

Data Mining

DAPT 631

Vishal Patel

Spring 2022



- Vishal
- Founder of DERIVE, LLC
- MS in Computer Science (IIT, Chicago), and MS in Decision Sciences (VCU, Richmond)
- Mining data since 2003





○ **Introduction**

○ **History**

○ **Course Structure**

○ Introduction

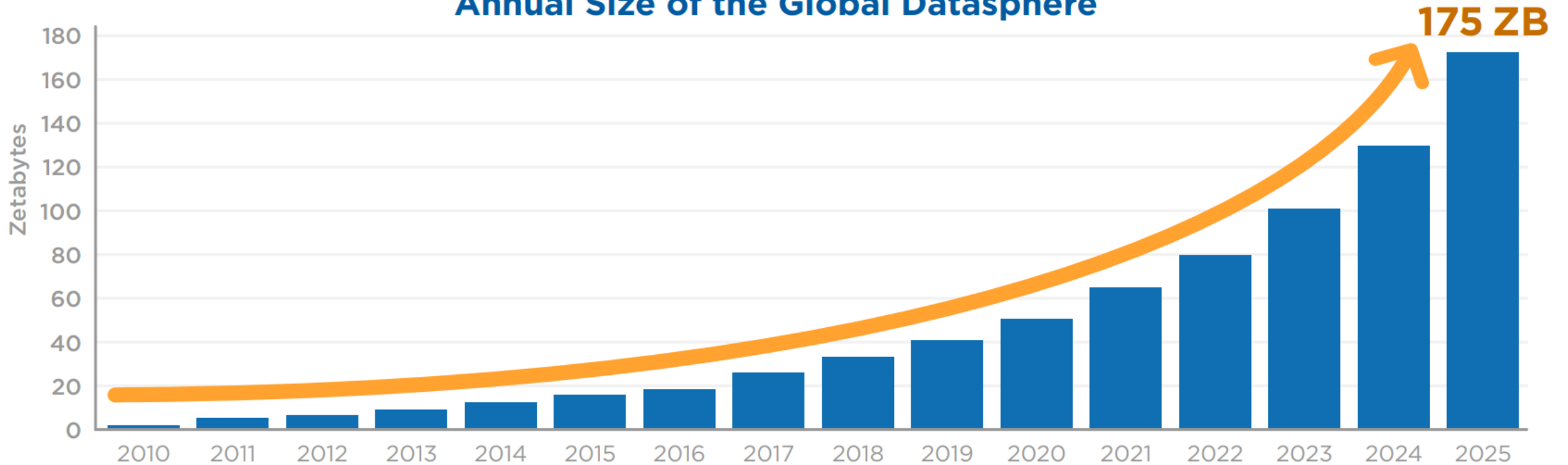
○ History

○ Course Structure

Cambrian Era



Annual Size of the Global Datasphere



Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

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

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

INFRASTRUCTURE

STORAGE

HADOOP

DATA LAKES

DATA WAREHOUSES

STREAMING / IN-MEMORY

NoSQL DATABASES

NewSQL DATABASES

GRAPH DBs

MPP DBs

SERVER-LESS

CLUSTER SVCS

ETL / DATA TRANSFORMATION

DATA INTEGRATION

DATA GOVERNANCE

DATA QUALITY

MGMT / MONITORING

DATA GENERATION & LABELLING

AI OPS

GPU DBs & CLOUD

AI HARDWARE

ANALYTICS & MACHINE INTELLIGENCE

BI PLATFORMS

VISUALIZATION

DATA ANALYST PLATFORMS

DATA SCIENCE NOTEBOOKS

DATA SCIENCE PLATFORMS

MACHINE LEARNING

COMPUTER VISION

HORIZONTAL AI

SPEECH & NLP

SEARCH

LOG ANALYTICS

SOCIAL ANALYTICS

WEB / MOBILE / COMMERCE ANALYTICS

APPLICATIONS - ENTERPRISE

SALES

MARKETING - B2B

MARKETING - B2C

CUSTOMER EXPERIENCE / SERVICE

HUMAN CAPITAL

LEGAL

REGTECH & COMPLIANCE

FINANCE

AUTOMATION & RPA

SECURITY

APPLICATIONS - INDUSTRY

ADVERTISING

EDUCATION

REAL ESTATE

GOVT & INTELLIGENCE

COMMERCE

FINANCE - LENDING

INSURANCE

HEALTHCARE

LIFE SCIENCES

TRANSPORTATION

AGRICULTURE

INDUSTRIAL

OTHER

OPEN SOURCE

FRAMEWORKS

QUERY / DATA FLOW

DATA ACCESS & DATABASES

ORCHESTRATION & PIPELINES

STREAMING & MESSAGING

STAT TOOLS & LANGUAGES

AI OPS & INFRA

AI / MACHINE LEARNING / DEEP LEARNING

SEARCH

LOGGING & MONITORING

VISUALIZATION

COLLABORATION

SECURITY

DATA SOURCES & APIs

DATA MARKETPLACES & DISCOVERY

FINANCIAL & ECONOMIC DATA

AIR / SPACE / SEA

PEOPLE / ENTITIES

LOCATION INTELLIGENCE

OTHER

DATA RESOURCES

DATA SERVICES

INCUBATORS & SCHOOLS

RESEARCH



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”.



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

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.

