# Lecture 1

# Introduction to Data Mining

**Golnaz Taheri, PhD**
Senior Lecturer, Stockholm University

Stockholms universitet

# Course Logistics

o https://ilearn.dsv.su.se/course/view.php?id=1678
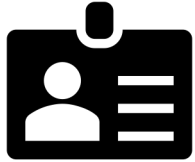
- Course activities:
  - o 7 weeks: 35 - 41
  - o 12 Lectures: Aug 28 – Oct 13
  - o 6 Labs &  Q&A sessions
  - o Written Exam: Oct 20

- Instructor and Responsible teachers:
  - o Golnaz Taheri          golnaz.taheri@dsv.su.se
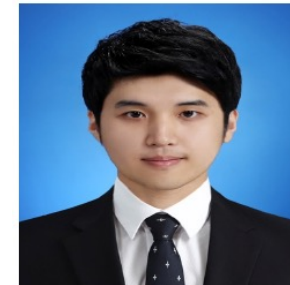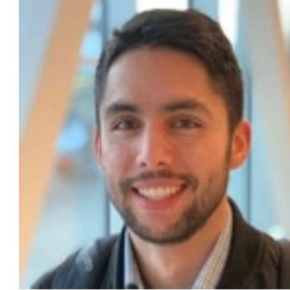  - o Ioanna Miliou          ioanna.miliou@dsv.su.se

Stockholms universitet

# Course Logistics

- **Course Assistants:**

  o Luis Quintero    luis-eduardo@dsv.su.se
  o Maria Bampa    maria.bampa@dsv.su.se
  o Zed Lee    zed.lee@dsv.su.se

Stockholms universitet

# Course page on ilearn

# Course Syllabus

| Week 35 | Week 36 | Week 37 | Week 39 | Week 40 | Week 41 |
|---------|---------|---------|---------|---------|---------|
| Introduction to Data Mining 08/28 | Dimensionality reduction 09/04 | Clustering \|\| 09/11 | Classification \|\| 09/25 | Model evaluation 10/02 | Advanced Topics II Graph Mining 10/09 |
| Introduction to Python 08/30 | Data Preparation using Python 09/06 | Clustering using Python 09/13 | Classification using Python 09/26 | Advanced Topic \| Neural Network 10/03 | Deployment 10/10 |
| Association Rules 08/31 | Clustering \| 09/07 | Classification \| 09/15 | Classification \|\|\| 09/29 | Model Evaluation 10/05 | Exam Review 10/13 |

Stockholms universitet

5

# Course workload

Assignments             3 hp

- Three programming assignments (Python)
- Online quizzes

Written Exam            4.5 hp

Stockholms universitet

# Homework Assignments

- To be done <span style="color:red">individually (strictly)</span>
- Will involve programming in <span style="color:red">Python</span>
- Each corresponding to a lab session

<span style="color:red">Plagiarism</span> is not acceptable, such as:

- Borrowing code from the internet (chatGPT) and submitting it as is or with minor changes
- Borrowing code from each other and submitting it as is or with minor changes
- Borrowing code from previous years and submitting it as is or with minor changes

Stockholms universitet

# Homework Assignments

Submissions:

- Before a given <span style="color:red">deadline</span>
- <span style="color:red">Late submissions:</span> Not Allowed
- <span style="color:red">Second deadline:</span> November 15th
    - <span style="color:blue">OBS: penalty of 50% off the obtained grade</span>

Stockholms universitet

# Quizzes

- 6 weekly online quizzes (lowest quiz grade to be dropped, and the best five will count)

- Questions on previous lectures (1 to 3 lectures)

- Only one attempt per quiz

- All quizzes will be timed!

- No make-up quizzes possible!

