# Machine Learning & Data Science

ImpactDeal 2022
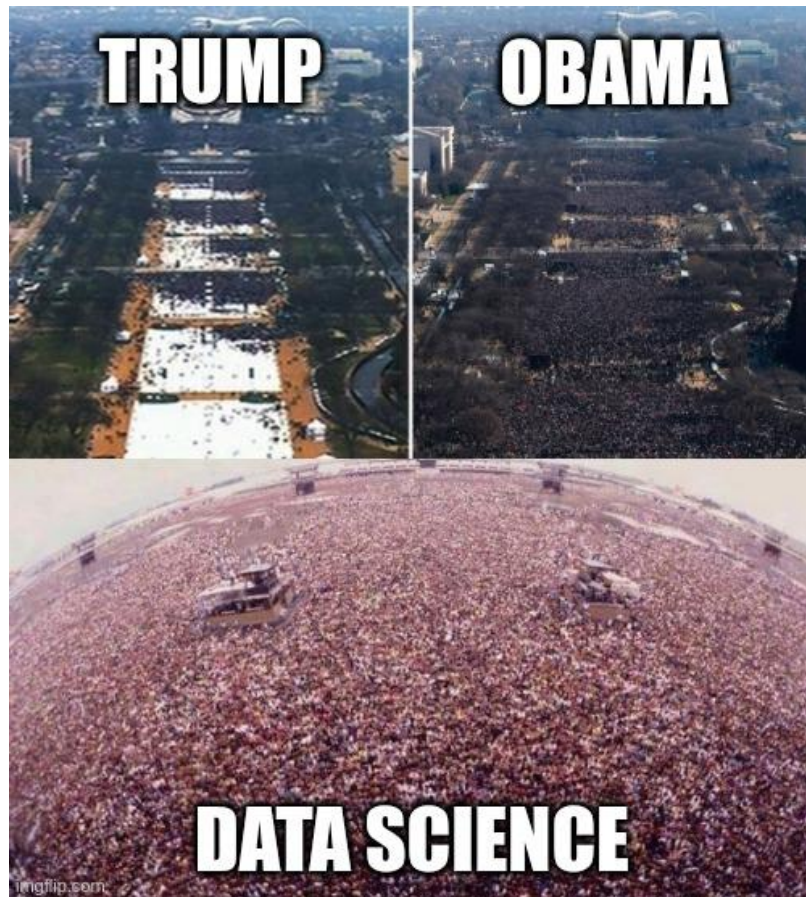
# Yet another online course…?



"How was the online class? What did you get from it?"

Me:

- Books,

- Papers,

- Bootcamps, academies, masters,

- Online courses,

- Online platforms,

- Conferences,

- Videos,

- Blogs,

- More blogs,

- …

# But we can do something different!



- **Hands on**: we will see and write a lot of code.

- **Collaborative**: we will work together and/or in groups.

This course is designed around the idea of participation:

- Ask questions!
- Give feedback!
- Turn your camera on!
- Communicate your ideas!

