

# Machine Learning: Data Foundation

## Day 1

# Objective

- Get to know each other
- Level set
- Establish the overview of the course



# Who am I?

# Who are you?

Name, Role at Salesforce, and Machine Learning Experience

"A change in perspective is worth 80 IQ points."

*Alan Kay*

# What is Machine Learning?

- It's the rocketship by which we travel to Planet AI
  - By the way, what's the difference between Machine Learning and AI?
  - Oh, and what's fueling that rocketship?
- automating automation
- getting computers to program themselves
- letting the data do the work
- how is ML different from traditional software development?
  - computers produce output from what input?
  - ML is the reverse: data + output = programs
  - ML is also...

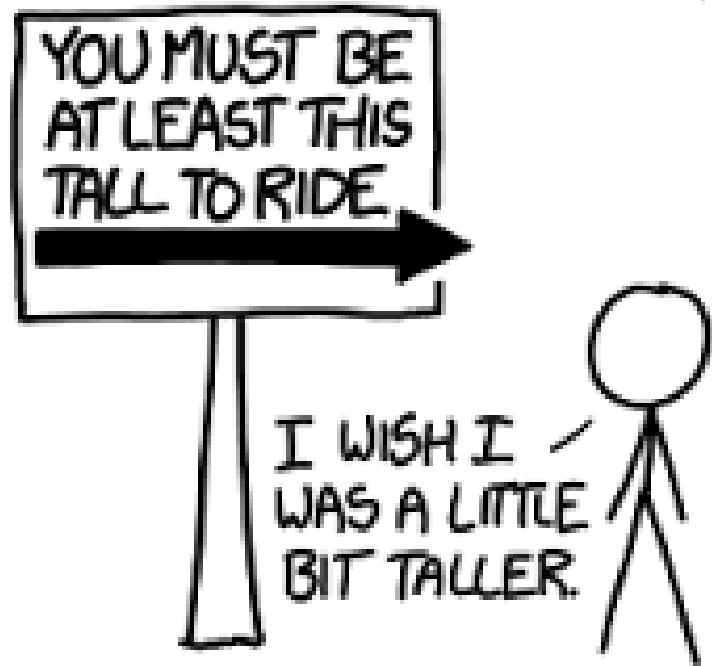
# The Rockstar Job of the 21st Century...



CC BY 2.0, <https://www.flickr.com/photos/onepointfour/12314110693>

# ...Requires a Lot of Sweeping!





Machine Learning will NOT succeed without...

- understanding our data
- cleaning up our data
- visualizing our data

# Isn't This Someone Else's Problem?

"Poor data quality is enemy number one to the widespread, profitable use of machine learning."

"The quality demands of machine learning are steep, and bad data can rear its ugly head twice — first in the historical data used to train the predictive model and second in the new data used by that model to make future decisions."

"...today, most data fails to meet basic “data are right” standards. Reasons range from data creators not understanding what is expected, to poorly calibrated measurement gear, to overly complex processes, to human error. To compensate, data scientists cleanse the data before training the predictive model. It is time-consuming, tedious work (taking up to 80% of data scientists' time), and it's the problem data scientists complain about most."

# A Day in the Life of salesforce



# Five Elements of Trust



Security



Availability



Scalability



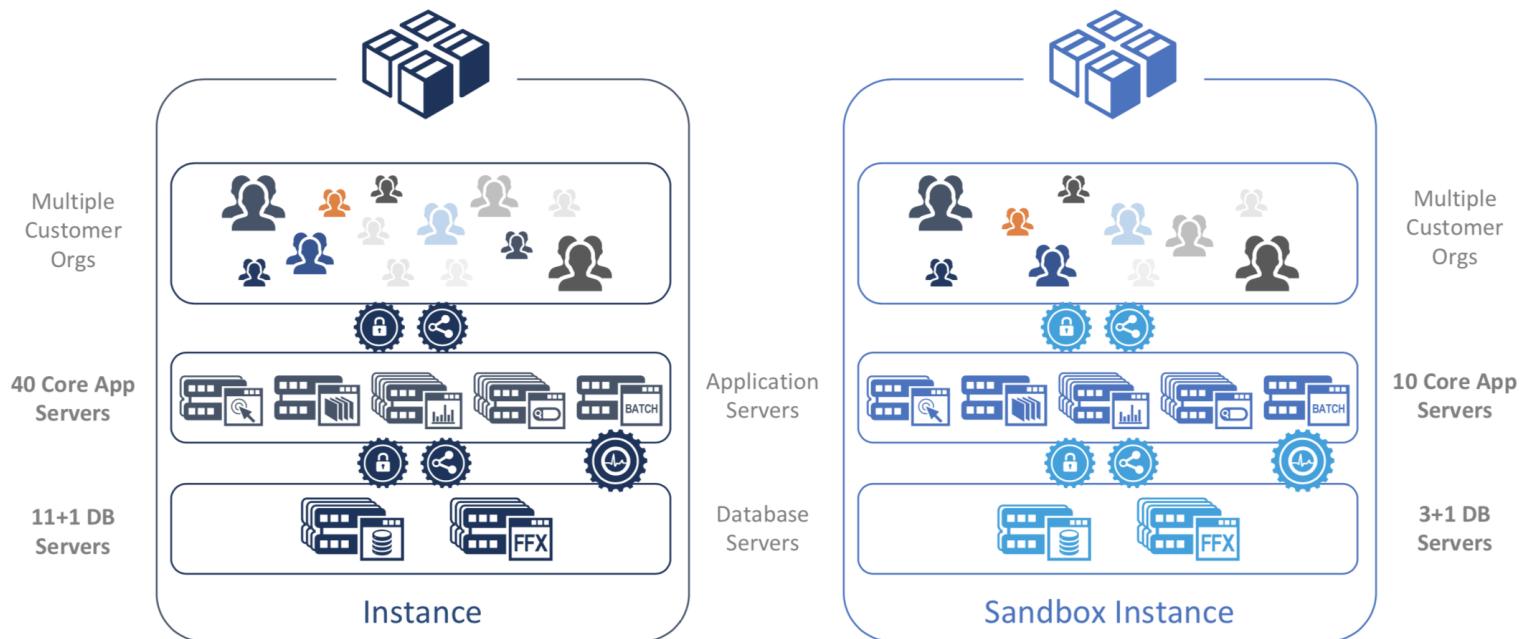
Multi-Tenant



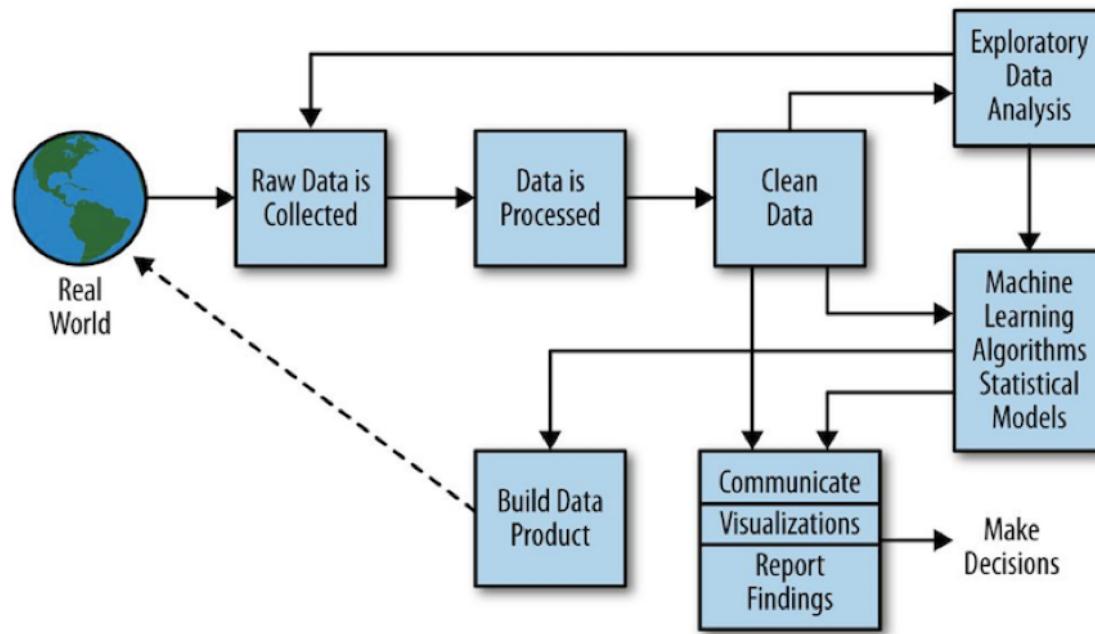
Continuous  
Innovation

What information from all of the data that we collect from our systems will alert us that things are about to go sideways?

# Core and Sandbox Instances



# Our Roadmap: Data Science Pipeline

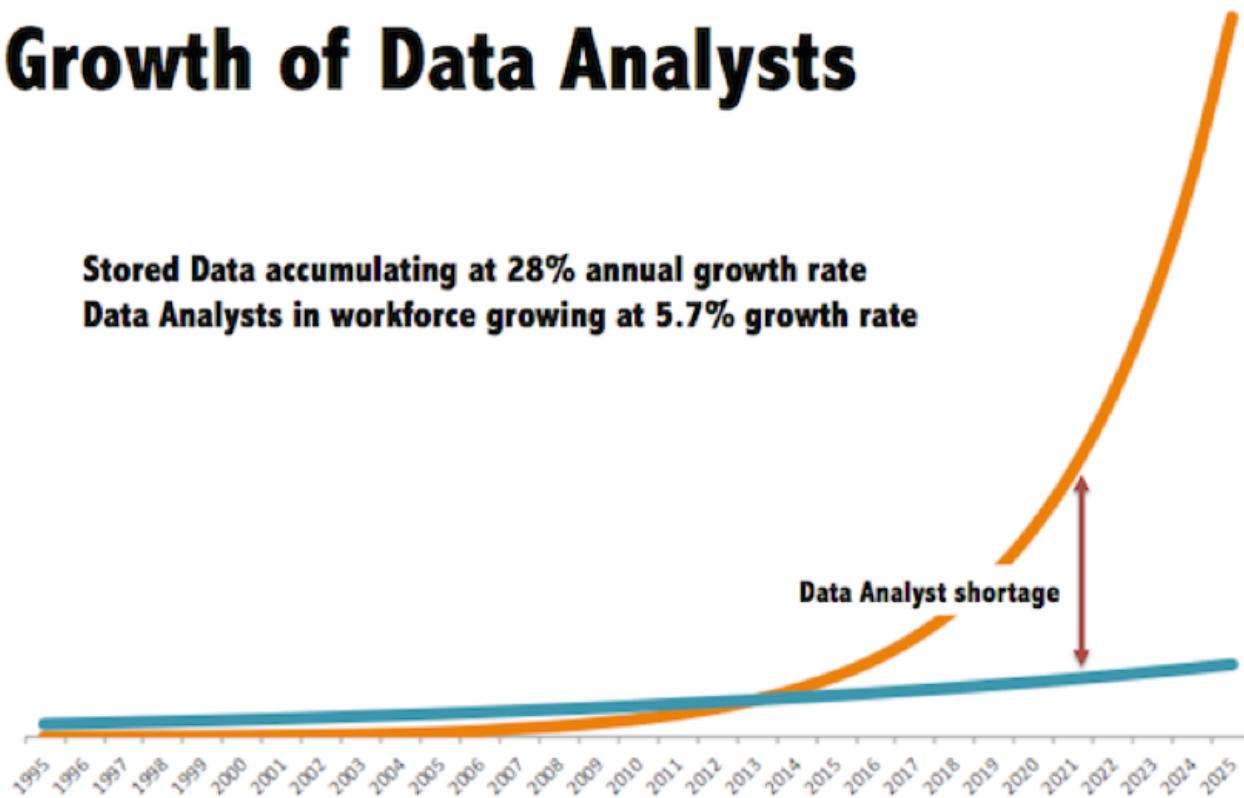


# Data and Data Processing

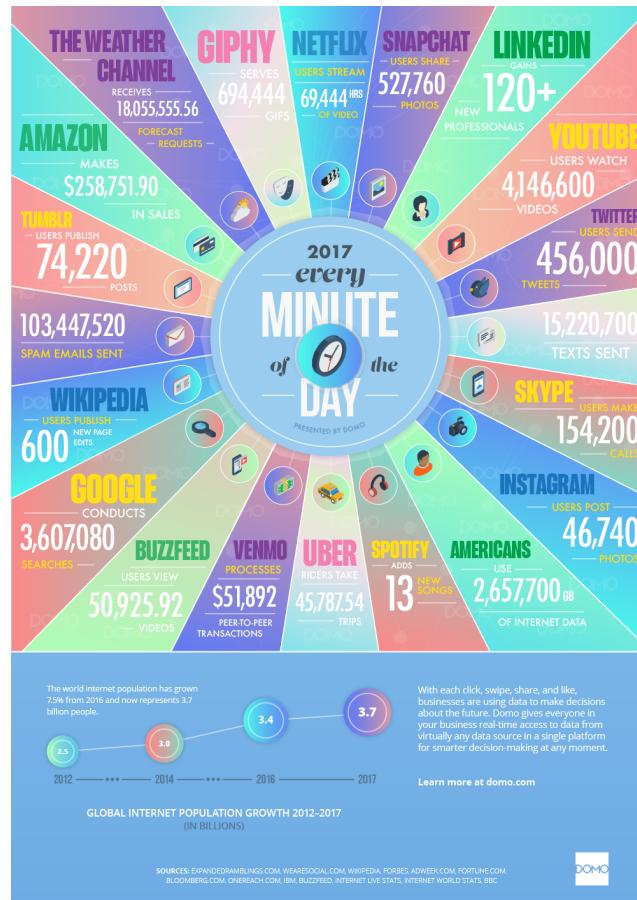
# Objective

- Connect Data Science with other Information Technology (IT) activities
- Explore how we model our domains
- Understand some of the obstacles we face

# Growth of Data vs. Growth of Data Analysts



# Humans Generate a LOT of Data!



<http://bit.ly/2f12JCT>

# Demand for Data Scientists Will Soar 28% by 2020

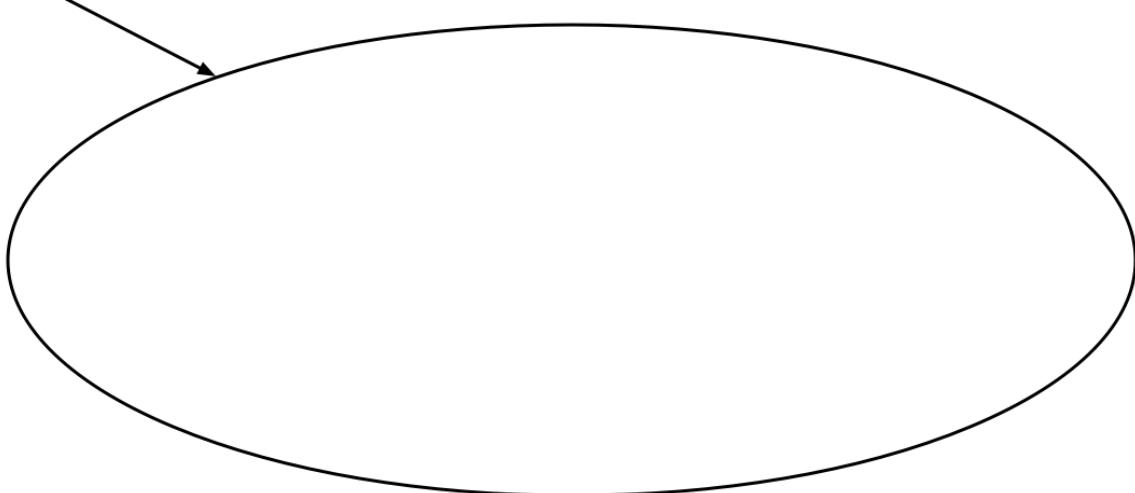
"Annual demand for the fast-growing new roles of data scientist, data developers, and data engineers will reach nearly 700,000 openings by 2020."

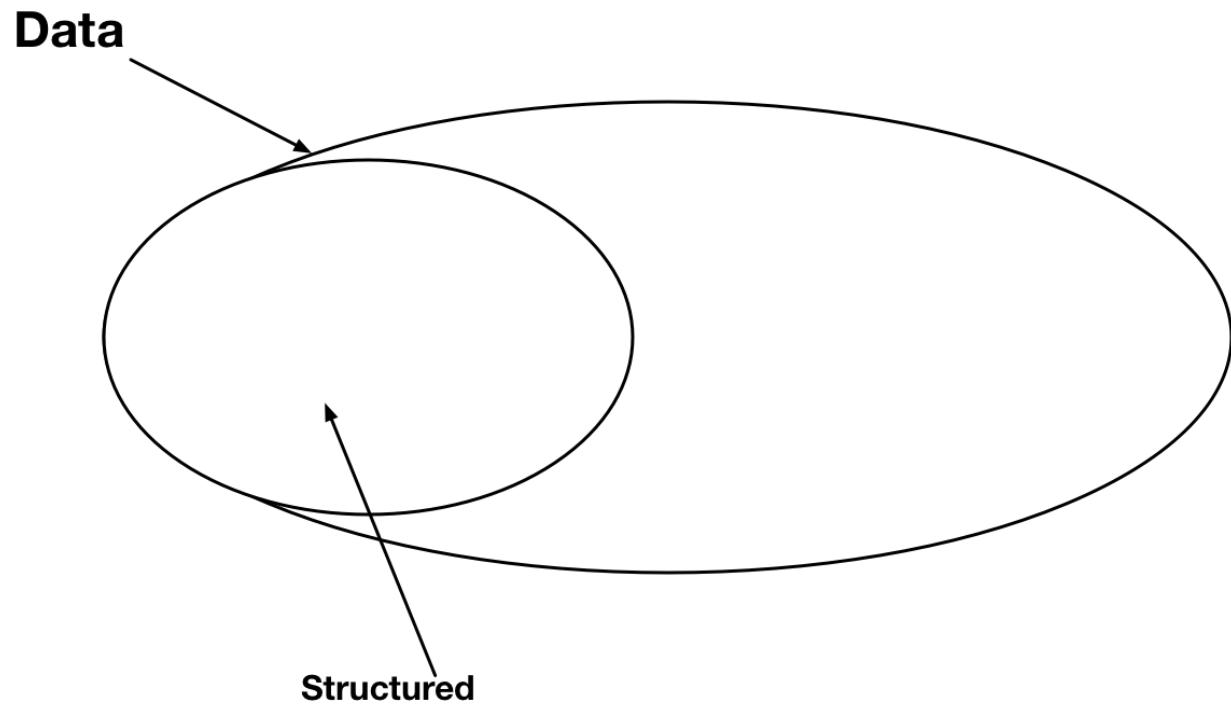
"By 2020, the number of jobs for all US data professionals will increase by 364,000 openings to 2,720,000 according to IBM."