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



WISDOM

KNOWLEDGE

INFORMATION

DATA

Data Mining



Data Mining Tasks



A horizontal sequence of five colored circles, each containing a data mining task. From left to right, the circles are teal, light green, orange, red-orange, and grey. Each circle has a subtle drop shadow.

Description

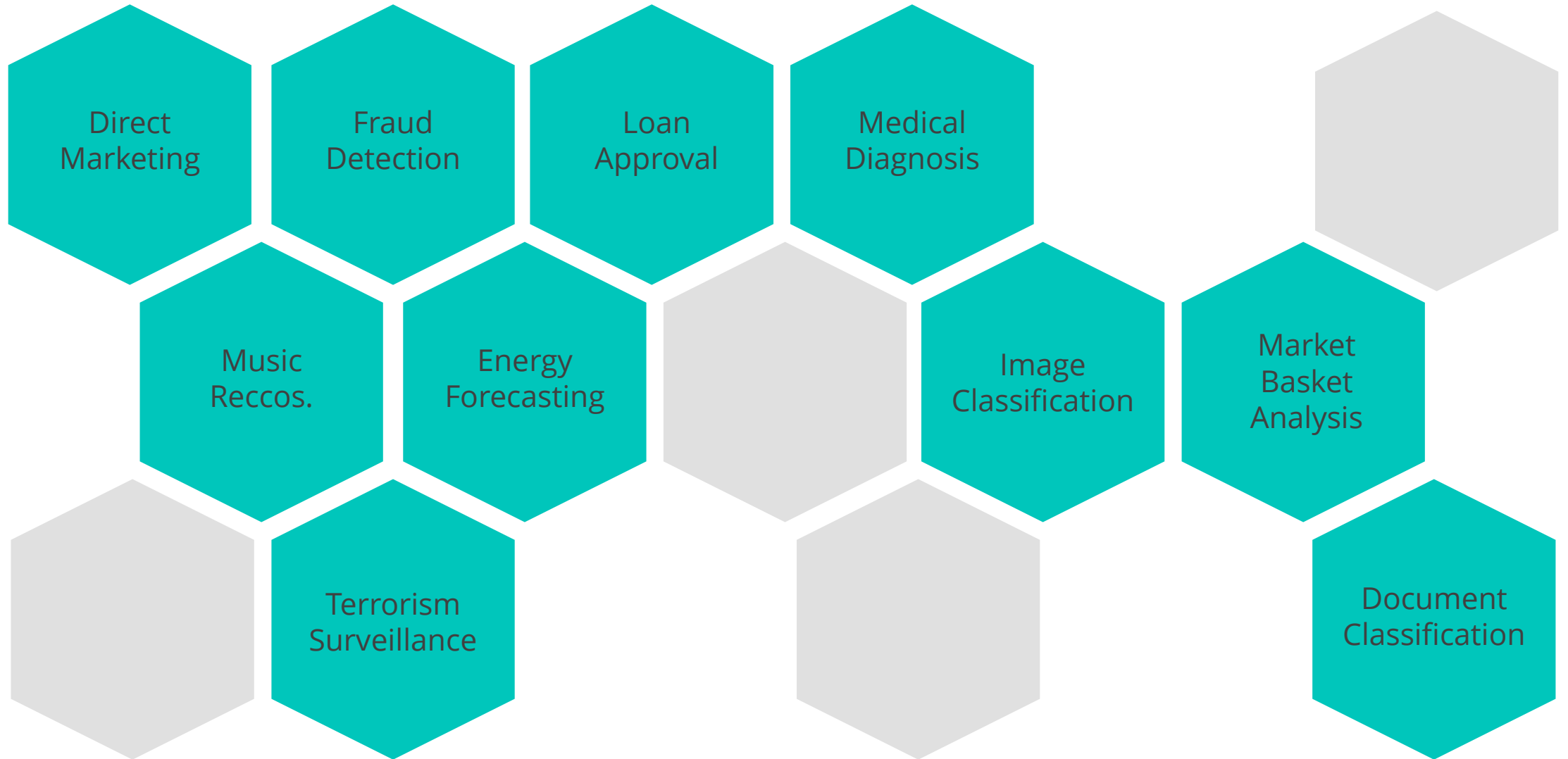
Estimation

Prediction

Classification

Clustering

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



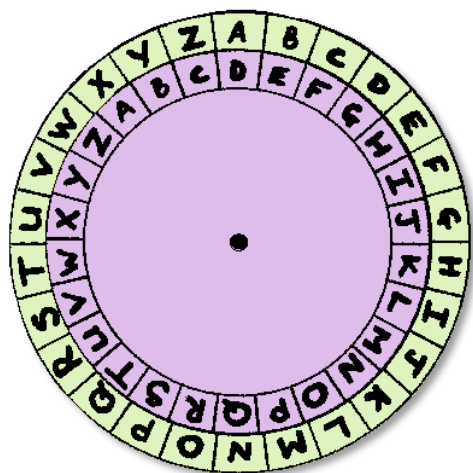
Siege of Plataea (5th Century BCE)