# Learning Objectives

## Knowledge

- Structure of a data project

- Tools for data science in Python

- Fundamentals of Machine Learning

## Skills

- Techniques for data exploration

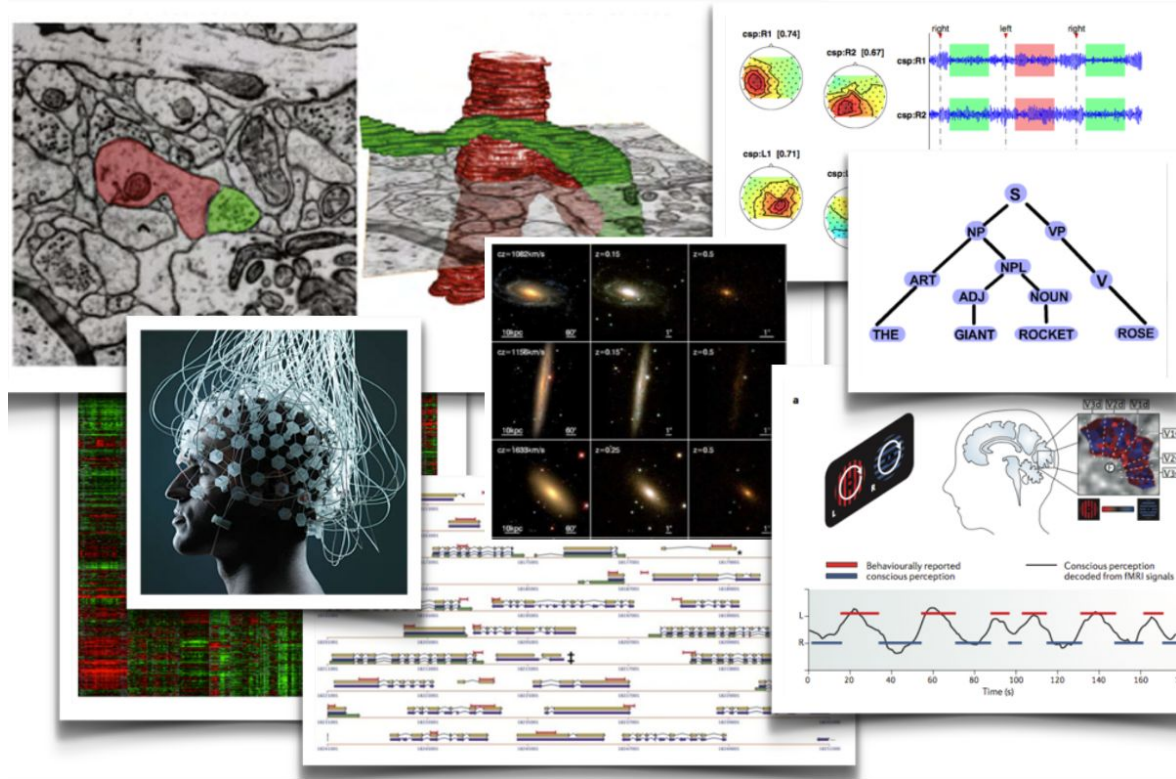- Training ML models

- Ability to deal with complex data

## Attitude

- Critical thinking about data

- Creativity with data analytics
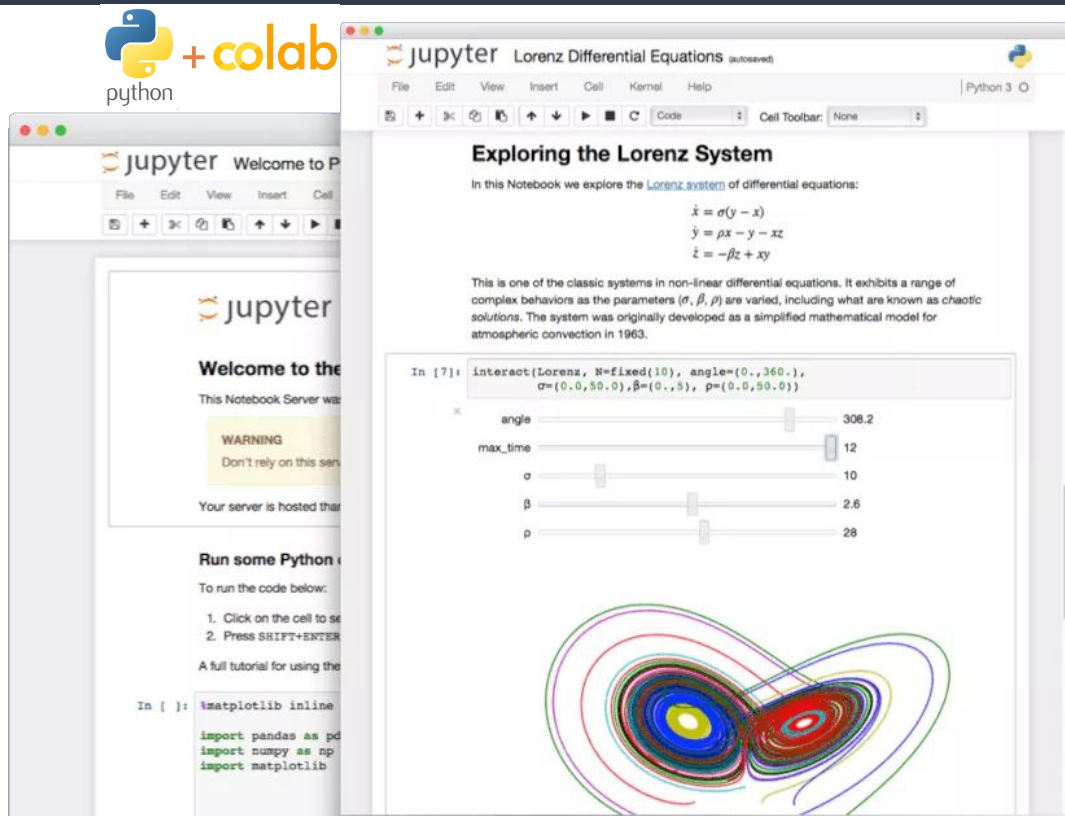
# Learning Tools

**Slides (PDF):**

- Introduction to the topics,

- Theoretical concepts,

- No-code examples.

# Learning Tools

**Jupyter notebooks (Colab):**

- Example code,

- Exercises,

# Learning Tools

## Quizzes (Google Forms):

- Short and simple questions,

- Useful to self-assess learning path,

- Helpful for Q&A session,

- No grades.

