# Data Science
# in the Business World

*BUSI488 & COMP488*

*Daniel M. Ringel*

UNC Kenan-Flagler Business School
Spring 2022

January 12th, 2023

**Class 2:** Data Formats, Data Types, Data Quality,
and Outlier Detection

Sections 001 and 002

UNC
KENAN-FLAGLER
BUSINESS SCHOOL

# Today's Agenda

**1**  What are Data?

**2**  Data Sources of Structured vs. Unstructured Data

**3**  Data Storage

**4**  Data Types and Structure

**5**  Data Quality

**6**  Cleaning Data

**7**  Anomalies

**8**  Working with Data in Cloud Computing (Google CoLab)

**9**  Learning to Code in Python in this Course (DataCamp)

***Prep-Check:***

✓ Read Syllabus
✓ Set-up CoLab
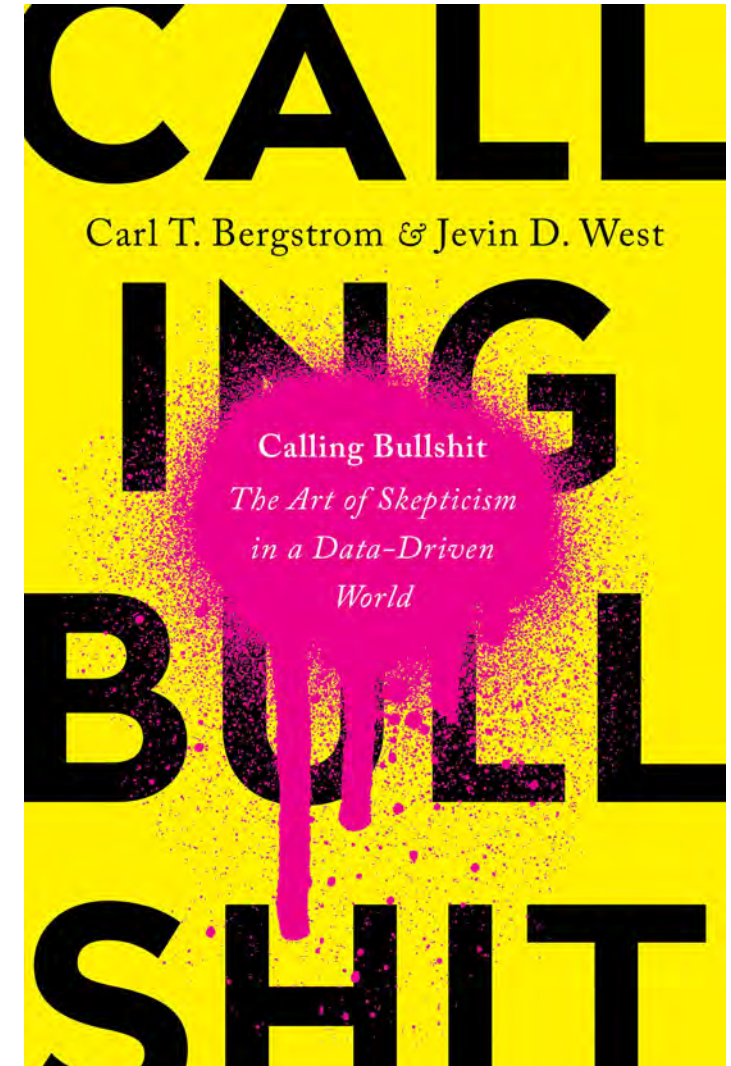✓ Started DataCamp HW1

# What are Data?

- Data (singular datum) are individual **units of information**

- A datum describes a **single quality or quantity** of some object or phenomenon

- In analytical processes, data are **represented by variables**

- Data are measured, collected, reported, and analyzed

- Data is **not equivalent** to **insight** or **knowledge**

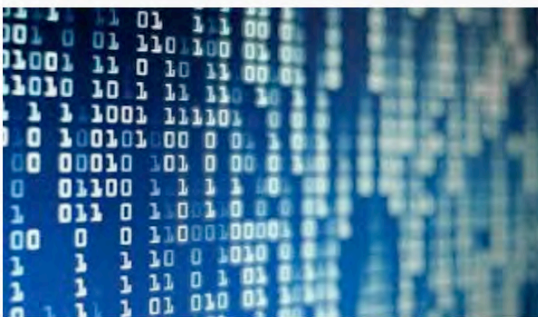- Data is the **least abstract concept**; Knowledge the most abstract

**Sources:**
Shannon, C.E. (1948), A mathematical theory of communication. *Bell system technical journal*, 27(3), pp.379-423.
Wikipedia (2020), https://en.wikipedia.org/wiki/Data, accessed January 15th 2020

# What do Data look like?

Data - Wik ... 
en.wikipedia ...

... big ...