Poll 1

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

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

EXAMPLE #2 FREQUENCY ANALYSIS, CRYPTOANALYSIS



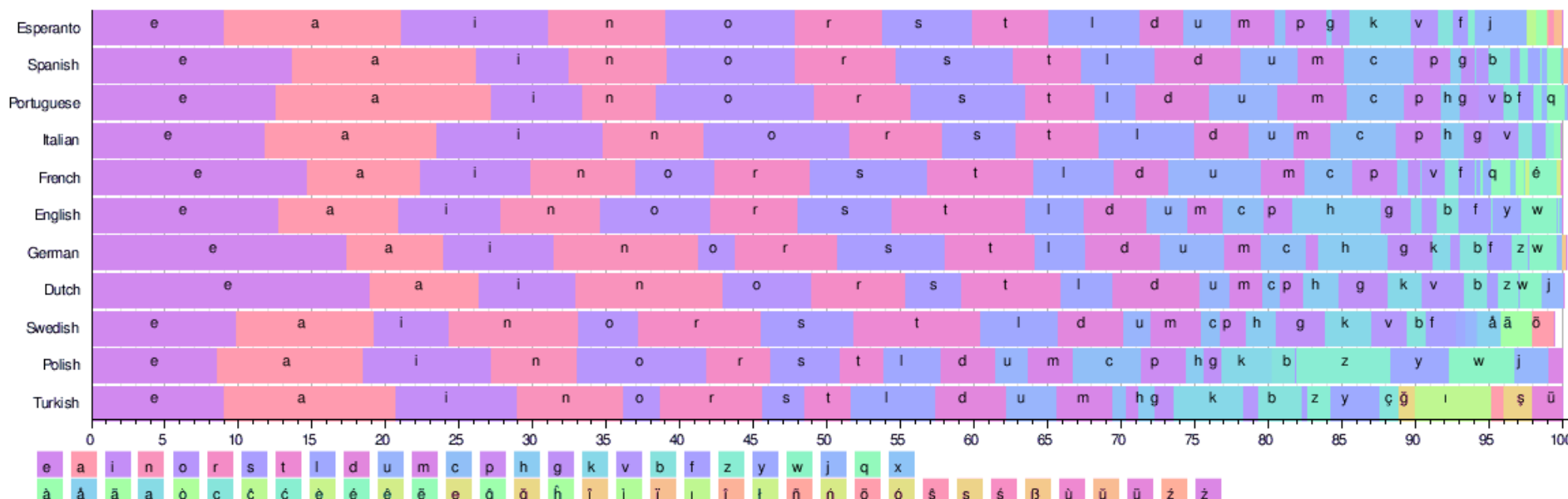
Original message: Et tu, Brute?

Encrypted message: Hw wx, Euxwh?



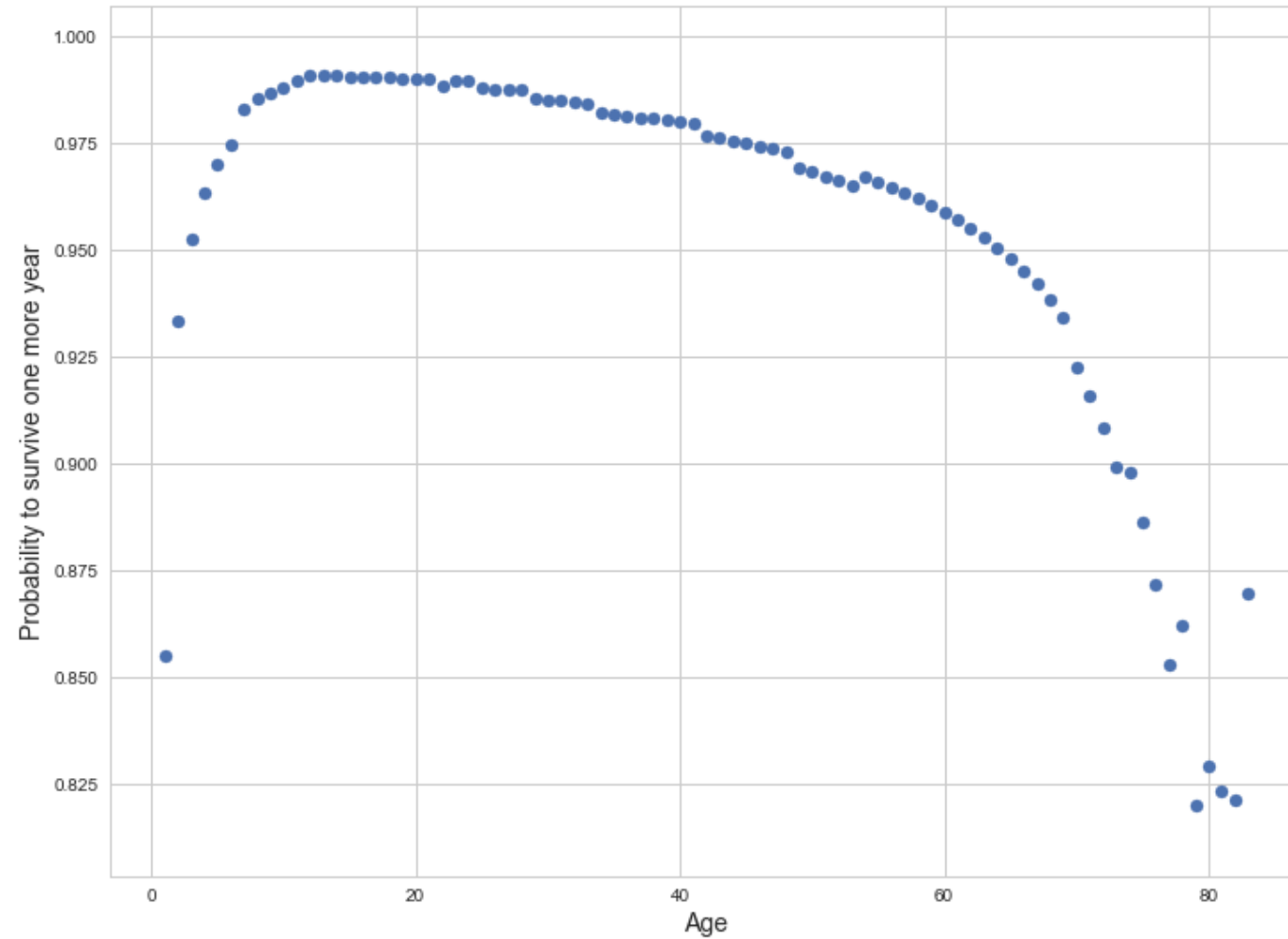
Al-Kindi (801–873 AD)

