

Data Mining

DAPT 631

Vishal Patel

Spring 2025



- Vishal
- I run my own Data Science practice at **DERIVE, LLC** since **2016**
- **MS in Computer Science** 2003 (IIT, Chicago), and **MS in Decision Sciences** 2012 (VCU, Richmond)
- Mining data since **2003**





ActiveCampaign >



JAGUAR



○ Introduction

○ History

○ Course Structure

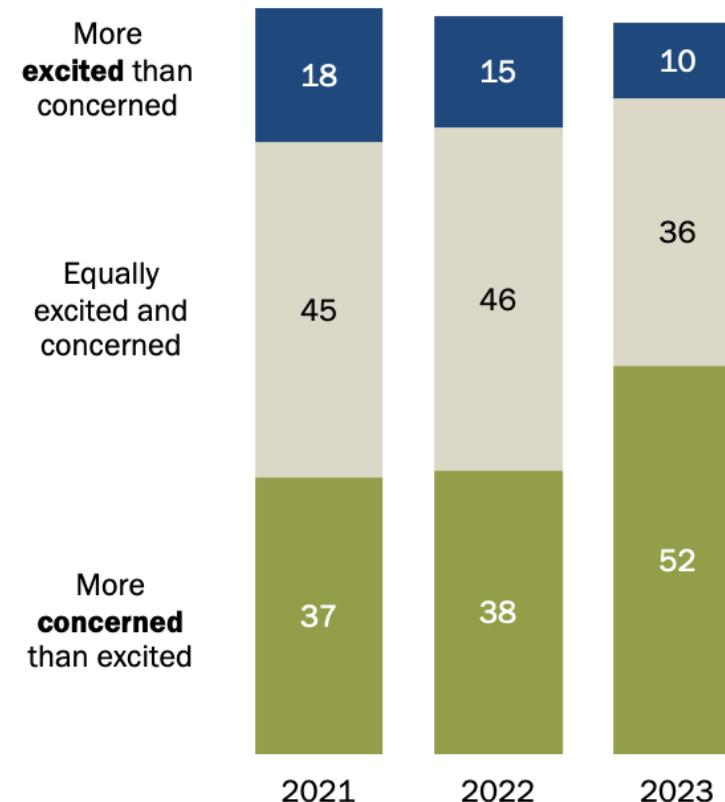
AI: Hype or Tripe?

How do you feel about AI?

- A. Excited
- B. Concerned
- C. Other

Concern about artificial intelligence in daily life far outweighs excitement

% of U.S. adults who say the increased use of artificial intelligence in daily life makes them feel ...



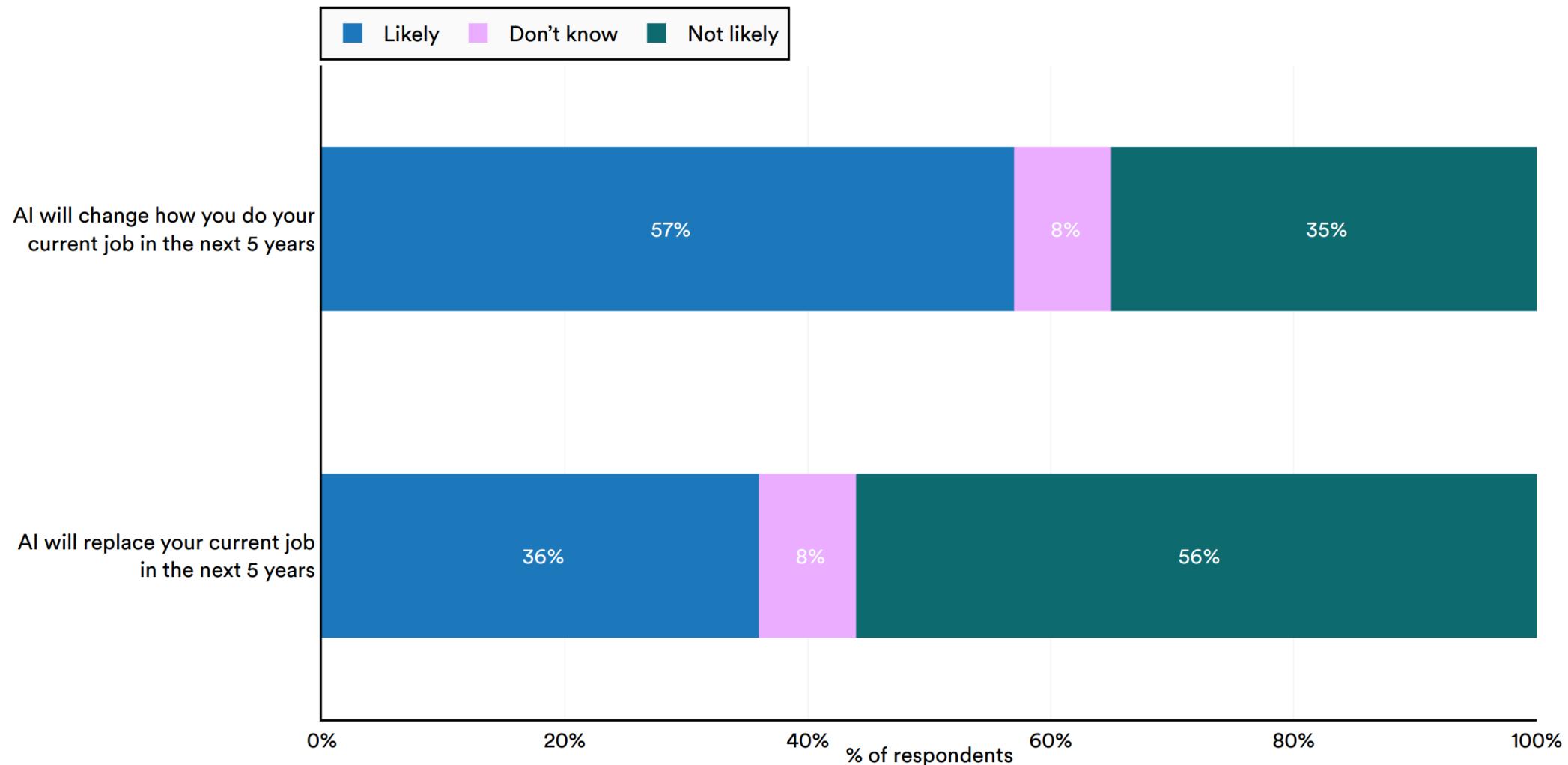
Note: Respondents who did not give an answer are not shown.

Source: Survey conducted July 31-Aug. 6, 2023.