Big Data Brings Challenges Beyond the ...
cpomagazine.com

How can big data turn into improved ...
atlas-network.com

Data Analysis: What, How, and Why to Do ...
import.io

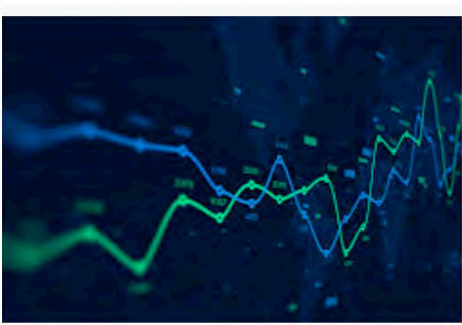What is Data: Types of Data, and How To...
simplilearn.com

The Future of Work | Transforming Data ...
tdwi.org

Why Is Big Data So Important ...
europeanbusinessmagazine.com

Big Data Analytics Startups Which Are ...
startup-buzz.com

Data Analytics Overtakes Big Data ...
flextrade.com

What is the Shelf Life of Data? | 201...
pobonline.com

The challenges of using data lakes in ...
seleritysas.com

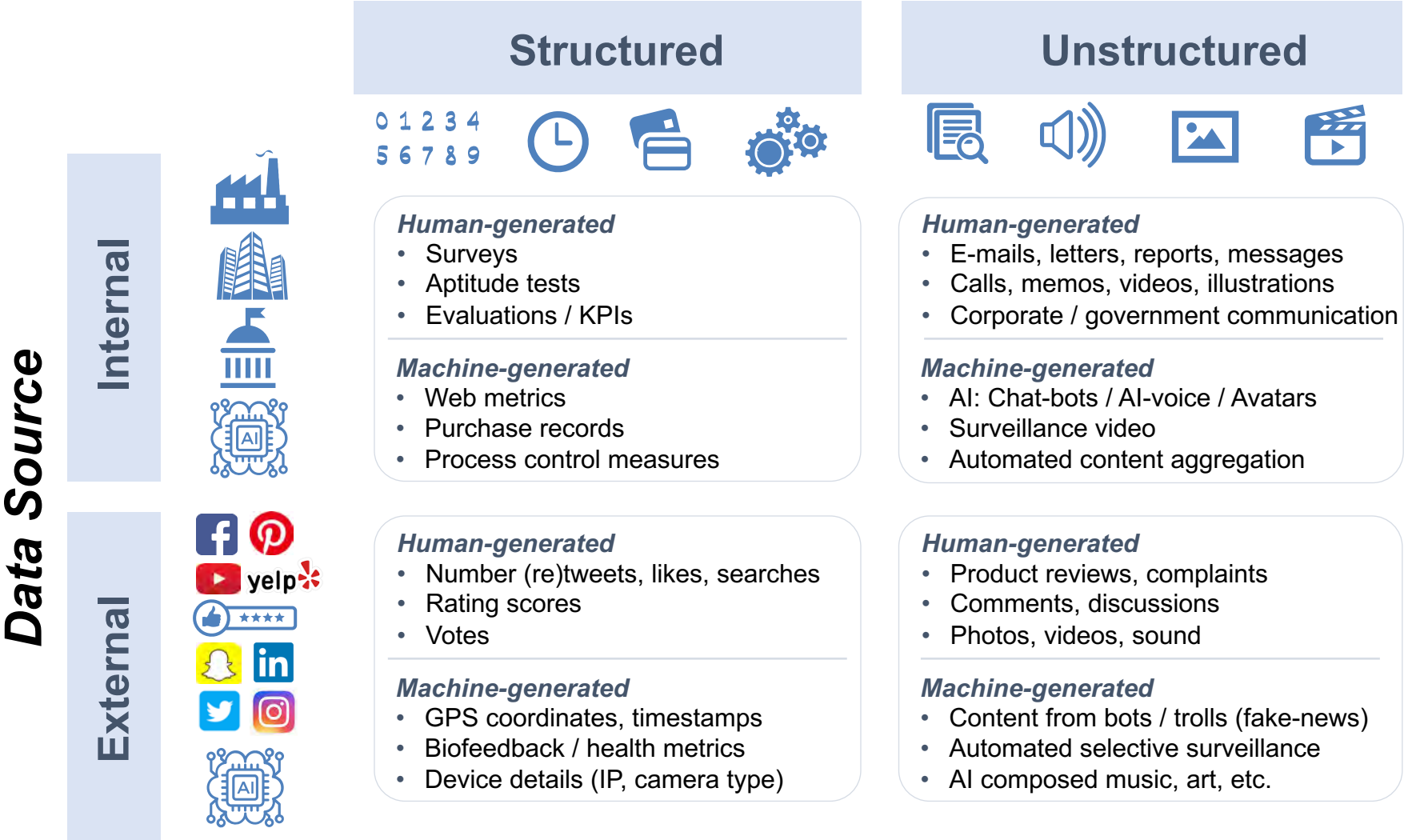Enable big data analytics and enhance ...
druva.com