---

### Sample Quiz

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam ornare turpis in nulla varius efficitur. Sed varius est eget ullamcorper lacinia. Sed hendrerit augue non iaculis consectetur. Mauris venenatis, urna non faucibus maximus, magna sapien feugiat ante, eu gravida urna dolor mollis metus. Aliquam molestie nulla sed purus varius tristique. Praesent vitae iaculis est. Donec malesuada tempus turpis in facilisis. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas.

Sign in to Google to save your progress. Learn more

*Required

Email *

Your email address

Sample question

○ Option 1

○ Option 2

○ Option 3

○ Option 4

Submit                              Clear form

# Learning Tools

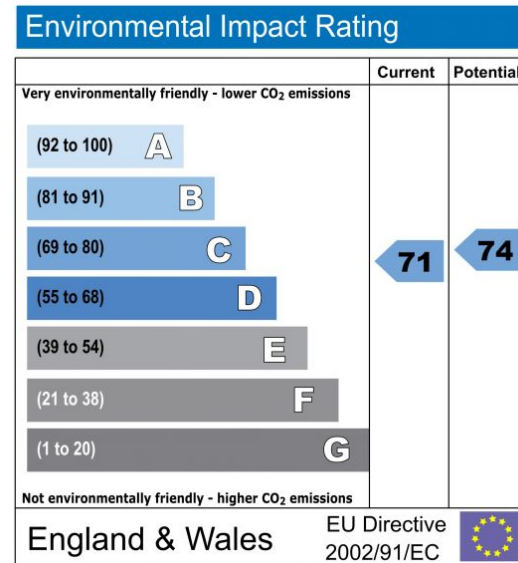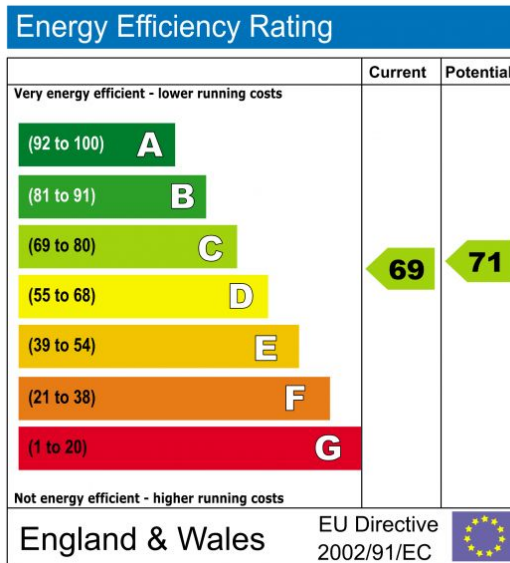**Collaboration:**

- GitHub (https://github.com/darioka/impactdeal-2022),

- Teams,

- Telegram.

# Project

## Estimation of Building Energy Efficiency

We will try to build a machine learning model able to predict the Energy Efficiency Rating (EER) of a dwelling, using historical data of Energy Performance Certificates of England and Wales.
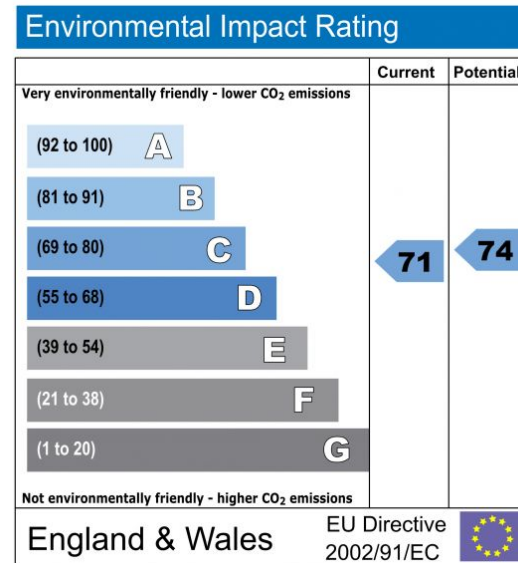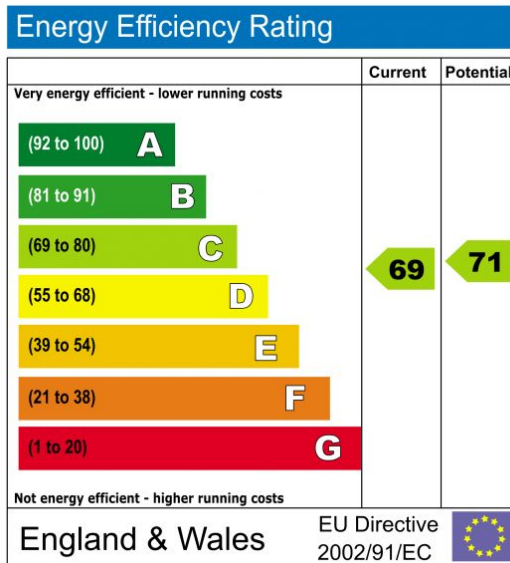
# Project

## Estimation of Building Energy Efficiency

**What is a EPC rating?**
EPC is a review of the energy efficiency of a property, which is labelled from A (very efficient) to G (inefficient). EPC are valid for 10 years and are needed whenever a property is sold or rented.