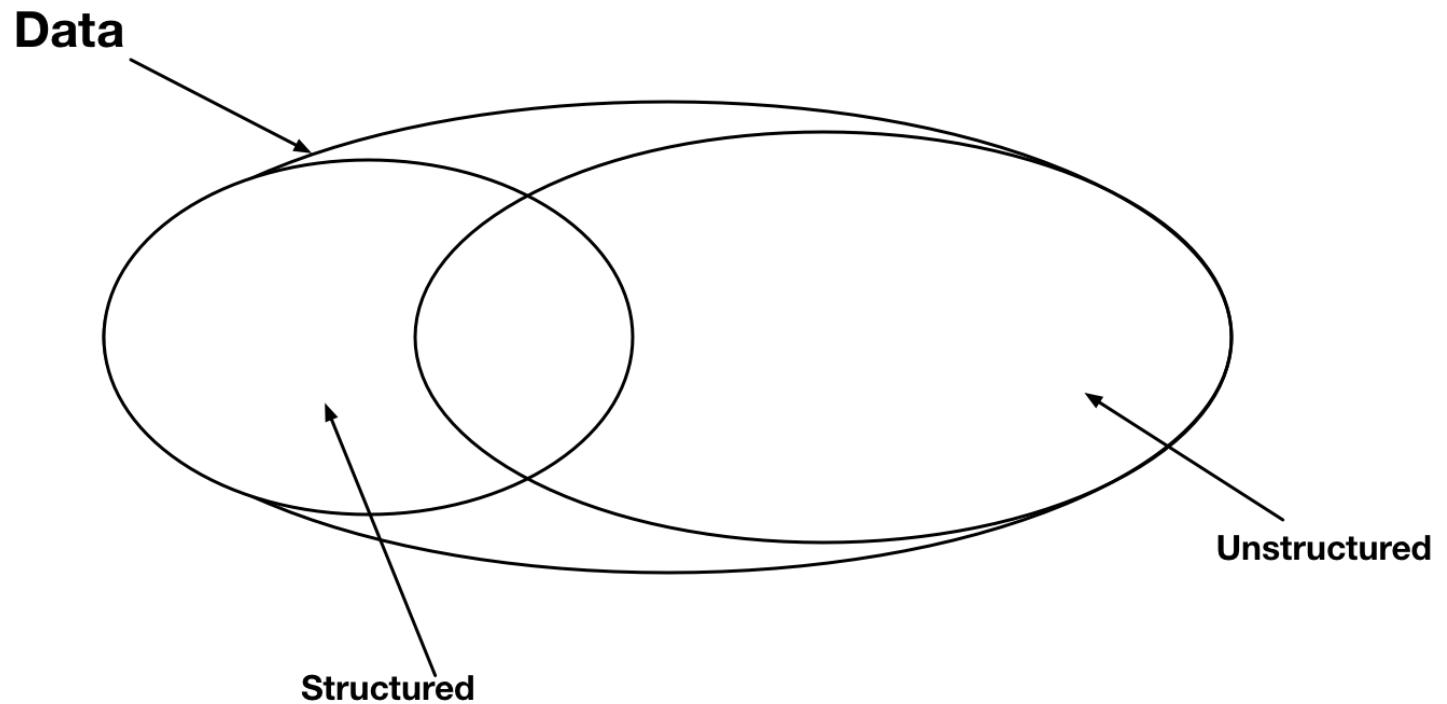
"Machine learning, big data, and data science skills are the most challenging to recruit for and potentially can create the greatest disruption to ongoing product development and go-to-market strategies if not filled."

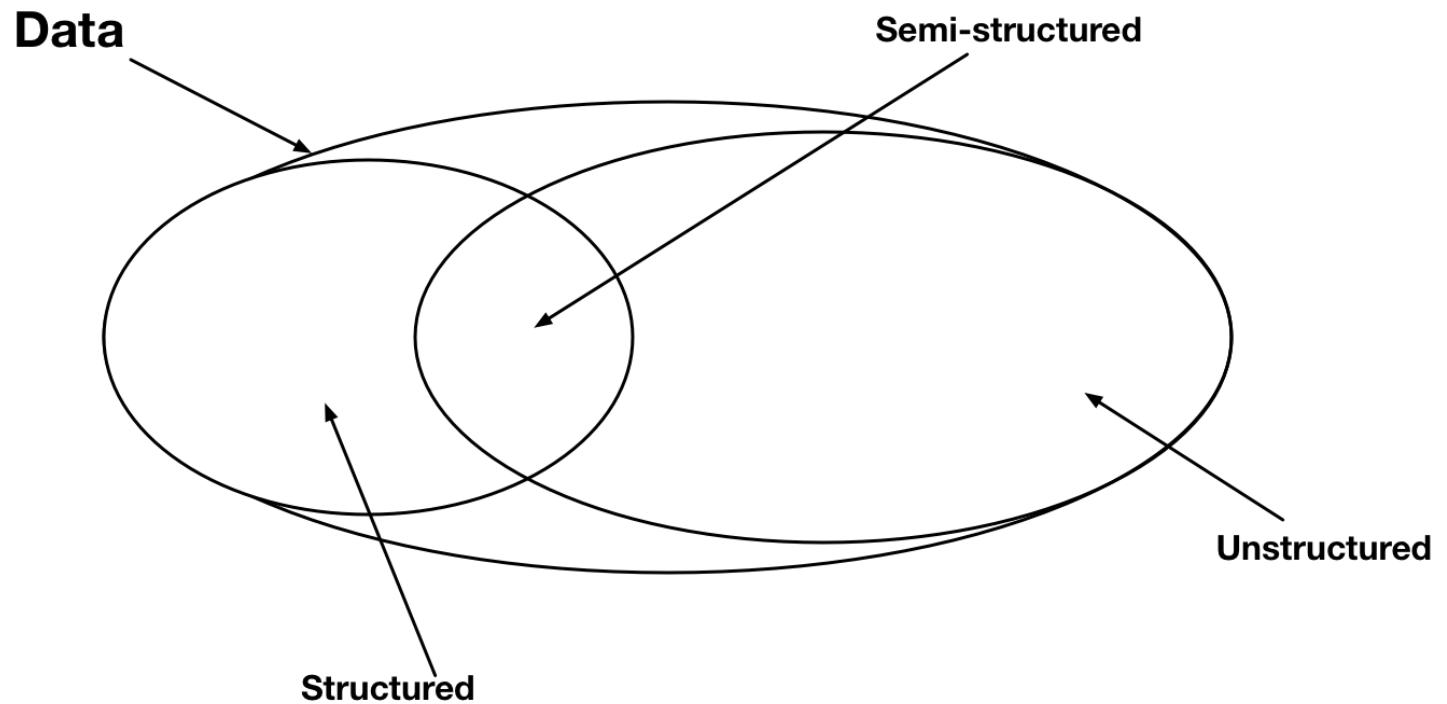
# What is Data?