PEW RESEARCH CENTER

Global opinions on the impact of AI on current jobs, 2023

Source: Ipsos, 2023 | Chart: 2024 AI Index report



The survey consisted interviews with 22,816 adults in 31 countries between May and June 2023.

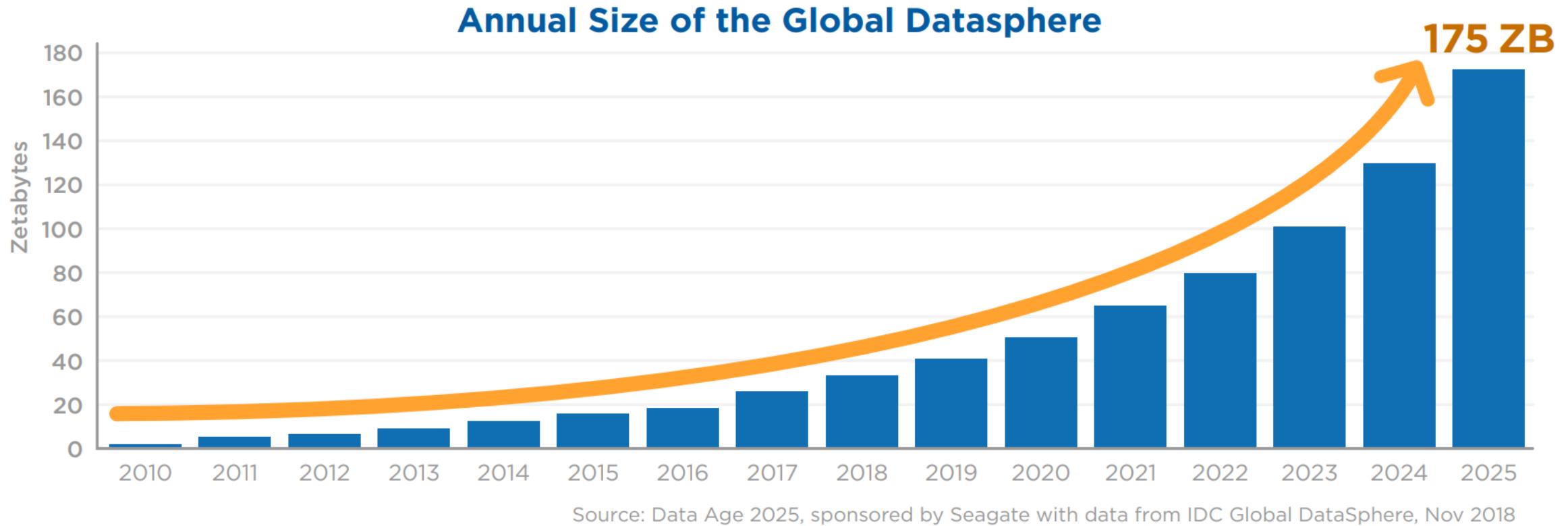
○ Introduction

○ History

○ Course Structure

Cambrian Era



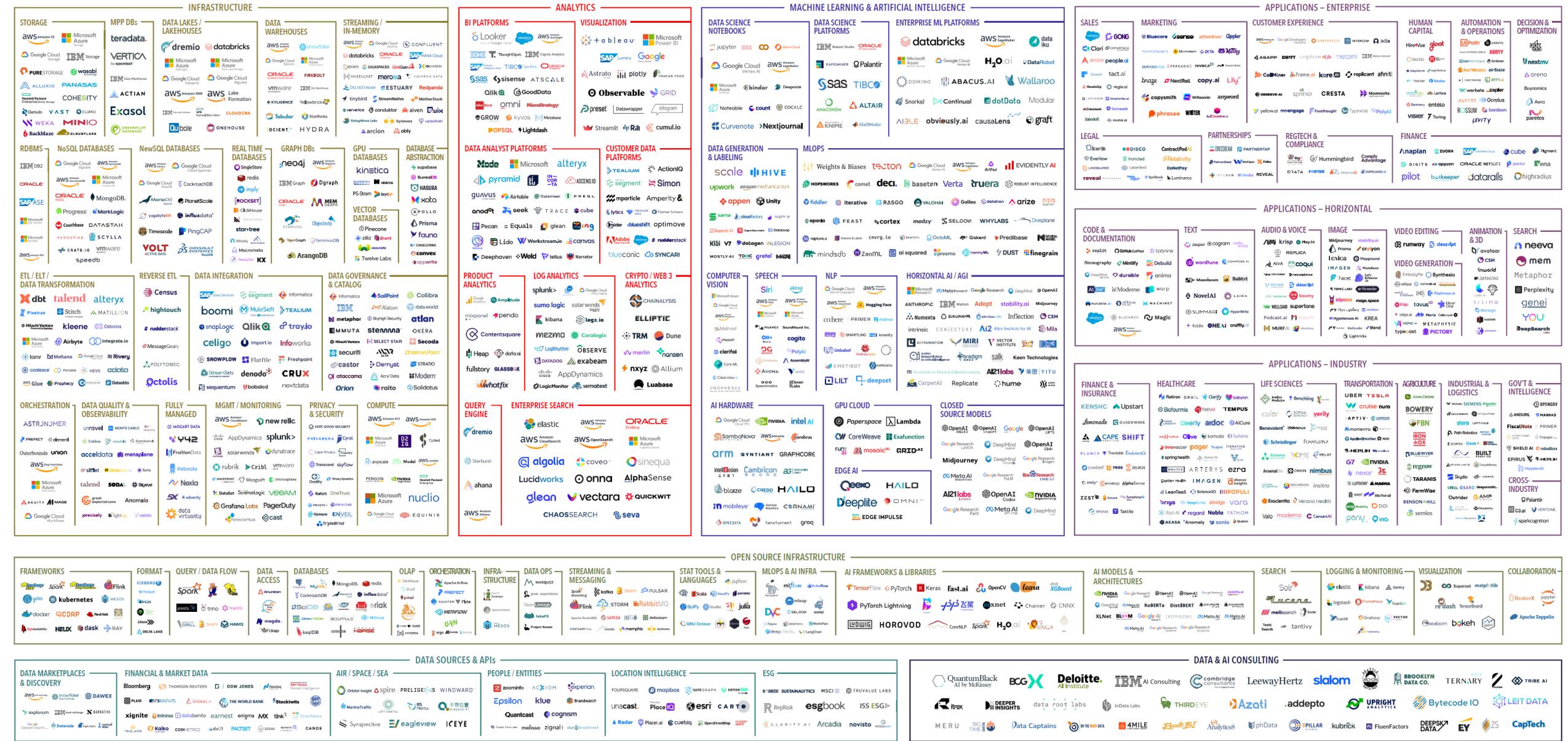


One zettabyte = One trillion gigabytes = One sextillion (10^{21}) bytes

For comparison, the universe is 4×10^{17} seconds old.



