# Data Science/AI Capabilities



*Word cloud of this document*

# Professor Prem Saggar

Crypto

DARPA

NATIONAL SECURITY AGENCY
UNITED STATES OF AMERICA

FINRA

Microsoft

# Twins, Triplets, and Crypto

- Twins & Triplets all have the same birthday 4 years apart
- All "6" boys
- 1/33*1/60,000*1/365 = 1/722 million <=  chances < 1/33*1/200,000,000*1/365 = 1/2.4 Trillion
- Co-occurrence? → 3,805,024,570 → 3,805,024,570 /722,700,000 = 5 others are possible in the world → 1 of 5 on Earth ☺ or 1/1 in maybe the universe

# Data Science Defined

- **Data science** is a multi-disciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from structured and unstructured data.
- The multiple disciplines are:

  - Mathematics

    - Calculus and Partial Differential Equations, …

  - Statistics

    - Distributions, missing data, regression techniques, predictions about a population based on a sample …

  - Computer Science

    - Programming, algorithms, complexity analysis, BIG DATA …

  - Information Science

    - Librarians, cataloging, information discovery, information retrieval, …

# Basics of Statistics

- Statistics is the science of making a prediction about an entire population based on a sample.

  - E.g.,

    - Predicting when the next attack will occur

    - Predicting who will win the next Presidential Election

  - How is this done? Through sampling, probabilities, and hypothesis testing.
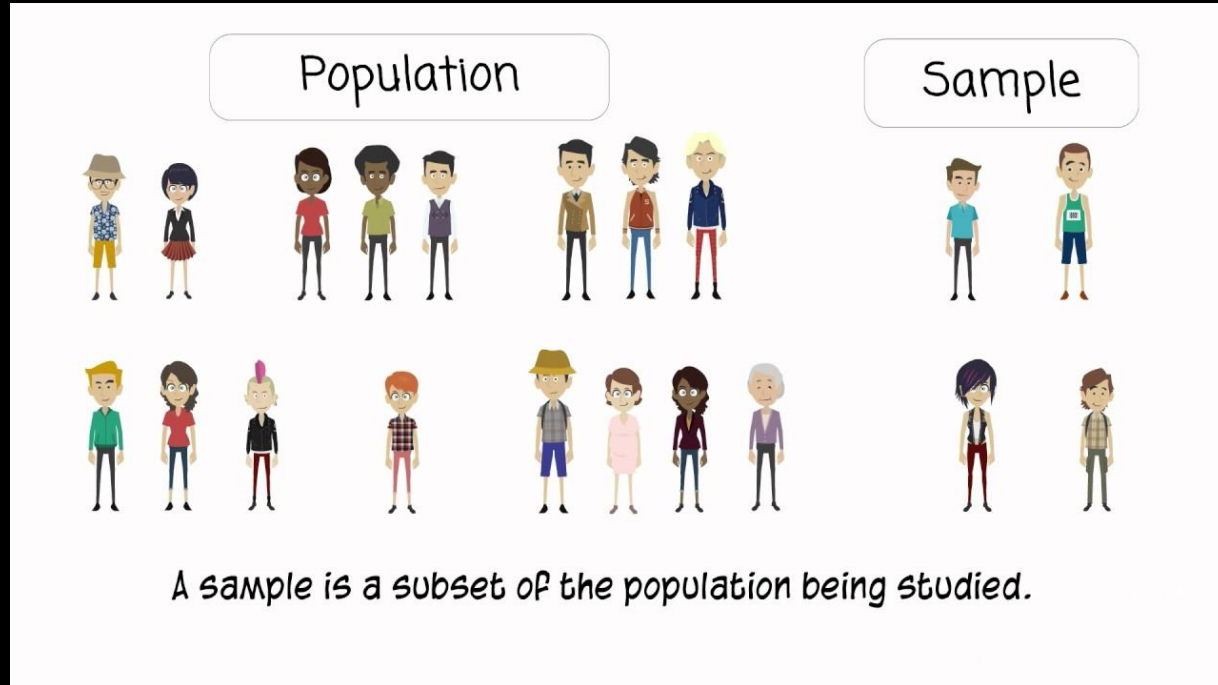
Statistics = Use random samples to make confident statements about entire populations. Intelligent guesses/speculation.
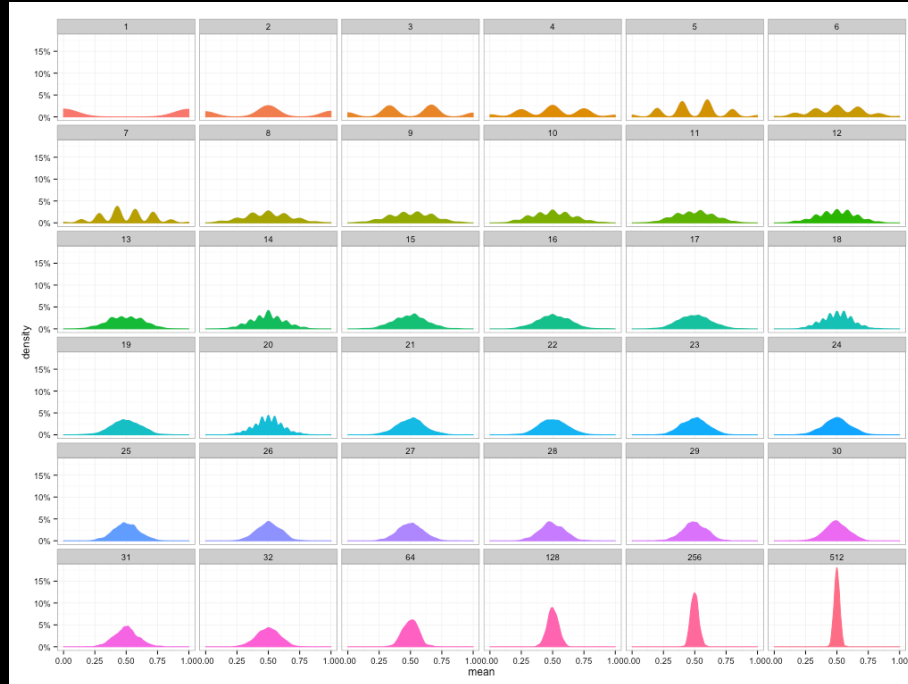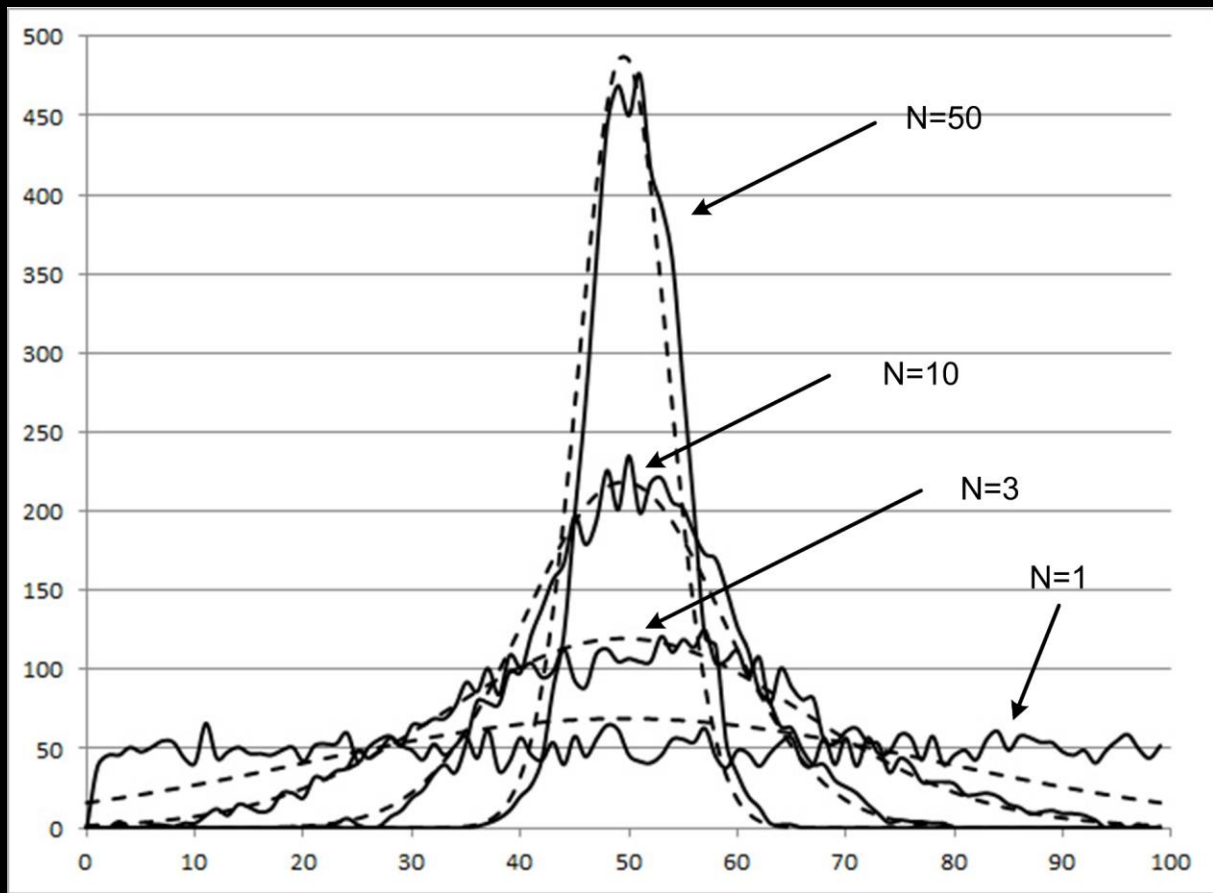
# Sample, shape, location, and spread