**Data**

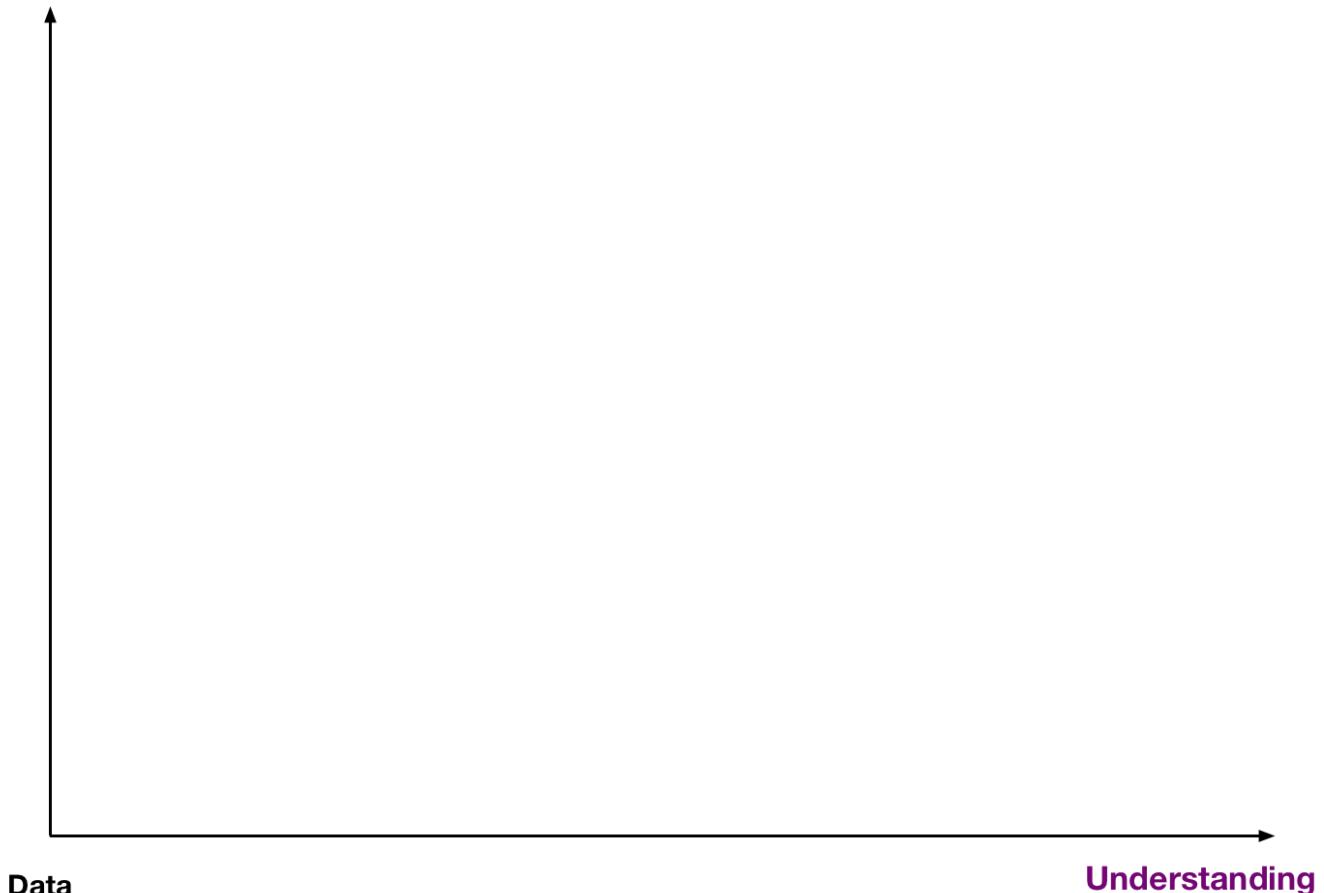




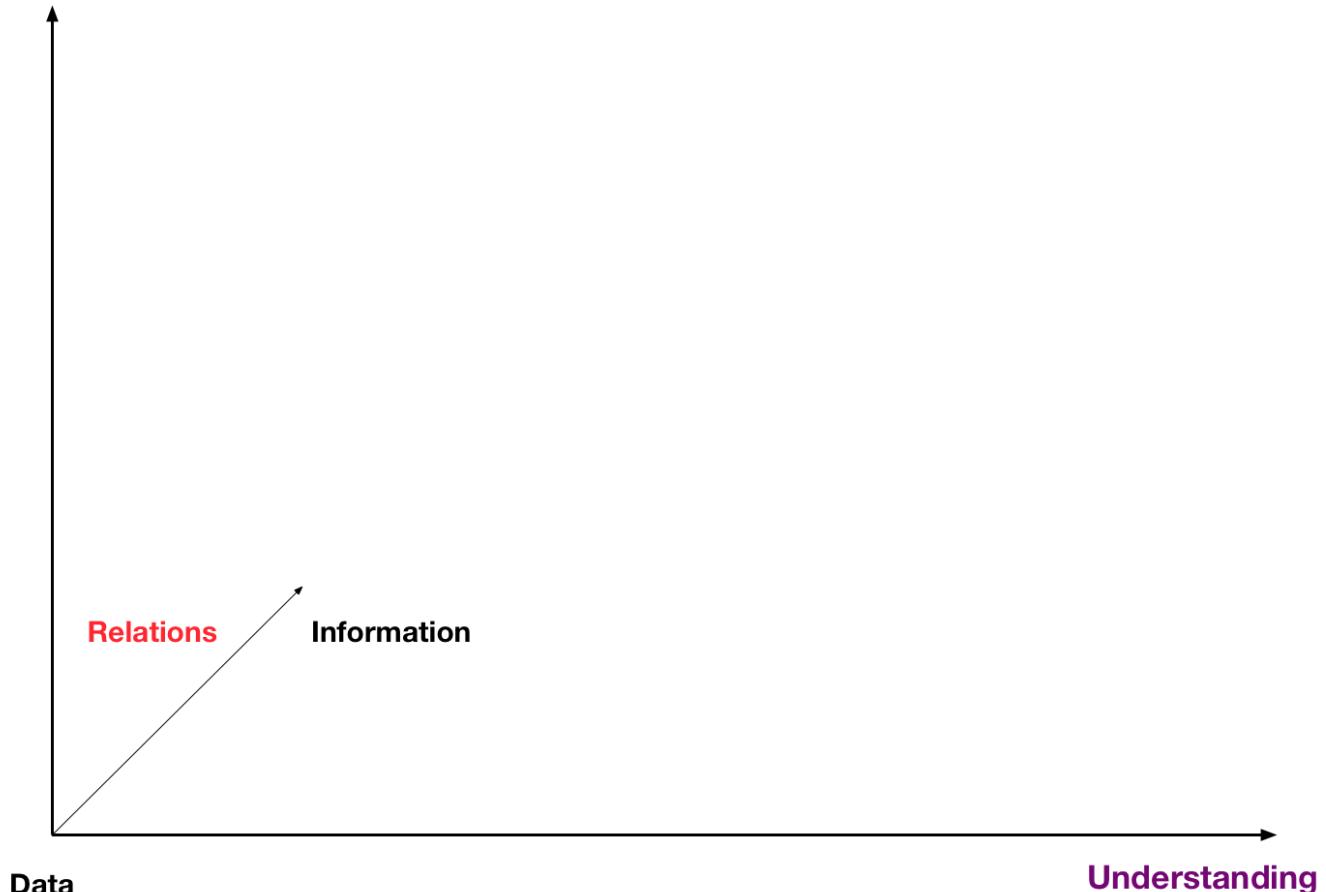




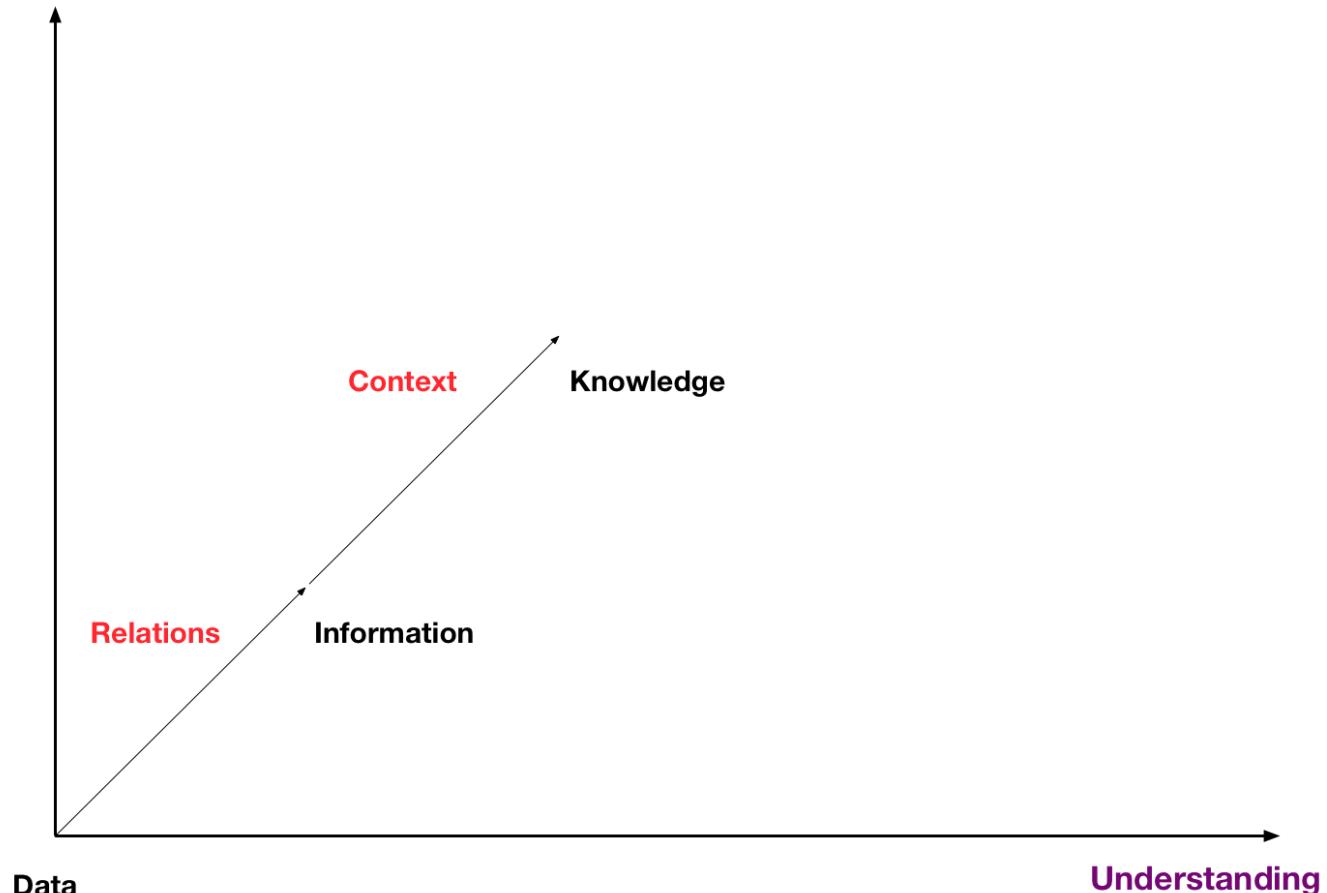
**Connectedness**



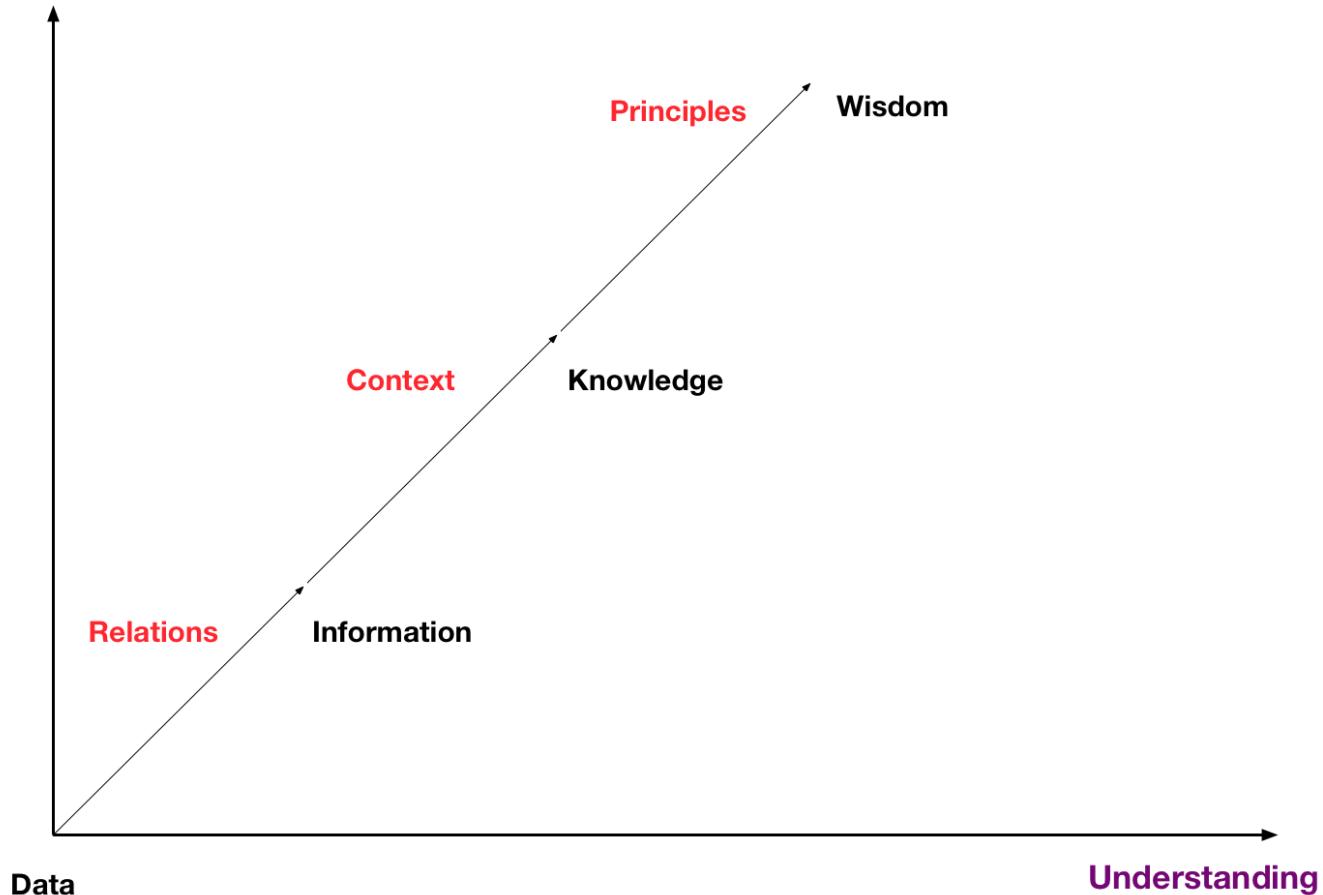
**Connectedness**



**Connectedness**



**Connectedness**



# What is the Problem?

# Wasting Time

- Data professionals spend 60% of their time "getting to insight"
  - 37% of this time is spent searching for data
  - 36% of this time is spent preparing data
  - Only 27% of that time is actual analysis
- 50% of the time is spent searching for or replicating existing data
- 30-50% of organizations are not where they want to be
- Costs to organizations annually
  - \$1.7 million/100 employees in U.S.
  - €1.1 million/100 employees in Europe

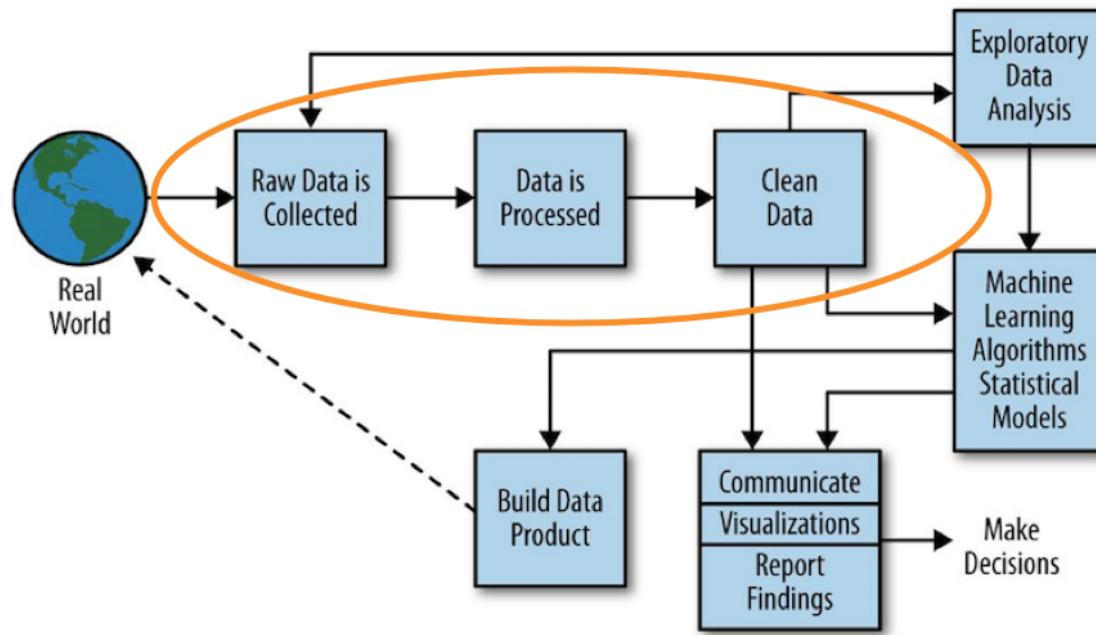
"It is evident that many professionals are not aware of what resources are available within data assets like data lakes, how to access the data, where it came from, or how to glean trusted insights..."

<https://tinyurl.com/ybwdq34u>

"Unless organizations make changes to their infrastructure now, and close the gaps on data discovery, integrity and cataloging, processes will only become more inefficient as data volume and variety continues to grow."

<https://tinyurl.com/ybwdq34u>

# Data Science Pipeline





# Metadata

- Data about data
- Types of metadata
  - Business metadata
  - Technical metadata
  - Media metadata
  - Semantic metadata



# Business Metadata

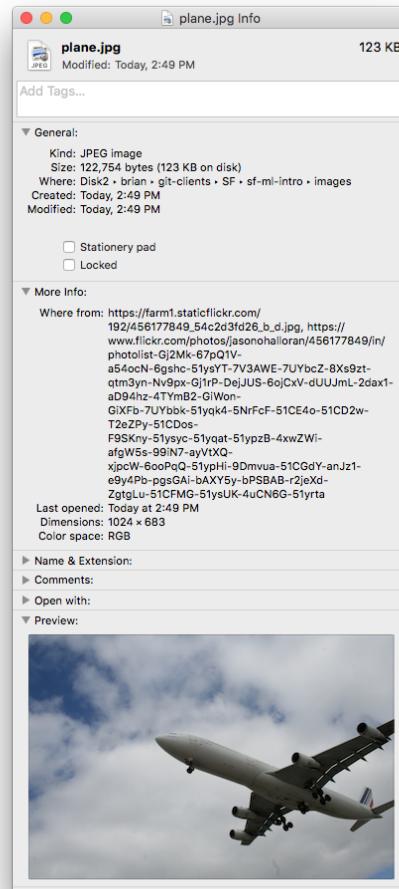
Term	Description
Airline	A company that owns one or more airplanes that are used to transport passengers from one airport to another.
Airport	A transportation destination at which airlines have routes to transport passengers.
Airplane	A vehicle that flies and lands at airports to pick up passengers.
Routes	A connection between two airports that an airline flies its passengers.

# Technical Metadata

```
sqlite> .tables
airlines  airports  routes
```

```
sqlite> .schema airports
CREATE TABLE airports (
[index] INTEGER,
[id] TEXT,
[name] TEXT,
[city] TEXT,
[country] TEXT,
[code] TEXT,
[icao] TEXT,
[latitude] TEXT,
[longitude] TEXT,
[altitude] TEXT,
[offset] TEXT,
[dst] TEXT,
[timezone] TEXT
);
CREATE INDEX ix_airports_index ON airports ([index]);
```

# Media Metadata



# Semantic Metadata

- One approach is JSON-LD (JSON + Linked Data)
- JSON-LD organizes and connects messy and disconnected data

```
{ "@context": "http://schema.org",
  "@type": "Organization",
  "url": "http://united.com",
  "contactPoint":
  [
    {
      "@type": "ContactPoint",
      "telephone": "+1-800-864-8331",
      "contactType": "customer service"
    }
  ]
}
```

```
@prefix schema: <http://schema.org/> .

[ a schema:Organization ;
  schema:contactPoint [ a schema:ContactPoint ;
    schema:contactType "customer service" ;
    schema:telephone "+1-800-864-8331"
  ] ;
  schema:url <http://united.com>
] .
```

<https://json-ld.org/>

<http://schema.org/Organization>

# Example of Semi-Structured

- Semantic Metadata added to an Unstructured HTML Document

```
<html>
  <head>
    <script type="application/ld+json">
      {"@context": "http://schema.org",
       "@type": "Organization",
       "url": "http://united.com",
       "contactPoint":
         [
           {
             "@type": "ContactPoint",
             "telephone": "+1-800-864-8331",
             "contactType": "customer service"
           }
         ]
       }
     </script>
   </head>
...
</html>
```

<https://developers.google.com/gmail/markup/>

<https://developers.google.com/search/docs/data-types/events>

# Data Sources

# Objective

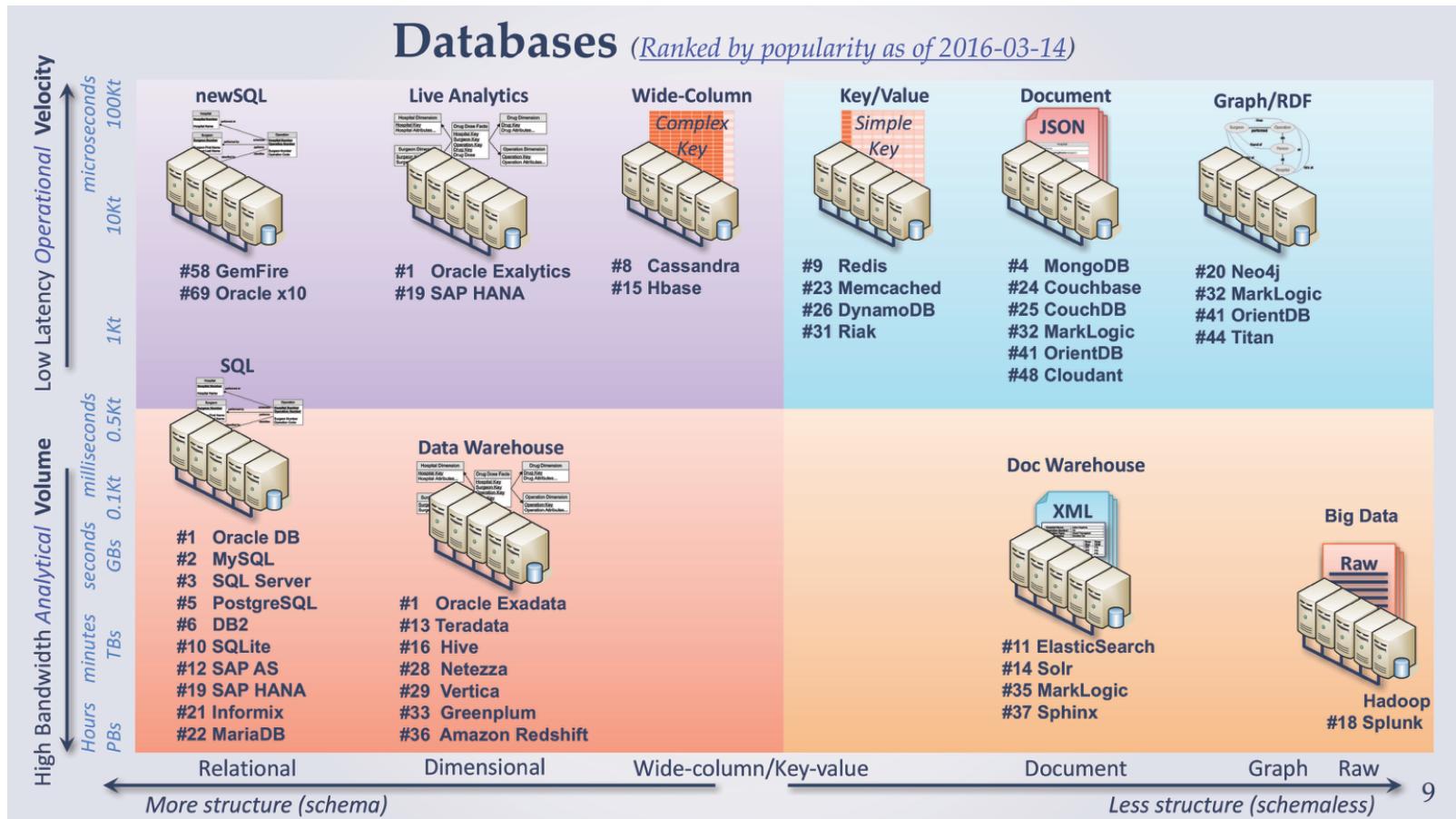
- Consider where data comes from
- Understand some of the integration issues we will face
- Consider our data storage options
- Read data into Python in a convenient format

What did you do Saturday night?

# Data Modeling Choices

- Relational
- Key-Value
- Column
- Document
- Graph

# Databases (Ranked by popularity as of 2016-03-14)



# RDBMS

- Established technology, well-understood
- Works best for stable domains

# Spreadsheets

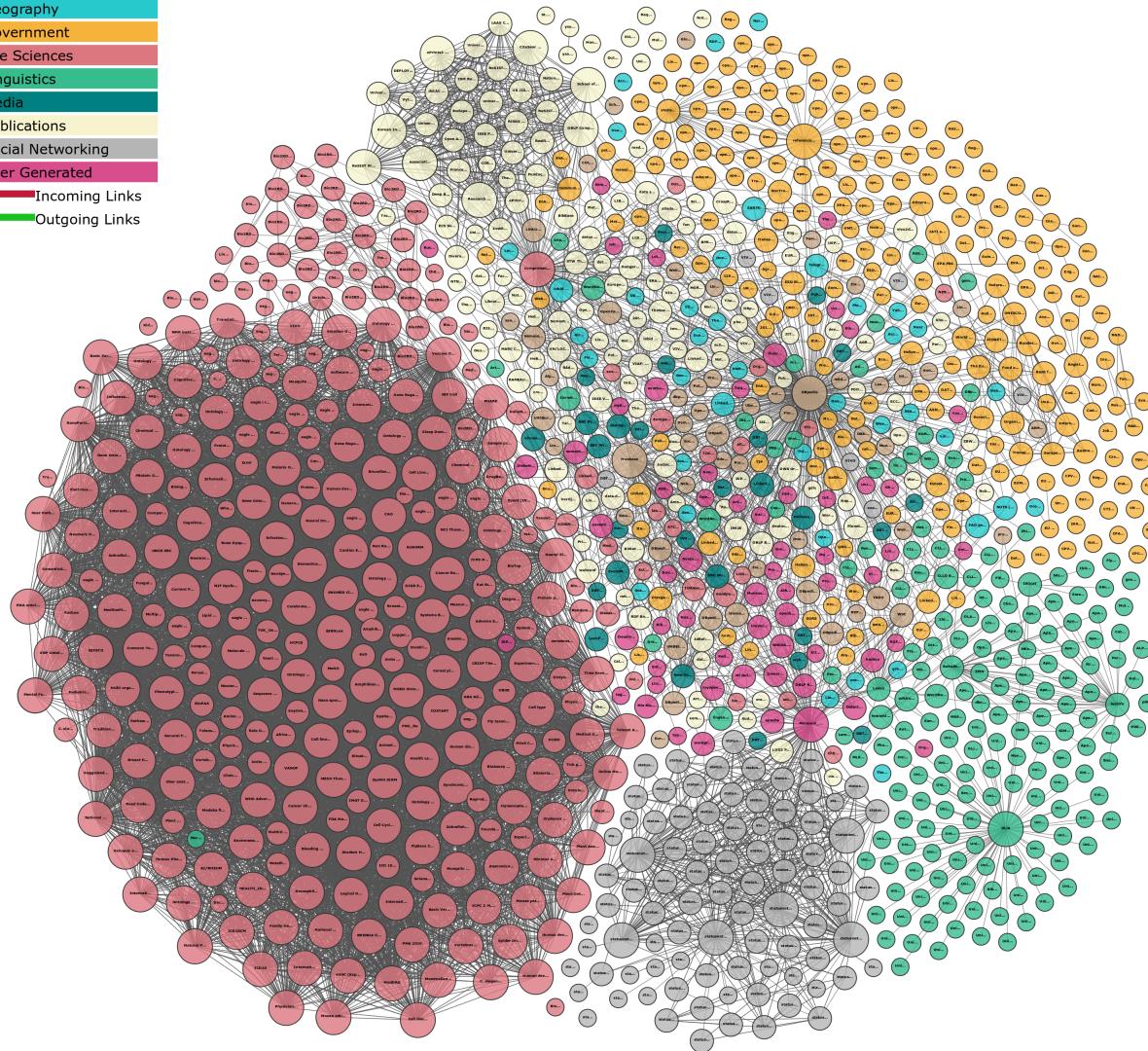
- Conceptually aligned
- Easy tooling
- Familiarity
- Rapid experimentation
- Customization

# NoSQL

- Volume
- Velocity
- Variety
- Veracity

Legend

Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networking
User Generated
Incoming Links
Outgoing Links



# Setting up a Data Science Environment

# Let's Install Some Essential Tools!

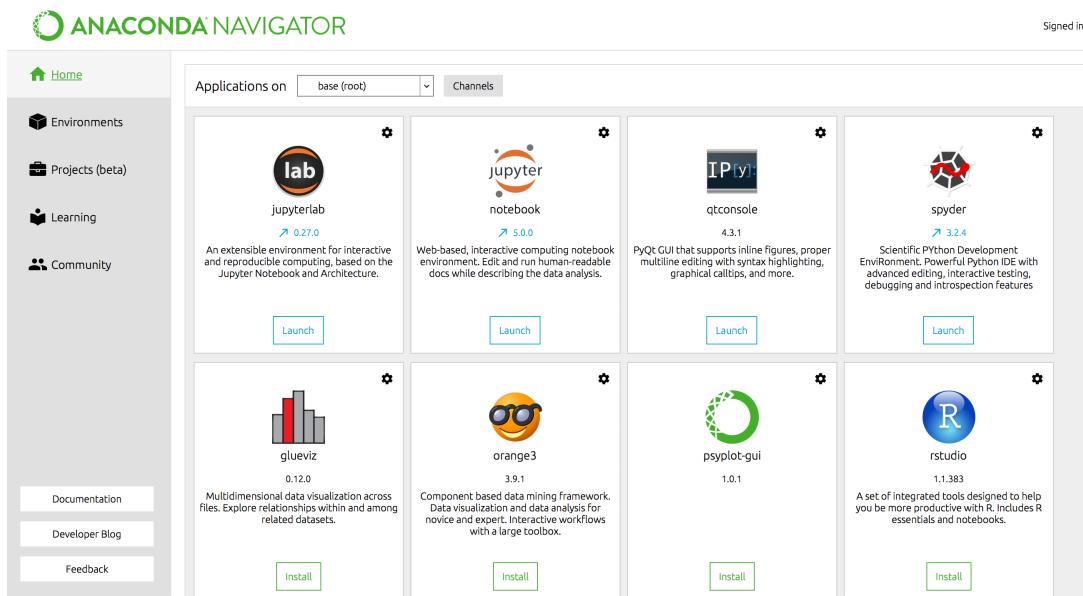
- Remember that one of our goals is to gain some familiarity with commonly-used data science tools
- We're going to use a tool called Jupyter notebooks to do our hands-on work
- Let's go...