THE 2023 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE





Watch later



Share

Era of Data Literacy

CONTINUUM[®]
ANALYTICS

- Data exploration and analysis are going to be a new kind of **literacy** that will be required to do great work in any field.
- Language is a human instinct and is a natural path to insight. We see this in our interaction with Python/PyData users, whose passion chiefly stems from this *expressiveness* and *agility*.
- An analytical language is “**thoughtware**”, not “software”.

ANAconda



PyData
DC 2016

What is Data Mining?

Data mining is the process of **discovering patterns** in large data sets involving methods at the intersection of machine learning, statistics, and database systems.

[Wikipedia]

Data mining is the **extraction** of implicit, previously-unknown, and potentially-useful **information** from data.

– Witten and Frank

Data mining is the process of **discovering** meaningful new correlations, patterns and trends by sifting through large amounts of data stored in repositories, using pattern recognition technologies as well as statistical and mathematical techniques.

– Gartner

Dictionary

data mining



data mining

noun COMPUTING

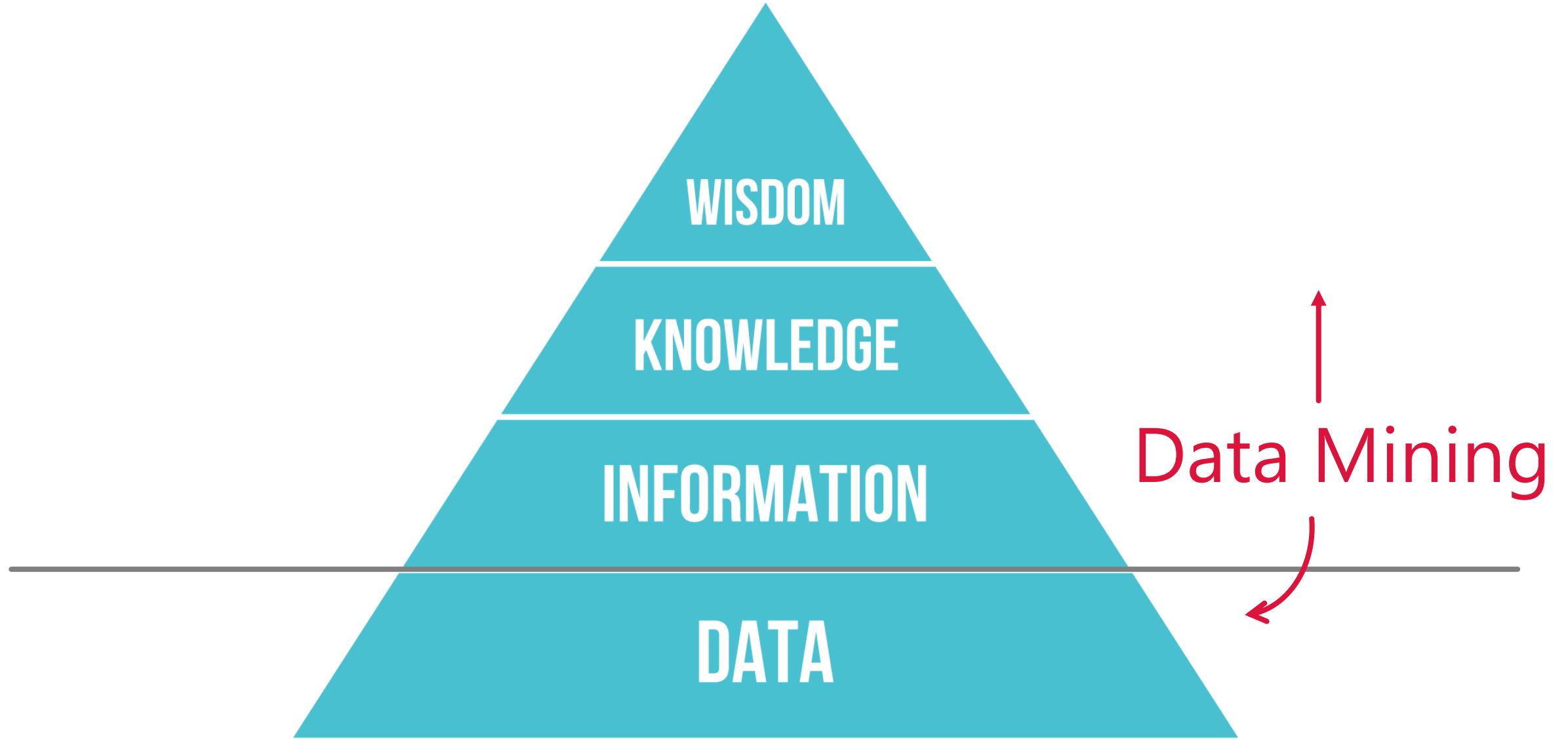
noun: data mining; noun: datamining

the practice of examining large pre-existing databases in order to generate new information.



ChatGPT

Data mining is a process of discovering patterns, trends, and valuable insights or knowledge from large volumes of data. It involves using various techniques and algorithms to analyze and extract meaningful information from datasets, often with the goal of making informed business decisions, identifying opportunities, or solving complex problems.





