \bar{x}}{s_x}$$

# Power Law

Characterized many frequency patterns when ordered from most to least:

County Populations [r-bloggers.com]

# links into webpages [Broader et al., 2000]

Sales of products [see book]

Frequency of words [Wikipedia, "Zipf's Law"]

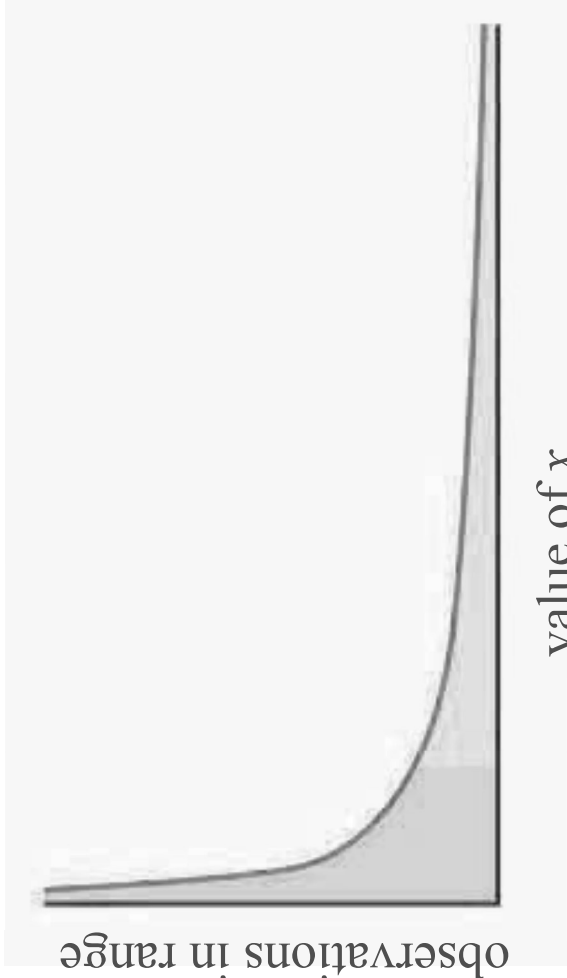("popularity" based statistics, especially without limits)

# Power Law

$$\log y = b + a \log x$$

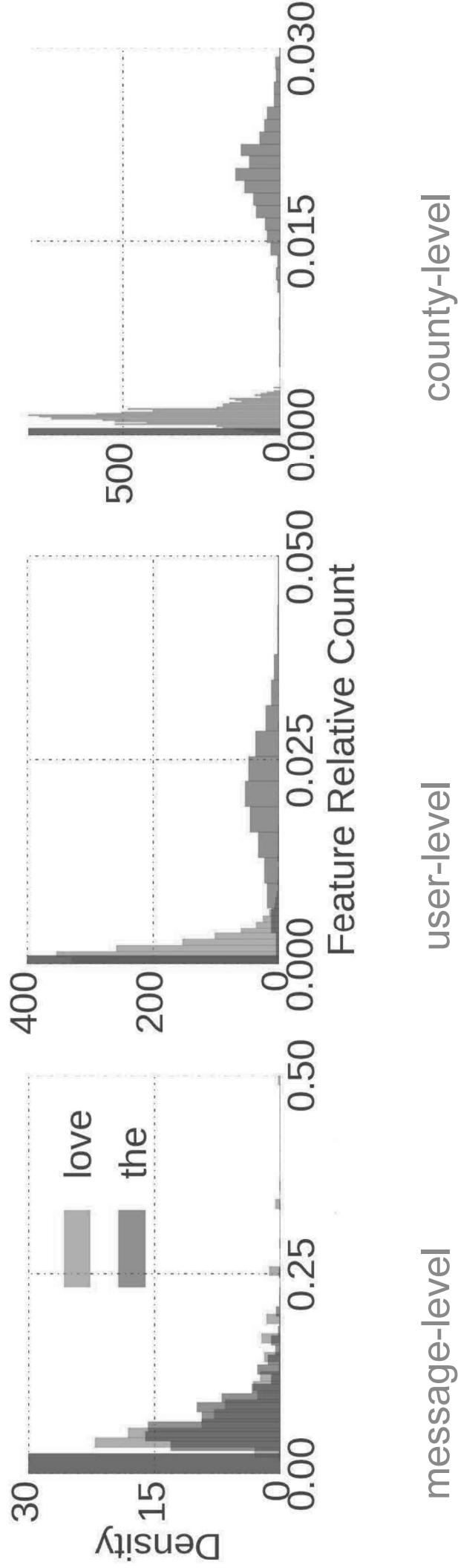raising to the natural log:

$$y = e^b e^{a \log x} = e^b x^a = cx^a$$

where c is just a constant

Characterizes "the Matthew Effect" -- the rich get richer



density: proportion of observations in range

value of $x$

# Power Law



message-level　　　　user-level　　　　county-level

Feature Relative Count

Legend: love, the

Almodaresi, F., Ungar, L., Kulkarni, V., Zakeri, M., Giorgi, S., & Schwartz, H. A. (2017). On the Distribution of Lexical Features at Multiple Levels of Analysis. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics* (pp. 79-84).