# Step 1: Install Anaconda

- Anaconda is an open source distribution of Python which includes the essential data science tools we want
- Go to <https://www.anaconda.com/download>
- Make sure you download the Python 3 version
- At the **Installation Type** phase:
  - Click **Change Install Location** and choose **Install for me only**
  - Choose **Customize** at the **Installation Type** phase, and disable the **Modify PATH** option

# Step 2: Launch the Anaconda Navigator

- Should have been installed in a directory named `anaconda3` underneath your home directory

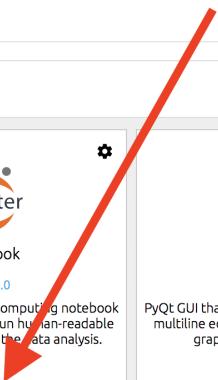


# Step 3: Launch Jupyter

 ANACONDA NAVIGATOR

Signed in as

Applications on base (root) Channels



 jupyterlab 0.27.0 <p>An extensible environment for interactive and reproducible computing, based on the Jupyter Notebook and Architecture.</p> <a href="#">Launch</a>	 notebook 5.0.0 <p>Web-based, interactive computing notebook environment. Edit and run human-readable docs while describing the data analysis.</p> <a href="#">Launch</a>	 qtconsole 4.3.1 <p>PyQt GUI that supports inline figures, proper multiline editing with syntax highlighting, graphical calltips, and more.</p> <a href="#">Launch</a>	 spyder 3.2.4 <p>Scientific Python Development EnvIRonment. Powerful Python IDE with advanced editing, interactive testing, debugging and introspection features</p> <a href="#">Launch</a>
 glueviz 0.12.0 <p>Multidimensional data visualization across files. Explore relationships within and among related datasets.</p> <a href="#">Install</a>	 orange3 3.9.1 <p>Component based data mining framework. Data visualization and data analysis for novice and expert. Interactive workflows with a large toolbox.</p> <a href="#">Install</a>	 psyplot-gui 1.0.1 <a href="#">Install</a>	 rstudio 1.1.383 <p>A set of integrated tools designed to help you be more productive with R. Includes R essentials and notebooks.</p> <a href="#">Install</a>

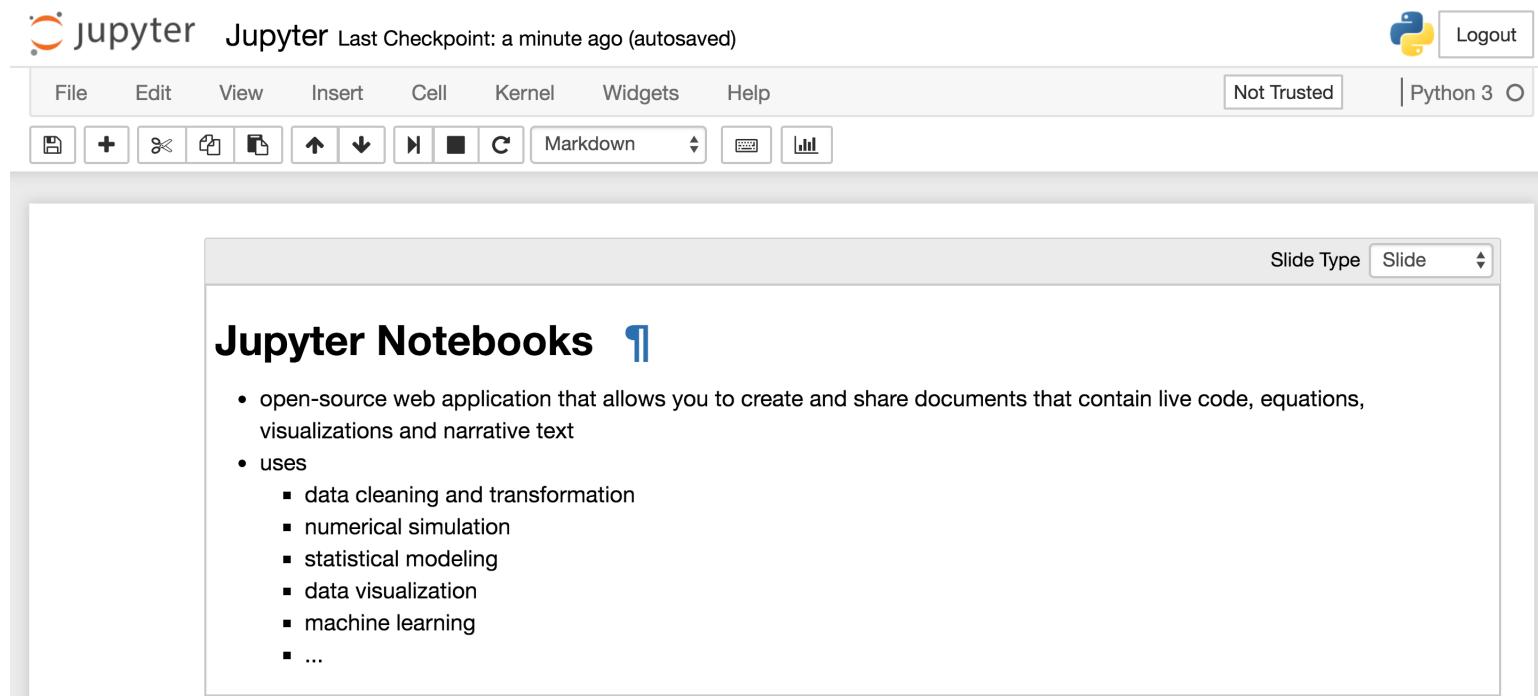
# Step 4: Copy/Paste URL into Your Browser

```
[I 16:25:50.469 NotebookApp] JupyterLab alpha preview extension loaded from /Users/dws  
/anaconda3/lib/python3.6/site-packages/jupyterlab  
JupyterLab v0.27.0  
Known labextensions:  
[I 16:25:50.472 NotebookApp] Running the core application with no additional extension  
s or settings  
[I 16:25:50.478 NotebookApp] Serving notebooks from local directory: /Users/dws  
[I 16:25:50.478 NotebookApp] 0 active kernels  
[I 16:25:50.478 NotebookApp] The Jupyter Notebook is running at: http://localhost:8888  
/?token=39a5ad5e45dbd4d6a7e78287693858ba145ea5dbaaf65856  
[I 16:25:50.478 NotebookApp] Use Control-C to stop this server and shut down all kerne  
ls (twice to skip confirmation)  
[C 16:25:50.478 NotebookApp]
```

Copy/paste this URL into your browser when you connect for the first time,  
to login with a token:

`http://localhost:8888/?token=39a5ad5e45dbd4d6a7e78287693858ba145ea5dbaaf65856`  
0:97: execution error: "http://localhost:8888/tree?token=f29cdf589b571450f2553e30ea03  
b6d5947fbe8a511170e" doesn't understand the "open location" message. (-1708)

# Step 5: Navigate to the file `Demo - Jupyter.ipynb`



The screenshot shows a Jupyter Notebook interface. At the top, there's a header bar with the Jupyter logo, the word "Jupyter", and a message "Last Checkpoint: a minute ago (autosaved)". On the right of the header are a Python 3 logo, a "Logout" button, and status indicators "Not Trusted" and "Python 3". Below the header is a toolbar with various icons for file operations like opening, saving, and creating new notebooks, along with buttons for kernel selection and cell execution modes (Code, Markdown, etc.). The main content area displays a slide titled "Jupyter Notebooks" with a small icon of a person looking at a screen. The slide contains the following text and list:

Jupyter Notebooks 

- open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text
- uses
  - data cleaning and transformation
  - numerical simulation
  - statistical modeling
  - data visualization
  - machine learning
  - ...

# Demo: Pandas

# Demo: Pandas

- contraction of "**panel data**"
- de facto data analysis toolkit (<https://pandas.pydata.org>)
- there are screenshots in this presentation to maintain continuity, but let's open the notebook named **Demo - Pandas.ipynb** and go through it together
- when done, click [here](#) to skip screenshots

# Pandas Series

```
>>> population_dict = {'California': 38332521,
                      'Texas': 26448193,
                      'New York': 19651127,
                      'Florida': 19552860,
                      'Illinois': 12882135}
>>> population = pd.Series(population_dict)
>>> population
California    38332521
Florida       19552860
Illinois      12882135
New York      19651127
Texas         26448193
dtype: int64
```

```
>>> area_dict = {'California': 423967, 'Texas': 695662, 'New York': 141297,
                  'Florida': 170312, 'Illinois': 149995}
>>> area = pd.Series(area_dict)
>>> area
California    423967
Florida       170312
Illinois      149995
New York      141297
Texas         695662
dtype: int64
```

# Pandas DataFrame

```
>>> states = pd.DataFrame({'population': population,
                           'area': area})
>>> states
      area  population
California  423967    38332521
Florida     170312    19552860
Illinois    149995    12882135
New York    141297    19651127
Texas       695662    26448193
```

```
>>> states.index  
Index(['California', 'Florida', 'Illinois', 'New York', 'Texas'], dtype='object')
```

```
>>> states.columns  
Index(['area', 'population'], dtype='object')
```

```
>>> states['area']  
California    423967  
Florida       170312  
Illinois      149995  
New York      141297  
Texas         695662  
Name: area, dtype: int64
```

```
>>> states[area > 200000]  
           area  population  
California  423967    38332521  
Texas        695662    26448193
```

```
>>> states.values  
array([[ 423967,  38332521],  
       [ 170312,  19552860],  
       [ 149995,  12882135],  
       [ 141297,  19651127],  
       [ 695662,  26448193]])
```

# Reading CSV Files into Pandas

```
In [1]: import pandas as pd  
       data = pd.read_csv('data/agg_database_daily.csv')
```

```
In [2]: data.head()
```

0	20161025	CHI	SP2	cs15	65.872505	69.13125	30.184166	1189676.4	1.6051193	14.272053	\N	\N	\N	\N	\N	59.996094	46.06433	13.931762	8.10202	51.894073	59.996094
1	20161031	CHI	SP4	cs40	30.010305	86	22.05	863386.7	1.1656047	9.44478	106.87152	87.29087	19.580648	15.005844	91.865668	106.8					
2	20161031	WAS	SP3	na23	6.484551	72	35.066666	1139595.2	1.7977492	13.583612	59.996094	32.447983	27.54811	20.739609	39.256485	59.99					
3	20161031	PHX	SP1	na45	4.7593827	50	17.758333	744373.94	1.1349427	9.120685	67.4989	26.249228	41.249672	30.199581	37.29932	67.					
4	20161031	CHI	SP1	gs0	10.942432	37	8.622222	1118318	1.1859665	4.268443	362.25	33.72686	328.52313	240.68625	121.56374	36					

```
In [3]: data.shape  
Out [3]: (55285, 20)
```

# Specifying No Headers

```
In [3]: import pandas as pd  
       data = pd.read_csv('data/agg_database_daily.csv', header=None)
```

```
In [4]: data.head()
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	
0	20161031	WAS	SP3	na17	12.881275	87	47.8	1329935.5	2.3832552	24.30115	59.996094	46.06433	13.931762	8.10202	51.894073	59.9960
1	20161025	CHI	SP2	cs15	65.872505	69.13125	30.184166	1189676.4	1.6051193	14.272053	\N	\N	\N	\N	\N	\N
2	20161031	CHI	SP4	cs40	30.010305	86	22.05	863386.7	1.1656047	9.44478	106.87152	87.29087	19.580648	15.005844	91.86568	106.871
3	20161031	WAS	SP3	na23	6.484551	72	35.066666	1139595.2	1.7977492	13.583612	59.996094	32.447983	27.54811	20.739609	39.256485	59.9960
4	20161031	PHX	SP1	na45	4.7593827	50	17.758333	744373.94	1.1349427	9.120685	67.4989	26.249228	41.249672	30.199581	37.29932	67.49

```
In [5]: data.shape  
Out [5]: (55286, 20)
```

# Specifying Column Names

```
In [6]: data.columns = ['date_key', 'datacenter', 'superpod', 'pod', 'max_redo_size',
'max_active_sessions', 'max_db_cpu_user', 'max_peak_buffer', 'avg_db_cpu_system',
'avg_db_cpu_user', 'total_db_size_in_tb', 'used_db_space_in_tb',
'free_db_space_in_tb', 'asm_free_db_space_in_tb', 'asm_used_db_space_in_tb',
'asm_total_db_size_in_tb', 'last_modified', 'asm_used_db_space_perc',
'oem_cpu_util', 'oem_read_io_latency']
```

```
In [7]: data.head()
```

	date_key	datacenter	superpod	pod	max_redo_size	max_active_sessions	max_db_cpu_user	max_peak_buffer	avg_db_cpu_system	avg_db_cpu_user	total
0	20161031	WAS	SP3	na17	12.881275		87	47.8	1329935.5	2.3832552	24.30115
1	20161025	CHI	SP2	cs15	65.872505		69.13125	30.184166	1189676.4	1.6051193	14.272053
2	20161031	CHI	SP4	cs40	30.010305		86	22.05	863386.7	1.1656047	9.44478
3	20161031	WAS	SP3	na23	6.484551		72	35.066666	1139595.2	1.7977492	13.583612
4	20161031	PHX	SP1	na45	4.7593827		50	17.758333	744373.94	1.1349427	9.120685

# Missing Values

```
In [7]: data['max_active_sessions']
Out[7]: 0           87
        1           69.13125
        2           86
        3           72
        4           50
        5           37
        6           38
        7           \N
        8           89
        9           94
        ...
        ...
```

# Handling Missing Values

```
In [8]: data = pd.read_csv('data/agg_database_daily.csv', header=None,
                           na_values=[r'\N'])
        data.columns = ['date_key', 'datacenter', 'superpod', 'pod', 'max_redo_size',
        'max_active_sessions', 'max_db_cpu_user', 'max_peak_buffer', 'avg_db_cpu_system',
        'avg_db_cpu_user', 'total_db_size_in_tb', 'used_db_space_in_tb',
        'free_db_space_in_tb', 'asm_free_db_space_in_tb', 'asm_used_db_space_in_tb',
        'asm_total_db_size_in_tb', 'last_modified', 'asm_used_db_space_perc',
        'oem_cpu_util', 'oem_read_io_latency']
```

```
In [7]: data['max_active_sessions']
Out [7]: 0          87
         1      69.13125
         2          86
         3          72
         4          50
         5          37
         6          38
         7          NaN
         8          89
         9          94
         ...
         ...
```

# Dropping Missing Values

```
In [9]: data['max_active_sessions'].dropna()
Out [9]: 0      87.000000
         1      69.131250
         2      86.000000
         3      72.000000
         4      50.000000
         5      37.000000
         6      38.000000
         8      89.000000
         9      94.000000
        10     35.542683
        11     57.950000
        12     44.000000
        14     46.000000
        15     47.000000
        16    107.000000
        17     83.000000
        18    77.470290
        19     89.000000
        ...
...
```

# Reading SQLite into Pandas

```
>>> import pandas as pd
>>> import sqlite3
>>> conn = sqlite3.connect("data/flights.db")
>>> data = pd.read_sql_query("select * from airlines", conn)
>>> data.shape
(6048, 9)
```

```
>>> data.head()
   index  id
0      0  1
1      1  2
2      2  3
3      3  4  2 Sqn No 1 Elementary Flying Training School
4      4  5

   callsign      country active
0    None        None      Y
1  GENERAL  United States      N
2  NEXTIME  South Africa      Y
3    None  United Kingdom      N
4    None        Russia      N
```

```
>>> conn.close()
```

# Exercise: Pandas

(open the notebook named **Exercise 1 - Pandas.ipynb**)

# Exercise: Reading TTL Data

(open the notebook named **Exercise 2 - Reading TTL Data.ipynb**)

# Recap: Data Sources Recap

- What are we trying to do?
  - we're reading data from relational databases, spreadsheets, etc.
- Why would we choose to put data into these different types of storage?
  - differently-shaped data
  - high-volume data
  - high-velocity data
  - in other words, there are compelling business and technical reasons
- We also talked about *linked data*
  - e.g., [DBpedia](#) is a dataset drawn from the editorial process of [Wikipedia](#)
  - for anything in Wikipedia, such as a city, there is structural information (e.g., when city was founded, geographic location, etc.)
  - the entities in DBpedia are generated from the pages of Wikipedia
  - they're resolvable through URLs, and the relationships themselves are also resolvable
    - node ids/links are like global, disambiguable, resolvable DB keys
    - relationship ids/links are like global, disambiguable, resolvable column names

# Demo: Reading Linked Data

# Demo: Reading Linked Data

- there are screenshots in this presentation to maintain continuity, but let's open the notebook named **Demo - Reading Linked Data.ipynb** and go through it together
- when done, click [here](#) to skip screenshots

```
import pandas as pd
import json
from SPARQLWrapper import SPARQLWrapper, JSON
```

```
def get_sparql_dataframe(service, query):
    """
    Helper function to convert SPARQL results into a Pandas data frame.
    """
    sparql = SPARQLWrapper(service)
    sparql.setQuery(query)
    sparql.setReturnFormat(JSON)
    result = sparql.query()

    processed_results = json.load(result.response)
    cols = processed_results['head']['vars']

    out = []
    for row in processed_results['results']['bindings']:
        item = []
        for c in cols:
            item.append(row.get(c, {}).get('value'))
        out.append(item)

    return pd.DataFrame(out, columns=cols)
```

```
wds = "https://query.wikidata.org/sparql"
```

```

# This SPARQL query will be sent to the SPARQL endpoint defined in the previous
# step. It mixes three vocabularies that each have their own definitions but is
# ultimately a selection from the Wikidata graph. We're looking for distinct rows
# of individuals who have an orcid (https://orcid.org) and any English description
# and labels we might also have about them. We're matching a pattern in the graph
# for any node that is connected to other nodes with these relationships.

# To understand what __`wdt:P496 means`__, expand it into its full URL by applying
# the prefix for wdt and then issuing an HTTP request to
# http://www.wikidata.org/prop/direct/P496.