Stockholms
universitet

# Homework Assignments

- HW1:           4 pts

- HW2:           5 pts

- HW3:           6 pts

- Quizzes:     5 pts

- Total points:   20 pts
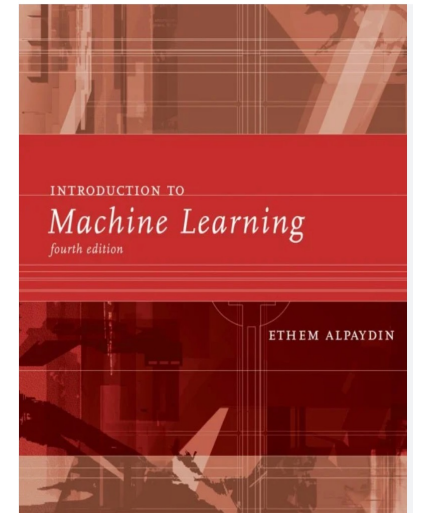
- To pass you need 12 pts

- Grading scheme: A – F

# Exam

- Two versions:
  - **DAMI**: on-campus
  - **DAMI-DIST**: online

- Two parts:

  Part A: multiple-choice questions

  Part B: free text questions

- This will examine your ability on what you have learned

- To pass you need at least 60% of the points

- Grade scheme: A – F

Stockholms universitet
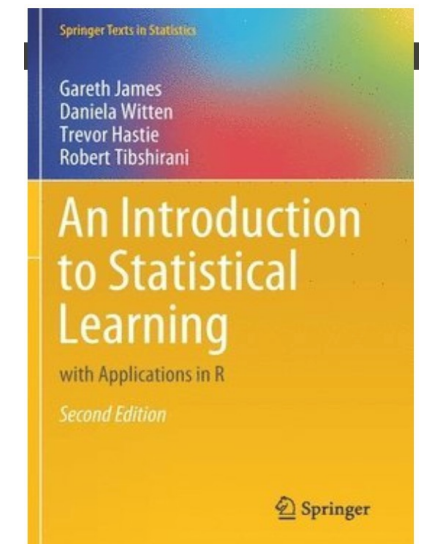
# Course textbook

**Main:**

- Introduction to Machine Learning (fourth edition)
      Publisher: MIT Press
      Year: 2022
      ISBN: 978-0-2620-4379-3

**Additional:**

- An Introduction to Statistical Learning with applications in R
      Publisher: Springer
      Year: 2013
      ISBN: 978-1-4614-7138-7
      URL: http://www-bcf.usc.edu/~gareth/ISL/

# Learning Objectives

- Become familiar with data science and its algorithms

- Be able to identify a correct algorithmic solution to a given problem

- Be able to apply these algorithmic solutions to solve practical problems

- Be able to perform basic data analysis on real data using Python

Stockholms universitet

# Introduction

- Why we need Data Analysis?

- What is Data Science (DS) ?

- What is Data Mining (DM) ?

- What is Artificial Intelligence (AI)?

- What is Machine Learning (ML) ?