Caesar Cipher



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 of more year(s) at a certain age.

Modern Statistics

Normal distribution

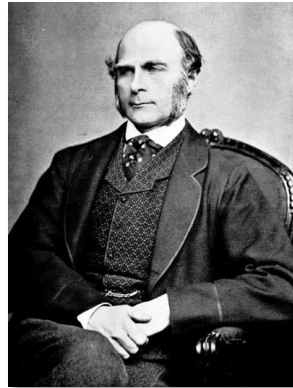
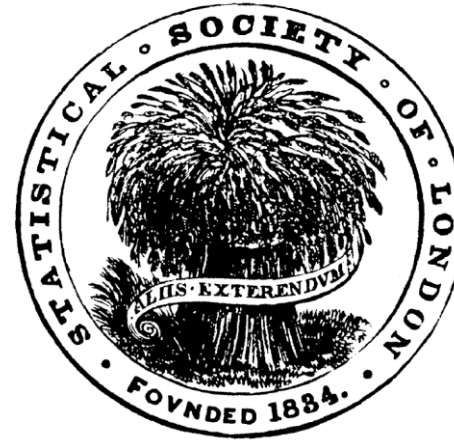
t distribution

Random sampling

Design of-
experiments

Bayesian Statistics

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



Sir Francis Galton
(1822–1911)
Correlation, regression



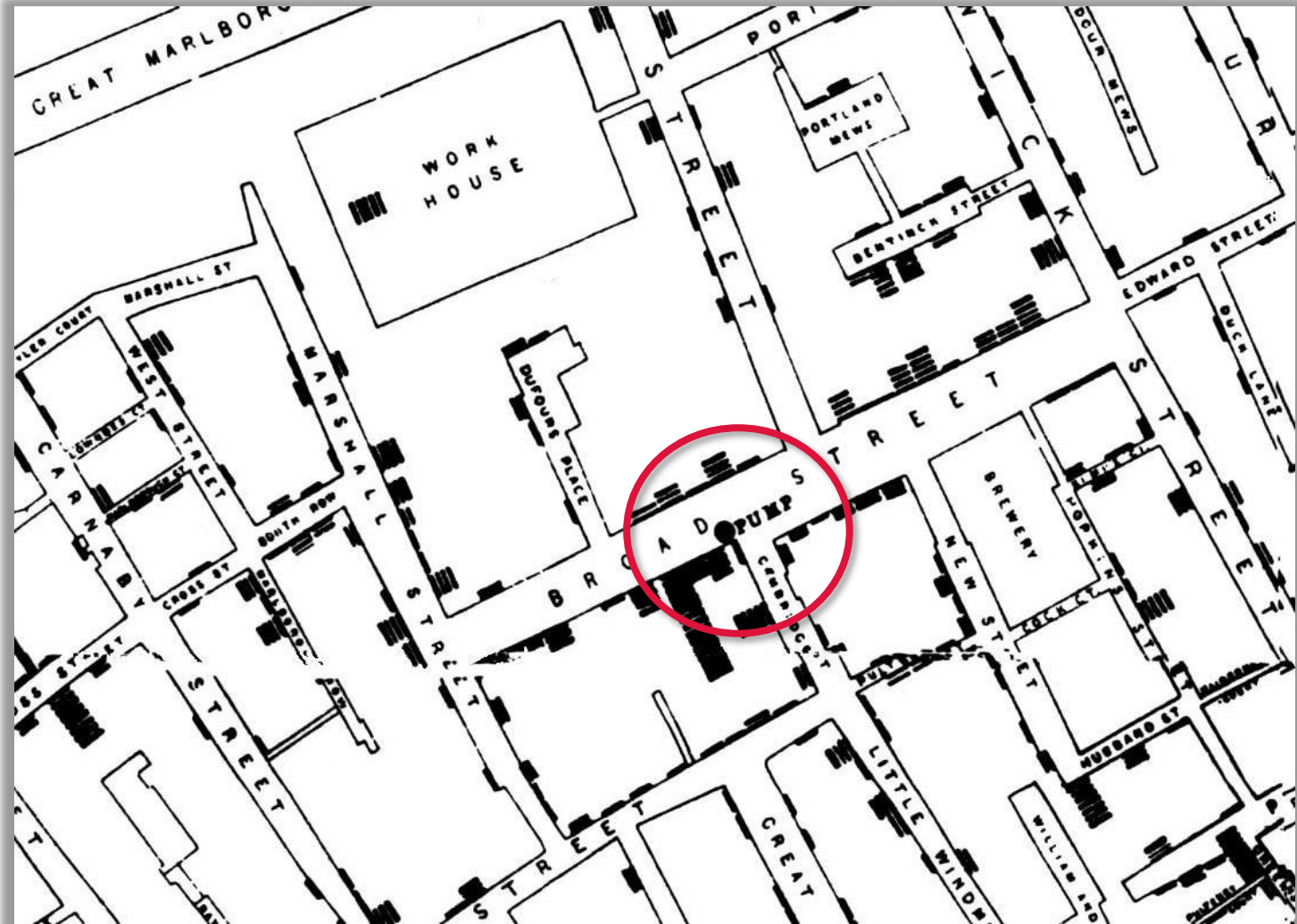
Carl Pearson
(1857–1936)
Founder of mathematical statistics