Data Accountability and Trust Act ...
rush.house.gov

Data Preprocessing : Concepts ...
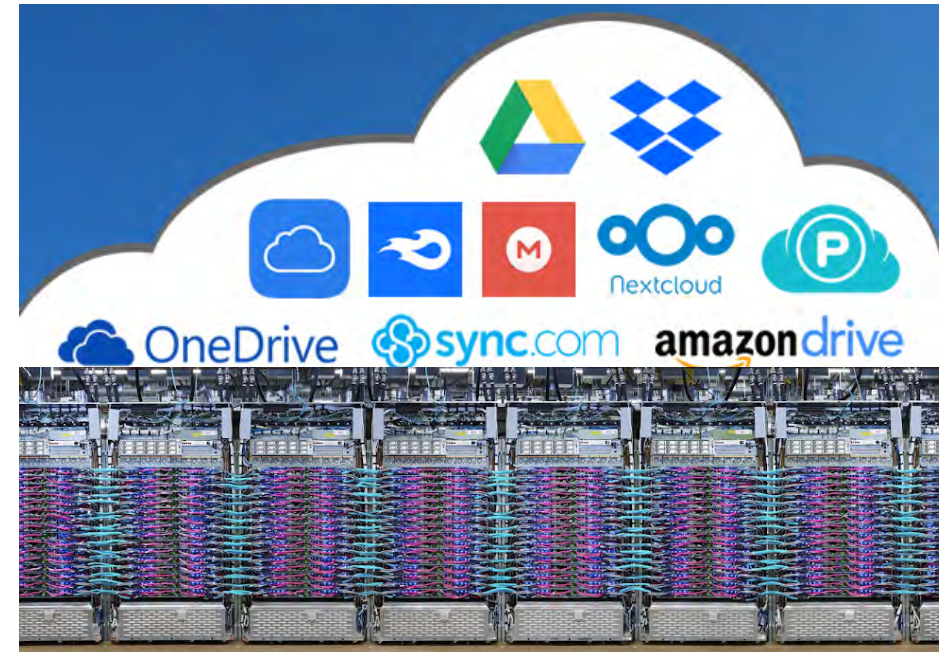towardsdatascience.com

4

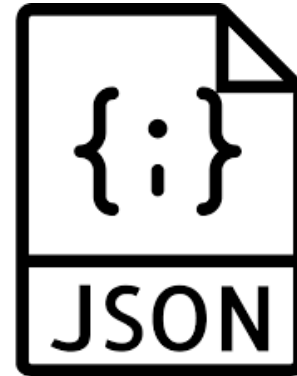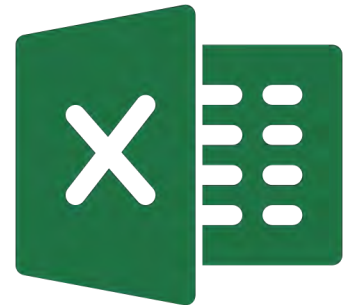# Data Sources and Data Structure

## Data Structure

| Structured | Unstructured |
|---|---|

### Internal

**Structured**

*Human-generated*
- Surveys
- Aptitude tests
- Evaluations / KPIs

*Machine-generated*
- Web metrics
- Purchase records
- Process control measures

**Unstructured**

*Human-generated*
- E-mails, letters, reports, messages
- Calls, memos, videos, illustrations
- Corporate / government communication

*Machine-generated*
- AI: Chat-bots / AI-voice / Avatars
- Surveillance video
- Automated content aggregation

### External

**Structured**

*Human-generated*
- Number (re)tweets, likes, searches
- Rating scores
- Votes

*Machine-generated*
- GPS coordinates, timestamps
- Biofeedback / health metrics
- Device details (IP, camera type)

**Unstructured**

*Human-generated*
- Product reviews, complaints
- Comments, discussions
- Photos, videos, sound

*Machine-generated*
- Content from bots / trolls (fake-news)
- Automated selective surveillance
- AI composed music, art, etc.

**Data Source**

# How Data are Stored 1.0

# How Data are Stored 2.0

**Data Formats for Storage
Some Examples**

CSV

JSON

.TXT

SQL

PICKLE

# Data in Flat Text Files

## Text Editor View



## Raw Text View

'\r\nProject Gutenberg's The Complete Works of William Shakespeare, by William\r\nShakespeare\r\n\r\nThis eBook is for the use of anyone anywhere in the United States and\r\nmost other parts of the world at no cost and with almost no restrictions\r\nwhatsoever. You may copy it, give it away or re-use it under the terms\r\nof the Project Gutenberg License included with this eBook or online at\r\nwww.gutenberg.org. If you are not located in the United States, you'll\r\nhave to check the laws of the country where you are located before using\r\nthis ebook.\r\n\r\n\r\nTitle: The Complete Works of William Shakespeare\r\n\r\nAuthor: William Shakespeare\r\n\r\nRelease Date: January 1994 [EBook #100]\r\nLast Updated: November 7, 2019\r\n\r\nLanguage: English\r\n\r\nCharacter set encoding: UTF-8\r\n\r\n*** START OF THIS PROJECT GUTENBERG EBOOK THE COMPLETE WORKS OF WILLIAM SHAKESPEARE ***\r\n\r\n\r\n\r\n\r\nThe Complete Works of William Shakespeare\r\n\r\n\r\n\r\nby William Shakespeare\r\n\r\n\r\n\r\n Contents\r\n\r\n\r\n\r\n THE SONNETS\r\n\r\n\r\n

| \r | carriage return |
| \n | new line |

# Data in Microsoft Excel



*Raw File view*

# Data in a Comma Separated Value (CSV) File



Lines
(with line break at end of each line often indicated by "\n" or by ↵ )

Columns are separated by commas—other separators possible(!!!)

Note that the line numbering is not part of the csv file and was added for better readability by the text editor used (Sublime Text—get it at https://www.sublimetext.com/ )

# Data in a JavaScript Object Notation (JSON) File



**Note:** originally, there were no line breaks—everything was on a single long line. I added breaks to better illustrate how information was organized

- JSON is a lightweight data-interchange format
- JSON is "self-describing" and easy to understand
- JSON is language independent *
- JSON is text only: it can easily be sent to and from a server, and used as a data format by any programming language

*JSON uses JavaScript syntax, but the JSON format is text only. Text can be read and used as a data format by any programming language.*

# Data in a Pickle File

The **Python Pickle Module** implements binary protocols for serializing and de-serializing a Python object structure.

