# Introduction to Data Science

Introduction to Data Science

https://sherbold.github.io/intro-to-data-science

# Outline

- Introduction to Big Data
- Data Science definition
- The Skillset of Data Scientists
- AI
- ML
- Summary

# What is „Big Data"?!?

Is this really about size?

# Naive Definition

- Naive definition:
  - Big data only depends on the data size
  - 1 Gigabyte? 1 Terabyte? 1 Petabyte?

- Naive interpretation misses important aspects
  - Time:
    - Analyzing 1 Gigabyte of data per day is different from analyzing 1 Gigabyte of data per second
  - Diversity:
    - Analyzing spread sheets with numeric data is different from analyzing Web pages that contain a mixture of text and images
  - Distribution:
    - Analyzing data from a single source is different from analyzing data from multiple sources

# Definition of Big Data

- Following Gartner's IT Glossary:
  - Big data is high-**volume**, high-**velocity** and/or high-**variety** information assets that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, and process automation.
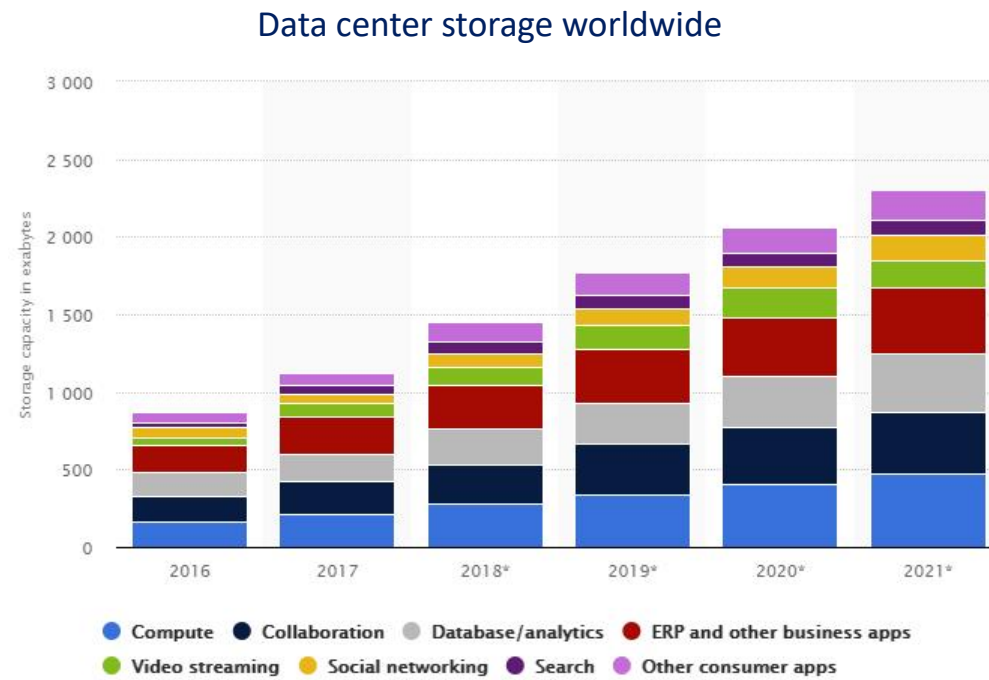
- The three Vs
  - Volume
  - Velocity
  - Variety

Some people actually use 10 Vs to define big data!
- Variability
- Veracity
- Validity
- Vulnerability
- Volatility
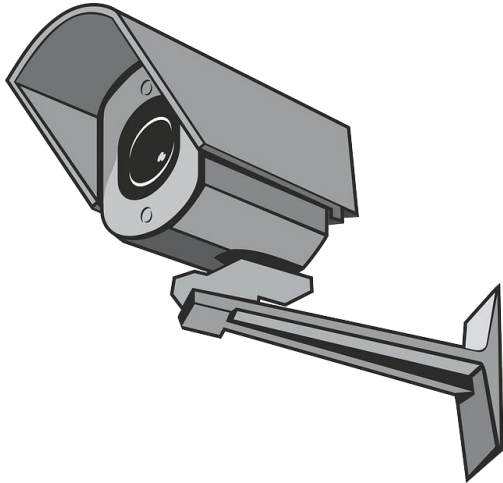- Visualization
- Value

# The 3 Vs: Volume

- Scale of the data must be „big"
  - No clear definition
  - „that demand […] innovative forms of information processing" (Gartner)

Data center storage worldwide
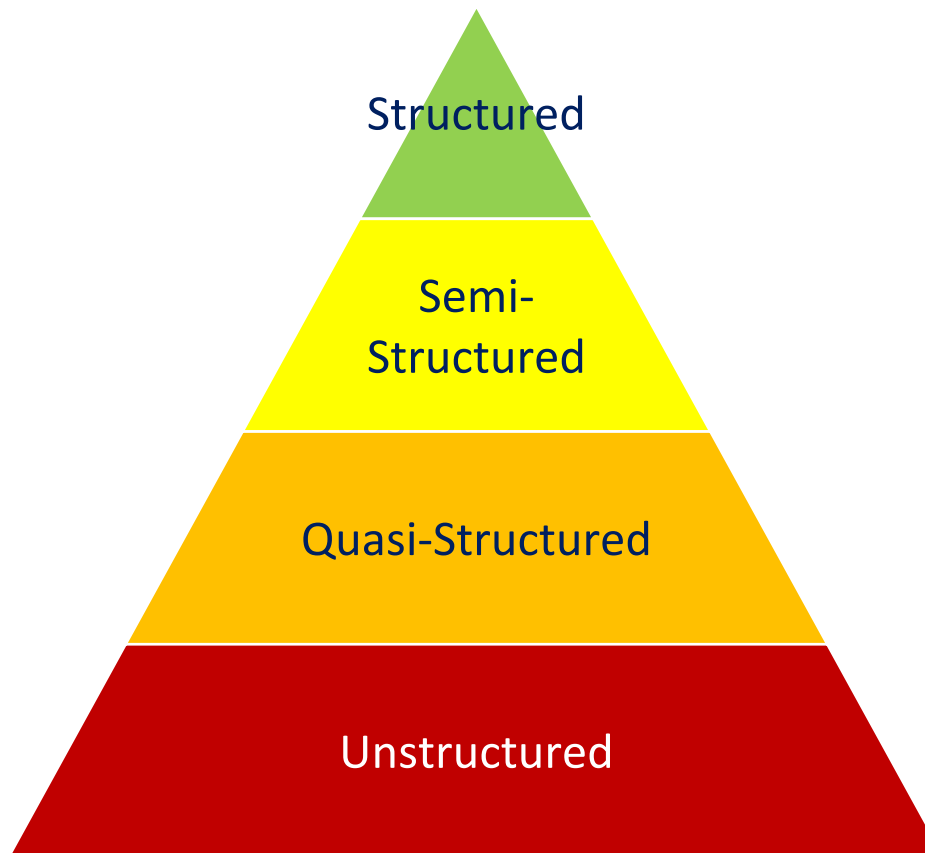


© Statista 2018

# The 3 Vs: Velocity

- Speed at which new data is created

- Speed at which data must be processed and analyzed
  - Often close to real-time