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

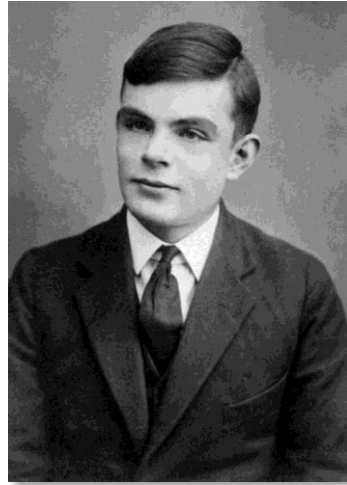
Computer Science

Neural Networks

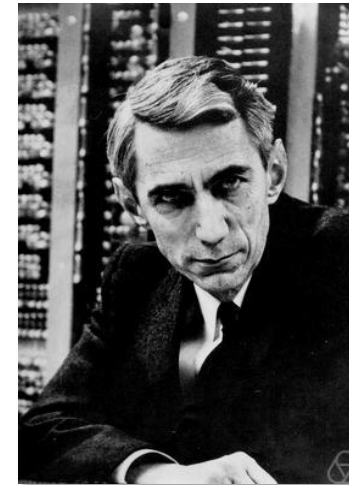
Decision Trees

Genetic Algorithms

Relational Databases



Alan Turing
(1912 –1954)
Theoretical Computer Science



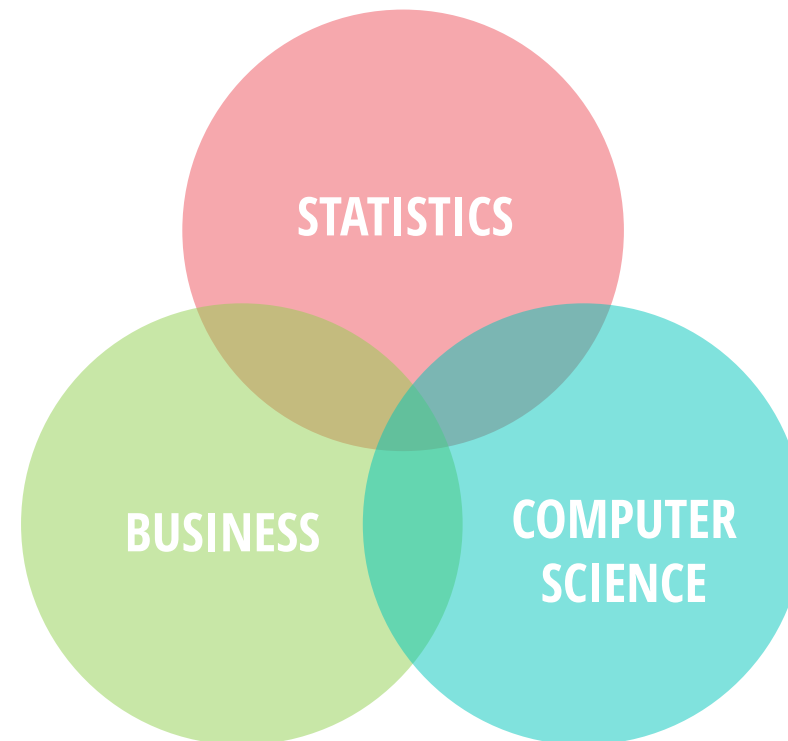
Claude Shannon
(1916 –2001)
Information Theory