Data Mining Tasks

Description

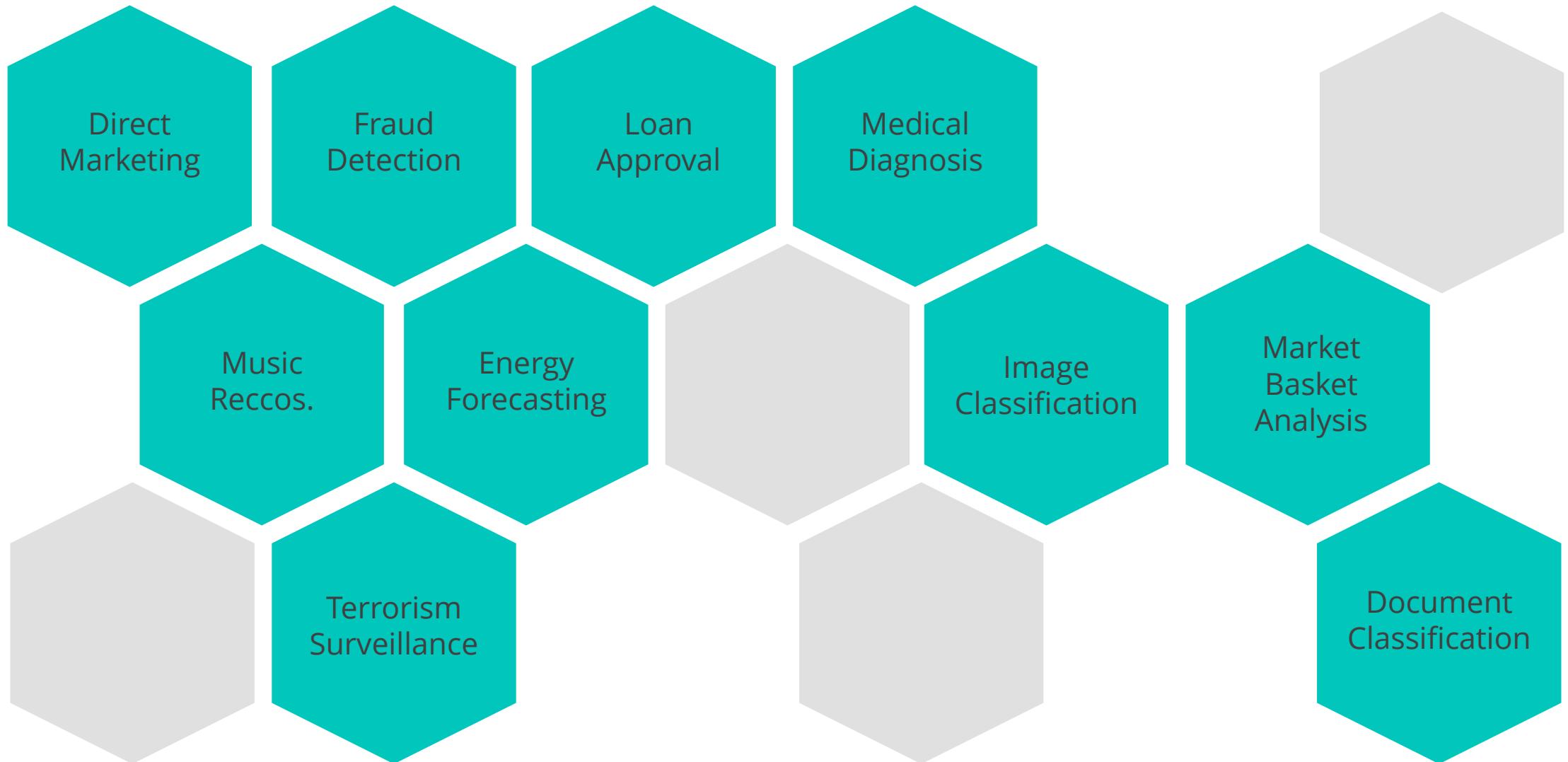
Estimation

Clustering

Classification

Prediction

Applications of Data Mining



○ Introduction

○ History

○ Course Structure

Statistics

Census
Mortality tables
Accounting

From Latin: *status* state

... teaches us what is the political arrangement
of all modern states of the world.

W Hooper, 1770

DATA COLLECTIONS + ANALYSIS + DECISION MAKING

Statistics

EXAMPLE #1: UNCERTAINTY



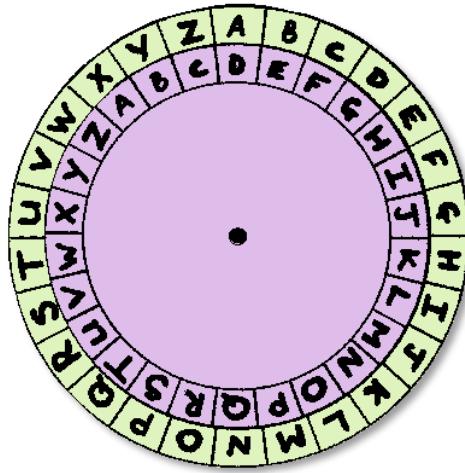
Siege of Plataea (5th Century BCE)

What do you think they used as the best estimate for the height of the wall?

- A. Mean
- B. Median
- C. Mode
- D. Max

Statistics

EXAMPLE #2 FREQUENCY ANALYSIS, CRYPTOANALYSIS



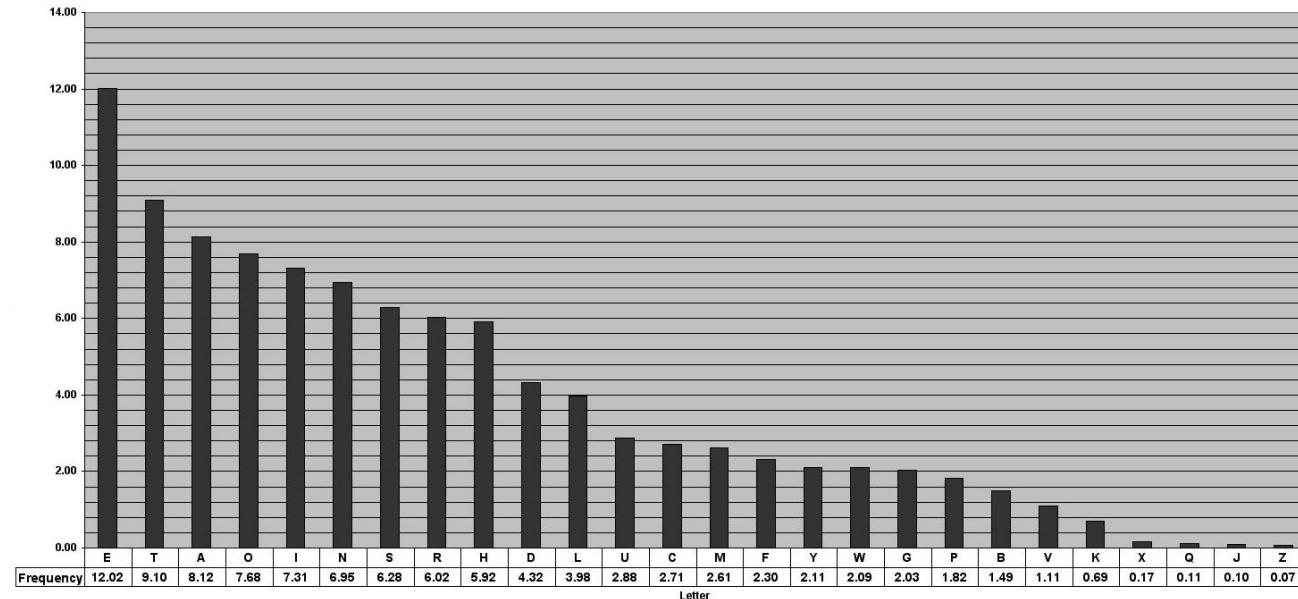
Caesar Cipher

Original message: Et tu, Brute?

Encrypted message: Hw wx, Euxwh?