**"Pickling"** is the process whereby a Python object hierarchy is converted into a byte stream, and **"unpickling"** is the inverse operation, whereby a byte stream (from a binary file or bytes-like object) is converted back into an object hierarchy.

Pickling (and unpickling) is alternatively known as "serialization", "marshalling," or "flattening"; however, to avoid confusion, the terms used in this course are "pickling" and "unpickling".

**When Not To Use pickle**
If you want to use data across different programming languages, pickle is not recommended. In contrast, JSON is standardized and language-independent. This is a serious advantage over pickle. It's also much faster than pickle.



*Check-out DataCamp's Tutorial*
https://www.datacamp.com/community/tutorials/pickle-python-tutorial

# Data in a Pandas Series and DataFrame

## Series

A one-dimensional labeled array capable of holding any data type

Column

Rows

| | |
|---|---|
| 0 | Albania |
| 1 | Algeria |
| 2 | Andorra |
| 3 | Anguilla |
| 4 | Antigua and Barbuda |

Index

series1 = pd.Series(['Albania', 'Algeria', 'Andorra', 'Anguilla', 'Antigua and Barbuda'], index = [0, 1, 2, 3, 4])

## DataFrame

A two-dimensional labeled data structure with columns of potentially different types

Columns

| | Country or Area | 1990 | 1995 |
|---|---|---|---|
| 0 | Albania | 28385.000000 | 40311.000000 |
| 1 | Algeria | 76160.000000 | 90270.000000 |
| 2 | Andorra | 539.947998 | 510.673004 |
| 3 | Anguilla | 93.099998 | 100.730003 |
| 4 | Antigua and Barbuda | 300.299988 | 374.500000 |

Rows

Index

data = {'Country or Area': ['Albania', 'Algeria', 'Andorra', 'Anguilla', 'Antigua and Barbuda'],
      '1990': [28385, 76160, 539.947998, 93.099998, 300.299988],
      '1995': [40311, 90270, 510.673004, 100.730003, 374.500000]}

df = pd.DataFrame(data, columns = ['Country or Area', '1990', '1995'])

# Many Names for (often) the same Thing

Dataset, File, Table, Sheet, Data Frame

Rows
Instances
Cases
Examples
Observations
Tuples

Cell,
Value

| | Country or Area | 1990 | 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Albania | 28385.000000 | 40311.000000 | 0.000000 | 0.000000 | 0.0 | 38284.000000 | 30683.000000 | 30491.00 |
| 1 | Algeria | 76160.000000 | 90270.000000 | 53380.000000 | 74460.000000 | 66470.0 | 50150.000000 | 64430.000000 | 43840.00 |
| 2 | Andorra | 539.947998 | 510.673004 | 560.340027 | 434.475006 | 254.0 | 450.151001 | 518.666016 | 456.6260 |
| 3 | Anguilla | 93.099998 | 100.730003 | 0.000000 | 0.000000 | 0.0 | 0.000000 | 68.190002 | 70.7300 |
| 4 | Antigua and Barbuda | 300.299988 | 374.500000 | 323.299988 | 279.200012 | 384.5 | 426.799988 | 249.600006 | 238.0000 |

Columns, Features, Variables, Fields

- Columns generally have the same type of data, but rows can be heterogeneous
- Types tell the computer how big something is and what operations it supports
- Types help us avoid errors caused by applying the wrong operations to the data

# Tidy Data Concept

**"Tidy Data" are also known in statistics as a *model matrix* or *data matrix*:**

- Standard method of displaying a multivariate set of data
- Rows correspond to sample individuals and columns to variables
- Entry in the $i^{th}$ row and $j^{th}$ column gives value of the $j^{th}$ variate as measured or observed on the $i^{th}$ individual

| | Country or Area | 1990 | 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Albania | 28385.000000 | 40311.000000 | 0.000000 | 0.000000 | 0.0 | 38284.000000 | 30683.000000 | 30491.000000 | 35883.000000 | 2789 |
| 1 | Algeria | 76160.000000 | 90270.000000 | 53380.000000 | 74460.000000 | 66470.0 | 50150.000000 | 64430.000000 | 43840.000000 | 37317.000000 | 0.00 |
| 2 | Andorra | 539.947998 | 510.673004 | 560.340027 | 434.475006 | 254.0 | 450.151001 | 518.666016 | 456.626007 | 565.559021 | 566. |
| 3 | Anguilla | 93.099998 | 100.730003 | 0.000000 | 0.000000 | 0.0 | 0.000000 | 68.190002 | 70.730003 | 68.190002 | 108. |

**More recently, Wickham defined "Tidy Data" as data sets that are arranged such that each variable is a column, and each observation (or case) is a row**

- Each column should be common type of measurement/feature
- Each row has all observations from one measurement/experiment
- Table should unite all observations for the features in common

*Wickham, H., 2014. Tidy data. Journal of Statistical Software, 59(10), pp.1-23.*