# For more info, consult Bob DuCharme Learning SPARQL (2nd Edition)

rq = """
PREFIX bd: <http://www.bigdata.com/rdf#>
PREFIX wikibase: <http://wikiba.se/ontology#>
PREFIX wdt: <http://www.wikidata.org/prop/direct/>

select distinct
?item
?itemLabel
?orcid
?description
WHERE {
?item wdt:P496 ?orcid
OPTIONAL {
?item schema:description ?description filter (lang(?description) = "en")
}
SERVICE wikibase:label {
bd:serviceParam wikibase:language "en" .
}
"""

```

# Exercise: Querying Linked Data

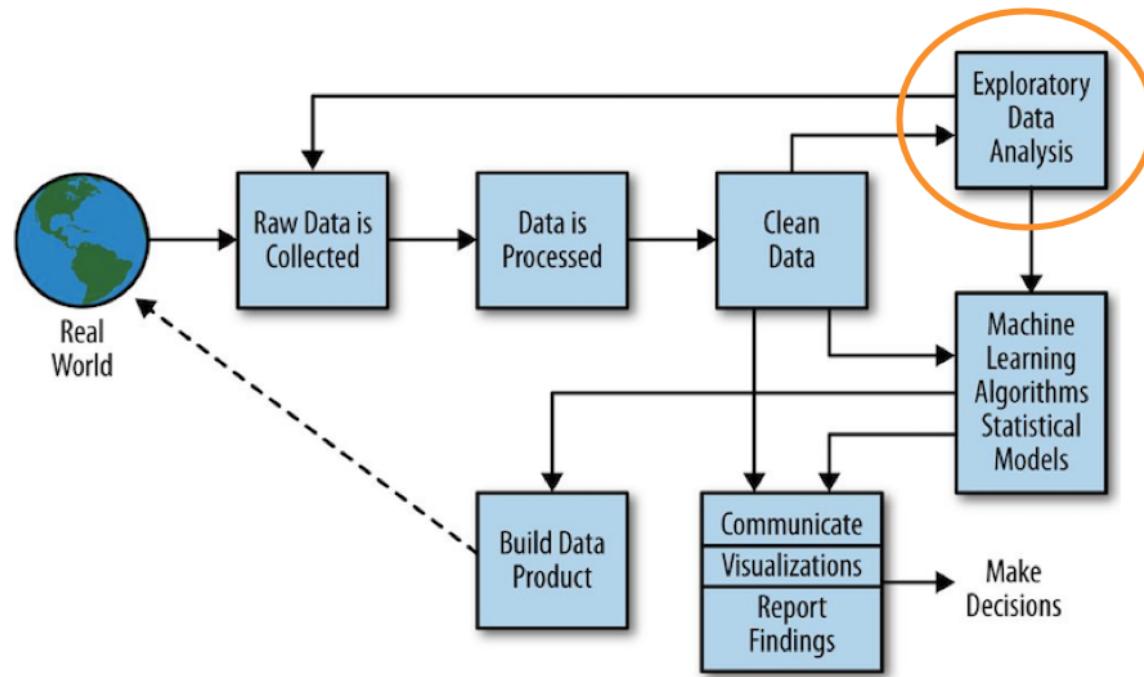
(open the notebook named `Exercise 3 - Querying Linked Data.ipynb`)

# Data-Driven

# Objective

- Consider the mental shift in letting the data drive our activities
- Explore the discipline needed to understand our data
- Friendly statistics intro/re-familiarization
- Understand some of the challenges that emerge with being data driven

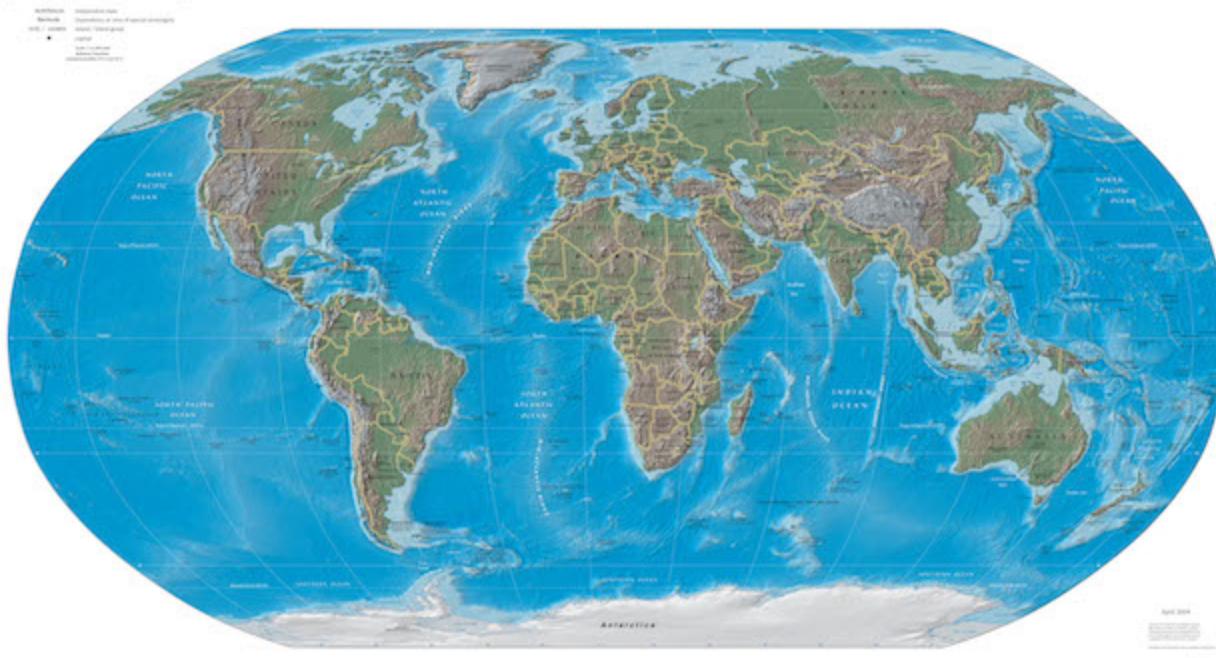
# Data Science Pipeline



What was the degree with the highest average starting salary at the University of North Carolina in the 1980's?

# Geography!

Physical Map of the World, April 2004



# Intro to (or Review of) Statistics

# Demo: Let's Make Some Data

- there are screenshots in this presentation to maintain continuity, but let's open the notebook named **Demo - Let's Make Some Data.ipynb** and go through it together
- when done, click [here](#) to skip screen shots

```
>>> # Generate a list of numbers 20 to 59  
>>> ages = range(20, 60)
```

```
>>> # Select 100 numbers at random from the set of numbers 20 to 59  
>>> random_ages = [random.choice(ages) for _ in range(100)]
```

```
>>> random_ages  
[27, 21, 20, 22, 21, 27, 25, 52, 56, 43, 33, 28, 31, 30, 41, 29, 38, 52,  
 48, 24, 50, 31, 29, 56, 45, 54, 37, 33, 22, 33, 23, 40, 37, 23, 28, 38,  
 41, 34, 35, 40, 25, 57, 32, 48, 56, 54, 39, 25, 57, 24, 26, 52, 26, 35,  
 57, 39, 26, 32, 50, 52, 54, 53, 30, 53, 56, 25, 43, 38, 53, 46, 51, 51,  
 35, 38, 20, 52, 47, 53, 45, 27, 22, 26, 25, 47, 36, 52, 54, 58, 31, 42,  
 46, 46, 43, 33, 35, 48, 49, 55, 36, 36]
```

```
>>> max(random_ages)  
58
```

```
>>> min(random_ages)  
20
```

# Range

```
>>> def calc_range(x):
...     return max(x) - min(x)
...
```

```
>>> calc_range(random_ages)
38
```

```
>>> nums = [10, 10, 100, 100]
>>> calc_range(nums)
90
```

```
>>> nums = [10, 50, 50, 50, 50, 100]
>>> calc_range(nums)
90
```

```
# Peak to Peak - returns the range, the min is a "peak" and the max is a "peak"
>>> numpy.ptp(random_ages)
38
```

# Mean

```
>>> def mean(x):
...     return sum(x) / len(x)
...
```

```
>>> mean(random_ages)
38.99
```

```
>>> import numpy
>>> numpy.mean(random_ages)
38.99000000000002
```

# Median

```
def median(x) :  
    n = len(x)  
    sorted_x = sorted(x)  
    mid = n // 2  
  
    if n % 2 == 0:  
        return (sorted_x[mid - 1] + sorted_x[mid]) / 2  
    else:  
        return (sorted_x[mid])
```

```
>>> median(random_ages)  
38
```

```
>>> numpy.median(random_ages)  
38.0
```

# Percentile

```
>>> numpy.percentile(random_ages, 33)  
32.0
```

```
>>> numpy.percentile(random_ages, 80)  
52.0
```

```
>>> numpy.percentile(random_ages, 50)  
38.0
```

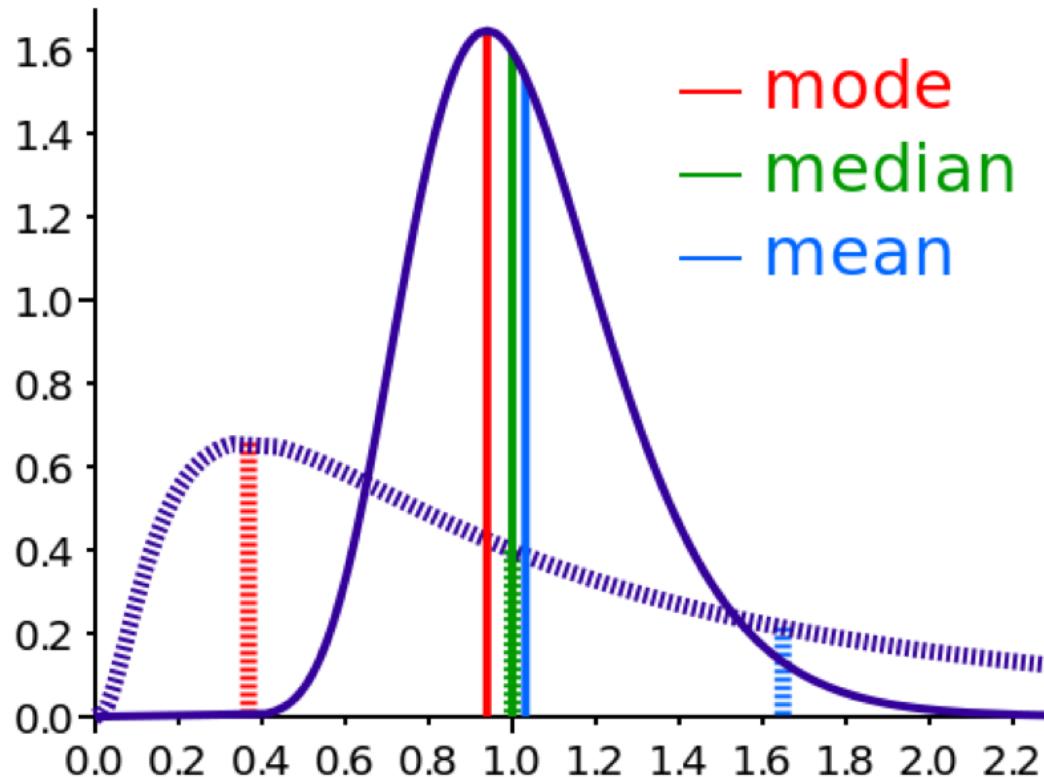
# Interquartile Range (IQR)

$$IQR = Q_3 - Q_1$$

```
>>> from scipy import stats  
>>> stats.iqr(random_ages)  
21.5
```

# Mode

```
>>> from scipy import stats  
>>> stats.mode(random_ages)  
ModeResult(mode=array([52]), count=array([6]))
```



- consider the spread of data in two hypothetical distributions
- how can we identify/quantify different spreads?
- focus on the **mean**, **median**, and **mode**

# Variance

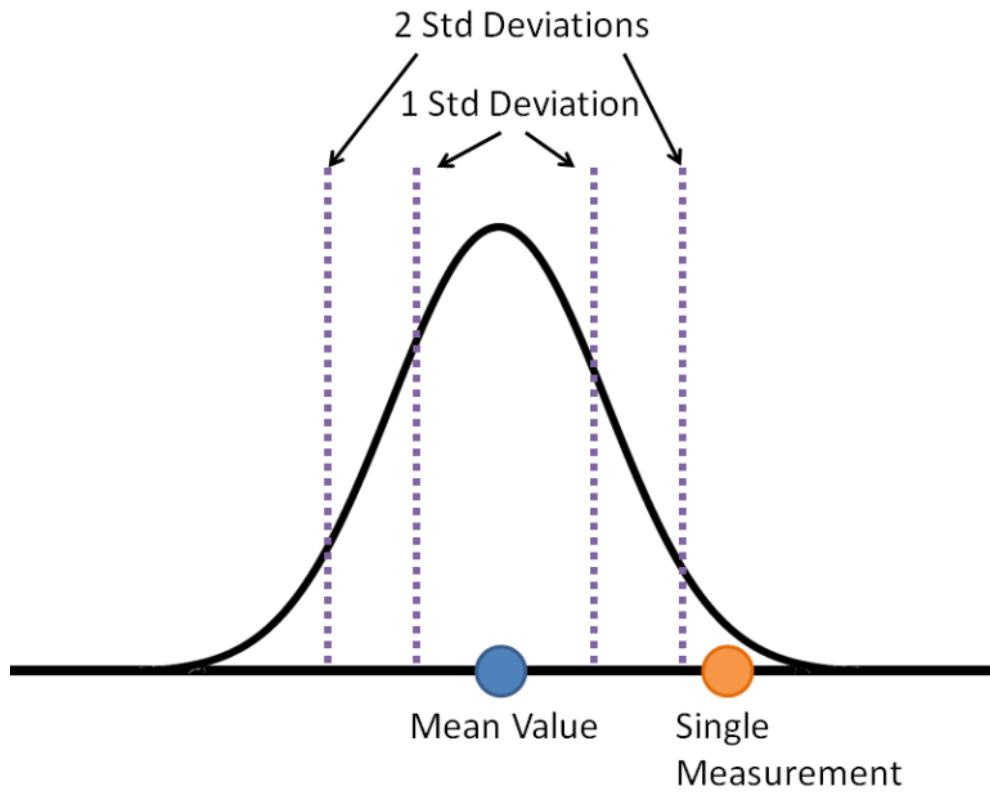
$$Var(X) = \frac{1}{n} \sum_{i=1}^n (a_i - \bar{x})^2$$

```
>>> numpy.var(random_ages)  
131.8298999999998
```

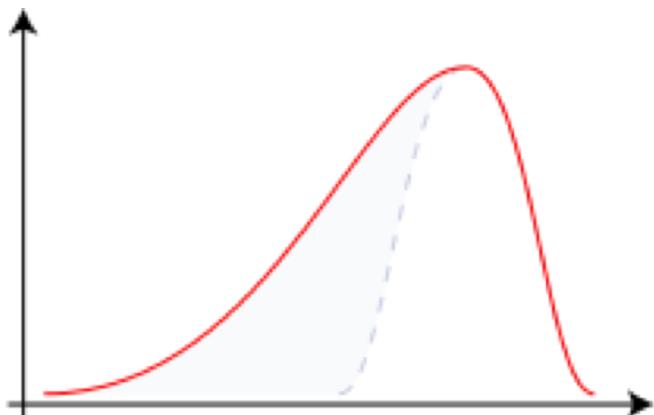
# Standard Deviation

$$\sqrt{Var(X)}$$

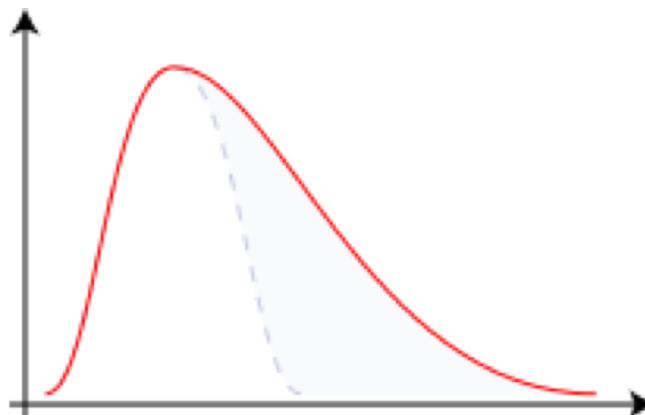
```
>>> numpy.std(random_ages)  
11.481720254386969
```



# Skewness



Negative Skew



Positive Skew