- Sample = make sure it's random, handle missing data (mcar, mar, nmar), imputation methods. NMAR!
- Shape = Is the data skewed, normal, or flat? If normal then we can use statistical analysis for normal distributions
- Location = Where does the data accumulate? Is it skewed, if so the median tells a better story.
- Spread = How much does the data differ? Standard deviations!

# Samples give statistics, and populations (which we may never know give parameters)



A sample is a subset of the population being studied.

# The Central Limit Theorem

We're all ██████████ until proven ████████████.
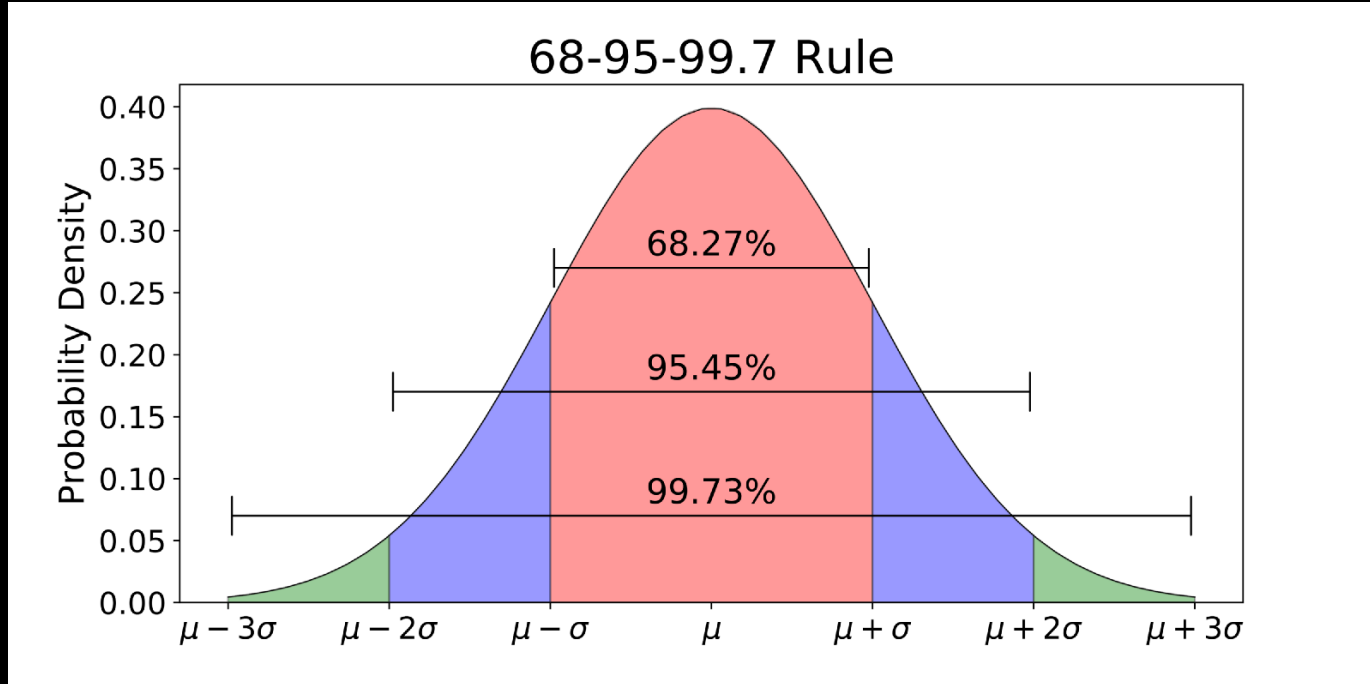
# Hypothesis Testing
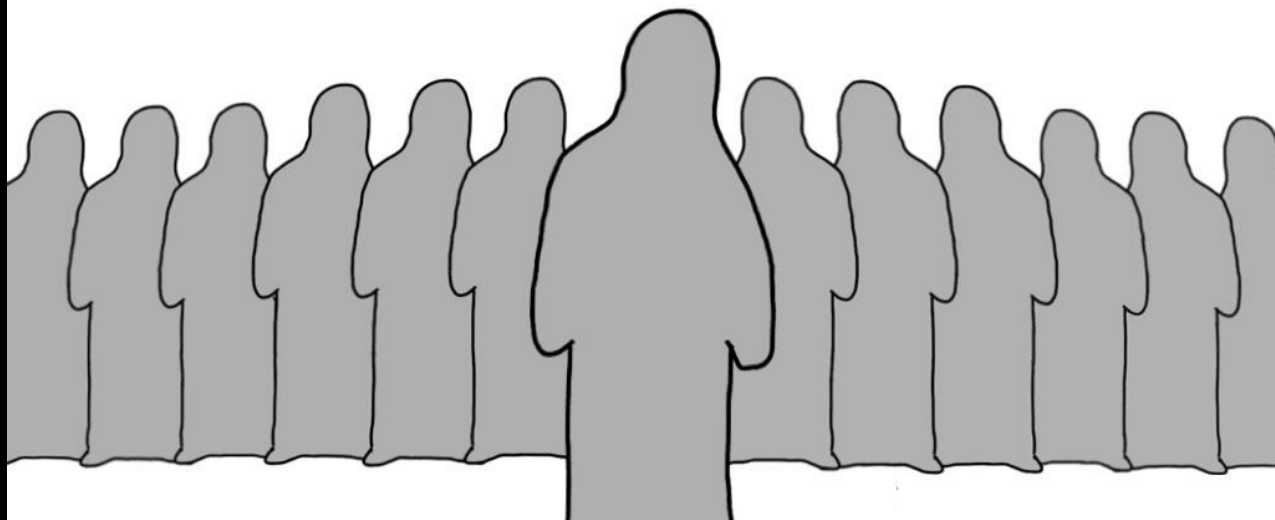
# How do we prove our statistical claims?