Stockholms universitet

# Why we need Data Analysis?

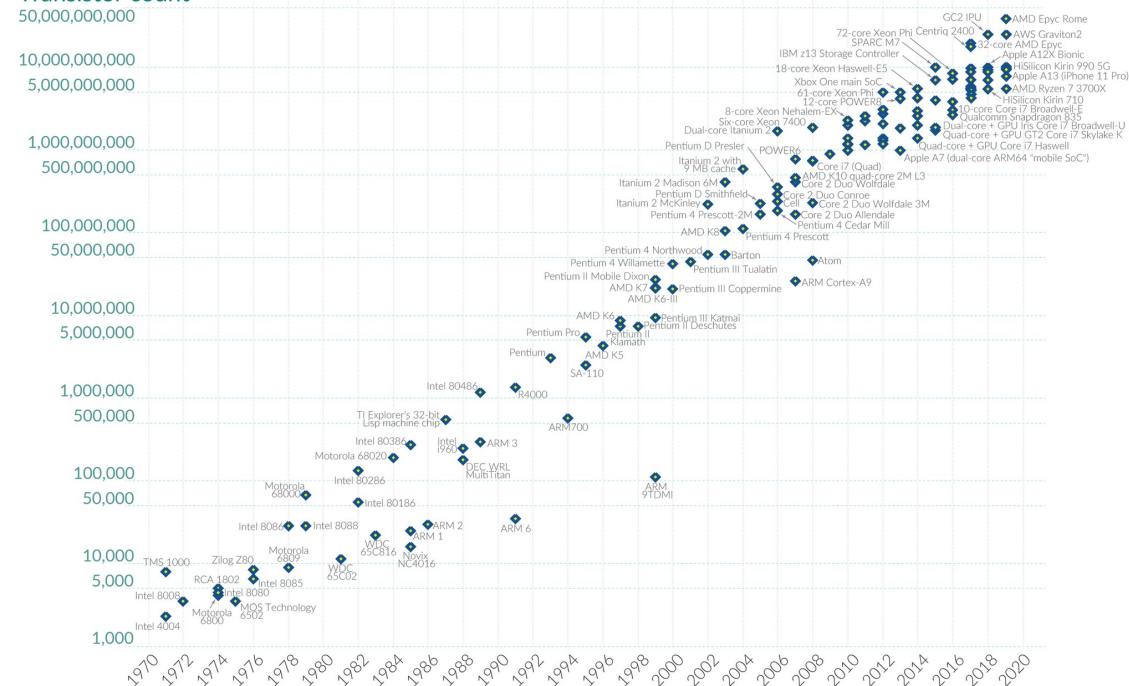- Computational power

  - More efficient <span style="color:red">processors,</span> larger <span style="color:red">memories</span>

Moore's Law: The number of transistors on microchips doubles every two years

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years.
This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

Our World in Data

Stockholms universitet

# Why we need Data Analysis?

- Data collection and transfer
  - <span style="color:red">Communication</span> and <span style="color:red">measurement</span> technologies have improved
- Data storage
  - Huge <span style="color:red">hard disks</span>
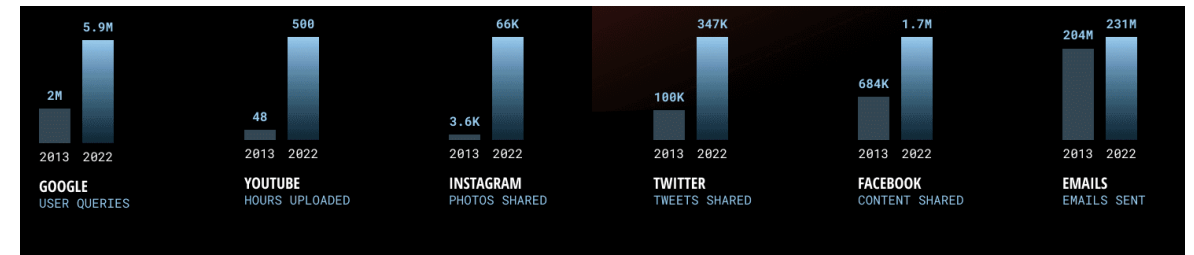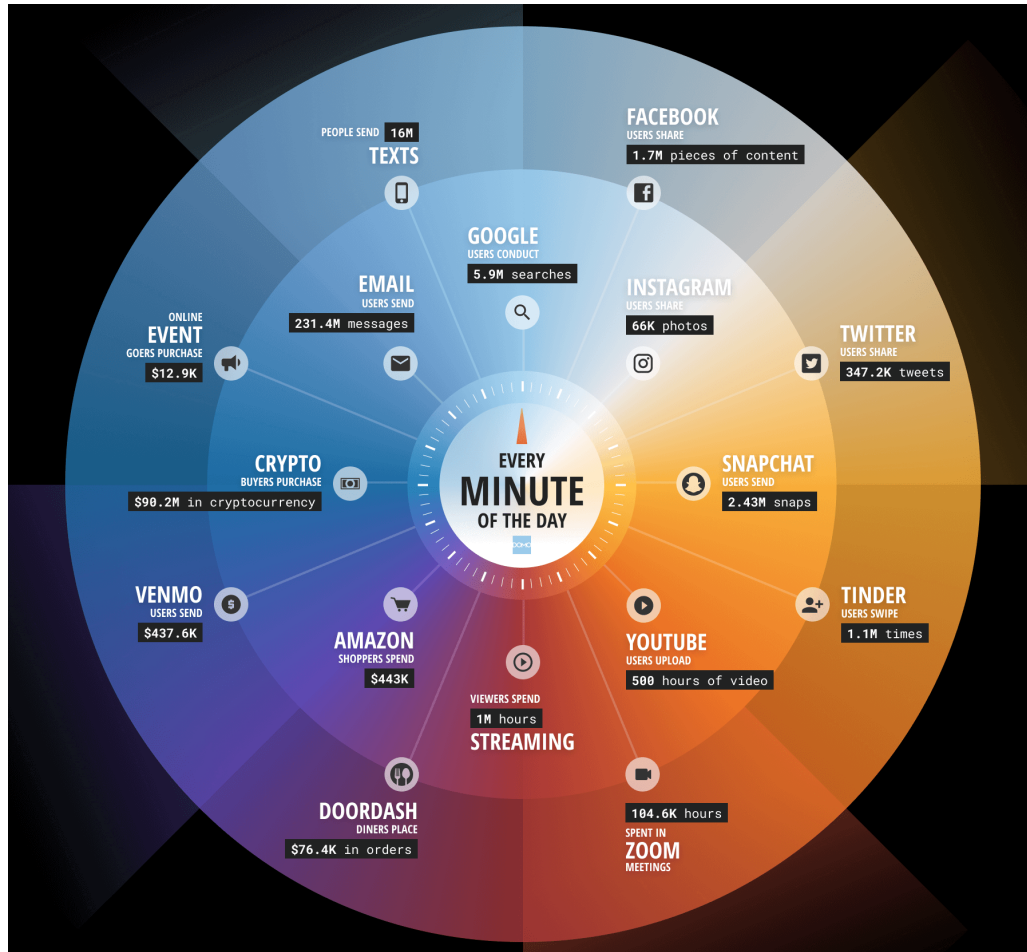  - Data on the <span style="color:red">cloud</span>