# Knowledge Check: Types of Data

**Categorical**          category names may be unrelated – no median or mean

**Boolean/binary**       two categories

**Ordinal**              Can order, so median makes sense.

**Interval**             Evenly spaced, so mean makes sense.  No zero.

**Time**                 Interval with (daily or seasonal) patterns

**Ratio**                Even spacing, well-defined zero

**Spatial**              Multidimensional ratio coordinates

# Data Structures and Types in Python



**Check-out DataCamp's Tutorial**

https://www.datacamp.com/community/tutorials/data-structures-python#files

# Imperfect Data

# Key Insight

**Rule #1**    **Garbage in,
garbage out**

**Rule #2**    **Quality data beats
fancy algorithms**

# What can go Wrong? Everything!

Each combination of
{table, records, fields}
$\times$
$\left\{\begin{array}{l}\text{names, inconsistent types,}\\\text{duplicates, missing, merged,}\\\text{disaggregated, pivoted, ...}\end{array}\right\}$

# Data Quality

# Validity

**Data-Type Constraints:** values in a column must be of same datatype (e.g., boolean, numeric, date, etc.)

**Range Constraints:** typically, numbers or dates should fall within a certain range

**Mandatory Constraints:** certain columns cannot be empty

**Unique Constraints:** a field, or a combination of fields, must be unique across a dataset (test for duplicates!)

**Set-Membership Constraints:** values of a column come from a set of discrete values (e.g., gender)

**Regular Expression Patterns:** text fields that have to be in a certain pattern (e.g., phone numbers)

**Cross-field Validation:** certain conditions that span across multiple fields must hold. For example, a patient's date of discharge from the hospital cannot be earlier than the date of admission.

**Foreign-key Constraints:** as in relational databases, a foreign key column can't have a value that does not exist in the referenced primary key.

# Accuracy

**The degree to which the data is close to the true values**

First defining all possible *valid* values helps easily spot invalid values

BUT, this does not mean they are *accurate*

A *valid* ROI is 22%, but this might not be *accurate* for your firm

**Note** that *accuracy* is not equivalent to *precision:* Saying that Horizon Hobby (a U.S. hobby firm that makes remote control cars) is a manufacturer of electric cars is true, but not very precise.

# Completeness

**The degree to which all required data are known**

**Missing data** is going to happen for various reasons
- → Error in collection process
- → Accidental deletion
- → No response
- → Not applicable
- → etc.

Mitigate the problem by collecting it again, imputing it, or assigning a meaning to it (e.g., not applicable)

# Consistency

**The degree to which the data is consistent, within the same data set or across multiple data sets**

Inconsistency occurs when two values in the data set contradict each other

A record in a customer database may indicate that a customer's lifetime value (CLV) is high even though they have not purchased anything in the past 5 years and their previous purchases were only clearance items with a total value of $99.

Similarly, a client's credit score may be 520 in one database and 789 in another.

# Uniformity

**The degree to which the data is specified using the same unit of measure**

Revenue may be recorded either in U.S. Dollars or in Euros

Dates might be in U.S. or in European formats

→ Data should always be converted to a common measure unit

# Data Cleaning Workflow

**Inspect**  Identify unexpected, incorrect, and inconsistent data

**Clean**  Resolve data quality violations, and fix or remove identified anomalies

**Verify**  Re-inspect to verify correctness

**Report**  Create report that details the changes made
and the quality of the currently stored data

# ANOMALY

# Anomaly!?



Fremont Bridge Bike Counter Time Series, Oct 2012 – Oct 2014
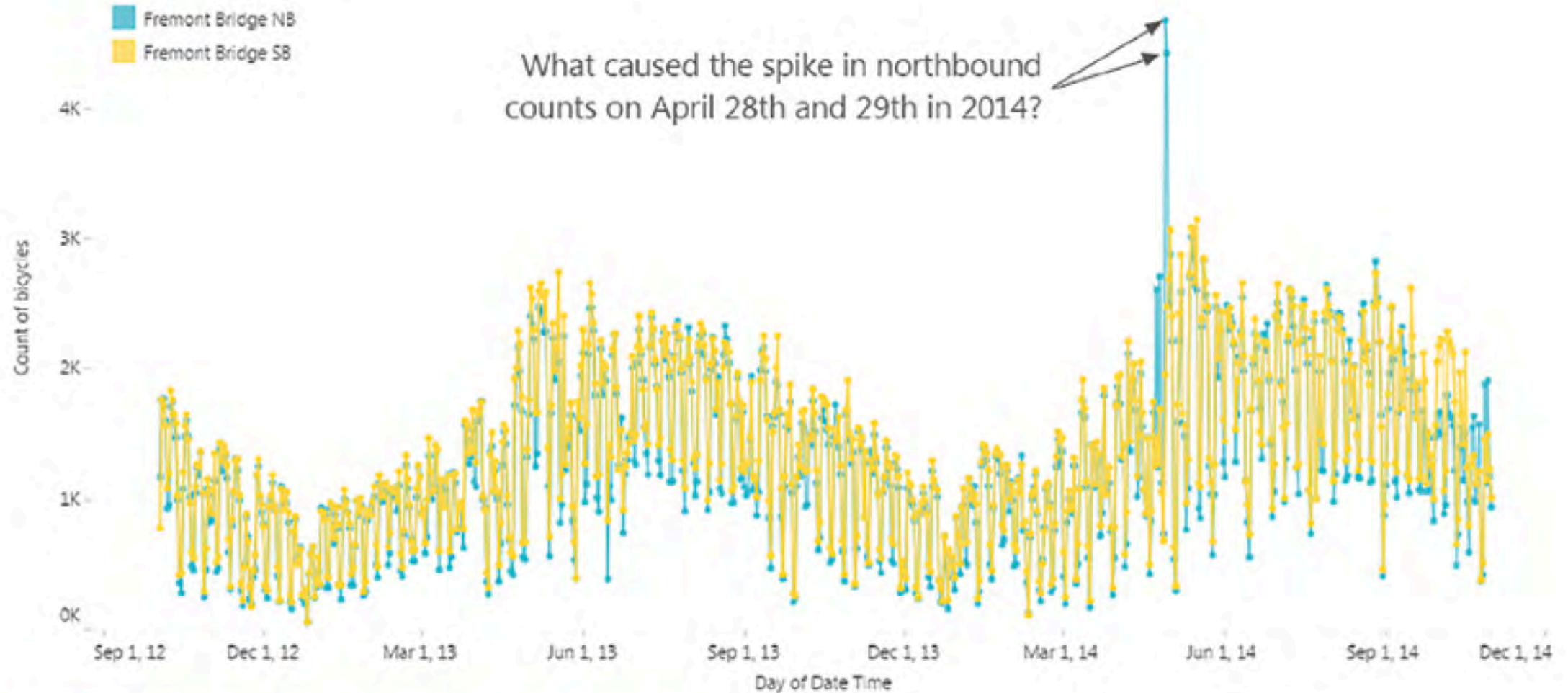
What caused the spike in northbound counts on April 28th and 29th in 2014?

