# Machine Learning

## Lecture 1: Introduction

Bojan Božić & Bujar Raufi
School of Computer Science
TU Dublin, Grangegorman

bojan.bozic@tudublin.ie; bujar.raufi@tudublin.ie

**OLLSCOIL TEICNEOLAÍOCHTA
BHAILE ÁTHA CLIATH**

**TU DUBLIN**

**TECHNOLOGICAL
UNIVERSITY DUBLIN**

**Slides adapted from Sarah Jane Delany and book slides from: Fundamentals of Machine Learning for Predictive Data Analytics. Kelleher, Mac Namee and D'Arcy**
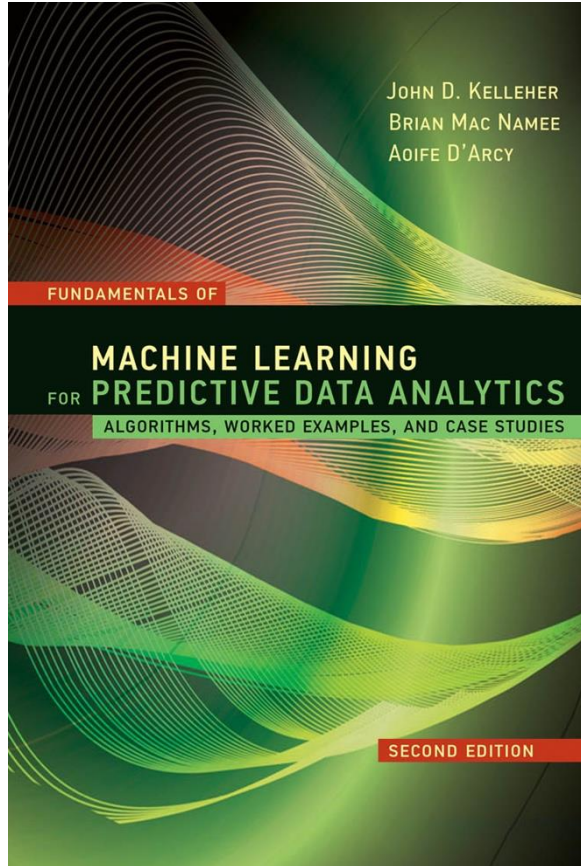
# Overview

- Administrivia

- Module Outline

- Machine Learning Today - the good, the bad and the ugly

- Machine Learning Introduction

- What is (supervised) ML?

- Supervised v Unsupervised Learning

- Representing Data as Features

# Administrivia

- Lecture and Lab:
  - Thursday: 6pm to 10pm
  - Weeks 1, 7 and 13 will be in-person (but streamed), the rest will be online
  - No lecture on week 13 (07.05.2026)!
- All notes, lecture recordings, tutorials, lab work, and assignments will be available on Brightspace.

- For all module queries please contact: bojan.bozic@tudublin.ie

# Textbook

*Fundamentals of Machine Learning for*

*Predictive Data Analytics*: John D. Kelleher, Brian

Mac Namee, Aoife D'Arcy

# ML with Python

- Tutorials and lab work require a laptop with Jupyter and scikit-learn.



https://scikit-learn.org/stable/

# Assessment

- Assessment is based on CA + final exam:

| Percentage | Activity |
|------------|----------|
| 30% | Assignment - due in week 13 |
| 20% | Lab Test - scheduled for week 9 |
| 50% | End of Semester Exam |

# Machine Learning - Overview

- Supervised Learning
- Classification: KNNs, Decision Trees, Naive Bayes
  - Neural Networks
  - Linear regression, Logistic Regression
- Dimensionality Reduction
  - Feature Selection, PCA
- The ML Process
  - Data Preprocessing, Missing Values, Scaling
  - Model Selection, Hyperparameters
- Evaluation
- Working with Text
- Unsupervised Learning

# Relevance of ML

- Explosion in rich, complex data to analyse - online and offline.

- Significant recent progress in algorithms and theory.

- Computational power is now available.