- Hypothesis Testing! The point of hypothesis testing is to make sure we don't jump to bad conclusions. Conclusions can be confusing (xylitol vs. fluoride). We are inherently speculating albeit rigorously. So we try to control the guessing by being conservative and use innocent until proven guilty.

# Standard deviations and probability of the population mean

# The alternative hypothesis tries to nullify the ghosts!

# Warning

- It's important to note that statistics and inferences about "populations" are always approximations. We aren't using probabilities when we do have the entire population e.g., the entire Data Science class is our population. Here I can get the population mean etc... No guess work is necessary. Now what if we said all data science students in the world currently? That's much harder if not impossible. Inference time!
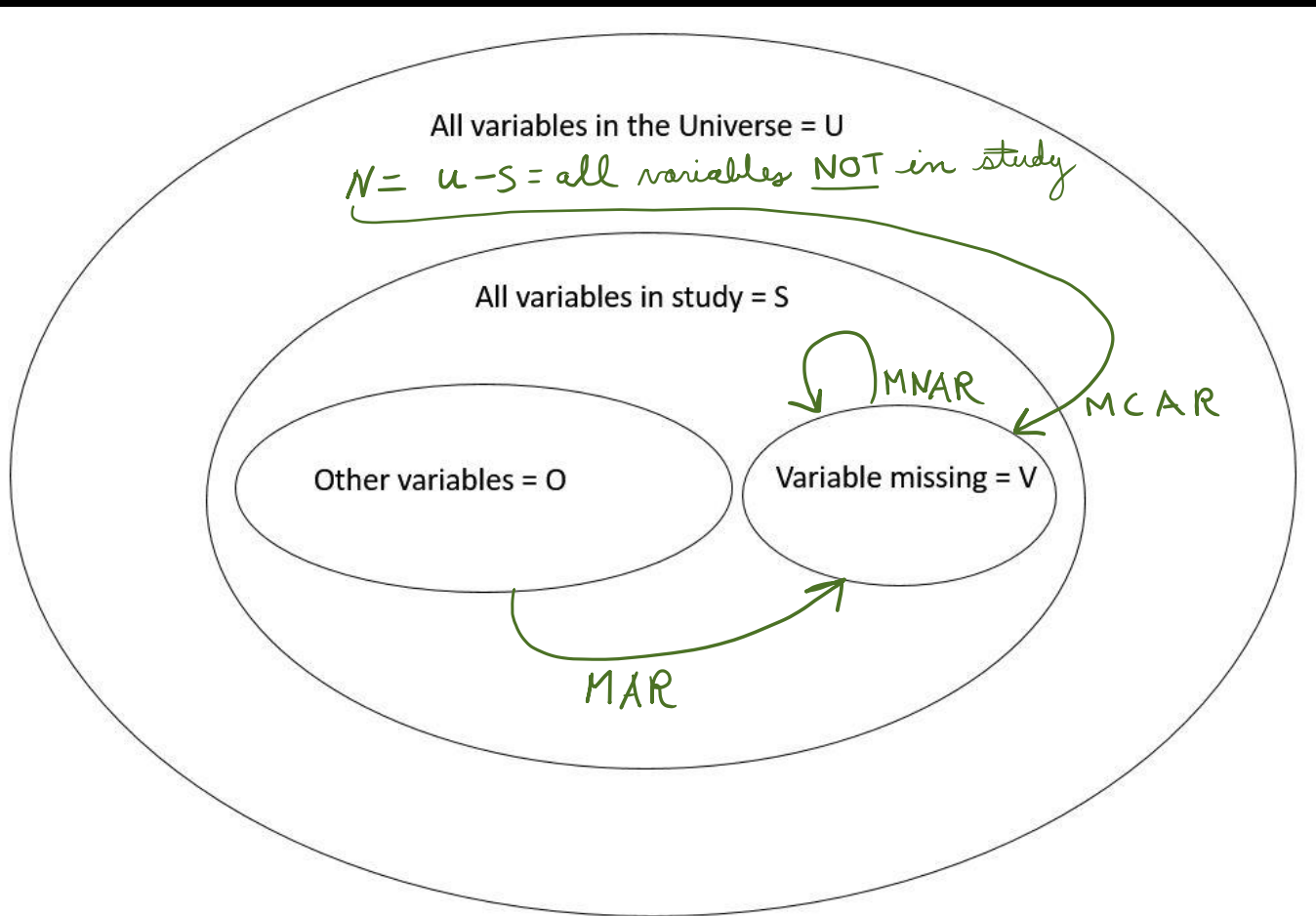
# Survey

➢ Assumption: student will find Waldo given enough time.

| Student | Time to find Waldo | Maximum Time |
|---|---|---|
| | | 10 |
| | | 10 |
| | | 10 |
| | | 10 |
| | | 10 |

# Missing data

| Missing Data | Cause of missing | Example | Type | Impute outcome | Imputable | Action |
|---|---|---|---|---|---|---|
| V | N | Variable missing because of freak event, or any variable you did not model | MCAR | freak event is random, so imputation is not biased | TRUE | Rerun if possible |
| V | O | Variable missing because of variable(s) you DID model, e.g., blood not drawn by clinician so verdict not possible or unable to find something based on time | MAR | 50/50 chances for missing data is modelled so imputation is ok | TRUE | Rerun if possible |
| V | V | Variable missing because of the variable! I'm not telling you, e.g., what is your income | MNAR | skewed incomes everyone looks middle class | FALSE | Change survey |

# Unsupervised vs. Supervised Learning

Supervised learning is aided by training data and human correction.

Here's some training data. Learn the patterns. Make your best guess at what the patterns are. We'll feed you test data to figure out if you've understood it. If you stray off course, we'll correct you and retrain.

Examples include Decision Trees and Neural Networks.

Unsupervised learning is uncorrected and runs on data. It can't classify things "yet". But is very good at clustering and anomaly detection.

CAT    DOG

IS THIS A
**CAT** or **DOG**?

OUTPUT
LAYER

What's
"deep"
about it?