- 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

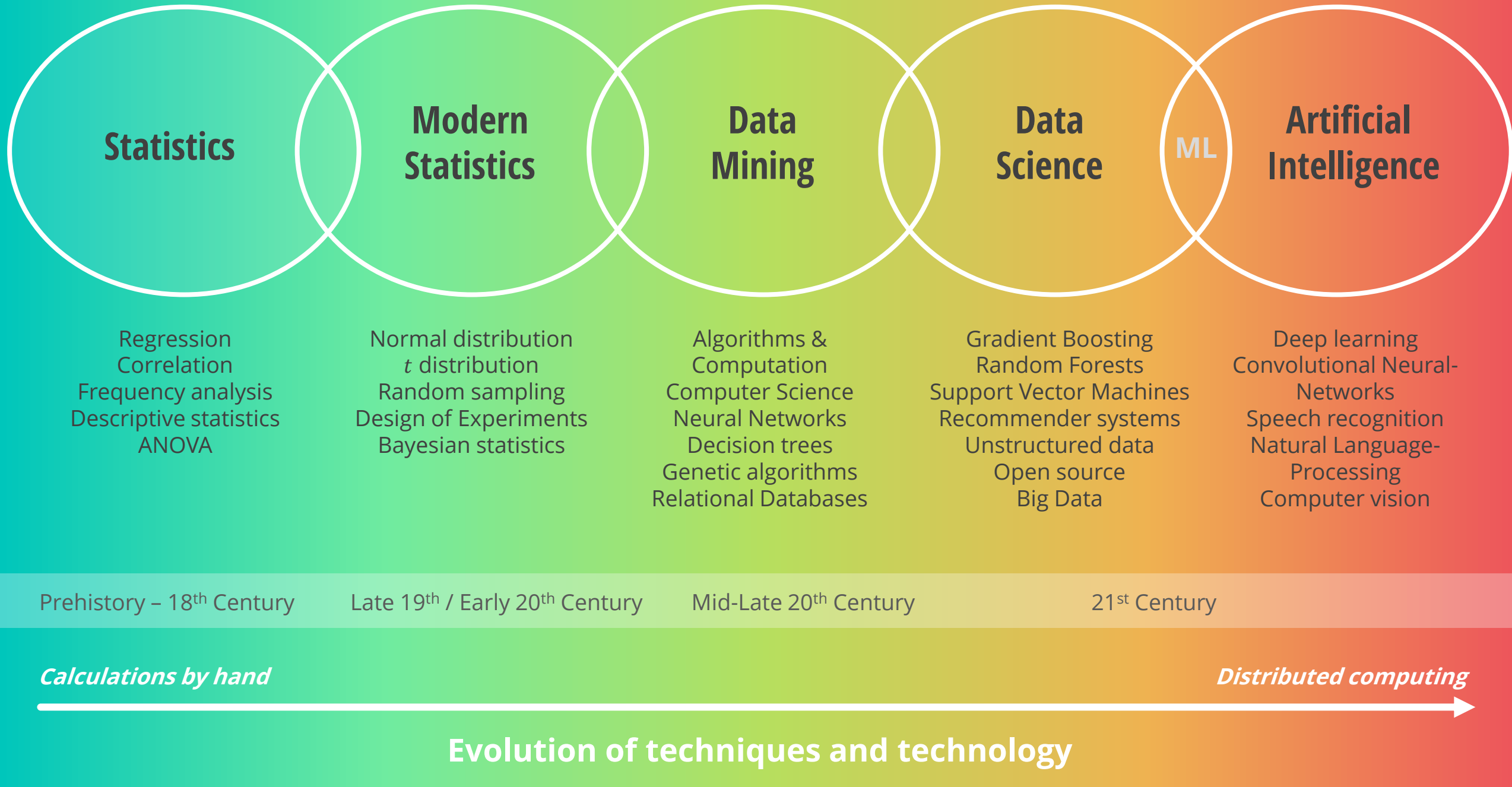
Data science is an **interdisciplinary** field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from data in various forms, both structured and unstructured, **similar to data mining.**[†]



Artificial Intelligence

Deep learning
Convolutional-
Neural Networks
Speech recognition
Natural Language-
Processing
Computer vision



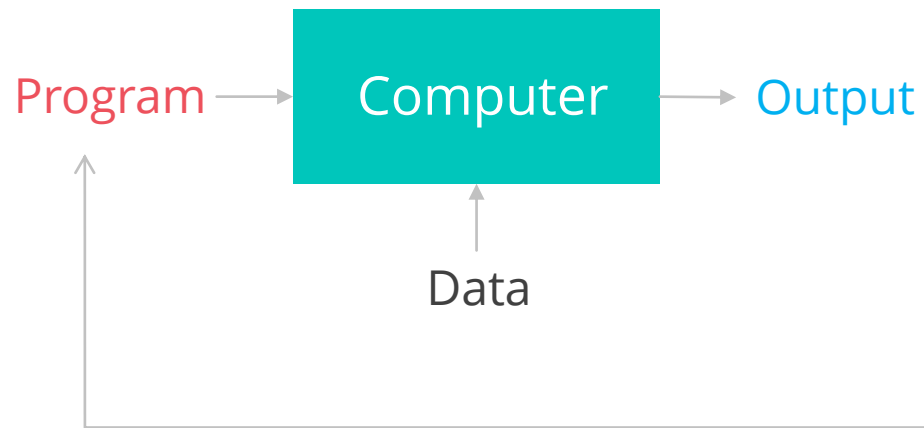


Machine Learning

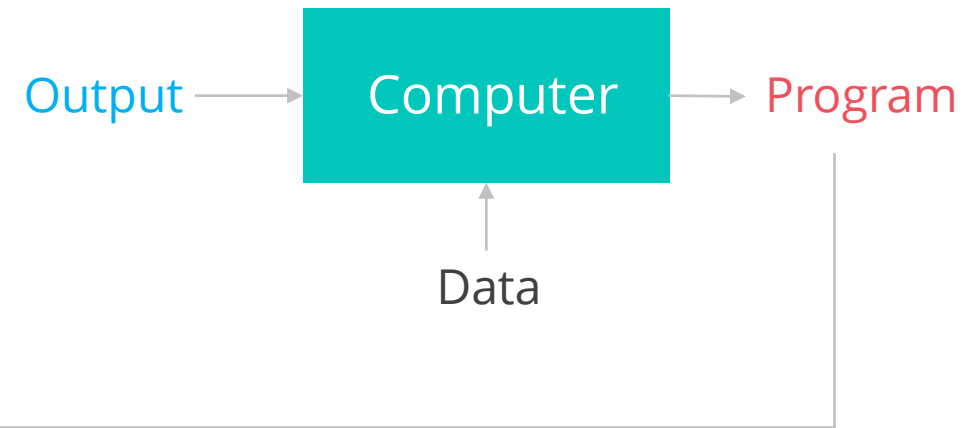
Field of study that gives computers the ability to learn without being explicitly programmed.