Stockholms universitet

# Why we need Data Analysis?

- It is possible to collect and store lots of raw data

- But…data analysis methods are lagging behind

- Need to analyze the raw data to extract knowledge

Stockholms universitet
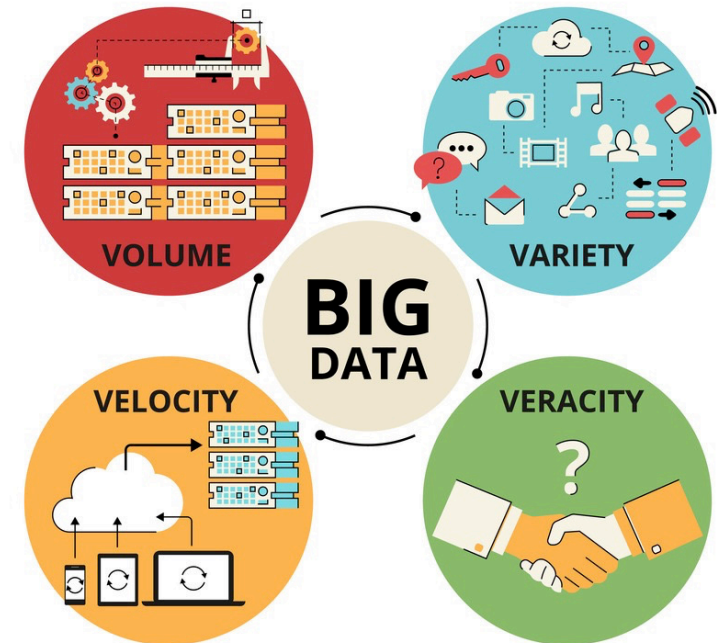
# Data Never Sleeps!

# The Four V's of Big Data

Volume: The first V of big data is all about the amount of data.

Velocity: The second V of big data, is all about the speed new data is generated and moves around.
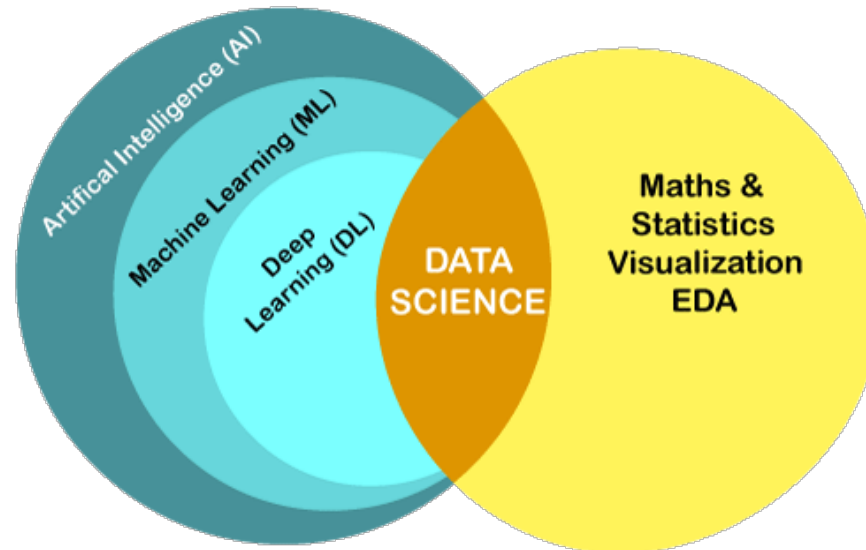
Variety: Data is generally one of three types: unstructured, semi-structured and structured and algorithms required to process the variety of data generated varies based on the type of data to be processed.