```
>>> stats.skew(random_ages)  
0.037925216234465986
```

# Descriptive Statistics

- Descriptive Statistics *describe* our data in various ways
- Measures of center
  - mean
  - median
  - mode
- Measures of dispersion
  - range
  - percentiles/quartiles
  - IQR
  - variance
  - standard deviation
  - skewness

# Demo: Summary Statistics of TTL Data

# Demo: Summary Statistics of TTL Data

- let's open the notebook named **Demo - Summary Statistics of TTL Data.ipynb** and go through it together

# Exercise: Summary Statistics of TTL Data

(open the notebook named `Exercise 4 - Summary Statistics of TTL Data.ipynb`)

# Understanding Relationships between Variables

# Covariance

- Measure of joint variability
  - i.e., how does one variable change as the other changes?
  - this draws our attention to a connection
- Non-zero covariance implies?
  - there is some connection
- Zero covariance implies?
  - independence

"Correlation is not causation."

*Anyone who has ever taken statistics*

"Empirically observed covariation is a necessary but not sufficient condition for causality."

*Edward Tufte*

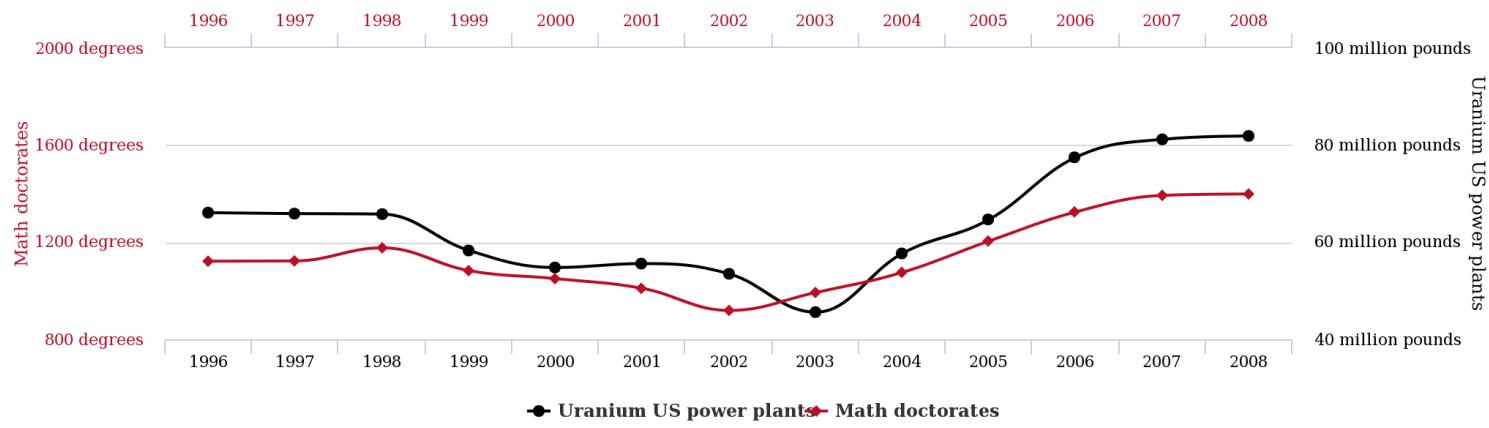
"Correlation is not causation but it sure is a hint."

*Edward Tufte*

# Correlation/Causation

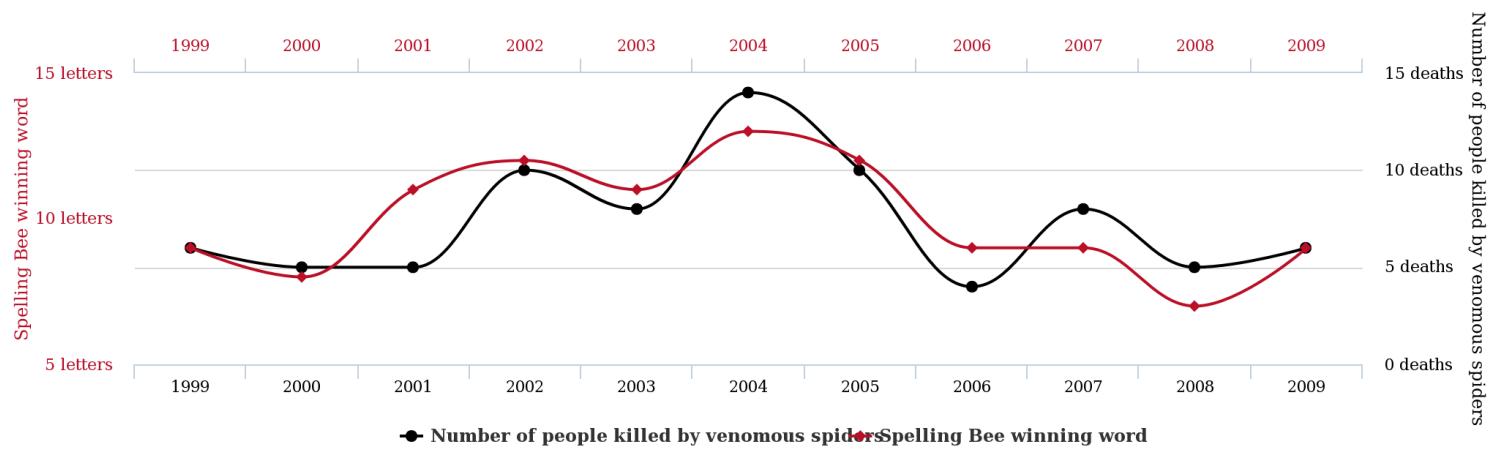
- What are the possible causal relationships?
  - could be unrelated, i.e., coincidence
  - reverse causation
    - every time windmills are spinning it's really windy
  - missed variable, i.e., some other factor causes both
    - whenever I go to sleep with my shoes on I wake up with a headache
    - when it rains, every time I see a flash I hear a boom
  - bi-directional relationships
    - temperature and pressure
  - or they are actually the same thing
    - °F goes up as °C goes up

**Math doctorates awarded**  
correlates with  
**Uranium stored at US nuclear power plants**



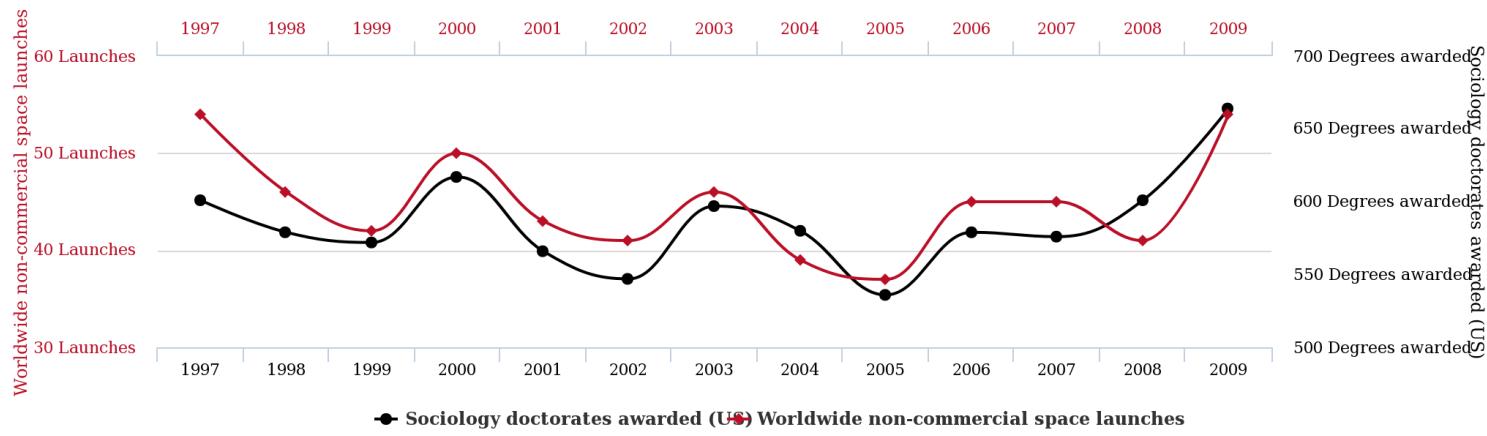
tylervigen.com

## **Letters in Winning Word of Scripps National Spelling Bee** correlates with **Number of people killed by venomous spiders**

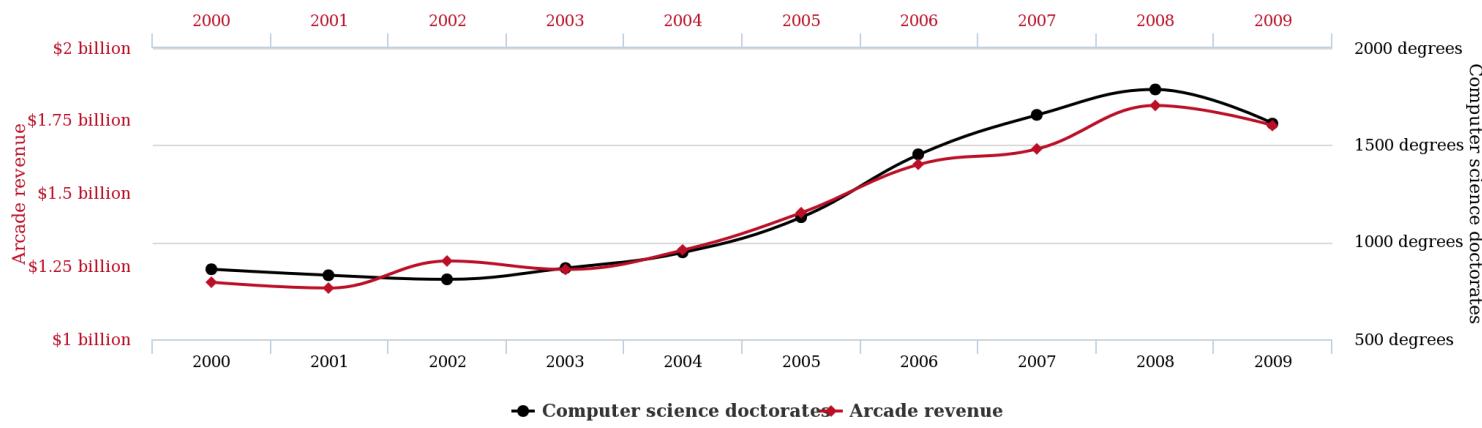


tylervigen.com

## Worldwide non-commercial space launches correlates with **Sociology doctorates awarded (US)**



**Total revenue generated by arcades**  
correlates with  
**Computer science doctorates awarded in the US**



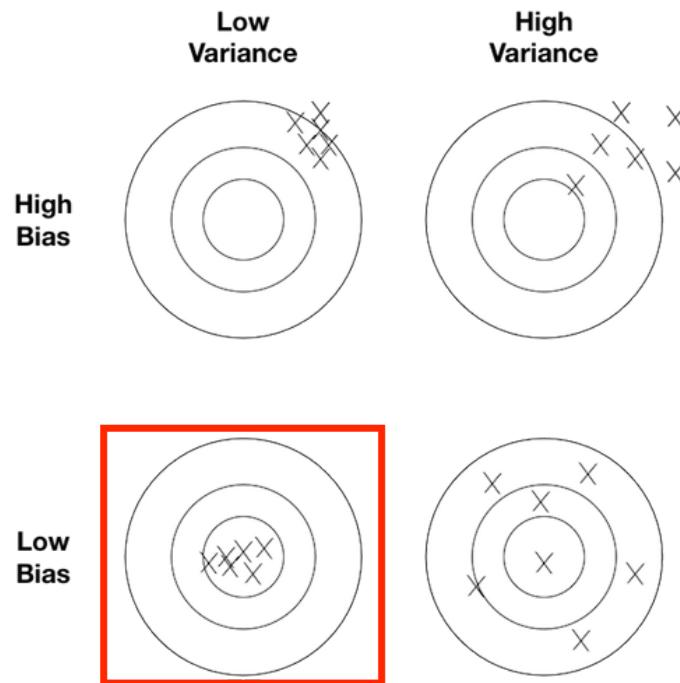
tylervigen.com

# Choosing Algorithms: Bias/Variance Tradeoff

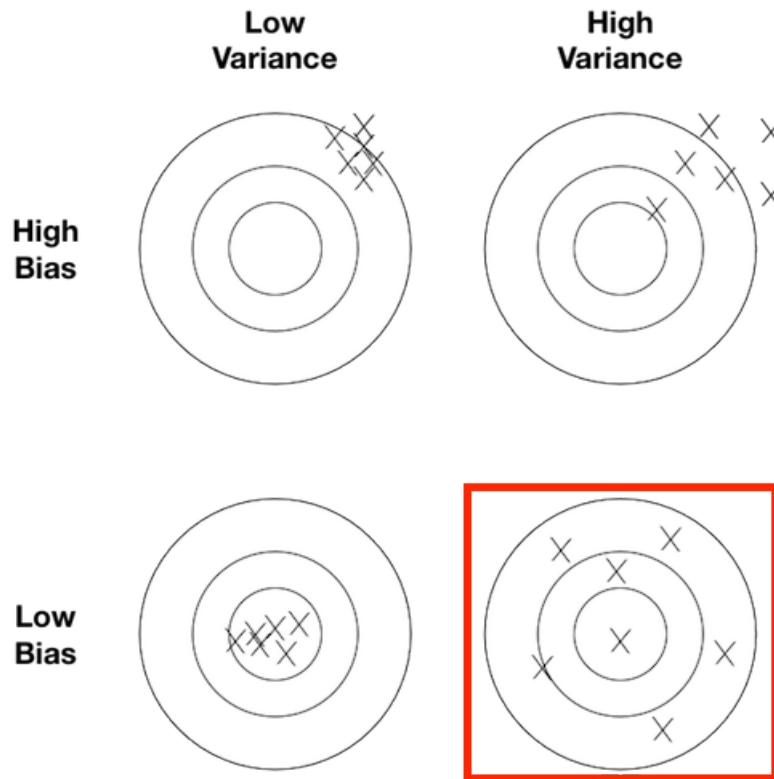
# Choosing an Algorithm

- from an exploratory data analysis standpoint, we might try things and see what works
  - but if we can't map it back to the business function, it doesn't matter
  - if it does work, but doesn't work from an operational standpoint, you can't use it either
    - e.g., can't run on a phone, or breaks email
  - therefore, not everything that works is available to us
  - sometimes things work well in one aspect don't work well in another aspect
- A **model** is the "function" which is produced from a machine learning algorithm
- **bias** refers to errors from incorrect assumptions in the learning algorithm
  - high bias means our model is pushing us in the wrong direction (underfitting)
- **variance** refers to errors from sensitivity to small fluctuations in the training set
  - High variance could result in modeling random noise in the training data, rather than the intended outputs (overfitting)

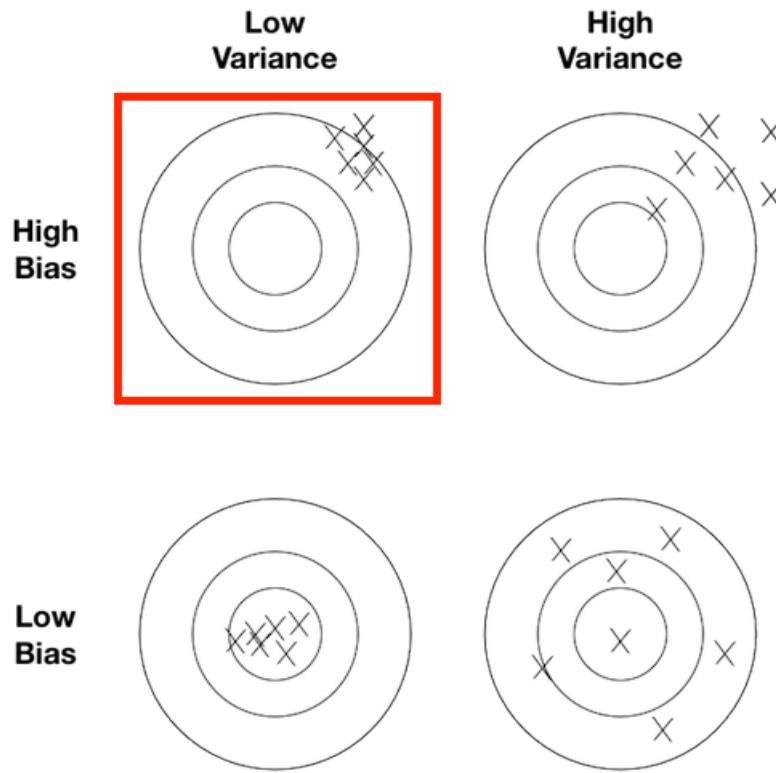
# Understanding Model Bias and Variance



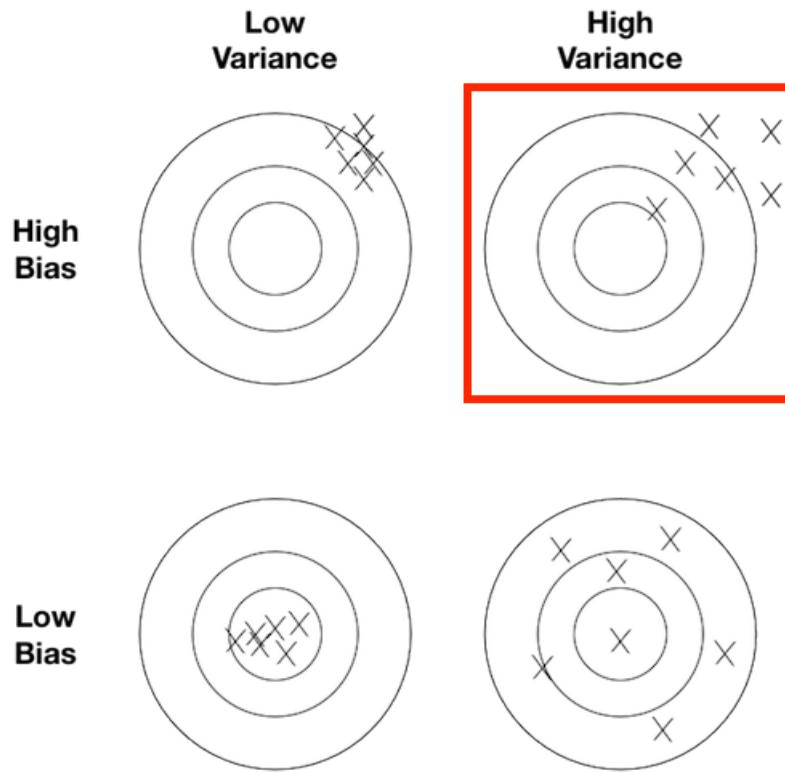
- *low bias* means our model is not pushing us in the wrong direction
- *low variance* means we're tightly clustered
- like someone who is good at playing darts, has natural hand-eye coordination, and is sober—they can put the darts wherever they want
- imagine data set that reflects the population, the features that we are trying to train off of, the method we are using to train (model we are choosing)—that should work



- let's say our dart thrower is drunk
  - darts will be all over the place, which maps to noise in the system
  - *high variance* is due to the model being too tightly tied to the particular data we are trying to train off of
  - if we pick different data, we'll get different results
  - we may need to choose a simpler model (throws a bunch of noise out, not as responsive to everything going on as a more complex model)



- simpler models are more like our friends who have some natural ability in dart throwing but no formal training
- we'll be wrong, but maybe right enough
- e.g., if we assume there is a linear relationship when there isn't, we'll be off from the real results



- our friend who has no skill, no training, and is drunk will have no hope of doing well
- imagine trying to build a model of credit worthiness based on your favorite color
- good algorithms + bad data = garbage

# Exploring Probability and Hypothesis Testing

# Introduction to Probability

# Flipping a Coin

- How many outcomes are there? (...assuming it doesn't land on its edge)
- What is the likelihood of each of those outcomes?
- We intuitively understand this when it comes to coin flipping, but let's formalize it