Artur Samuel, 1959

Traditional Programming



Machine Learning



○ Introduction

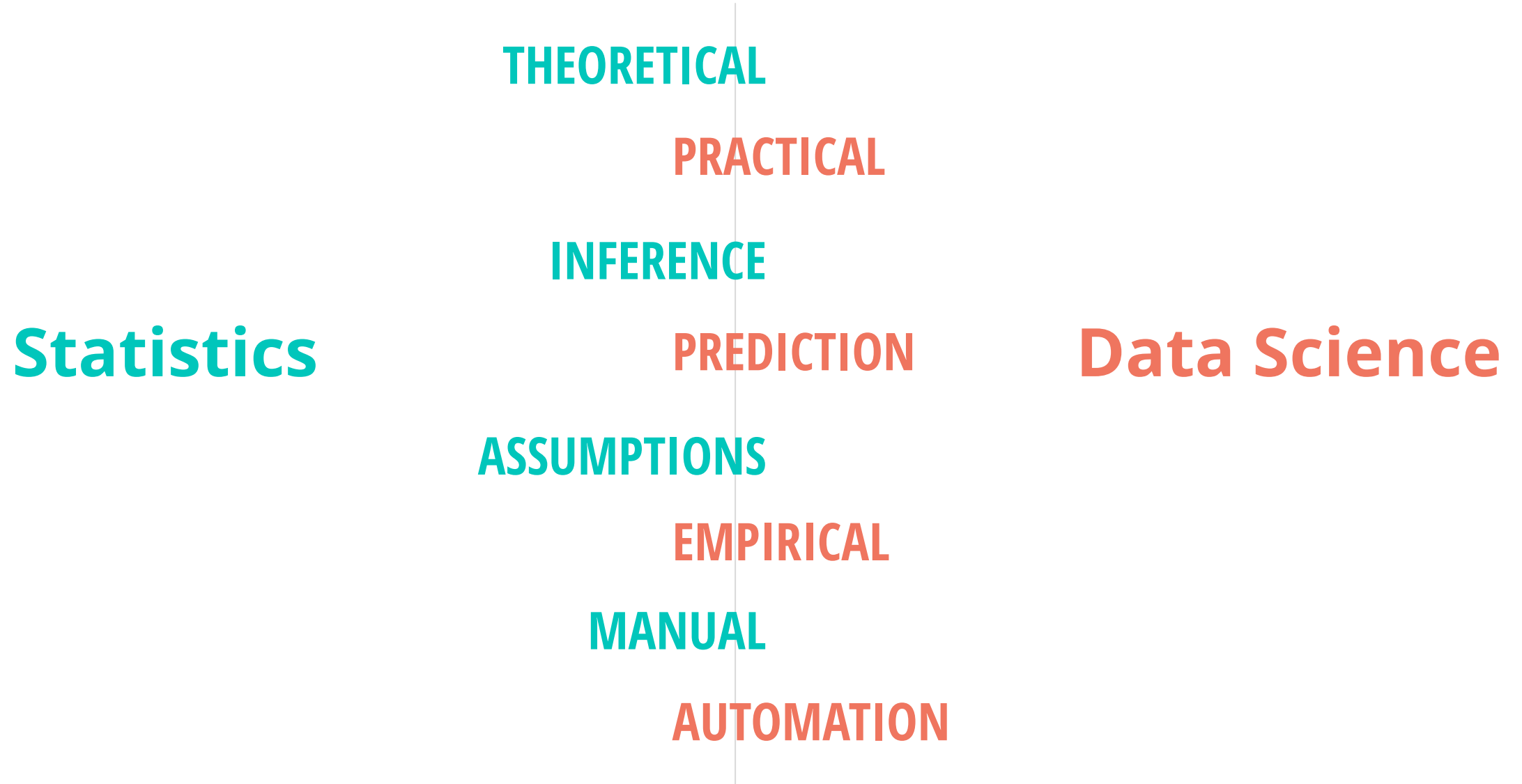
○ History

○ **Course Structure**

Data Science \approx 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



Course Outline

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

Class Structure

1. Ask **questions** at any time!
2. **Collaboration** is encouraged.
3. All content (course material) will be available on Blackboard (and on a git repository).
4. Data Mining + Python
5. Homework assignments in **Python**