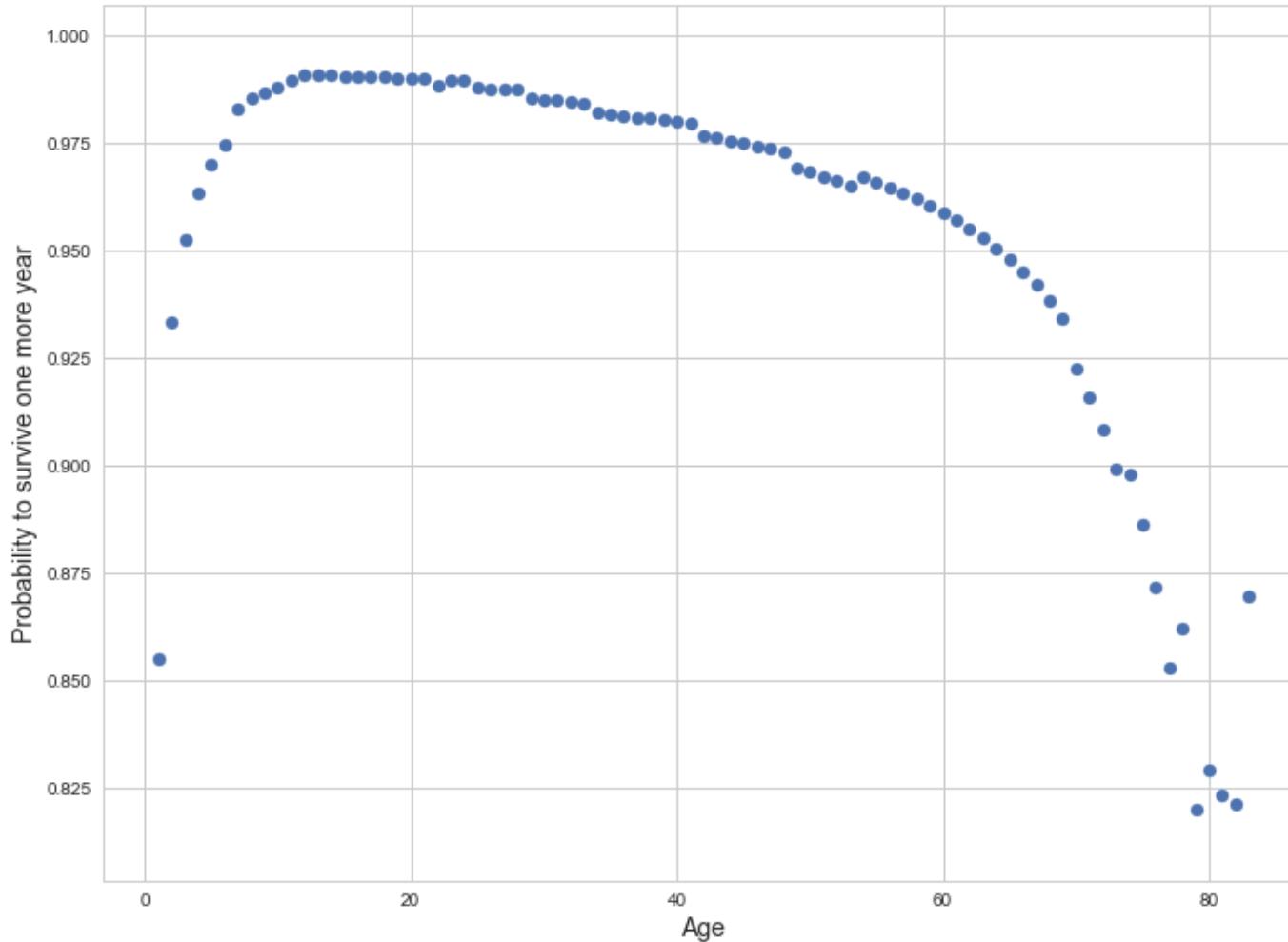


Al-Kindi
(801–873 AD)



Statistics

EXAMPLE #3 MORTALITY TABLES, DEMOGRAPHY



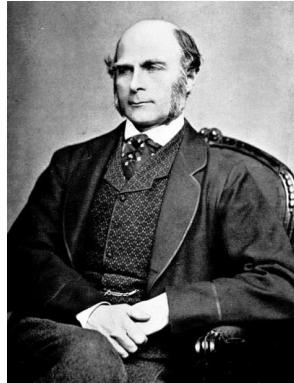
Data from Edmond Halley's *An Estimate of the Degrees of Mortality of Mankind* (1693), table p.600.

The graph shows the probability of surviving one or more year(s) at a certain age.

Modern Statistics

Normal distribution
 t distribution
Random sampling
Design of experiments
Bayesian Statistics