# The 3 Vs: Variety

- Diversity in data types and data sources

**Structured**
- Data with defined types and structure
- Example: comma separated values

**Semi-Structured**
- Textual data with parseable pattern
- Example: XML files with schema

**Quasi-Structured**
- Textual data with erratic formats that can be formated with effort
- Example: Clickstream data

**Unstructured**
- Data that has no inherent structure, often with multiple formats
- Example: Web site, videos

# Examples for data types



**Structured**

**Quasi-Structured**

**Semi-Structured**
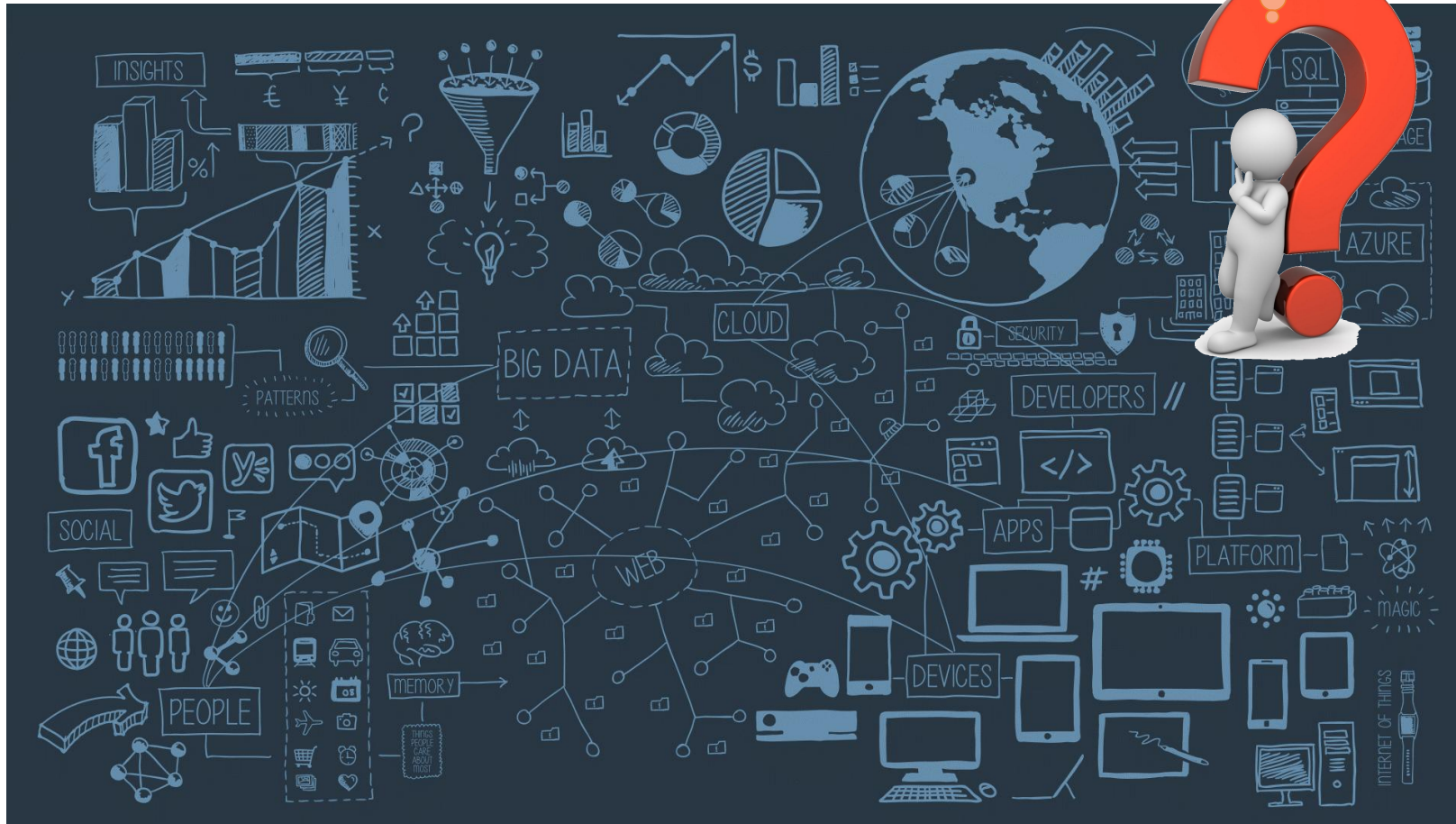
**Unstructured**

9

# Outline

- Introduction to Big Data

- Data Science definition

- The Skillset of Data Scientists
- AI
- ML

- Summary

# Defining Data Science

- Unfortunately, there is no clear definition (yet?)

- Goal is the extraction of knowledge from data

- Combination of techniques from different disciplines

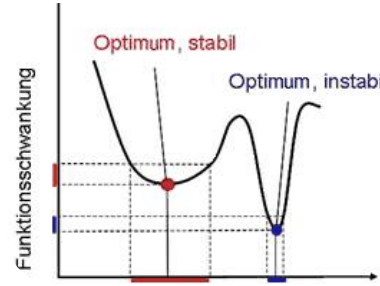- Scientific principles guide the data analysis

# What is „Data Science"?!?

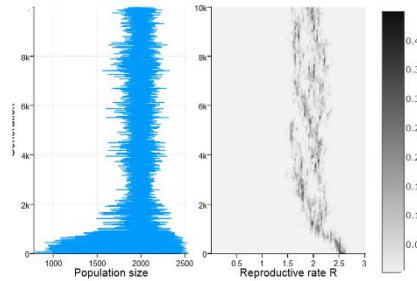# Mathematical Aspects



**Computational Geometry**



**Optimization**



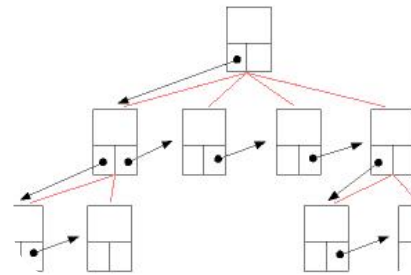**Stochastics**
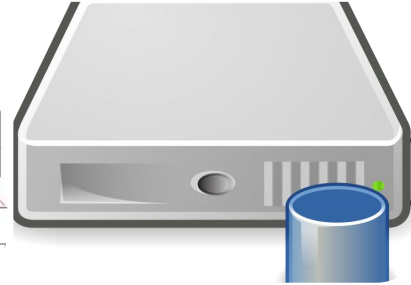


**Scientific Computing**



**Machine Learning**

# Computer Science Aspects
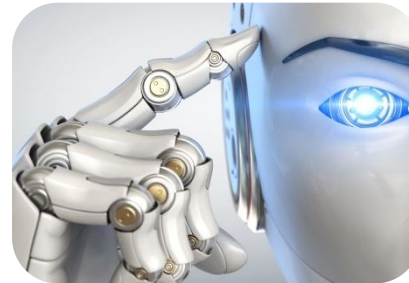


Data Structures and Algorithms

Databases

Distributed Computing

Software Engineering
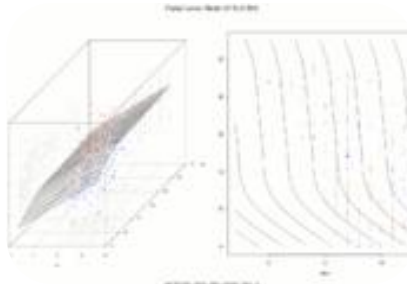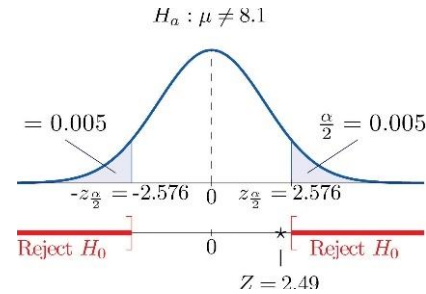
Artificial Intelligence

Machine Learning

# Statistical Aspects



Linear Models



Statistical Tests



Inference



Time Series Analysis



Machine Learning

# Applications



Intelligent Systems



Robotics



Marketing



Medicine



Autonomous Driving



Social Networks

# Outline

- Introduction to Big Data
- Data Science definition
- The Skillset of Data Scientists
- AI
- ML
- Summary

# What are Data Scientists?

- Not computer scientists
  - But should know about databases, data structures, algorithms, etc.