Data source: http://www.seattle.gov/transportation/bikecounter_fremont.htm

# Anomalies in Data

**Definition of Anomaly**
1. Something different, abnormal, peculiar, or not easily classified
2. Deviation from the common rule **:** irregularity
Source: Merriam-Webster

**Anomaly detection** is a step in data mining that identifies data points, events, and/or observations that deviate from a dataset's normal behavior.

Anomalies can **indicate** critical incidents, errors in the data collection, or potential business opportunities (e.g., a change in consumer behavior).

Download the Python Notebook for Anomaly Detection from CANVAS and put it on your Google Drive

# Anomalies: Novel vs. Outlier

> *Humans are relatively good at finding patterns,*
> *and they're also quite good at finding things that don't fit a pattern.*

**Novelty detection** is the mechanism by which an intelligent organism can identify an incoming sensory pattern as being hitherto unknown.

- Consider $N$ observations from the same distribution described by $p$ features
- Let's add one more observation
- Is the new observation so different from the others that we can doubt it is regular?
  (i.e., does it come from the same distribution?)

**Outlier detection** is the identification of rare items, events or observations which raise suspicions by differing significantly from most of the data.

- Similar to novelty detection
- Goal is to separate a core of regular observations from some polluting ones
- BUT: we don't have a clean data set representing the population of regular observations that can be used as basis

# Anomalies in Financial Markets



1. What is going on?
2. Novelty, Outlier, neither or both?
3. Opportunity, Error, something else?
5. What to do with it?

# Outlier Detection

***Threats and Opportunities***

**Outliers might be**

- items that are so far outside the norm that they need not be considered
- errors in the data
- novel items that are worth exploring

**Outliers can**

- alert you to a problem or and unknown opportunity
- cause potentially severe errors in your models that lead to incorrect conclusions

**Important to detect all outliers to**

- Analyze them to know why you had them there in the first place
- Eliminate them if appropriate

Based on Will Badr's 2019 post: "5 Ways to Detect Outliers/Anomalies That Every Data Scientist Should Know (Python Code)"

# Outlier Detection: Visual Inspection

## Easy

[2.99,1.99,4.49,6.99,59.90,0.99,3.29,4.89,5.79]

## Obvious



## Harder

| | |
|---|---|
| 1.00 | 0.95 |
| 1.75 | 1.68 |
| 2.17 | 2.08 |
| 2.97 | 2.95 |
| 3.51 | 3.51 |
| 4.06 | 3.98 |
| 4.35 | 4.26 |
| 4.99 | 4.95 |
| 5.63 | 5.60 |
| 6.57 | 6.48 |
| 7.32 | 7.23 |
| 8.01 | 8.00 |
| 5.16 | 9.90 |
| 5.31 | 5.26 |
| 6.26 | 6.18 |
| 6.81 | 6.72 |
| 7.31 | 7.27 |
| 7.78 | 7.69 |
| 8.06 | 8.02 |
| 8.78 | 8.77 |
| 9.03 | 8.95 |
| 9.42 | 9.39 |
| 9.79 | 9.75 |
| 10.57 | 10.52 |
| 10.95 | 10.95 |

## Obvious

# Outlier Detection: Standard Deviation

**Normal distribution:**
About 68% of the data values lie within one standard deviation of the mean
About 95% are within two standard deviations
About 99.7% lie within three standard deviations