A rigorous mathematical discipline
for analysis, decision making, and inference



Sir Francis Galton
(1822–1911)
Correlation, regression



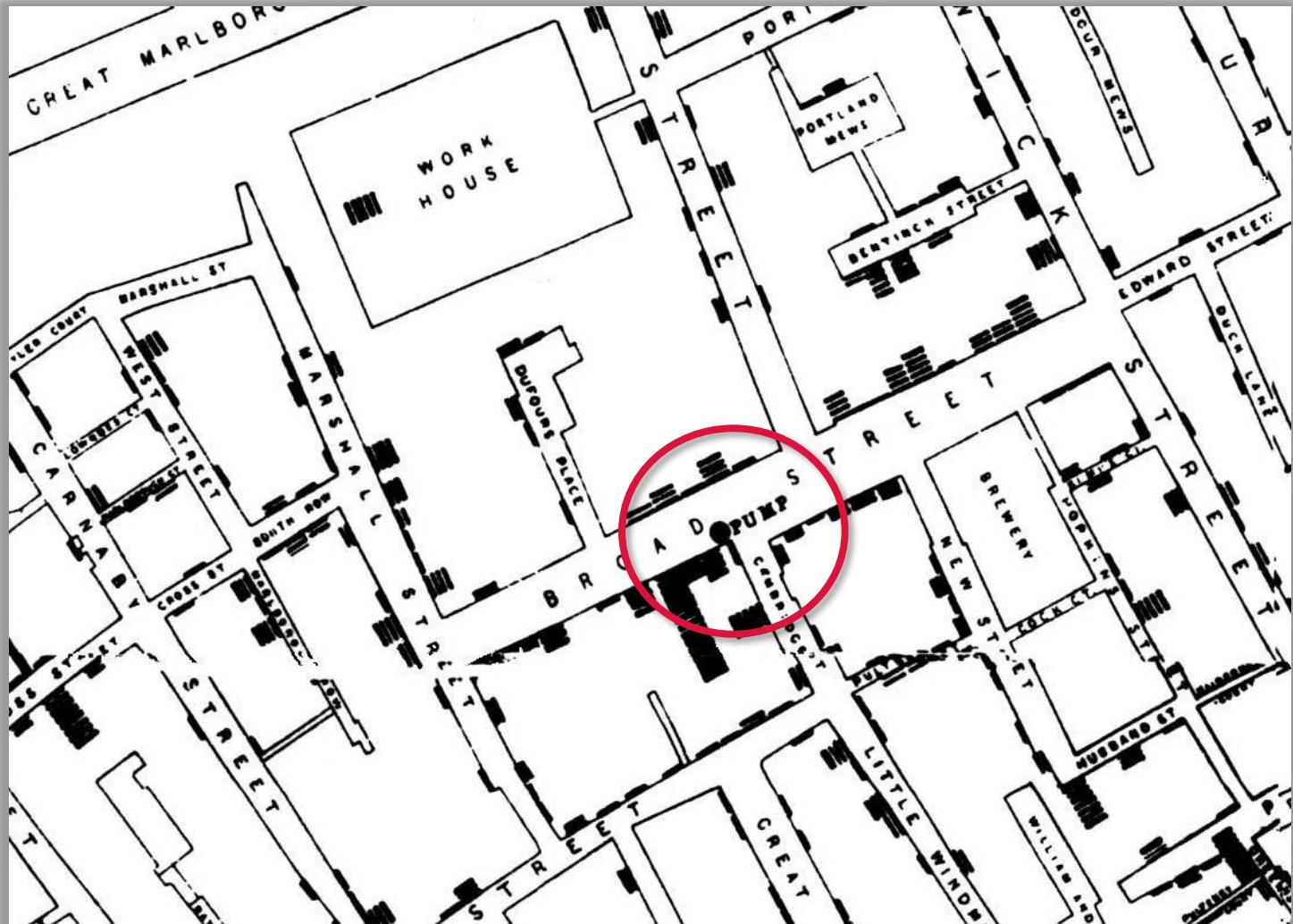
Carl Pearson
(1857–1936)
Founder of mathematical statistics



R A Fisher
(1890–1962)
ANOVA, Maximum Likelihood, DOE

Modern Statistics

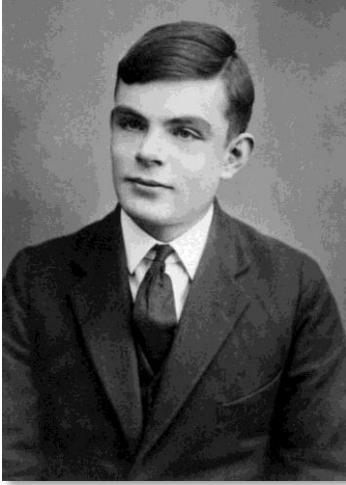
EXAMPLE: DATA VISUALIZATION



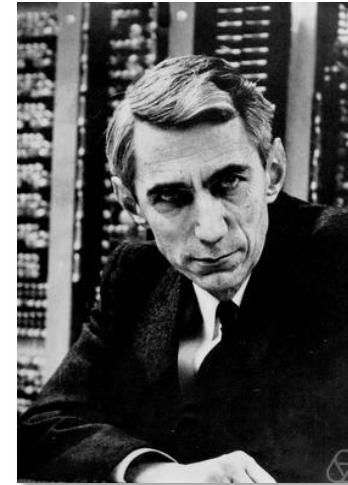
Original map by **John Snow** showing the clusters of cholera cases
in the **London epidemic of 1854** [[Source](#)]

Data Mining

Algorithms &
Computation
Computer Science
Neural Networks
Decision Trees
Genetic Algorithms
Relational Databases



Alan Turing
(1912 –1954)
Theoretical Computer Science



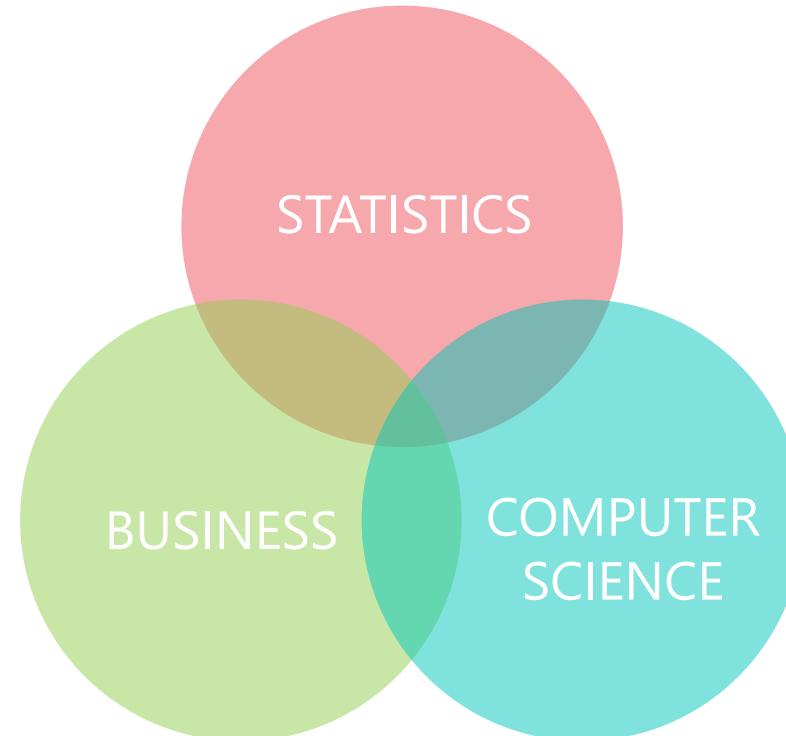
Claude Shannon
(1916 –2001)
Information Theory

- Warren McCulloch and Walter Pitts created a computational model for **neural networks**. (1943)
- John Holland introduced **Genetic Algorithm** based on the concept of Darwin's theory of evolution. (1960)
- E. F. Codd published an important paper to propose the use of a **relational database** model. (1970)

Data Science

Gradient Boosting
Random Forests
Support Vector-
Machines
Recommender-
systems
Unstructured data
Open source
Big Data

Data science is an **interdisciplinary** field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from data in various forms, both structured and unstructured, similar to data mining.[†]