**How is EPC rating calculated?**
A trained professional conducts an inspection of the property and assesses the energy efficiency of walls, windows, heating and water systems, etc.

# Project
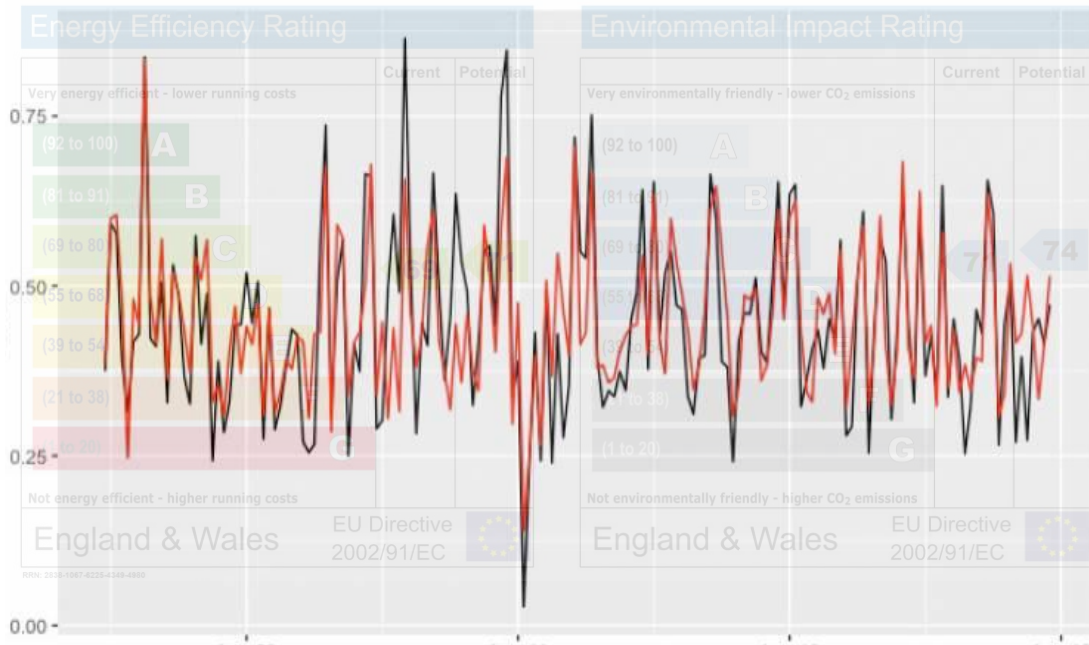
## Estimation of Building Energy Efficiency

**What is the goal of this project?**
We want to build a machine learning model to predict the EPC band of a property. The model will be trained on historical EPC data, available on opendatacommunities.org.

**How do we do it?**
We will follow the project throughout the course, applying the techniques we will be learning and discussing their implications on the EPC prediction problem. Analyses and training will be performed on EPC data from three major UK cities (*) and can be done entirely on Colab notebooks.

(*) the data has been downloaded, subsampled, pseudonymized (address and postcode) and uploaded to the course's Github repository.