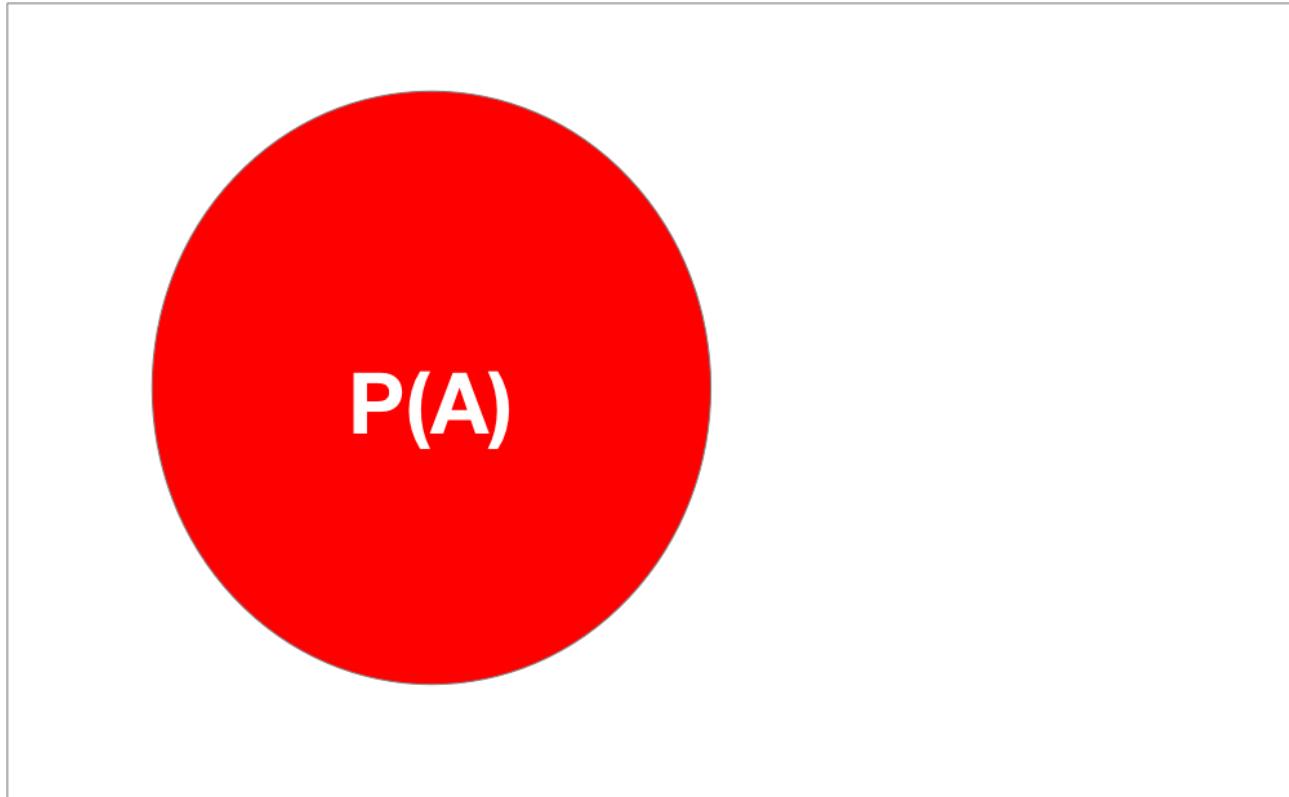


# Definitions

- **Random Experiment** - Observing something uncertain
- **Outcome** - Result of an experiment
  - heads, tails, rolling a 3, picking an even number
- **Sample space** - Set of all possible outcomes
  - Coin: heads, tails
  - A di(c)e: 1, 2, 3, 4, 5, 6



# Probability that Event A Happens



# Probability of Event A: Roll a 1

$$P(A) = \frac{|\text{Outcome}|}{|\text{Sample Space}|}$$

$$P(A) = \frac{|\{1\}|}{|\{1, 2, 3, 4, 5, 6\}|}$$

$$P(A) = \frac{1}{6}$$

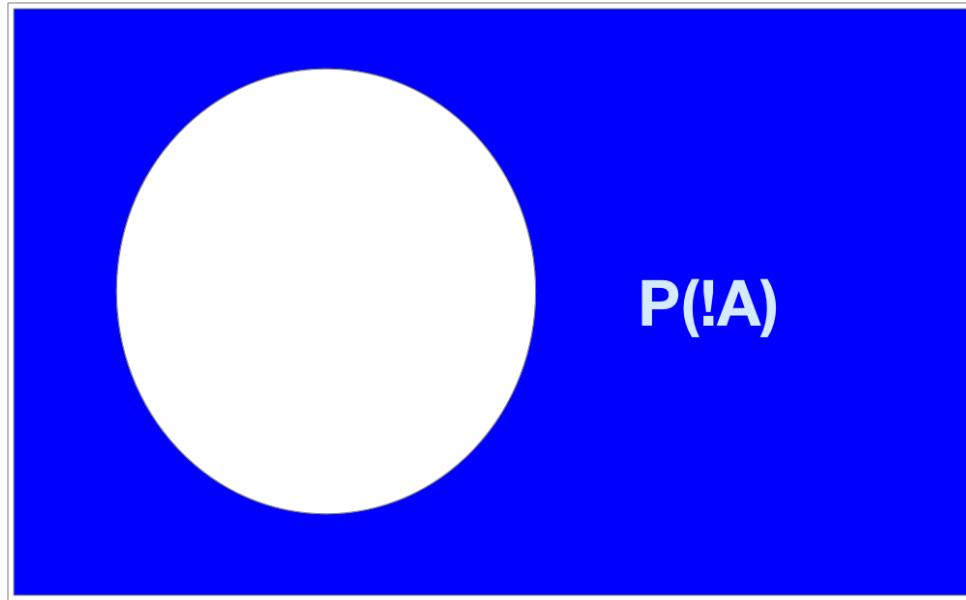
# Probability of Event A: Roll an Even Number

$$P(A) = \frac{|\text{Outcome}|}{|\text{Sample Space}|}$$

$$P(A) = \frac{|\{2, 4, 6\}|}{|\{1, 2, 3, 4, 5, 6\}|}$$

$$P(A) = \frac{3}{6} = \frac{1}{2}$$

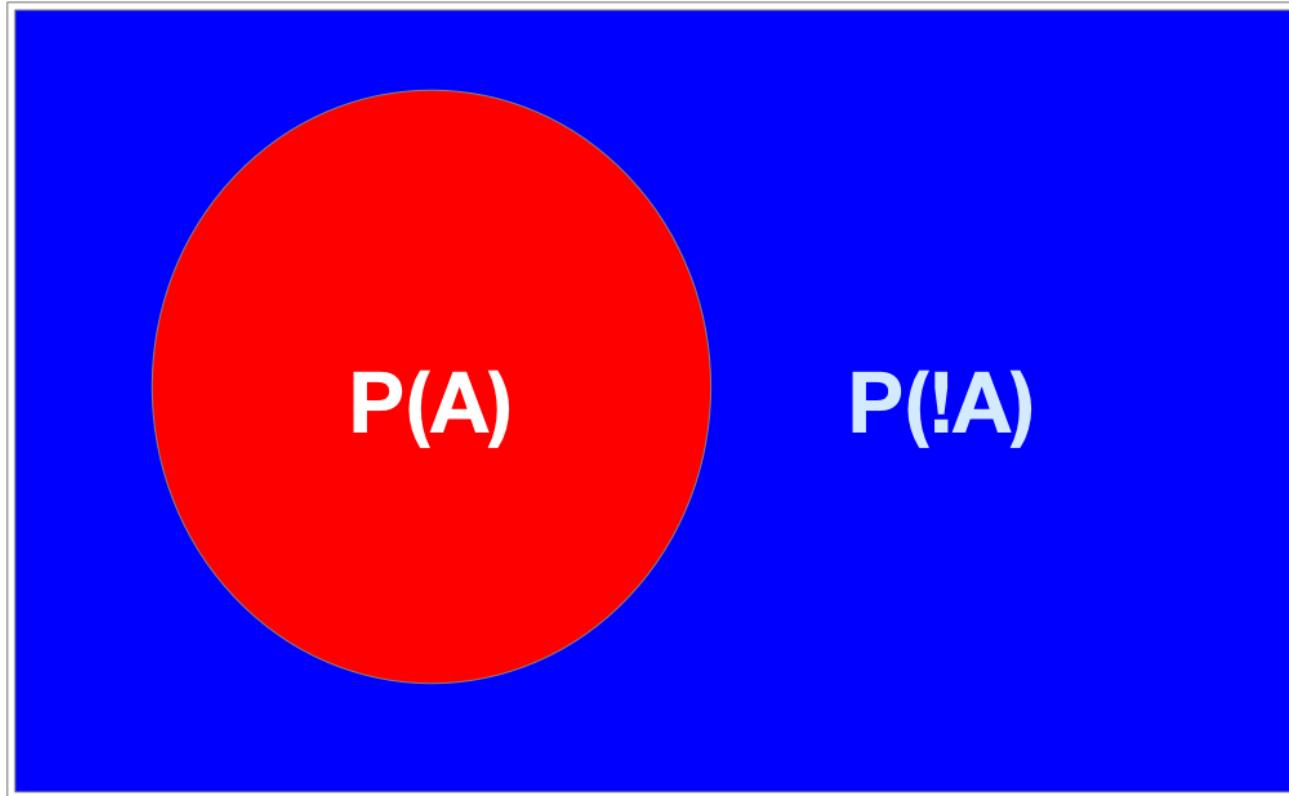
# Probability that Event A Doesn't Happen



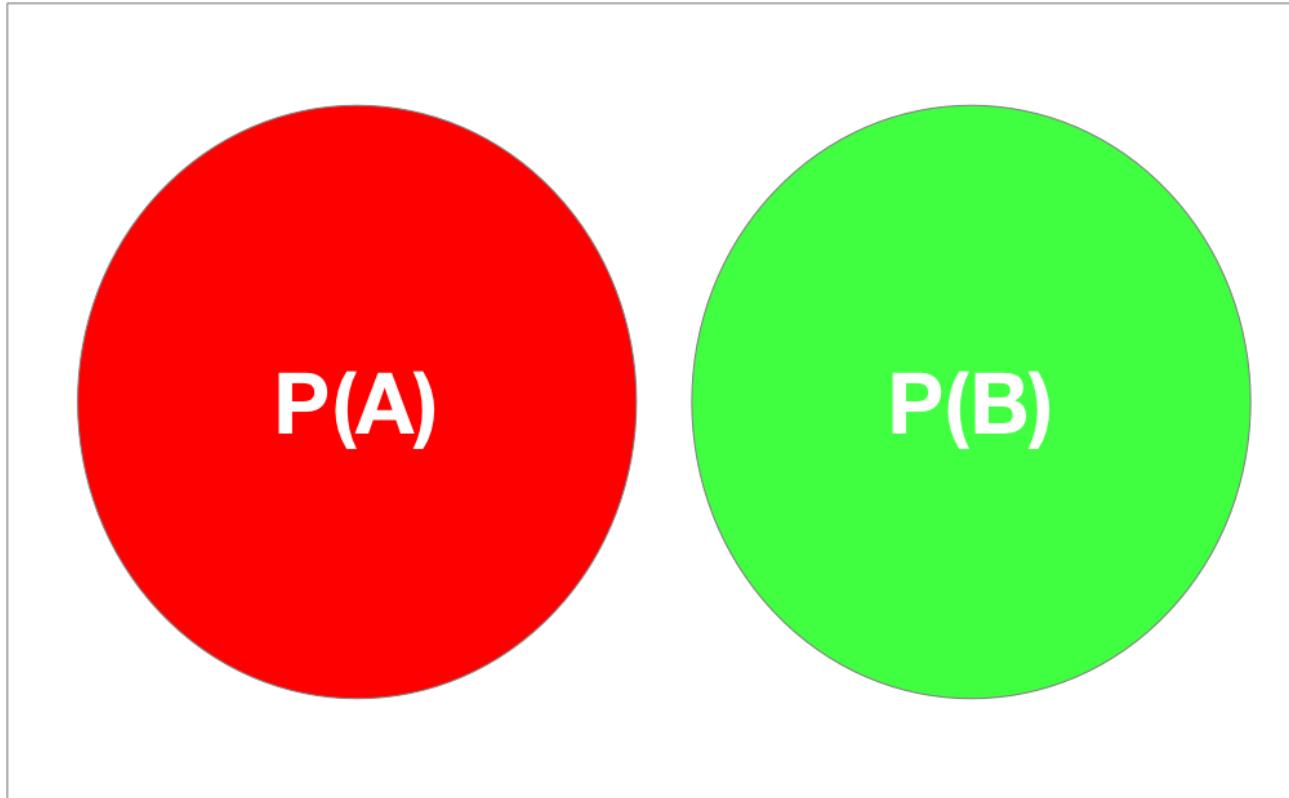
$$P(\text{!}A) = ?$$

$$P(\text{!}A) = 1 - P(A)$$

# Whole Sample Space



# Consider Two Events A and B



Event A = roll a 3

Event B = roll a 4

$$P(A) = P(\text{roll a 3})$$

$$P(B) = P(\text{roll a 4})$$

$$P(A) = P(\text{roll a 3}) = \frac{1}{6} = 0.167$$

$$P(B) = P(\text{roll a 4}) = \frac{1}{6} = 0.167$$

$$P(\neg A) = 1 - P(A) = \frac{5}{6} = 0.833$$

$$P(\neg B) = 1 - P(B) = \frac{5}{6} = 0.833$$

# Probability of A and B Happening (Two Rolls)

$$P(A \cap B) = ?$$

# Assuming Independent Events

$$P(A \cap B) = P(A) \cdot P(B)$$

$$= 0.167 \cdot 0.167$$

$$= 0.02789$$

# Probability of A Given B

(i.e., probability of A happening, given that B has happened)

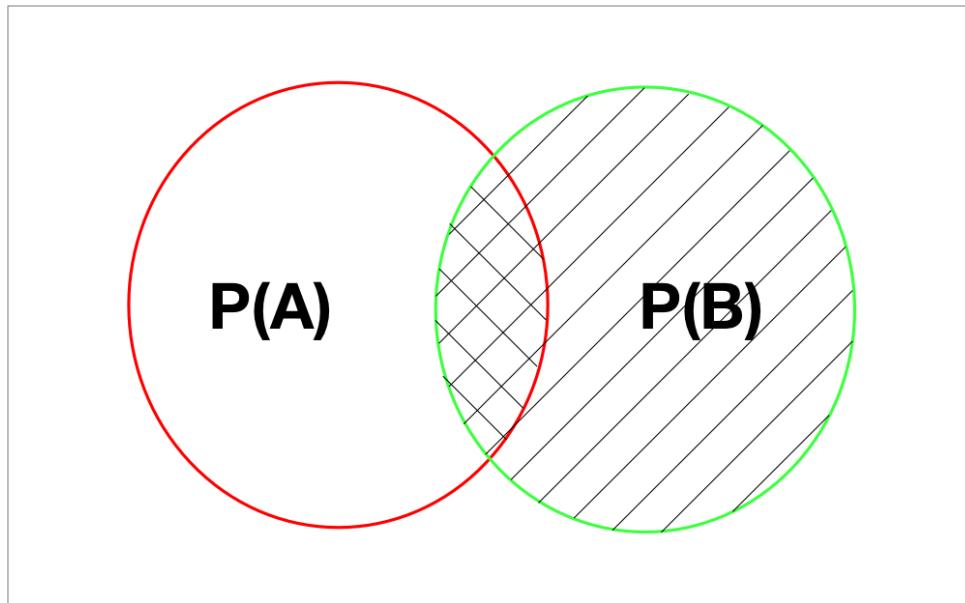
$$P(A|B) = ?$$

$$P(A) = P(\text{roll a 3})$$

$$P(B) = P(\text{roll a 4})$$

$$P(A|B) = P(A)$$

# Probability of A Given B



$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

# $P(A|B)$ for Independent Events

$$P(A) = \text{roll a 3}$$

$$P(B) = \text{roll a 4}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) = \frac{P(A) \cdot P(B)}{P(B)}$$

$$= P(A)$$

# $P(A|B)$ for Non-Independent Events

$A$  = roll an odd number

$B$  = roll a number  $\leq 3$

$$P(A) = \frac{1}{2}$$

$$P(B) = \frac{1}{2}$$

$$P(event) = \frac{|Outcome|}{|Sample Space|}$$

This time, the event in question is  $A | B$   
(i.e., die shows an odd number given that the roll was  $\leq 3$ )

$$P(A | B) = \frac{|A \cap B|}{|B|}$$

What is the numerator, i.e., number of ways that both A and B occur?

What is the denominator, i.e., number of ways that B can occur?

$$P(A|B) = \frac{|A \cap B|}{|B|}$$

Now divide both top and bottom by | Sample Space |

$$= \frac{\frac{|A \cap B|}{|Sample\ Space|}}{\frac{|B|}{|Sample\ Space|}}$$

$$= \frac{P(A \cap B)}{P(B)}$$

A = "Odd" and B = "Rolled 1, 2, or 3"

$$B = \{1, 2, 3\}$$

$$A \cap B = \{1, 3\}$$

$$P(A|B) = \frac{|A \cap B|}{|B|}$$

$$= \frac{2}{3}$$

# Exercise

1. Go to [setosa.io](https://setosa.io)
2. Change the *perspective* from **world** to  $P(B|A)$  and  $P(A|B)$  to be sure you understand what is happening in the simulation
3. Adjusting the *drop frequency* can also make it easier to see what's going on
4. Adjust  $P(A)$  so that the red shelf completely covers the blue shelf and be sure you understand how that affects  $P(A|B)$
5. Reload the page and perform a similar adjustment for  $P(B)$
6. Adjust the  $P(A \cap B)$  slider up and down and be sure you understand the effect on the probabilities

## 2. Solve the following problem

- Out of a sample of 1000 Tesla Model 3 pre-orders, 400 purchased the long range battery (LRB), 325 purchased the all-wheel drive (AWD) option, and 180 buyers purchased both options. If a random buyer from this group purchased the LRB, what is the probability that the buyer also purchased the AWD?



# Making Claims about a Population: Hypothesis Testing

# Hypothesis Testing

"Could we expect to see the result simply by chance?"

"Are we simply demonstrating a connection?"

"Are we demonstrating a specific result?"

# Example: Evaluating the Efficacy of Sales Training

- a bunch of sales people, you hear about a training that helps people get better at selling software
- How do you know it is worth it?
  - send a few people, see if it works
- How could software sold be different for the people you've sent?
  - sell more
  - sell the same
  - sell less
- When evaluating whether or not it's worth it maybe you need to delay measurement

# Example: Birth Weight of Babies



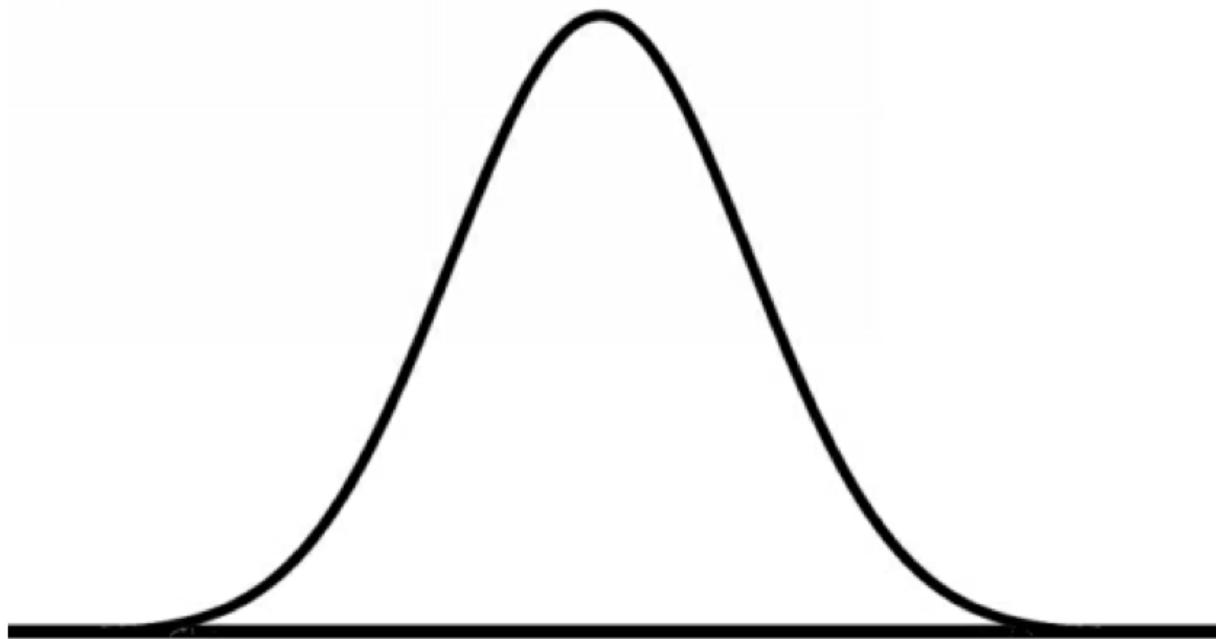
- Langston is one of 30 babies born in a particular hospital
- Suppose we want to know whether male babies born in that hospital are larger than average at birth?

Baby Number	Weight (lbs)
1	5.5
2	5.9
3	6
4	6
5	6.1
6	6.5
7	6.5
8	6.8
9	7
10	7.2
11	7.4
12	7.5
13	7.5
14	7.6
15	7.7
16	7.7
17	7.8
18	7.9
19	8
20	8
21	8.2
22	8.3
23	8.5
24	8.9
25	9
26	9.1
27	9.4
28	10
29	11
30	12
Average	7.8333

# Compare Our Hospital to Population Mean

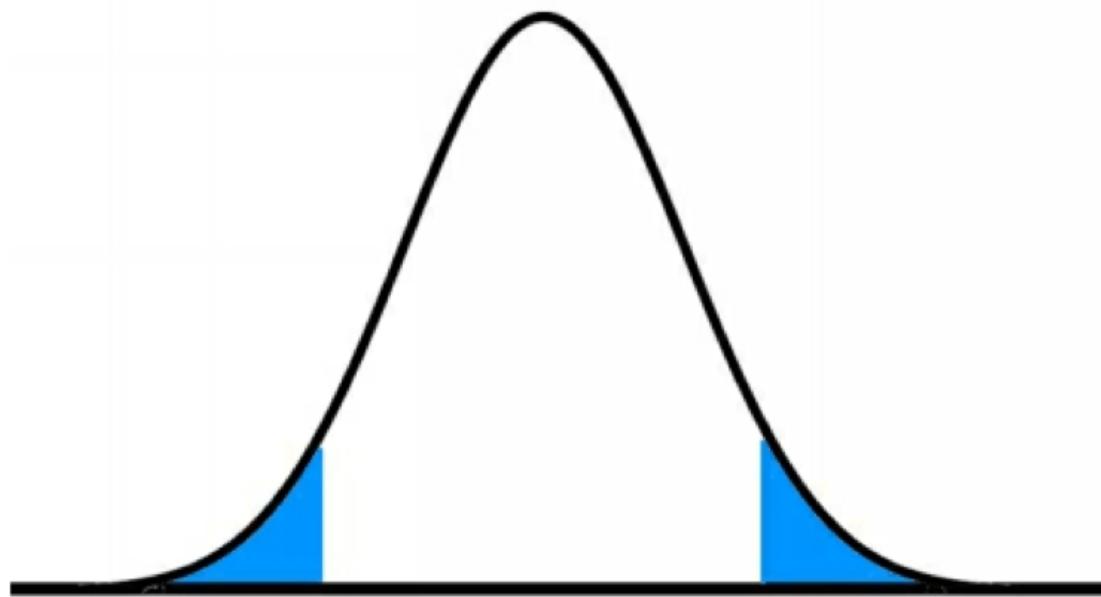
- avg. birth weight for a male baby in U.S. is 7.5 lbs. ( population mean,  $\mu$  )
- standard deviation of 1.25 lbs. ( population standard deviation,  $\sigma$  )