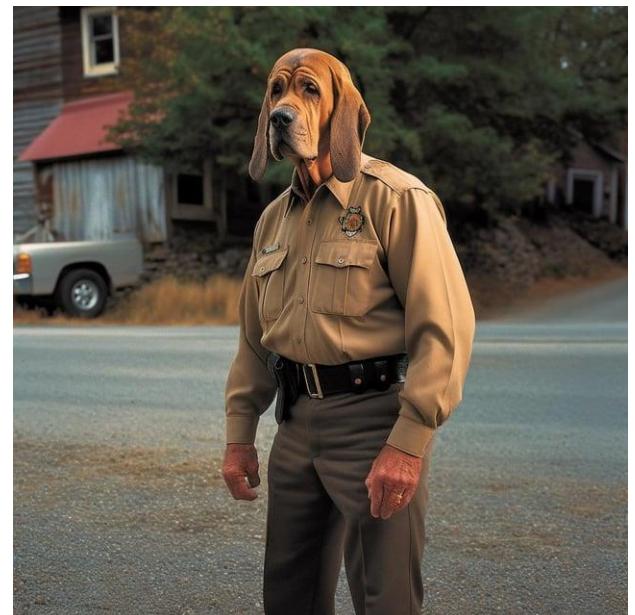
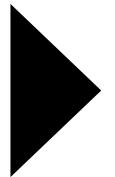
Artificial Intelligence

Deep learning
Reinforcement-
Learning
Speech recognition
Natural Language-
Processing
Computer vision



Everyone: AI art will make designers obsolete

AI accepting the job:



Jan 2023

Jun 2023

Statistics

Regression
Correlation
Frequency analysis
Descriptive statistics
ANOVA

Modern Statistics

Normal distribution
 t distribution
Random sampling
Design of Experiments
Bayesian statistics

Data Mining

Algorithms & Computation
Computer Science
Neural Networks
Decision trees
Genetic algorithms
Relational Databases

Data Science

Gradient Boosting
Random Forests
Support Vector Machines
Recommender systems
Unstructured data
Open source
Big Data

ML

Artificial Intelligence

Deep learning
Reinforcement Learning
Speech recognition
Natural Language Processing
Computer vision

Prehistory – 18th Century

Late 19th / Early 20th Century

Mid-Late 20th Century

21st Century

Calculations by hand

Distributed computing

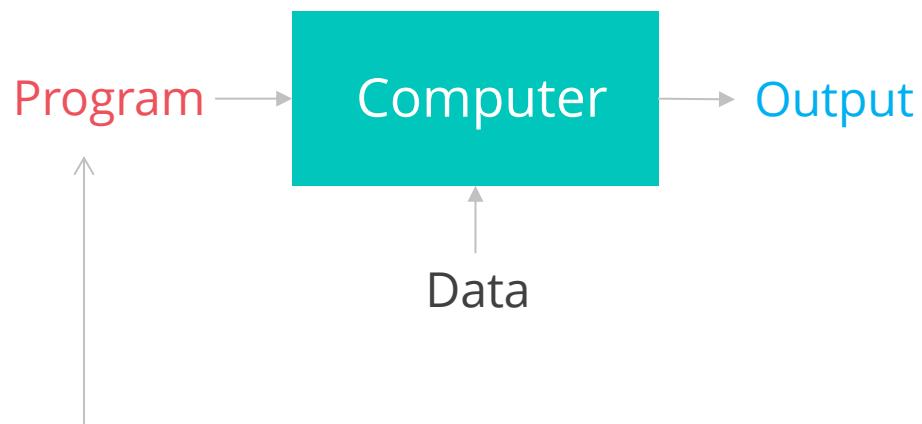
Evolution of techniques and technology

Machine Learning

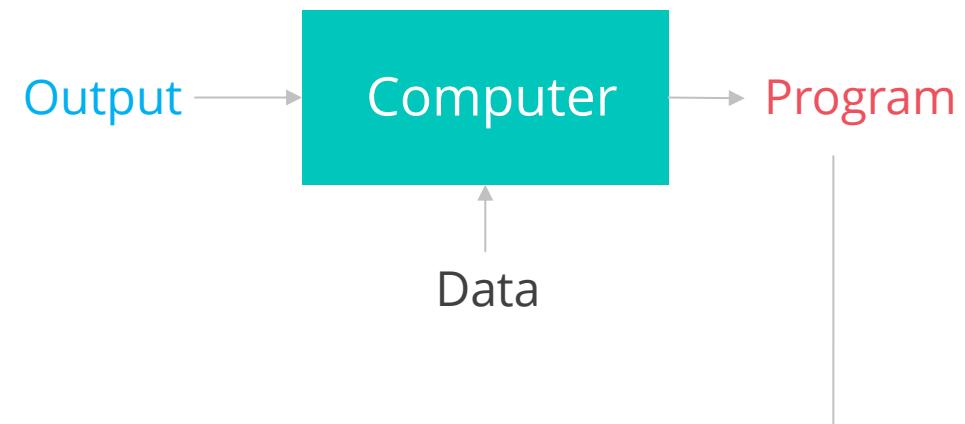
Field of study that gives computers the ability to learn
without being explicitly programmed.

Artur Samuel, 1959

Traditional Programming



Machine Learning



THIS IS YOUR MACHINE LEARNING SYSTEM?

| YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

| JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.



○ Introduction

○ History

○ Course Structure

Data Science ≈ Data Mining

- The specific definitions and boundaries between these disciplines remain fuzzy.
- For the purpose of this class, I will use the terms ‘Data Science’ and ‘Data Mining’ interchangeably (with a preference to the former).
- We will cover several Data Science techniques in this class, e.g., Gradient Boosting.

Two Cultures

Statistics

THEORETICAL

INFERENCE

ASSUMPTIONS

MANUAL

AUTOMATION

Data Science

PRACTICAL

PREDICTION

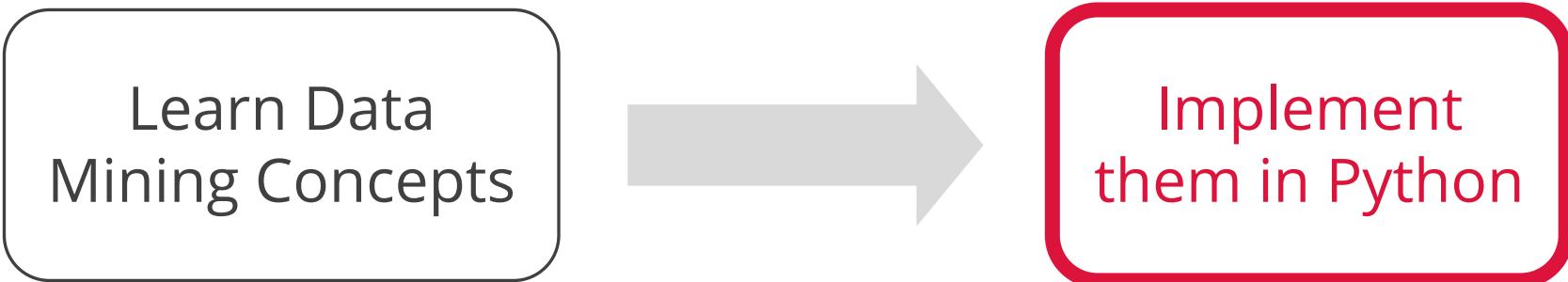
EMPIRICAL

Course Outline

1. Introduction
2. The Data Science Process
3. Supervised Learning
4. Unsupervised Learning
5. Wrap Up

1. Ask **questions** at any time!
2. **Collaboration** is encouraged.
3. All course content will be available on Canvas and GitHub.
4. Data Mining + Python
5. Homework assignments in **Python**