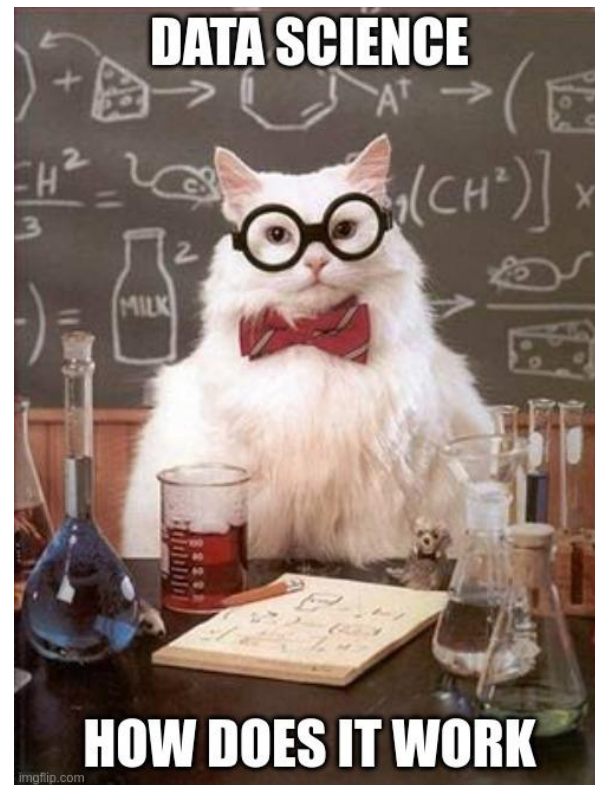# Data Science Fundamentals

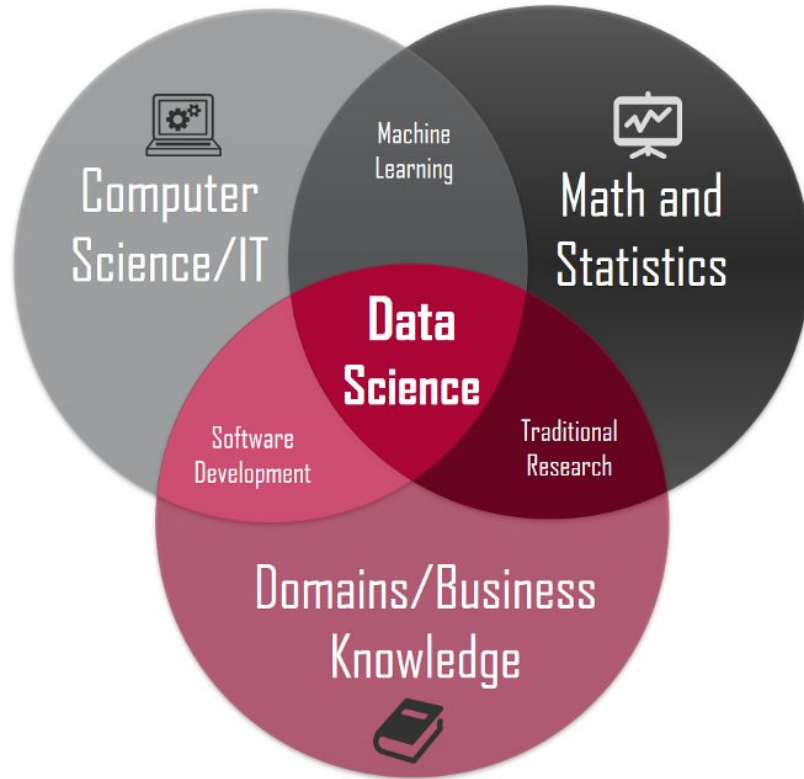# What is Data Science?

Not very well defined...

- ***Circular***:
  - *Data science is what data scientists do*
- ***General***:
  - *Data science is the science of learning from data*

Very broad field, but:

- it has to do with **science**,
- it has to do with **data**.
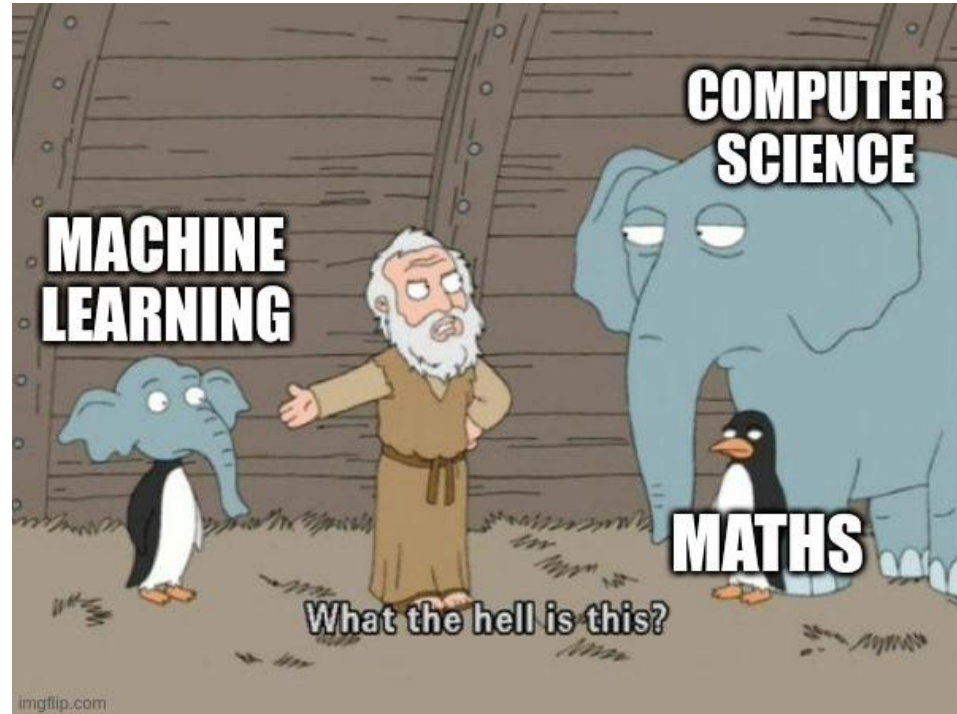
# What is Data Science?