Veracity: Denotes the trustworthiness of the data. Is the data accurate and high-quality?

Stockholms universitet

# What is Data Science

- **Data science** is an interdisciplinary field that uses statistics, scientific computing, scientific methods, algorithms and systems to extract knowledge and insights from structured, and unstructured data.
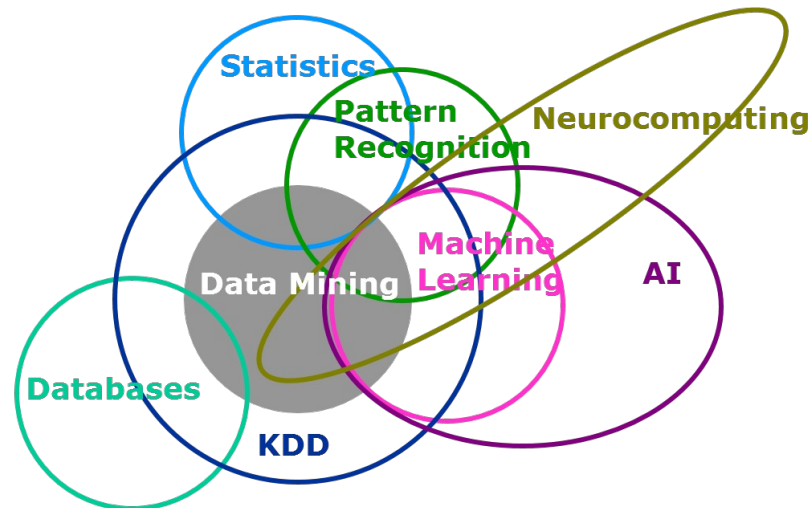
Stockholms universitet

# Why Data Science is important?

- Data science is revolutionizing the way companies operate. Many businesses, regardless of size, need a robust data science strategy to drive growth and maintain a competitive edge.
  - DS allows businesses to uncover new patterns and relationships that have the potential to transform the organization

  - DS can reveal unnoticed gaps and problems. Greater insight about purchase decisions, customer feedback, and business processes can drive innovation in internal operations and external solutions.

  - DS can help companies predict change and react optimally to different circumstances.