- Not mathematicians
  - But should know about optimization, stochastics, etc.

- Not statisticians
  - But should know about regression, statistical tests, etc.

- Not domain experts
  - But must work together with them

# Skills of Data Scientists

Quantitative
- Maths
- Algorithms
- Statistics

Collaborative
- Teamwork
- Communication skills

Data Scientists

Technical
- Programming
- Infrastructures

Skeptical
- Create hypotheses, but be skeptical about them

A bit of everything

… but actually as much as possible of everything

# Different types of Data Scientists

- **According to Microsoft Research:**

  - Polymath
    - „Do it all"

  - Data Evangelist
    - Data analysis, disseminating and acting on insights

  - Data Preparer
    - Querying existing data, preparing data for analysis

  - Data Shapers
    - Analyzing and preparing data

  - Data Analyzer
    - Analyzing data

  - Platform Builder
    - Collect data and create infrastructures

  - Moonlighters (50%/20%)
    - „Spare time" data scientists

  - Insight Actors
    - Use the outcome and act on insights.

Miyung Kim, Thomas Zimmermann, Robert DeLine, Andrew Begel: Data Scientists in Software Teams: State of the Art and Challenges, IEEE Transactions on Software Engineering (Online First)

# Data Science Definition

- Data science is a multidisciplinary field that uses scientific methods, algorithms, processes, and systems to extract valuable insights and knowledge from structured and unstructured data.

- It combines elements of statistics, computer science, domain knowledge, and data visualization to analyze large and complex datasets, uncover patterns, make predictions, and inform decision-making.

# Data Science vs. Business Intelligence

- Business Intelligence (Gartner IT Glossary)
  - […] best practices that enable access to and analysis of information to improve and optimize decisions and performance.



| | Business Intelligence | Data Science |
|---|---|---|
| Techniques | Dashboards, alerts, queries | Optimization, predictive modelling, forecasting |
| Data Types | Structured, data warehouses | Any kind, often unstructured |
| Common questions | What happened…? How much did…? When did…? | What if…? What will…? How can we…? |

# Dashboard examples

# HOTEL REVENUE MANAGEMENT

KPI'S  CUSTOMERS  AGENTS

## Key Performance Indicators

### REVENUE BY COUNTRIES



NORTH AMERICA

ASIA

Atlantic Ocean

AFRICA

SOUTH AMERICA

Indian Ocean

AUS

ANTARCTICA

Microsoft Bing

© 2023 Microsoft Corporation Terms

| $49.63M Revenue | $37.30M Revenue Apli Discount | $41.52M Net Revenue + Meals | 277.241K N° of Guests |
|---|---|---|---|

### REVENUE BY MONTH

$7.6M
$6.7M
$5.1M
$4.0M $4.1M $4.0M
$3.8M
$3.5M $3.3M
$2.7M $2.7M
$2.3M

ENE  FEB  MAR  ABR  MAY  JUN  JUL  AGO  SEP  OCT  NOV  DIC

### ADULTS AND YOUNGER

15K

● Adults
● Younger

262K

### NIGHTS

● WeekDay ● Weekend

353,504
131,662

0K  200K  400K

### RESERVATIONS BY DAY

2K

0K
Jan 2018     Jul 2018     Jan 2019     Jul 2019     Jan 2020     Jul 2020

### STATUS

Cancell...  ● Not  ● Yes

105K

172K

# Mall Analysis

## Mall

Big Bazaar
Reebok
Stanza
Van Heusen
Bagzone All
Puranmal
Louis Philipe Sports
Peter England
Nike
Wills Lifestyle
Afton
Delsey
Vip

### YTD Sales
**$23,270.31K**
Goal: $24,896.93K
Target

### YTD Footprint
**2231.15K**
Goal: 2389.64K
Target

### Avg Footprint
**8252**

### Sales Current Month
**$1,981.84K**
Goal: $1,916.86K
Previous Month

### Footprint Current Month
**189.6K**
Goal: 183.54K
Previous Month

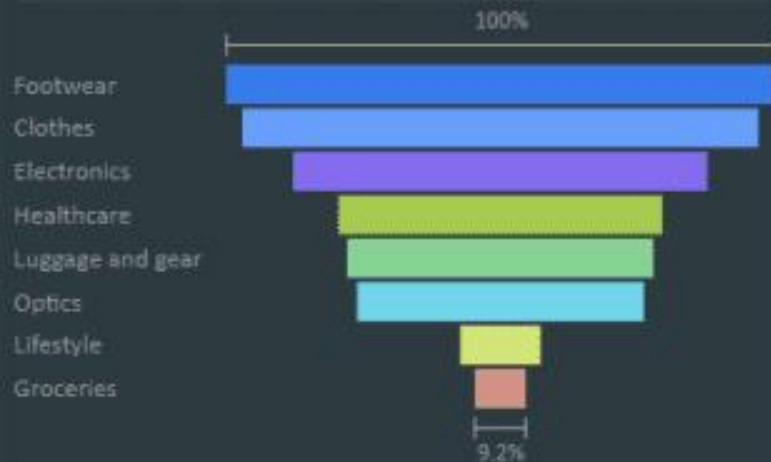### Avg sales per Footprint
**1043**

### Sales Conversion Rate%
**68%**

### Store Capture Rate %
**100%**

### Avg Stay Time (min)
**30**

## Shop by Category (Footprint)

100%

- Footwear
- Clothes
- Electronics
- Healthcare
- Luggage and gear
- Optics
- Lifestyle
- Groceries

9.2%

## Sales Comparison (Current vs Target)

● Sales   ● Target Sales        Footprint ─● Sales

- $2.4M
- $2.2M
- $2.0M
- $1.8M
- $1.6M

Sales and Target Sales

Month: January, February, March, April, May, June, July, August, September, October

## Footprint by Shop (Top 5)

| Shop | Footprint |
| --- | --- |
| Stanza | 224K |
| Wills Lifestyle | 198K |
| Reebok | 197K |
| Van Heusen | 194K |
| Bagzone | 188K |