# What is Machine Learning?

Better definitions:

***The study of computer algorithms that can "learn" from data to solve tasks, without being explicitly programmed to do so.***

Machine learning algorithms are based on **training data** and are often formulated as **minimization** of some **loss function** i.e. as **optimization** problems.
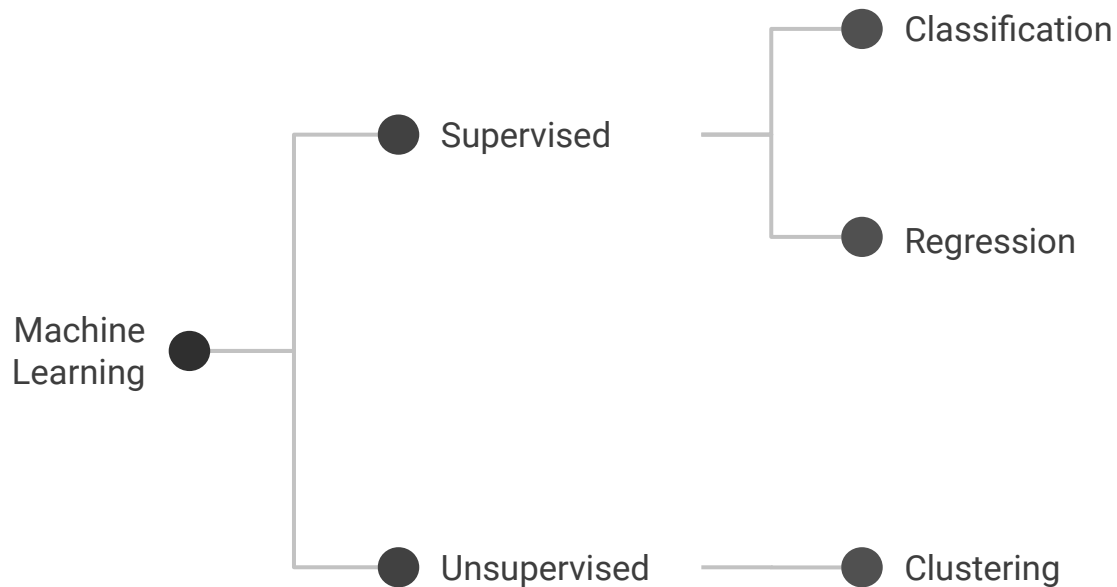
# What is Machine Learning?

Examples:

- **Recommenders**:
  - Not explicitly programmed to suggest object A if object B has been chosen. Not a set of rules.
  - Try to predict users' preferences, based on historical data.

- **Image recognition**:
  - Does not require manual feature extraction or expert computer vision knowledge,
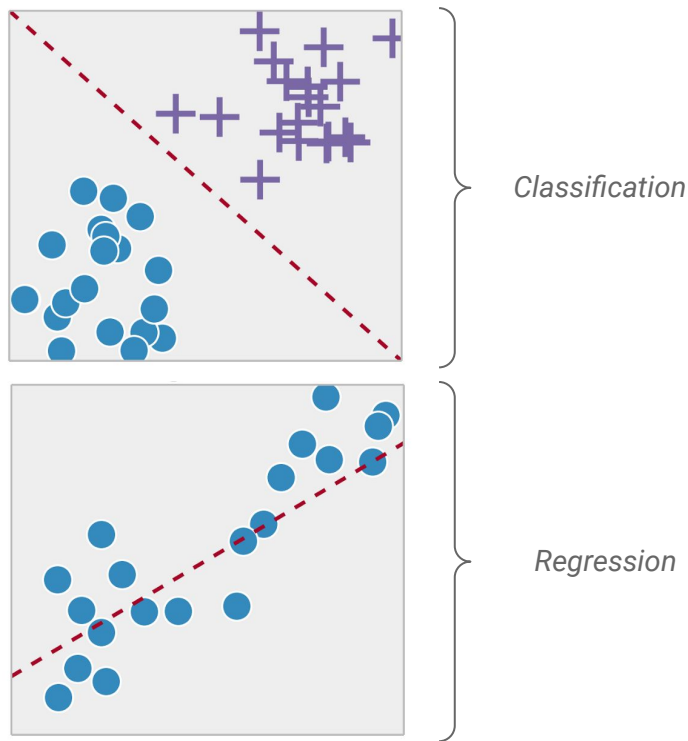  - Learn properties and relationships of pixels in images.