ACTIVATED
NEURONS

INPUT
LAYER

DEEP
neural
NETWORK

# Decision Tree Example: "BigTip"



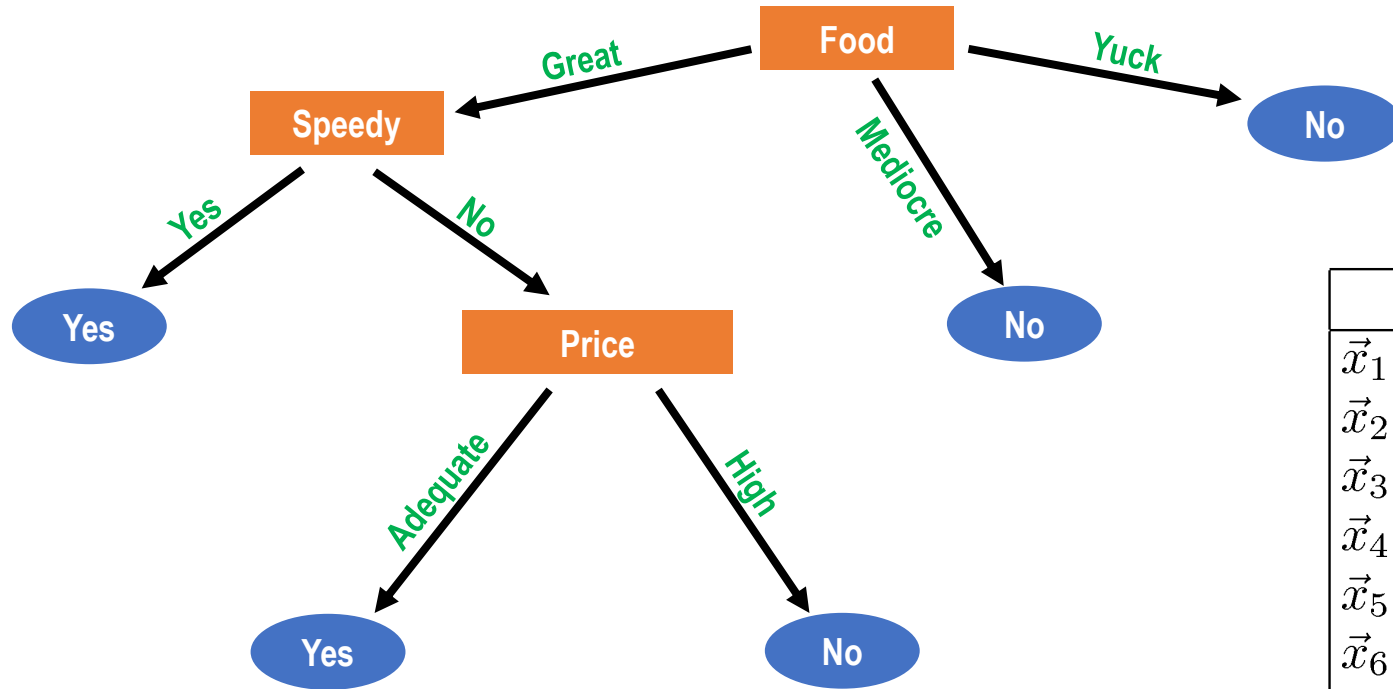| | F  S  P | BigTip |
|---|---|---|
| $\vec{x}_1 = ($ g, y, a $)$ | | $f(\vec{x}_1) = 1$ |
| $\vec{x}_2 = ($ g, n, h $)$ | | $f(\vec{x}_2) = 0$ |
| $\vec{x}_3 = ($ g, y, h $)$ | | $f(\vec{x}_3) = 1$ |
| $\vec{x}_4 = ($ g, n, a $)$ | | $f(\vec{x}_4) = 1$ |
| $\vec{x}_5 = ($ m, y, a $)$ | | $f(\vec{x}_5) = 0$ |
| $\vec{x}_6 = ($ y, y, a $)$ | | $f(\vec{x}_6) = 0$ |
| $\vec{x}_7 = ($ g, y, a $)$ | | $f(\vec{x}_7) = 1$ |
| $\vec{x}_8 = ($ g, y, h $)$ | | $f(\vec{x}_8) = 1$ |
| $\vec{x}_9 = ($ m, y, a $)$ | | $f(\vec{x}_9) = 0$ |
| $\vec{x}_{10} = ($ g, y, a $)$ | | $f(\vec{x}_{10}) = 1$ |