0K      100K      200K

Footprint

# Patient Record Details

## Filters

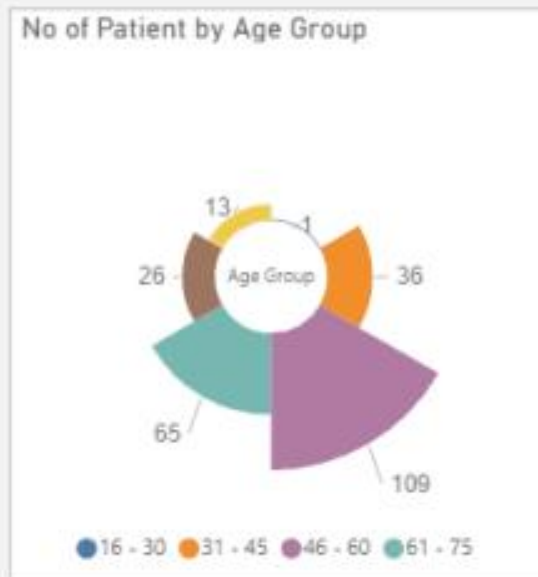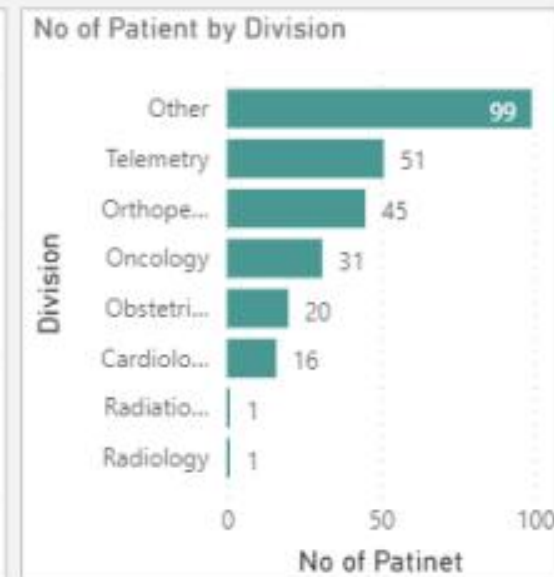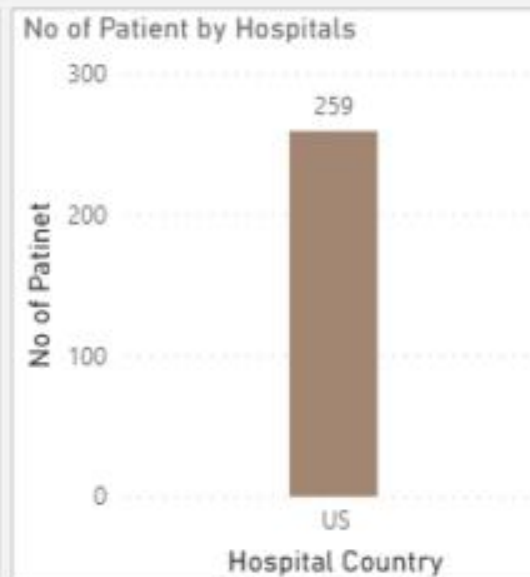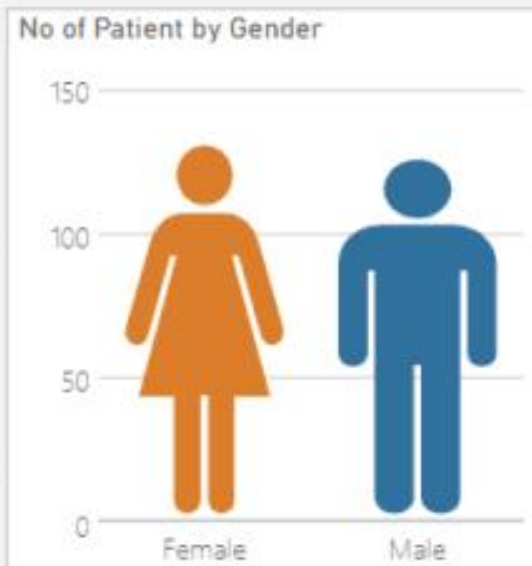| Date Period | Hospital Country, Hospit... | Division, Department Na... | Physicians | Patient Name | Surgical Specialty, Surgi... |
|---|---|---|---|---|---|
| Feb 2017 | All | All | All | All | All |

## Total Patient
**259**
Last Month: 302 (-14.24%)

## Patient in ICU
**15**
Last Month: 16 (-6.25%)

## Total Died Patient
**3**
Last Month: 6 (-50%)

## ReAdmit Patient
**18**
Last Month: 19 (-5.26%)

## Avg Days of Discharge
**8** ✓
Last Month: 5 (+60%)

## No of Patient by Gender



Female / Male

## No of Patient by Hospitals

259 — US (Hospital Country)

## No of Patient by Division

| Division | No of Patient |
|---|---|
| Other | 99 |
| Telemetry | 51 |
| Orthope... | 45 |
| Oncology | 31 |
| Obstetri... | 20 |
| Cardiolo... | 16 |
| Radiatio... | 1 |
| Radiology | 1 |

## No of Patient by Age Group

Age Group: 13, 1, 36, 26, 65, 109
● 16 - 30   ● 31 - 45   ● 46 - 60   ● 61 - 75

## No of Patient by LOS Bucket

| LOS Bucket | No of Patient |
|---|---|
| < 1 | 5 |
| 1 to 5 | 208 |
| 6 to 10 | 41 |
| 11 to 15 | 5 |
| 16 to 20 | 3 |
| 21 to 25 | 3 |
| 31+ | 1 |

## No of Patient by Discharge Type

- Clinical Advice / consent
- Discharged Themselves
- Died

No of Patient

## No of Patient by City



Allentown, Reading, Philadelphia, Camden, Wilmington
© 2022 Microsoft Corporation Terms

## Avg Waiting Time by Division

| Division | Avg Waiting Time |
|---|---|
| Cardiolo... | |
| Obstetri... | |
| Other | |
| Oncology | |
| Telemetry | |
| Orthope... | |

0.00   10.00   20.00

LinkedIn

Kamil.M.S (BI Consultant)

☐ 2019
☐ 2020

**Connections** 1235

**Companies** 1086

**Invitations Received** 504

**Invitations Sent** 322

**Reactions** 2008

## Total Connections by Month Name

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 16 | 13 | 55 | 145 | 47 | 153 | 118 | 213 | 149 | 119 | 91 | 116 |
| Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct | Nov | Dec |

## Total Connections by Company

- Microsoft
- Freelance
- IKEA Group
- Tata Consultancy ...
- Accenture
- Cognizant

(scale: 0, 5, 10, 15, 20)

## Total Connections by Position



## Messages Received and Messages Sent

604 (41.94%)
836 (58.06%)