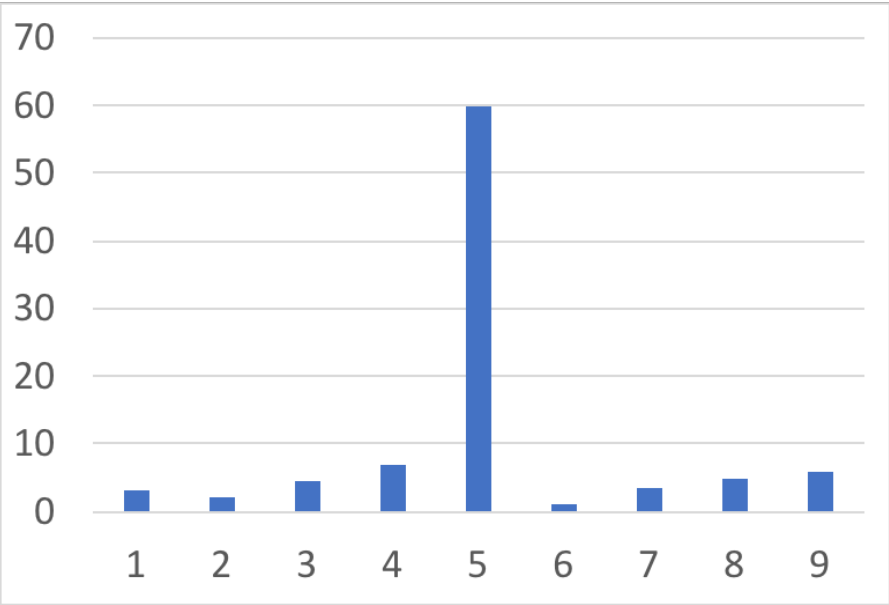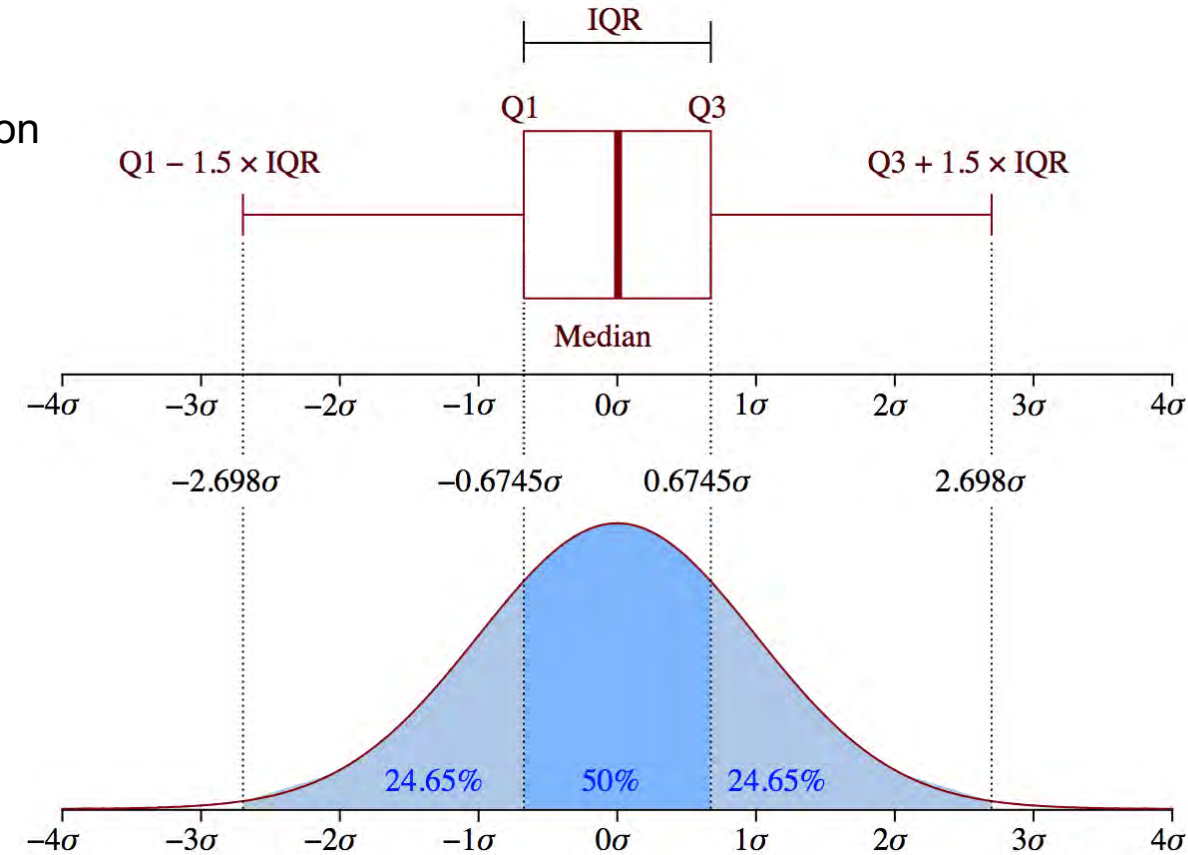Any data point that is more than 3 times the standard deviation might be an outlier candidate

*Based on Will Badr's 2019 post: "5 Ways to Detect Outliers/Anomalies That Every Data Scientist Should Know (Python Code)" and Wikipedia.org*

# Outlier Detection: Boxplots

- Graphical depiction of numerical data through their quantiles
- Lower and upper whiskers as the boundaries of the data distribution
- Data points outside of the whiskers can be considered outliers