# Machine Learning Approaches

# Machine Learning Approaches
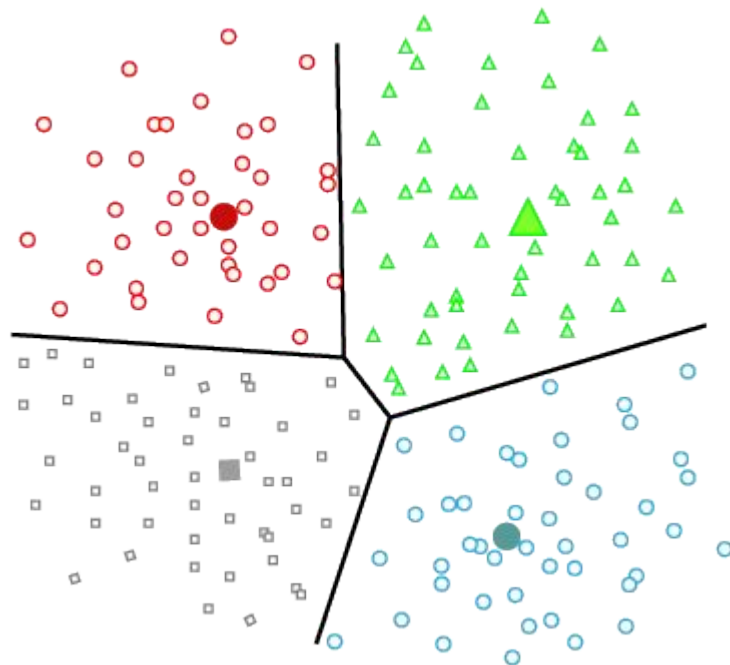
## Supervised

- **Predict the value of the target** for each sample,
- Models requires both **features** and **target** in the training data (i.e. labeled data).


- **Classification**:
  - Target is a discrete variable,
  - Model example: logisitc regression,
  - Application example: fraud detection.
- **Regression**:
  - Target is a continuous variable,
  - Model example: linear regression,
  - Application example: demand forecasting



*Classification*

*Regression*

# Machine Learning Approaches

## Unsupervised

- **Find patterns in the data**,
- There is not "target" value to predict or the target is absent from training data.


- **Clustering**:
    - Divide input samples into groups,
    - Model example: K-means,
    - Application example: customer segmentation
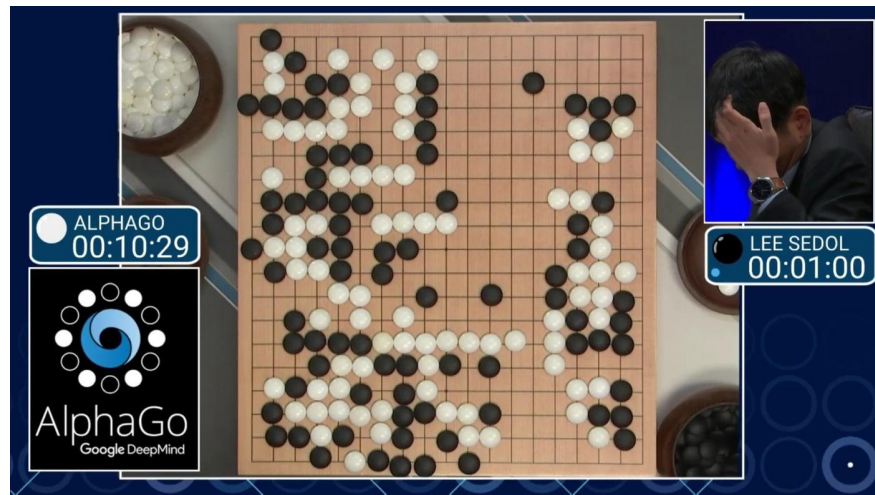
# Machine Learning Approaches

There are techniques and/or tasks that may be either supervised, unsupervised or both:

- **Semi-supervised learning**: training dataset has both labeled and unlabeled samples.

- **Dimensionality reduction** and **manifold learning**: transformations of data from a high-dimensional space to a low-dimensional representation.

- **Association rule learning**: methods for discovering relationships and strong rules between variables in large databases.

# Machine Learning Approaches

## Reinforcement Learning

- **Agents** interacting and taking **actions** in an environment.
- Solve a task trying to find an optimal strategy that maximizes a **reward** for the agent.
- Different approach: no need of large training data.
- Usually much harder than supervised/unsupervised. Active research area.
- Examples: AlphaGo

# Machine Learning Projects

# Project Workflow

1. Business Understanding

2. Data Preparation

3. Modeling

4. Deployment

5. Management



Machine Learning Life Cycle

01
02
03
04
05

# Project Workflow

1.  **Business Understanding**
    - Problem definition
    - Objectives
    - Expected outcomes
    - Assessments (risk, infrastructure, …)
    - Understand integration/deploy