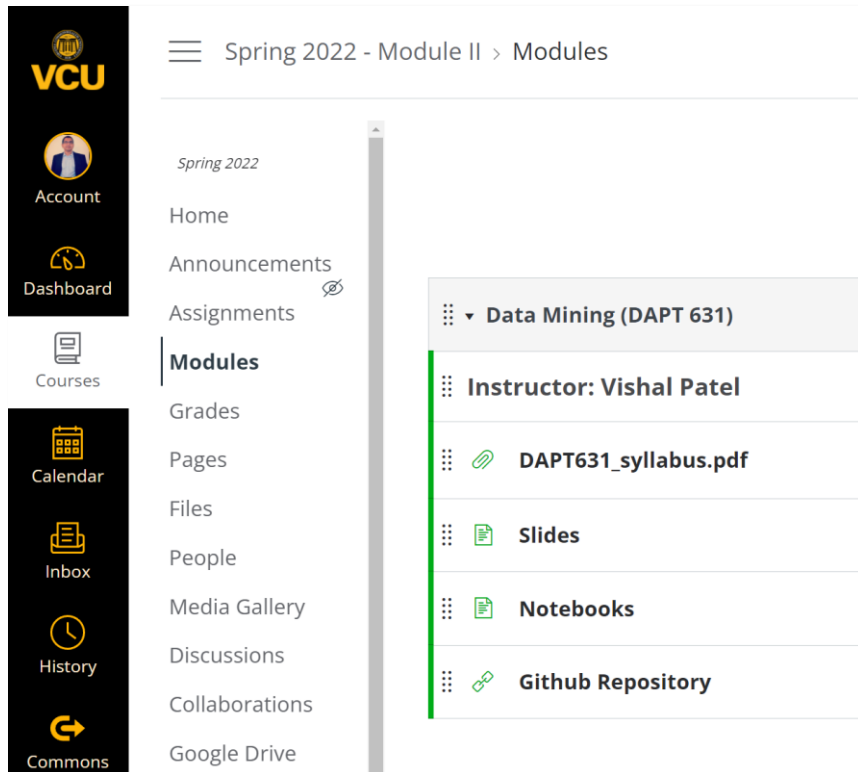
My Objectives

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.

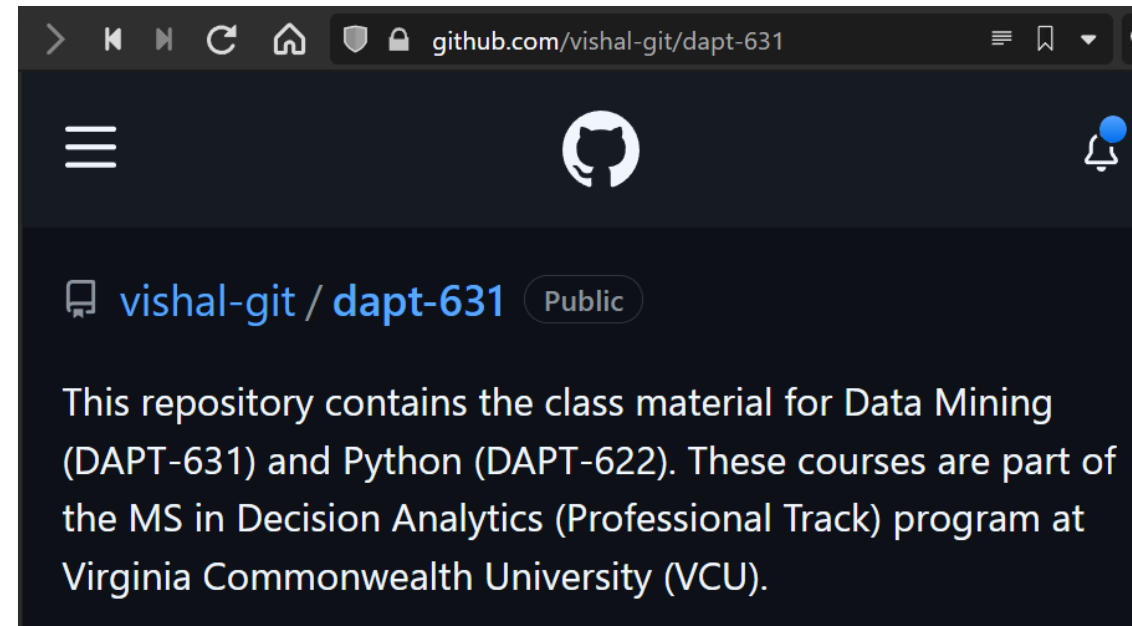
Data Mining + Python



Course Material



The screenshot shows the VCU Canvas LMS interface. On the left is a dark sidebar with navigation icons and labels: VCU logo, Account, Dashboard, Courses, Calendar, Inbox, History, and Commons. The main content area has a top header 'Spring 2022 - Module II > Modules'. Below this is a left-hand menu with options: Spring 2022, Home, Announcements, Assignments, Modules (highlighted), Grades, Pages, Files, People, Media Gallery, Discussions, Collaborations, and Google Drive. The main content area displays the 'Data Mining (DAPT 631)' module, which includes the following items: Instructor: Vishal Patel, DAPT631_syllabus.pdf, Slides, Notebooks, and Github Repository.



The screenshot shows the GitHub repository page for 'vishal-git / dapt-631'. The page is dark-themed and includes the GitHub logo and a notification bell in the top right. The repository name 'vishal-git / dapt-631' is displayed in blue, with a 'Public' label next to it. Below the repository name, a paragraph of text reads: '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).'

vishal@derive.io

www.linkedin.com/in/VishalJP

@derive_io