Stockholms universitet

21

# What is Data mining

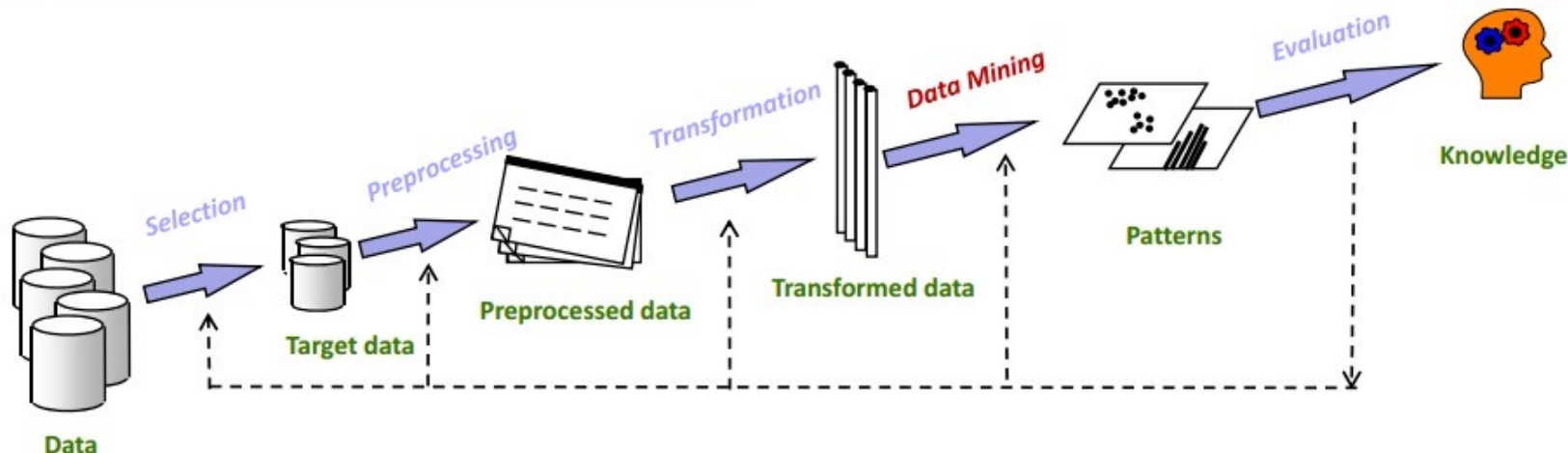- Data mining is the process of extracting and discovering patterns in large data sets.
- The overall goal of data mining is extracting information (with intelligent methods) from a data set and transforming the information into an understandable structure.
- Data mining is the analysis step of the "knowledge discovery in databases" process, or KDD.

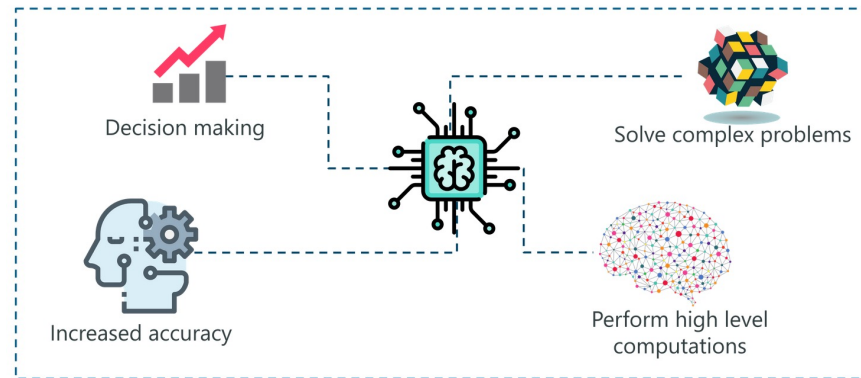Stockholms universitet

# What is KDD

- Knowledge Discovery in Databases (KDD) is a process that involves the extraction of useful, previously unknown, and potentially valuable information from large datasets.
- KDD is a multi-step process that encourages the conversion of data to useful information. Data mining is one of the steps of KDD which is the pattern extraction phase of KDD.

[Fayyad, Piatetsky-Shapiro & Smyth, 1996]

# What is Artificial Intelligence

- In 1956, the term Artificial Intelligence (AI) was defined by John McCarthy as:

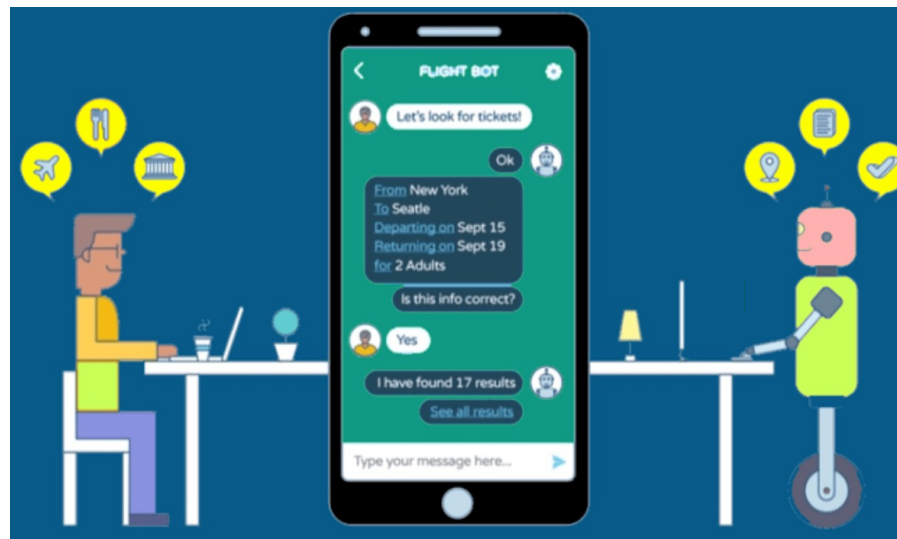  'The science and engineering of making intelligent machines.'