● Message Received ● Message Sent

## Reactions by Type

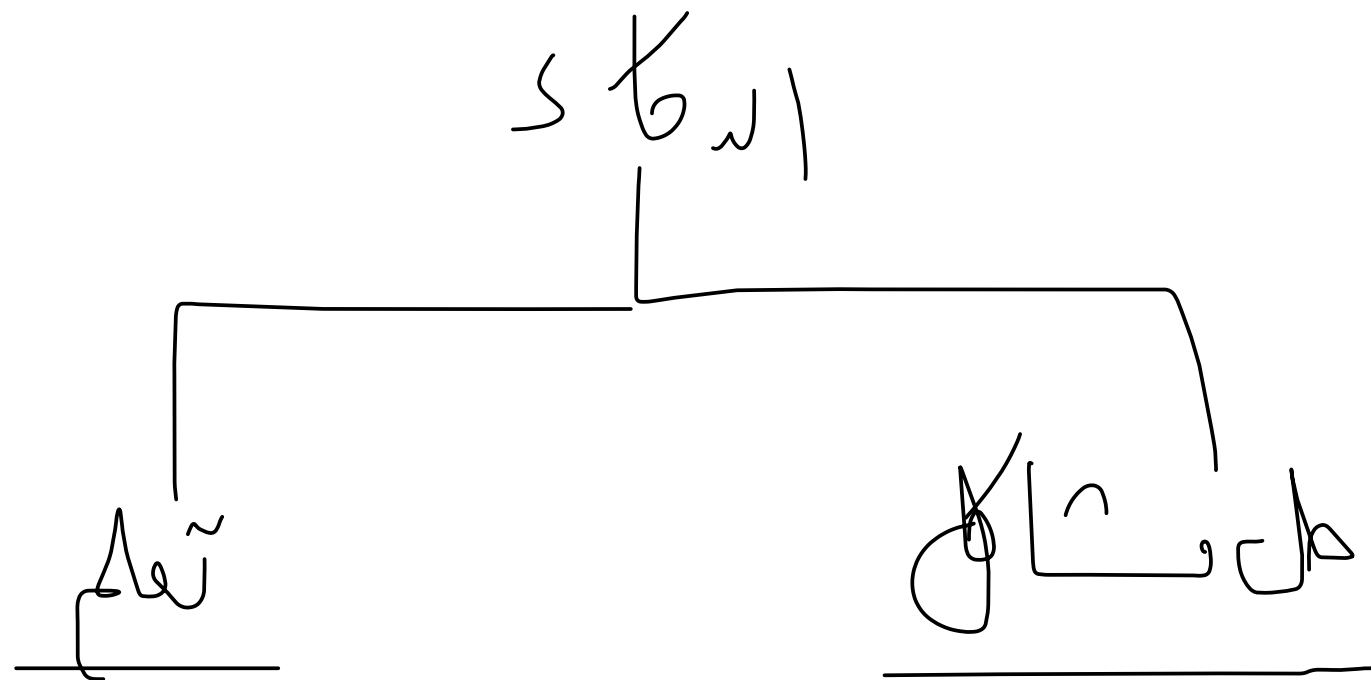| LIKE | EMPATHY | INTEREST | PRAISE | APPRECIATI... | MAYBE |
|---|---|---|---|---|---|
| 1907 | 81 | 9 | 6 | 4 | 1 |

# AI definition

- Artificial Intelligence (AI) refers to the simulation of human intelligence in machines or computer systems.

- It involves the development of algorithms, software, and hardware that enable computers to perform tasks that typically require human intelligence, such as understanding natural language, recognizing patterns, making decisions, and solving problems.
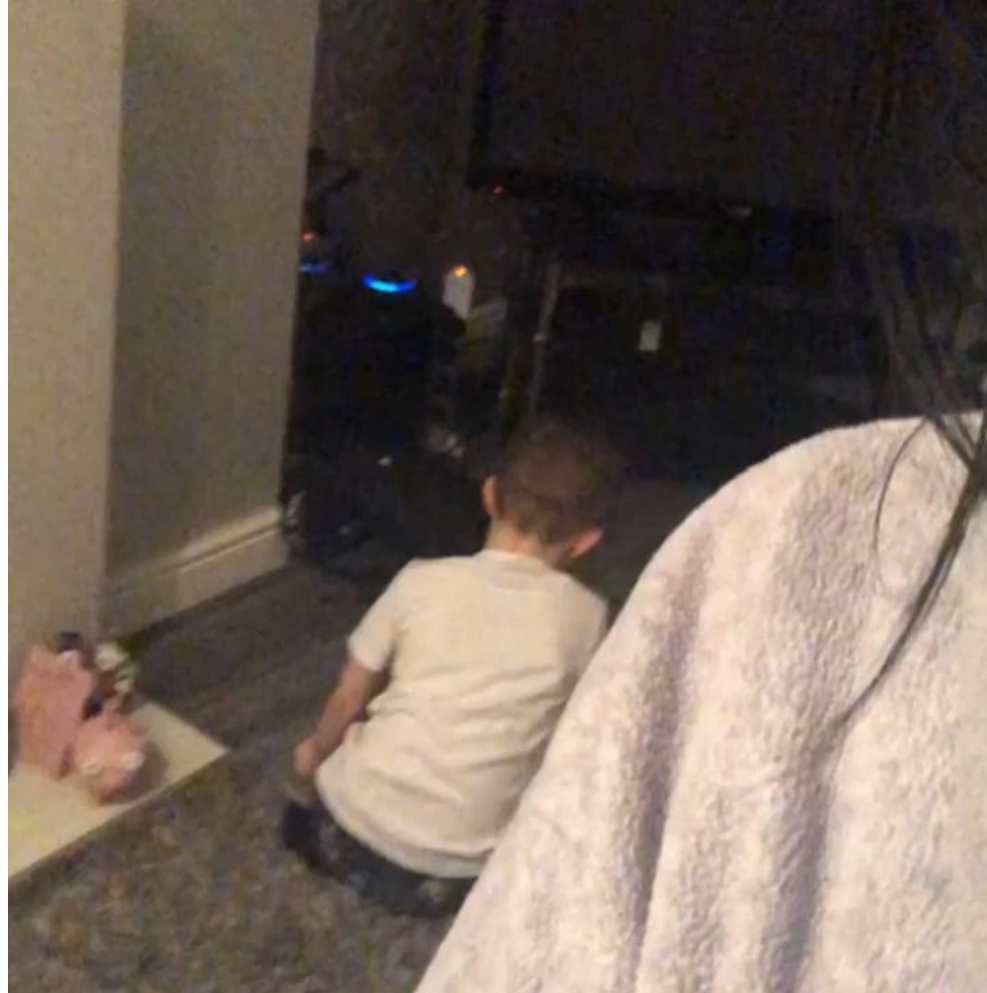
الانسكاب

طبيعي                    اصطناعي
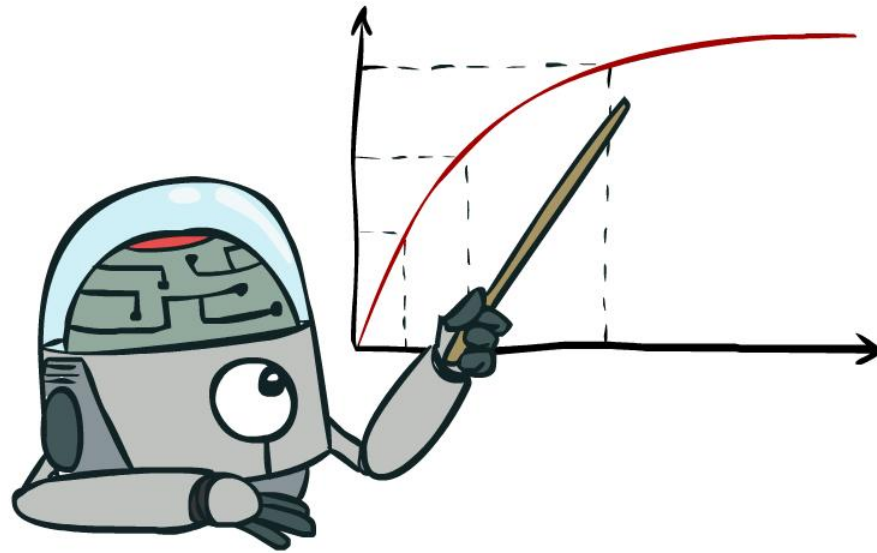
State of the Art: **Human** vs **Robot**

clideo.com

# Types of AI

1.  Narrow AI (Weak AI): This type of AI is designed for specific tasks or domains. It excels at performing functions, such as voice recognition, image classification, or playing board games like chess or Go. Narrow AI systems do not possess general intelligence or consciousness and are trained for specific applications.