Our Data

# Association Rules!!!

| Buying Pattern |
| --- |
| Beer, Orange Juice, Diapers |
| Beer, Q-Tips, Diapers |
| Beer, Chips |
| Beer, Advil |

Derived rules
Beer → Diaper 50% confidence
Diapers → Beer 100% confidence

# Association Rules example for financial transactions.

Please consider the following mock financial contract data. What patterns can we find for Great customer acceptance?

| | above 10 Million | delivery | QA | Acceptance |
|---|---|---|---|---|
| 0 | Contract above 10 million | no set delivery times | has Quality Assurance deliverables | Great customer acceptance |
| 1 | Contract above 10 million | deliveries bi-weekly | no Quality Assurance deliverables | Average customer acceptance |
| 2 | Contract above 10 million | deliveries bi-weekly | has Quality Assurance deliverables | Great customer acceptance |
| 3 | Contract above 10 million | no set delivery times | no Quality Assurance deliverables | Poor customer acceptance |
| 4 | Contract less than 10 million | deliveries bi-weekly | has Quality Assurance deliverables | Great customer acceptance |
| 5 | Contract less than 10 million | deliveries bi-weekly | no Quality Assurance deliverables | Average customer acceptance |

There are actually 190 rules we can derive and here are some key ones:
1) has Quality Assurance Deliverables → Great Customer acceptance @ **100%** confidence (happens every time)
2) no Quality Assurance Deliverables → Average customer acceptance @ **66.7%** confidence (happens 2 out of every 3 times)
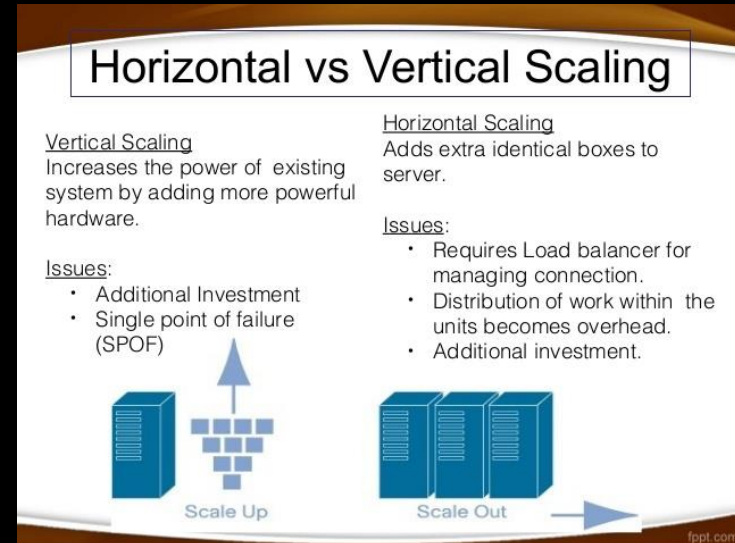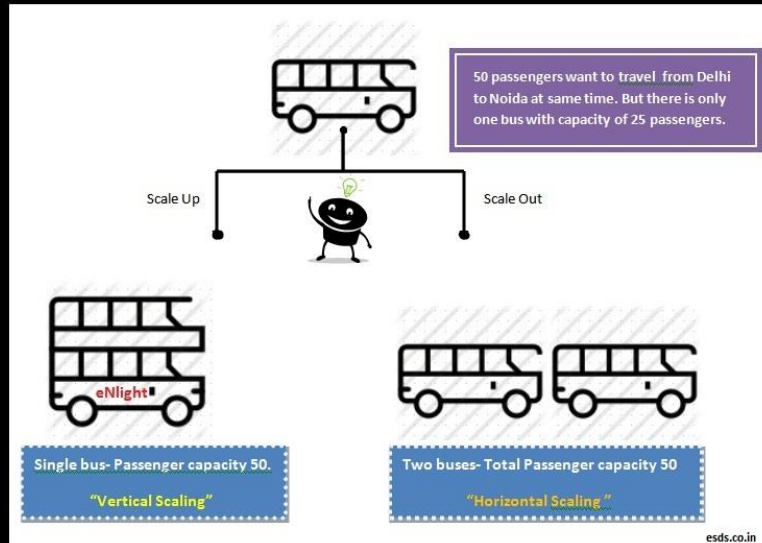
# What can we do with Association Rules?