- AI is a machine's ability to perform the cognitive functions we associate with human minds, such as reasoning, learning, interacting with an environment, problem solving, and even exercising creativity.
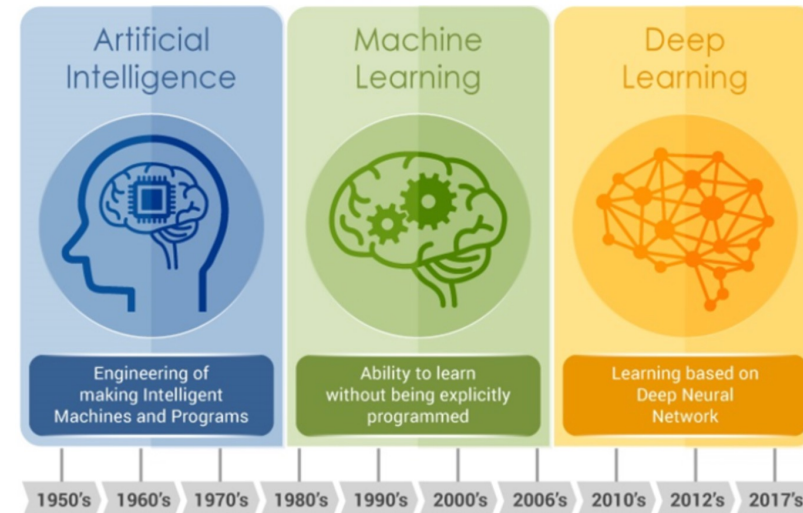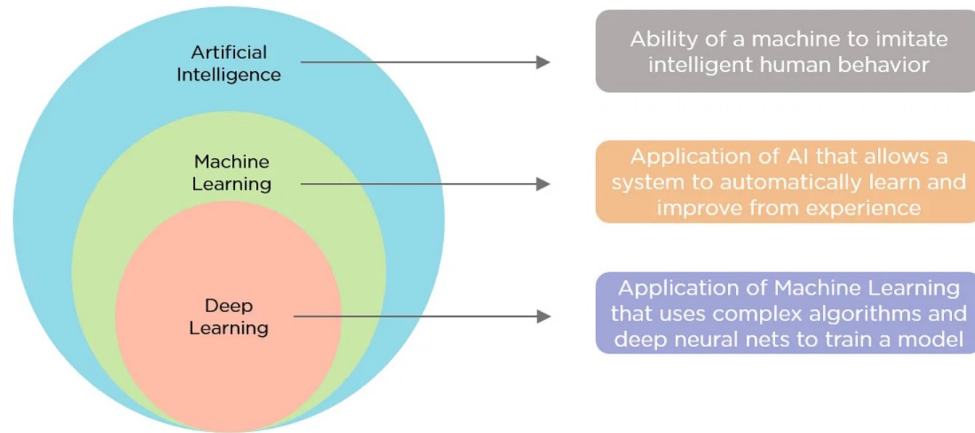
# AI example