2.  General AI (Strong AI): General AI represents machines or systems that possess human-like intelligence, including the ability to understand, learn, and apply knowledge across a wide range of tasks and domains. Achieving general AI is a long-term goal of AI research and development and is yet to be realized.
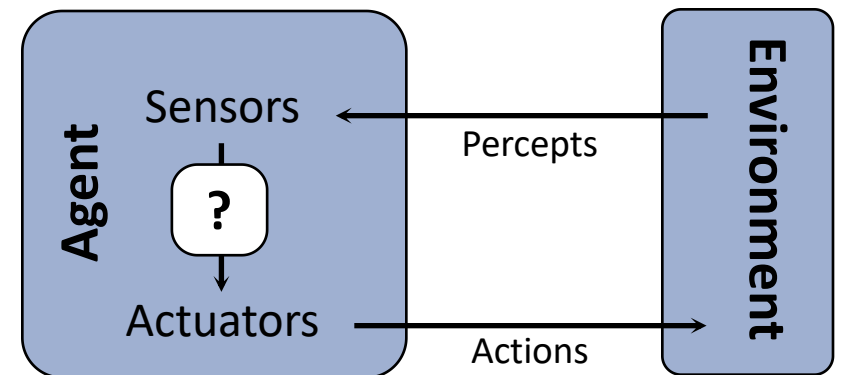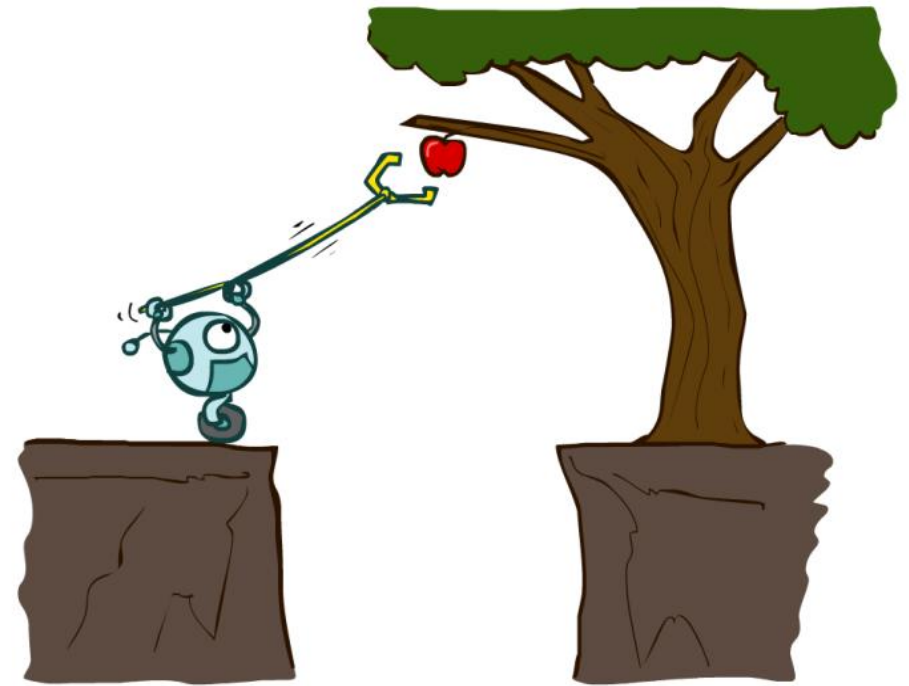
# Alexa helps to solve assignments ☺
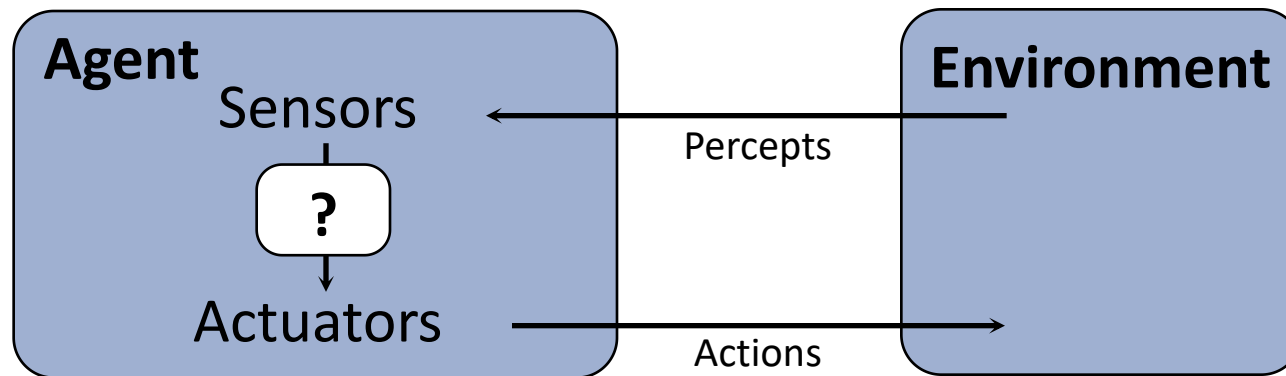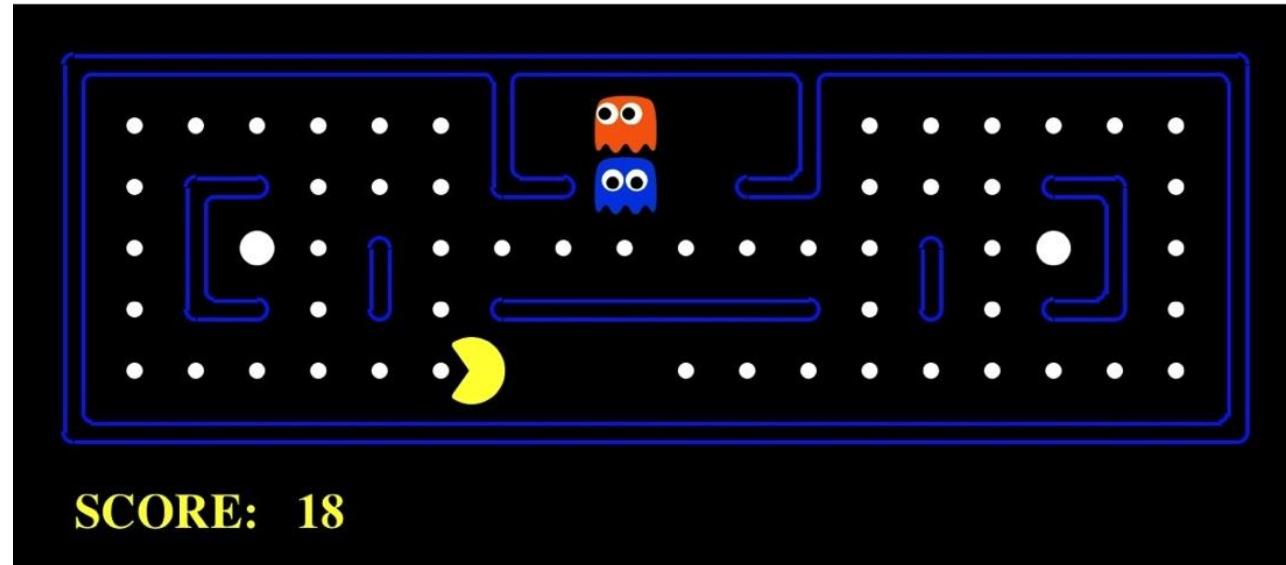
# Maximize Your Expected Utility

# Designing Rational Agents

- An **agent** is an entity that *perceives* and *acts*.