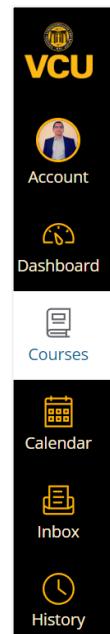
Data Mining + Python



Course Objectives

1. Provide a **practical** knowledge of data mining algorithms.
2. Give a **broader** perspective to help understand what role data mining plays in the decision-making process.
3. Help you develop an **appreciation** for the beauty of the theoretical foundations underlying data mining.
4. Help you **think** more like a Data Scientist.
5. (For myself) Continue learning.

Course Material



VCU Canvas

The image shows a GitHub repository page for 'dapt-631'. The URL in the address bar is 'github.com/vishal-git/dapt-631'. The repository is public, created by 'vishal-git'. It has 76 commits, 4 forks, and 0 stars. The repository description states: 'This repository contains the class material for Data Mining (DAPT-631) and Python (DAPT-622). These courses are part of the MS in Decision Analytics (Professional Track) program at Virginia Commonwealth University (VCU).'. The code tab is selected. The commit history shows: 'removed all slides' (cbfe86b · 24 minutes ago), 'added apriori, cf, and nbconvert notebo...' (data folder · 9 months ago), 'Added load_digits screenshot' (misc folder · 2 years ago), and 'added apriori, cf, and nbconvert notebo...' (notebooks folder · 9 months ago).

GitHub

		HyFlex		HyFlex		HyFlex		HyFlex		
	Week 1	Week 2	Week 3	Week 4	Week 5	Week 6	Week 7	Week 8	Week 9	
	Friday	10-Jan 11:30 - 12:30	24-Jan Lunch	7-Feb Lunch	21-Feb Lunch	7-Mar Lunch	21-Mar Lunch	4-Apr Lunch	18-Apr Lunch	May 2 Lunch
Session 1	12:30 - 2:15	Forecasting	Forecasting	Practicum Client Meetings	Risk Analysis	Forecasting	Risk Analysis	Forecasting	Practicum Work Session	Forecasting Presentations
Session 2	2:30 - 4:15	Forecasting	Forecasting	Forecasting	Tableau	Forecasting	Data Mining	Forecasting	Practicum Work Session	Practicum Status Report
Session 3	4:30 - 6:15	Practicum Introductions	Risk Analysis	Risk Analysis	Tableau	Tableau	Data Mining	Risk Analysis	Practicum Work Session	Practicum Status Report
Special Events	6:45 - 8:00	Social		Social		Social		Spring Gala		
	Saturday	11-Jan	25-Jan	8-Feb	22-Feb	8-Mar	22-Mar	5-Apr	19-Apr	3-May
	7:30 - 8:00	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast
Session 4	8:00 - 9:45	Data Mining	Data Mining	Python	Data Mining	Python	Python	Tableau	Data Mining	Data Mining
Session 5	10:00 - 11:45	Data Mining	Data Mining	Python	Data Mining	Python	Python	Tableau	Data Mining	Data Mining
	11:45 - 12:30	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch
Session 6	12:30 - 2:15	Statistics	Statistics	Statistics	Statistics	Statistics	Tableau	Statistics	Forecasting	Tableau Presentations
Session 7	2:30 - 3:45	Statistics	Statistics	Statistics	Statistics	Statistics	Tableau	Statistics	Tableau	CATME Feedback

31 ½ hours!

What We Will Cover

1. Introduction to Data Mining
2. The Data Science Process
3. Introduction to Regression
4. Linear Regression model
5. Build a model using partitioning
6. Decision trees
7. Classification Trees
8. Random Forests
9. Gradient Boosting Trees
10. Hyper-parameter optimization

11. Introduction to Neural Networks

12. Introduction to Clustering

13. Agglomerative Clustering

14. k-means Clustering

15. DBSCAN Clustering

16. PCA

17. Association Analysis (Apriori algorithm)

18. Collaborative Filtering

19. Data Wrangling

+ 20 Jupyter Notebooks

		HyFlex		HyFlex		HyFlex		HyFlex	
	Week 1	Week 2	Week 3	Week 4	Week 5	Week 6	Week 7	Week 8	Week 9
	Friday	10-Jan	24-Jan	7-Feb	21-Feb	7-Mar	21-Mar	4-Apr	18-Apr
1	1	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch
Session 1	1	Intro to Data Mining Intro to Python Intro to Jupyter Notebook	Forecasting	Practicum Client Meetings	Risk Analysis	Forecasting	Gradient Boosting Classifier Accuracy Neural Networks	Forecasting	Practicum Work Session Forecasting
Session 2	2:30 - 4:15	Forecasting	Intro to Python (Caesar Cipher) Data Science Process	Tableau	Forecasting	Data Mining	Forecasting	Practicum Work Session	Practicum Status Report
Session 3	4:30 - 6:15	Practicum Introductions	Risk Analysis	Intro to pandas	Tableau	Data Mining	Risk Analysis	Practicum Work Session	Practicum Status Report
Special Events	6:45 - 8:00	Social	Social	Decision Trees Random Forests	Social	Spring Gala	Neural Networks Clustering		
	Saturday	11-Jan	25-Jan	8-Feb	22-Feb	pandas EDA	22-Mar	5-Apr	19-Apr
	7:30 - 8:00	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast	Breakfast
Session 4	8:00 - 9:45	Data Mining	Data Mining	Python	Data Mining	Python	Python	Tableau	Data Mining
Session 5	10:00 - 11:45	Data Mining	Data Mining	Python	Data Mining	Python	Python	Tableau	Data Mining
	11:45 - 12:30	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch	Lunch
Session 6	12:30 - 2:15	Statistics	Python	Statistics	Statistics	Statistics	Tableau	Statistics	Forecasting
									Tableau Presentations
Session 7	2:30 - 3:45	Statistics	Python	Statistics	Statistics	Tableau	Statistics	Tableau	CATME Feedback

vishal@derive.io

www.linkedin.com/in/VishalJP

@derive_io