- we desire a 95% confidence level that our babies are bigger

Standard Normal Curve



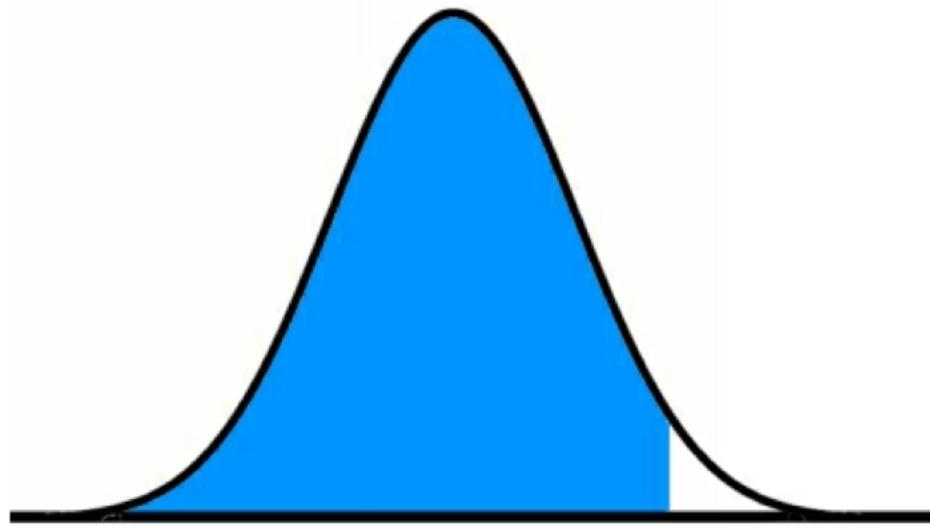
*How do our hospital's babies compare to the overall population mean?*

Standard Normal Curve  
Filled From Both Sides



- If we want to test that our babies do not weigh the same (i.e., more or less) as the average US male baby, then a two-tailed test would be appropriate

Standard Normal Curve  
Filled From Left  
i.e. 1 Tailed From Left



In our case, we are interested in seeing if our babies are *bigger* so we use a one-tailed test

# Making the Claim: The Hypothesis

- To perform a hypothesis test, a null ( $H_0$ ) and alternative ( $H_a$ ) hypothesis needed to be stated.
- $H_0$ : Our babies are not significantly heavier than the population mean
- $H_a$ : They are
- To determine if the difference is significant we use a Z-Test

# Z Test

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

Variable	Meaning
$\bar{x}$	sample mean
$\mu$	population mean
$\sigma$	population standard deviation
$n$	test sample size

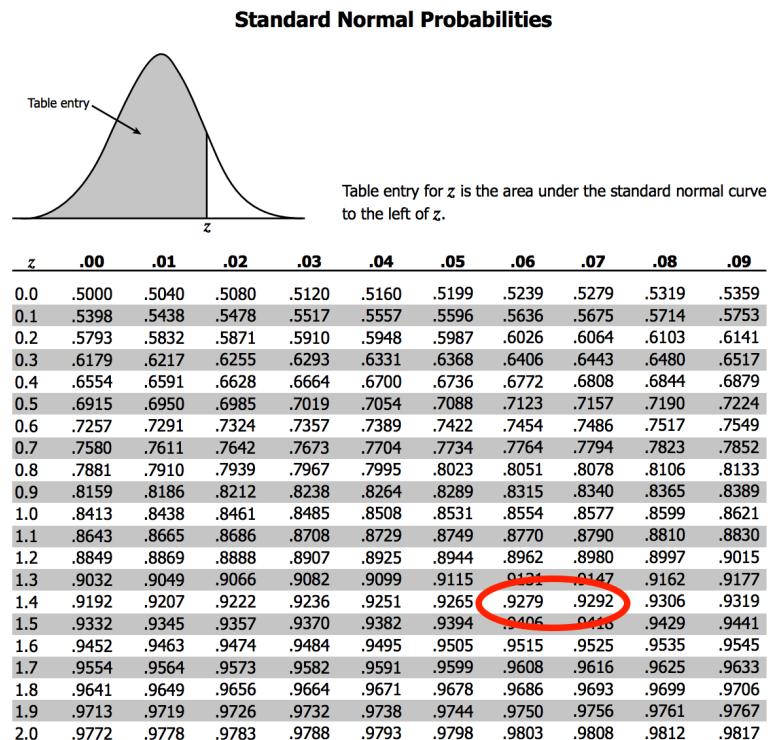
# Z Test (Rewritten)

$$Z = \frac{\bar{x} - \mu}{\sigma} \times \sqrt{n}$$

$$Z = \frac{7.8333 - 7.5}{1.25} \times \sqrt{30}$$

$$Z = 1.461$$

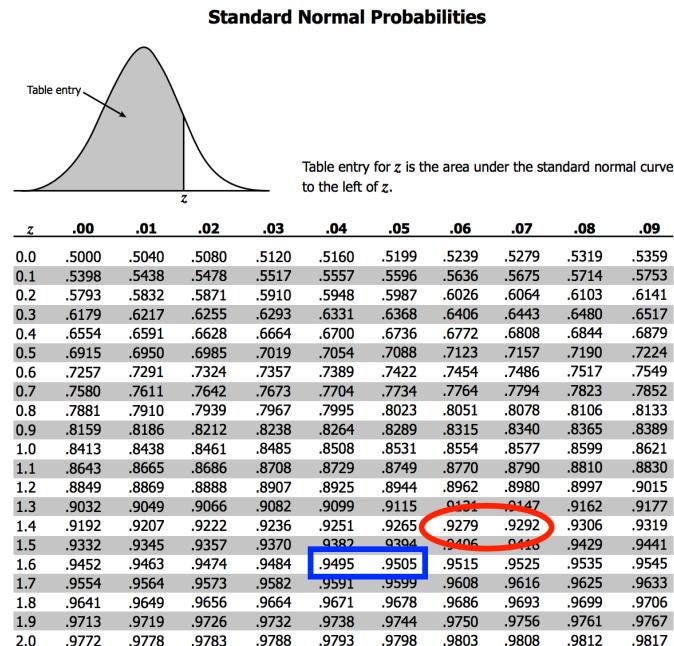
# What Does Our Z Score Mean?



- it means we believe—with 92% confidence—our babies are bigger
- therefore...our results are not statistically significant

# What Do We Do Now?

- We could get more data...
- Let's see go back to the Z table and see why...



- We need a Z score of 1.65—what sample size would be needed to reject the null hypothesis?

# Z Test (Rewritten)

$$Z = \frac{\bar{x} - \mu}{\sigma} \times \sqrt{n}$$

$$1.65 = \frac{7.8333 - 7.5}{1.25} \times \sqrt{n}$$

$$n = 38.3$$

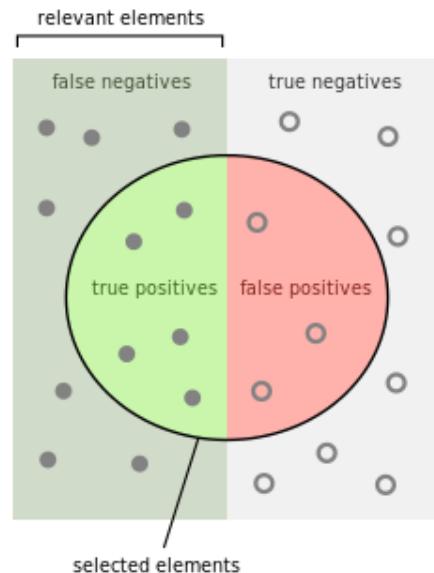
$$n = 39$$

- To be 95% confident our babies are bigger, we need a sample size of 39 (and of course the mean would have to be the same or larger)

# Understanding Classification Results

# Precision, Recall and Accuracy of a Model

How accurate are we? How many things might we have missed?



How many selected items are relevant?  
How many relevant items are selected?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

# Example: Hypothetical Terrorist Identification System

- It is guaranteed not to miss a single terrorist
- **You're all terrorists**
- We did not miss a single terrorist
- PROBLEM: Too many false positives
- Let's say there are 100 terrorists in the country, and 350 million people. Our system got them all, but wasn't very precise!
- Recall is a function of how many we got right—we got all 100 terrorists so recall was great!

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

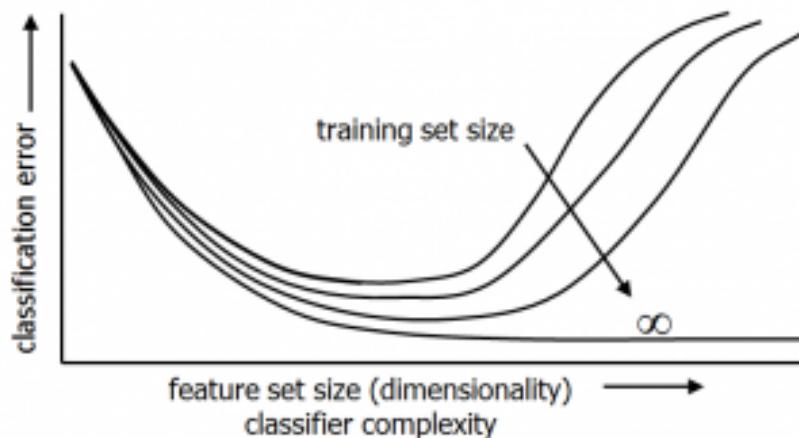
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

In what kind of real-world system would we prefer recall over precision?

# Curse of Dimensionality

# Curse of Dimensionality

- Our intuition is that more is always better, but that is not always true
- On the other hand, many things that were not possible before are now possible with Big Data
- The more features you look at, the more data you need to justify it, the harder it is to draw conclusions
- In other words, as the dimensionality increases, the volume of the space increases so fast that the available data become sparse
- Hughes Phenomenon: predictive power of a classifier/regressor increases as number of dimensions/features considered increases, then decreases



# Example Dataset

$X^T = \{x_1, \dots, x_N\}, N = 20$  (i.e., observations)

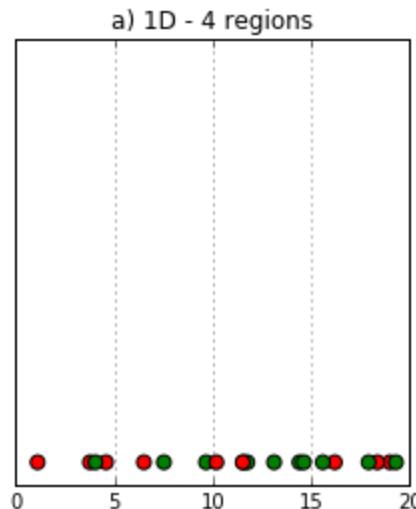
$T^T = \{t_1, \dots, t_N\}, t \in \{g, r\}$

(i.e., target class of each input, either red or green)

- we want a classifier to output the probability that each  $x$  is green

$$p(t_n = g | x_n)$$

- as a first approach, we'll partition X into four equal-sized regions



$R_1$ : 4 dots, 1 green, 3 red

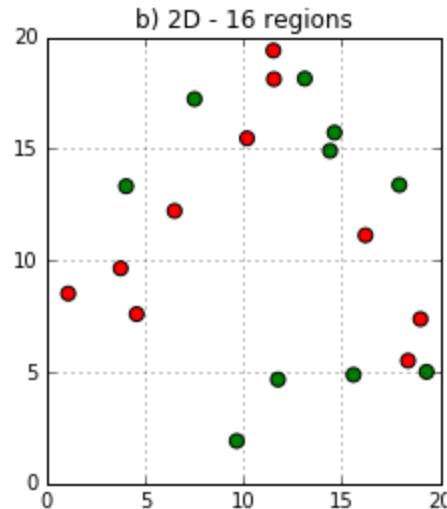
$R_3$ : 7 dots, 4 green, 3 red

$R_2$ : 3 dots, 2 green, 1 red

$R_4$ : 6 dots, 3 green, 3 red

$$p(t_n = g | R_1) = \frac{1}{4} = 0.25$$

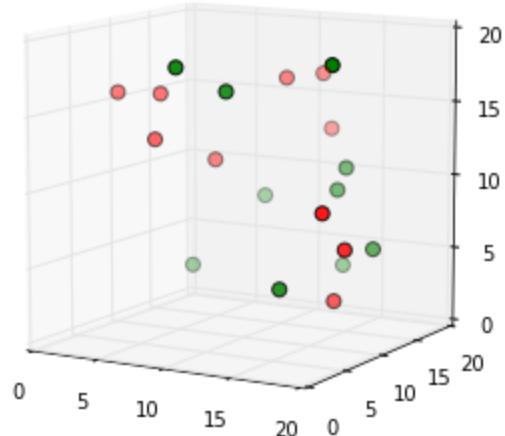
$$\frac{20}{4} = 5 \text{ obs per region on average}$$



$$X^T = \{(x_{11}, x_{12}), \dots, (x_{N1}, x_{N2})\}$$

$$\frac{20}{16} = 1.25 \text{ obs per region}$$

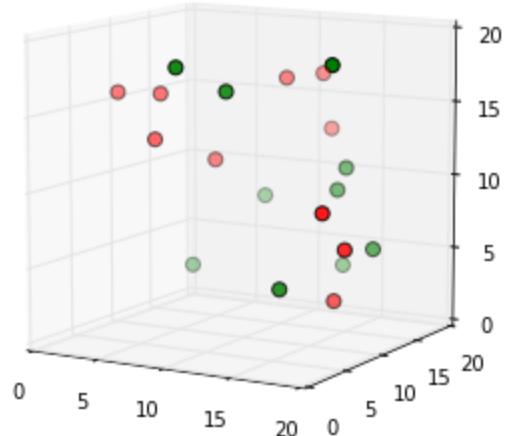
c) 3D - 64 regions



$$X^T = \{(x_{11}, x_{12}, x_{13}), \dots, (x_{N1}, x_{N2}, x_{N3})\}$$

$$\frac{20}{64} \approx 0.31 \text{ obs per region}$$

c) 3D - 64 regions

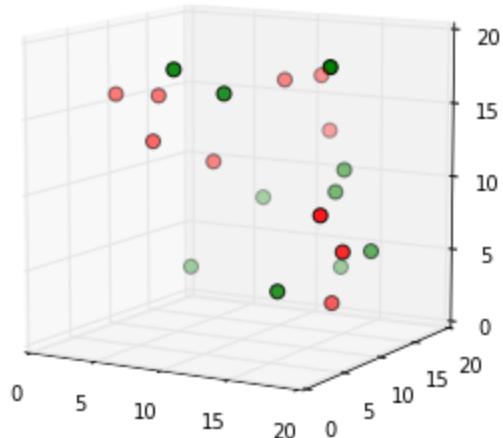


Sampling density is proportional to  $N^{\frac{1}{D}}$

$$20^{\frac{1}{1}} = 8000^{\frac{1}{3}}$$

As the number of dimensions increase the number of samples we need grows exponentially

c) 3D - 64 regions



Sampling density is proportional to  $N^{\frac{1}{D}}$

$$20^{\frac{1}{1}} = 8000^{\frac{1}{3}}$$

As the number of dimensions increase the number of samples we need grows exponentially

We need enough dimensions to have ample variation to detect the classification of each input–too many dimensions requires too much data, two few does not have enough feature variation to detect anything