- A **rational agent** selects actions that maximize its (expected) **utility**.

- Characteristics of the **percepts, environment,** and **action space** dictate techniques for selecting rational actions
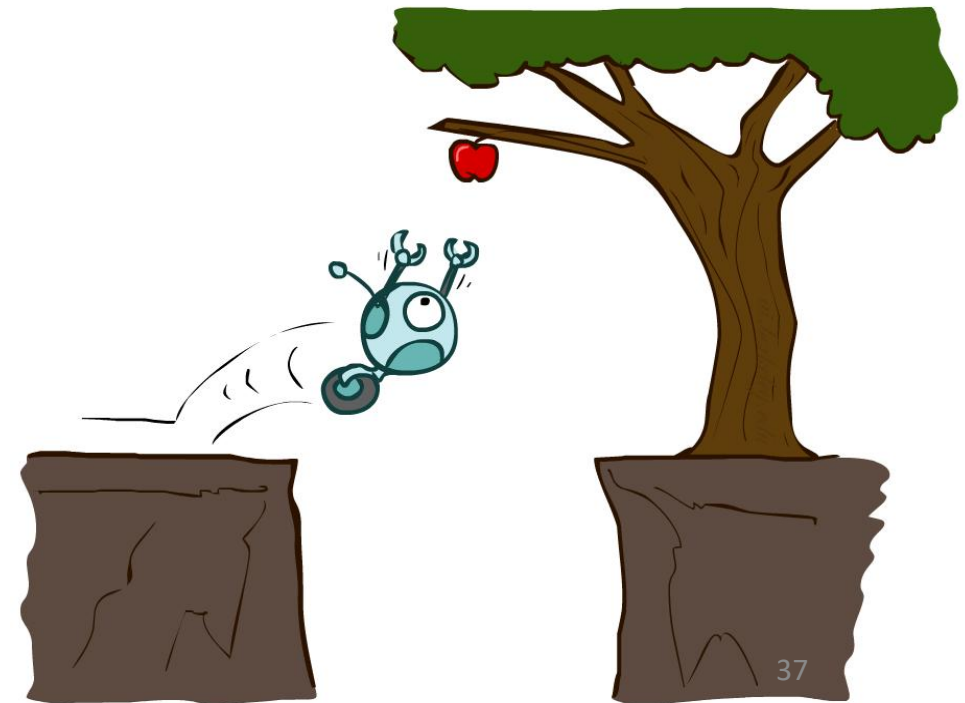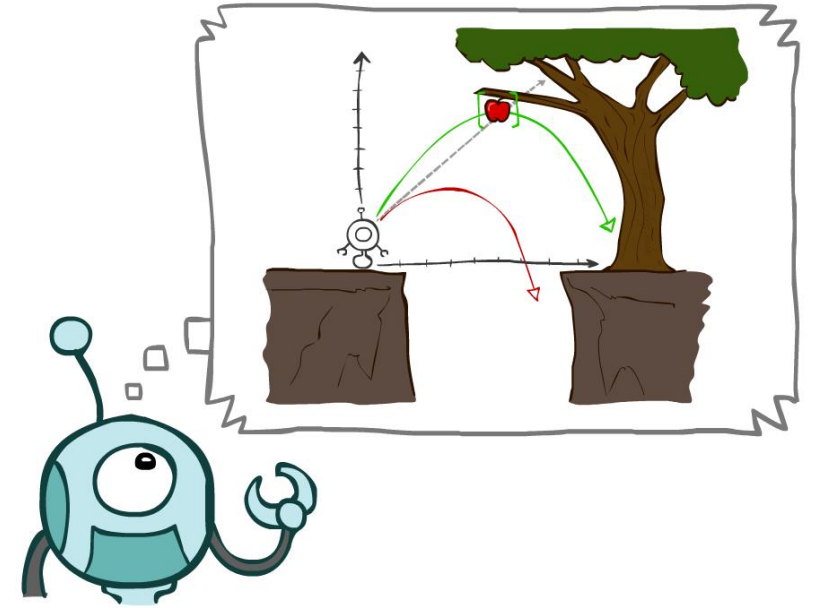


**Agent**

Sensors

**?**

Actuators

**Environment**

Percepts

Actions

# Pac-Man as an Agent



SCORE: 18

**Agent**
Sensors

**?**

Actuators

**Environment**

Percepts

Actions

Demo1: pacman-l1.mp4

# Search



How do we formulate a search problem?

# Machine Learning (ML) definition

Machine Learning (ML) is a subset of artificial intelligence (AI) that focuses on the development of algorithms and statistical models that enable computer systems to improve their performance on a specific task through learning from data, without being explicitly programmed
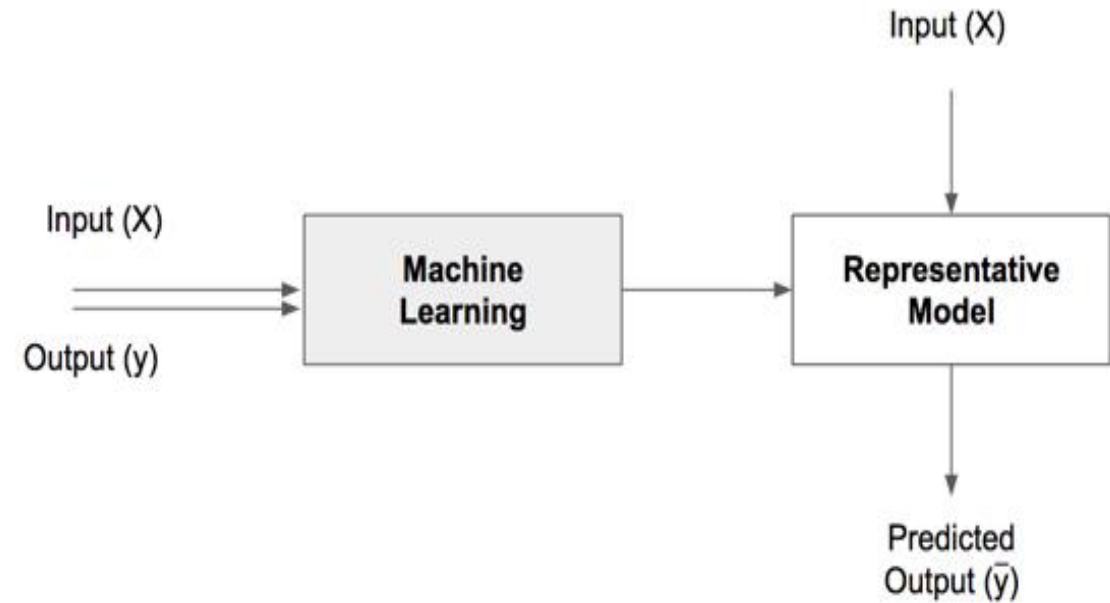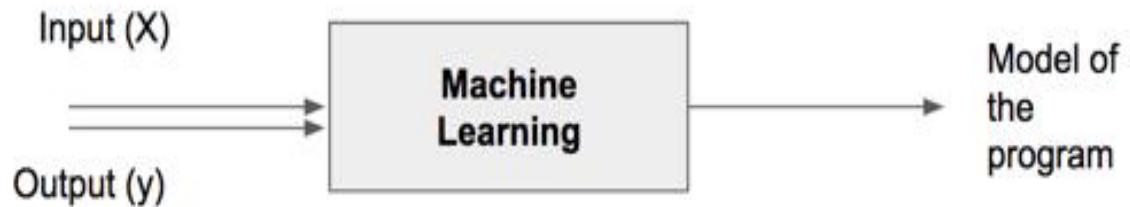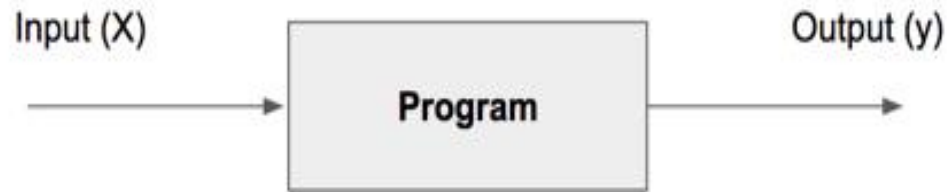
ML allows computers to automatically discover patterns, make predictions, or take actions based on past experiences or examples.