2.  Data Preparation
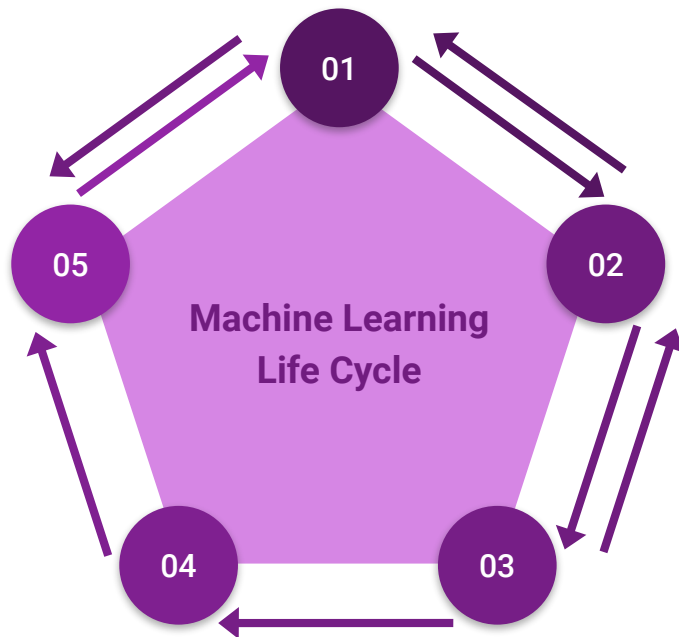
3.  Modeling

4.  Deployment

5.  Management

# Project Workflow

1. Business

2. **Data**
   - Data gathering
   - Data exploration
   - Data cleaning/preprocessing
   - Feature engineering

3. Modeling

4. Deployment

5. Management

# Project Workflow

1. Business Understanding

2. Data Preparation

3. **Modeling**
   - Design experiments
   - Train models
   - Evaluate and test models
   - Interpret results
   - Review and cross-check

4. Deployment

5. Management

# Project Workflow

1. Business Understanding

2. Data Preparation

3. Modeling

4. **Deployment**
   - **Create model artifacts**
   - **Write production-ready code**
   - **Integration**
   - **Plan maintenance**

5. Management



Machine Learning Life Cycle

# Project Workflow

1. Business Understanding

2. Data Preparation

3. Modeling

4. Deployment

5. **Management**
   - Monitoring input/output
   - Monitoring performances
   - Model updates

# Project Workflow

Data Exploration

Modeling

Deploy

# Data Exploration

It's the process of analyzing a dataset to understand its characteristics.

It is a mixture of:
- summary statistics,
- data visualization,
- hypothesis testing,
- manual drill-down.

# Data Exploration

Typical steps in data exploration are:

- Variable identification,
- Univariate analysis,
- Bivariate analysis,
- Missing data analysis.

# Data Exploration

## Variable identification

*What kind of information does the dataset contain?*

- Identify the target (if any),
- Understand the meaning of the variables,
- Identify data types.

# Data Exploration

## Univariate analysis

*How the data is distributed?*

For categorical variables:

| Methods | Visualization |
|---|---|
| Counts Frequencies | Bar plots |

# Data Exploration

## Univariate analysis

*How the data is distributed?*

For numerical variables:

| Methods | Visualizations |
|---|---|
| Central tendency (mean, median, …)<br><br>Dispersion (range, quartiles, variance, …) | Histograms<br>Box Plots<br>Violin plots |

# Data Exploration

## Bivariate analysis

*What are the relationships between the variables?*

| Methods | Visualizations |
|---------|----------------|
| Correlation | Scatter plots<br>Heatmap |
| Joint frequencies<br>Distribution difference<br>between groups | Histograms<br>Box Plots<br>Violin plots |

# Data Exploration

## Missing Values

*Is the data complete?*

Common situation with real-world data sets, where some variables contains no information, for example because of errors during the data collection process.

Missing data can be a problem for any analysis:

- Loss of information and statistical power,
- If systematic, introduce bias and distortions.



THE DATA IS MESSY AND FULL OF ERRORS

# Modeling



**Data Preparation**

- Cleaning
- Preprocessing
- Feature Engineering

**Model Training**

- Model fitting
- Hyperparameter tuning

**Machine Learning Life Cycle**

**Evaluation**

- Test performance
- Interpretation

# Modeling

*You may think this is how data scientists spend their time…*

*…but it is the other way around.*



Datasets
Model



Datasets
Model

# Deploy

**Deployment** is the integration of a machine learning model into a production environment.

Examples:

- A program that targets customer at risk of churning and send them personalized offers every month.

- An online bank website that recognizes uploaded documents.

- A virtual assistant device that answers to questions.

# Deploy

Challenges:

- Integration,

- Reproducibility,

- Scalability,

- Drift.

# Python for Machine Learning

# Python for Machine Learning

# Python for Machine Learning

*What do **you** think?*

# Python for Machine Learning