1. Allow us to objectively report these are the rules within our data
2. Discover new behaviors we weren't aware of
3. Finds rules with confidences that are not 100% but are still important
4. After examining the rules found we can say these are the target outcomes we're interested in e.g., Great Customer Acceptance and Contracts less than 10 million.
5. Once we have target outcomes, we can use Decision Trees to illustrate and predict rules for these target outcomes.

# Interview Problem for data scientists

- Given the following:

  - 3,6,4 = 18, 108, 648, 3888

  - 2,-2,5 = -4,8,-16,32,-64
- What is

  - 2,3,4 = ?
- Is there a formula for the last term?
- Can you write a program to get the last term?
- Can you prove your program works?
- What is the time complexity?

# Big Data

➢ What is Big Data? Is it Dangerous ☺?

➢ Well what is Big?

    ➢ Big is a relative term here and is assumed to be big by "today's" standards. However, it really means too "big" for you to handle via vertical scaling. By vertical we mean single node improvement vs multi node horizontal improvement.

# The 4 Vs of Big Data

- ➢ Volume

  - ➢ How will your system deal with storage as data approaches ∞?
- ➢ Velocity

  - ➢ How will your system deal with fast responses as requests approach ∞?
- ➢ Variety

  - ➢ How will your system deal with different types of data as the types approach ∞?
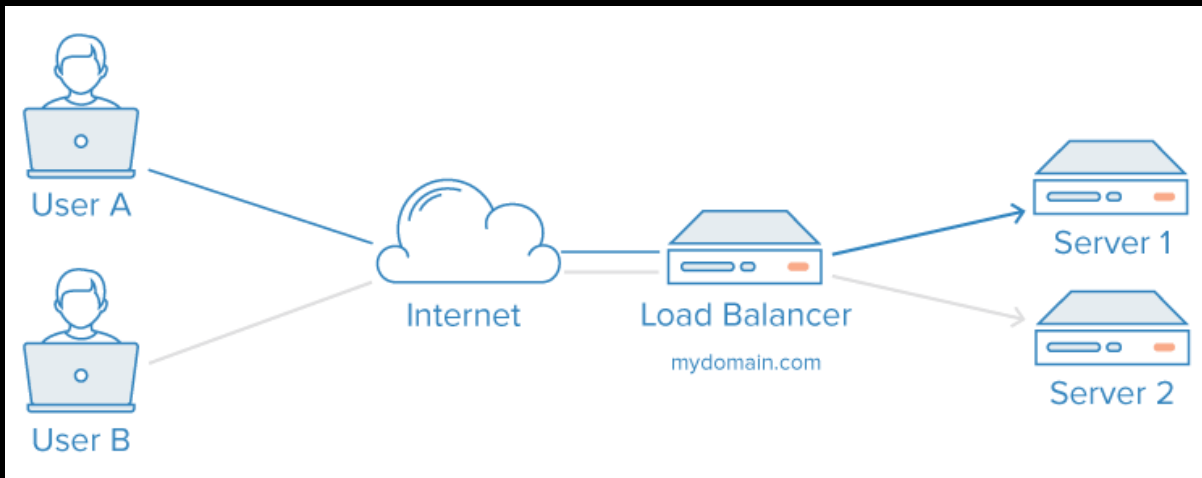- ➢ Veracity

  - ➢ How will your system ensure the accuracy of your Big Data is not crushed as the other 3 Vs approach ∞?

# Volume

➢ Traditional approaches were to add more disks or (Storage area networks) SANs in data centers. This did not scale and the SANs were not geographically fault tolerant.

➢ Current approach is to use Cloud storage that is practically limitless and is geographically fault tolerant. Many cloud storage providers exist making it financially practical to move storage to the cloud.

# Velocity

➢ How do we deal with millions → infinite concurrent requests?

    ➢ Traditionally we'd add load balancers and scale up the webserver.

    ➢ Now cloud providers have geographic load balancers (that are themselves load balanced) and webservers that are horizontally scaled

# Variety

➤ How do we deal with data changing and different data sets?

➤ Relational databases tend to crack with many changes (ultra structured so change is costly)

➤ Extraction, Transform, and Load (ETL) is cumbersome and expensive to run when data moves or is changed?

➤ Current solutions make data formatting easy and impose minimal structure (put this type of data in this folder).

➤ Data Lake → Folders ←→ Tables