# ML vs Traditional Programming
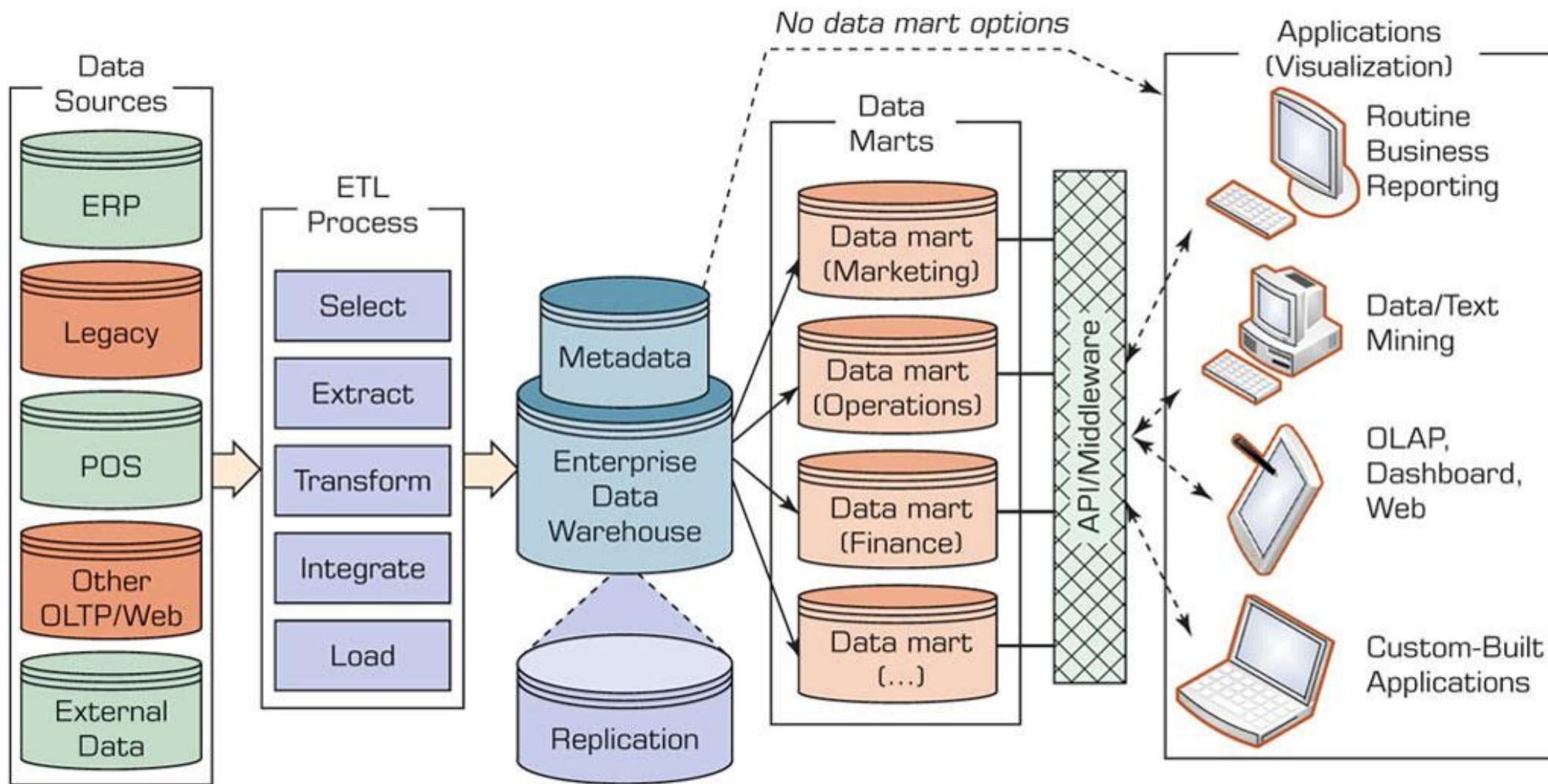
| Aspect | Machine Learning (ML) | Traditional Programming |
|---|---|---|
| Purpose | Learning from data to make predictions, recognize patterns, and automate decision-making. | Executing predefined instructions and algorithms to achieve specific tasks. |
| Data-driven | Relies on data for training and learning patterns. | Not data-driven; instructions are explicitly programmed. |
| Flexibility | Adaptable to changing data and can improve with more examples. | Less adaptable and requires manual code changes for modifications. |
| Problem Complexity | Suited for complex problems with large datasets or uncertain environments. | Effective for well-defined, deterministic tasks. |
| Expertise Required | Requires knowledge of data preprocessing, algorithm selection, and model tuning. | Requires expertise in programming languages, algorithms, and problem-solving. |
| Maintenance | May require periodic retraining and adjustment as data changes. | Maintenance involves debugging, updating, and code optimization. |

# ML models vs programs

# Data Pipeline

# Some Data Science Tasks

| Tasks | Description | Algorithms | Examples |
|---|---|---|---|
| Classification | Predict if a data point belongs to one of predefined classes. The prediction will be based on learning from known data set. | Decision Trees, Neural networks, Bayesian models, Induction rules, K nearest neighbors | Assigning voters into known buckets by political parties eg: soccer moms. Bucketing new customers into one of known customer groups. |
| Regression | Predict the numeric target label of a data point. The prediction will be based on learning from known data set. | Linear regression, Logistic regression | Predicting unemployment rate for next year. Estimating insurance premium. |
| Clustering | Identify natural clusters within the data set based on inherit properties within the data set. | K means, density-based clustering - DBSCAN | Finding customer segments in a company based on transaction, web and customer call data. |
| Association analysis | Identify relationships within an itemset based on transaction data. | FP Growth, Apriori | Find cross selling opportunities for a retailor based on transaction purchase history. |

# Summary

- Big data has a high volume, velocity, and variety

- Different data structures
  - Structured, semi-structured, quasi-structured, unstructured

- Data science is a very diverse discipline
  - Maths, computer science, statistics, applications

→ Data scientists require a diverse skillset