- # Chatbots

- Answering a customer's inquiries can take a long time.
- The use of algorithms to train machines to meet customer needs through chatbots is an AI solution.
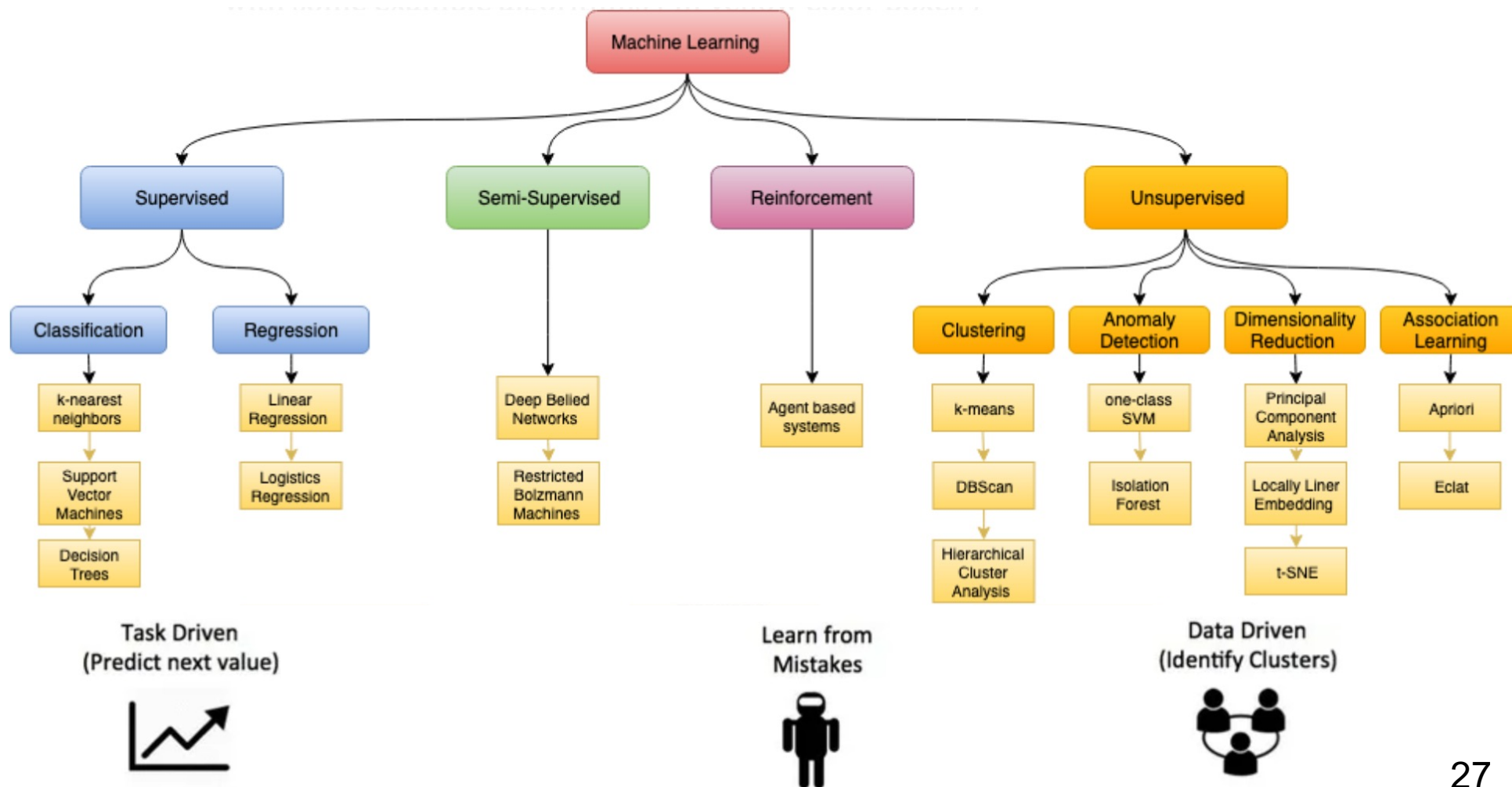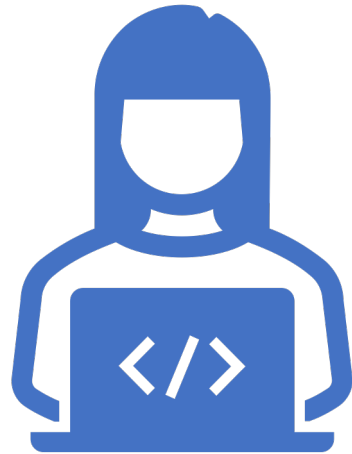- This allows machines to answer as well as take and track orders.

Stockholms universitet

# What is Machine Learning

- Machine learning is a branch of AI and computer science which focuses on the use of data and algorithms to imitate the way that humans learn, gradually improving its accuracy.

# Types of Machine Learning

Thanks!

golnaz.taheri@dsv.su.se

Stockholms universitet