*Based on Will Badr's 2019 post: "5 Ways to Detect Outliers/Anomalies That Every Data Scientist Should Know (Python Code)" and Wikipedia.org*
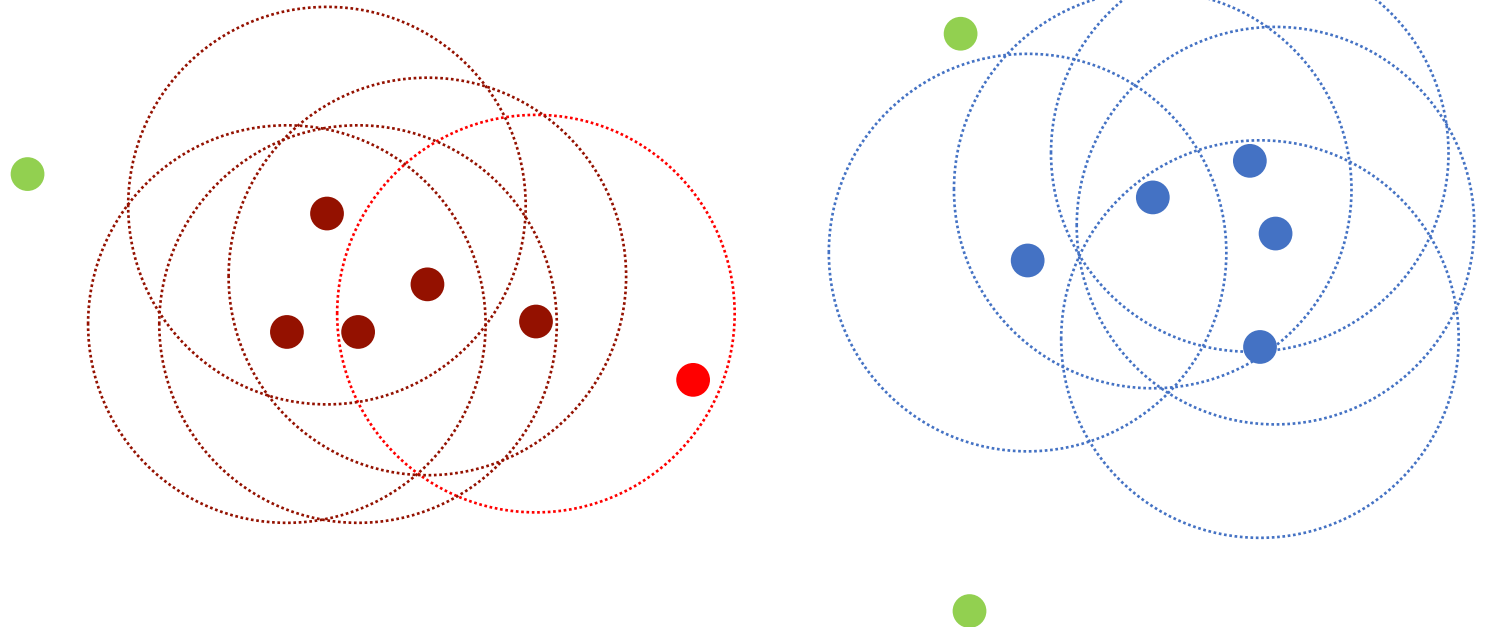
# Outlier Detection: Clustering

- **DBScan** can be used as density-based anomaly detection method

- Applicable to single or multi-dimensional data

- **Identifies different "Points"** (i.e., samples, observations, objects, items)
- 🔴🔵 **Core Points** are central samples (i.e., observations, objects, items) in a cluster
  - Set the min number to form a cluster using the *hyperparameter* **min_samples**
  - Set the maximum allowable distance between two samples to still be considered as being in same cluster with **eps** hyperparameter)
  - 🔴 **Border Points** are in same cluster as core points but further away from its center
  - 🟢 **Noise Points** are identified as not belonging to any cluster

The downside with this method is that the higher the dimension, the less accurate it becomes.

You also need to make a few assumptions like estimating the right value for eps which can be challenging.

Note that k-means and hierarchal clustering can also be used to detect outliers

# Working with Data in Cloud Computing (Google CoLab)

*Colaboratory, or "**Colab**" for short, **is** a product from **Google** Research. **CoLab** allows anybody to write and execute arbitrary python code through the browser, and **is** especially well suited to machine learning, data analysis and education.*

**CO** **Create, open, edit, run, save python (jupyter) notebooks**

**CO** **Upload Files and Datasets to Google Drive**

**CO** **Export Data/Files to Google Drive**

**CO** **Runtimes – use GPUs or even TPUs**

**CO** **Clone Repositories from Github**
Copy the clone link of the Github repository to CoLab

**CO** **Terminal Commands**
!pip install **library_name**

**CO** **Collaborate: Share your notebook**

> **This course requires YOU to use CoLab**
> - Follow instructions on Canvas: **Google CoLab**
> - Complete CoLab **Tutorial** <u>before</u> Class 02 (CANVAS>Files>CoLab_Tutorial)

Go to DataCamp:
www.datacamp.com

Instructions for DataCamp:
https://kenan-flagler.instructure.com/courses/3629908/pages/datacamp

# Looking Ahead

**Next Class**: Tuesday, January 17th, 2023

*Missing Data and Data Pitfalls*

**DataCamp Homework 1 due!**
- Introduction to Python
  *(approximately 4 hours)*
- *Due on January 17th by 11:59pm*

**Readings:**
- none