# Hash Functions and Indexes

Review:

*h: hash-key -> bucket-number*

Objective: uniformly distribute hash-keys across buckets.

Example: storing word counts.

# Hash Functions and Indexes

Review:

*h: hash-key -> bucket-number*

Objective: uniformly distribute hash-keys across buckets.

Example: storing word counts.

$$h(word) = \left( \sum_{char \in word} ascii(char) \right) \% \; \#buckets$$

# Hash Functions and Indexes

Review:

*h: hash-key -> bucket-number*

Objective: uniformly distribute hash-keys across buckets.

Example: storing word counts.

$$h(word) = \left( \sum_{char \in word} ascii(char) \right) \% \#buckets$$

Data structures utilizing hash-tables (i.e. O(1) lookup; dictionaries, sets in python) are a friend of big data algorithms! Review further if needed.

# Hash Functions and Indexes

Review:

*h: hash-key -> bucket-number*

Objective: uniformly distribute hash-keys across buckets.
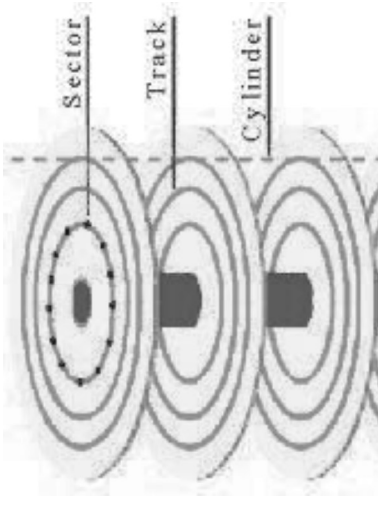
Example: storing word counts.

**Database Indexes:** Retrieve all records with a given *value*. (also review if unfamiliar / forgot)

Data structures utilizing hash-tables (i.e. O(1) lookup; dictionaries, sets in python) are a friend of big data algorithms! Review further if needed.

# IO Bounded

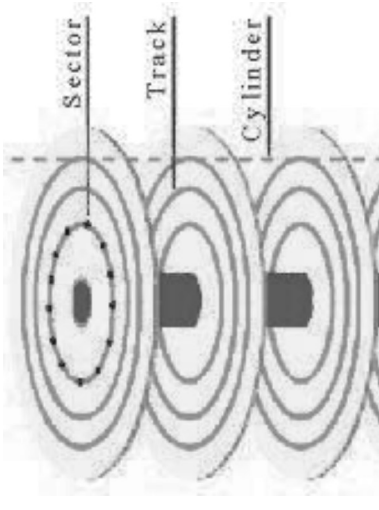Reading a word from disk versus main memory: $10^5$ slower!

Reading many contiguously stored words
is faster per word, but fast modern disks
still only reach 150MB/s for sequential reads.

Sector

Track

Cylinder

# IO Bounded

Reading a word from disk versus main memory: $10^5$ slower!

Reading many contiguously stored words
is faster per word, but fast modern disks
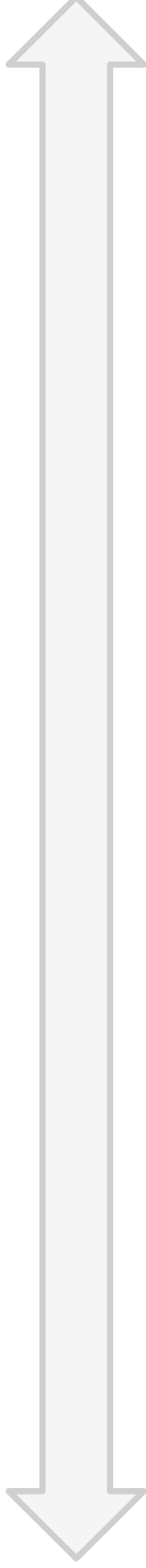still only reach 150MB/s for sequential reads.



IO Bound: biggest performance bottleneck is reading / writing to disk.

(starts around 100 GBs; ~10 minutes just to read).

# Data

**Structured** ⟷ **Unstructured**

- Unstructured ≈ requires processing to get what is of interest
- Feature extraction used to turn unstructured into structured
- Near infinite amounts of potential features in unstructured data

# Data

Structured ← → Unstructured

| Structured | | | Unstructured |
|---|---|---|---|
| mysql table | email header | satellite imagery | images |
| vectors matrices | facebook likes | | text (email body) |

- Unstructured ≈ requires processing to get what is of interest
- Feature extraction used to turn unstructured into structured
- Near infinite amounts of potential features in unstructured data