- Industry demand - Data scientists, Data engineers…

- New applications in many disciplines - medicine, engineering, humanities…
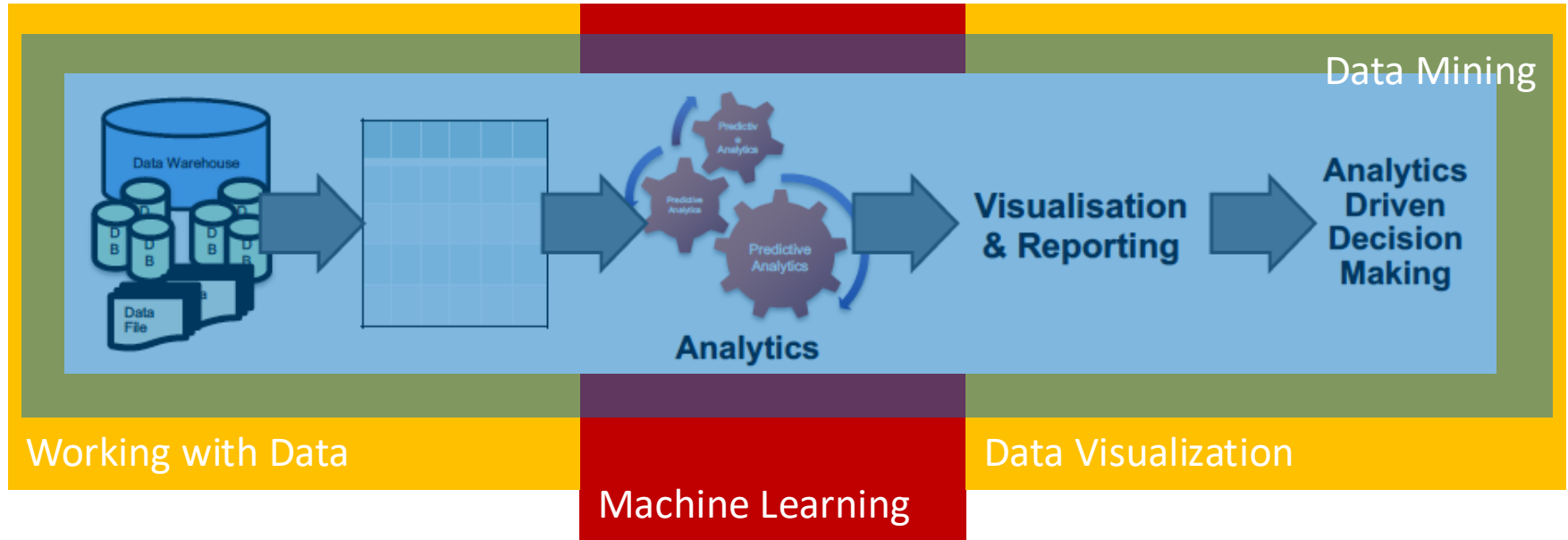
2.4 billion users end of 2024
Average Daily Usage of Instagram is 33 Minutes

more than 1.8 billion users worldwide. Around 22.22% of the global population uses Gmail. 121 billion emails daily (as of 2022 data)

# How ML fits in this program



Data Mining

Working with Data

Machine Learning

Data Visualization

# ML Today – Advances in AI

- Dall-E:  generating realistic images from a description in natural language https://openai.com/dall-e-2/



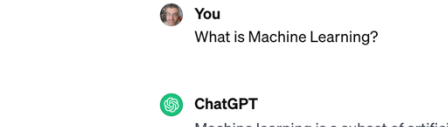"a robot using a computer in the desert digital art"

"children playing football in a school year abstract art"

"photograph of an astronaut riding a horse"

# Advances in AI

- ChatGPT:  A chatbot built on large language models that interacts in a conversational way https://openai.com/blog/chatgpt/

# Reactions

**Two Government departments confirm use of Chat GPT**

**ASIA PACIFIC**

**As Australian colleges crack down on ChatGPT, disabled students defend AI**

**Irish universities in race against time with ChatGPT to avoid widespread cheating scandals**

Education heads fear ChatGPT risk to 'academic integrity'

**THE SHIFT**

**How ChatGPT Kicked Off an A.I. Arms Race**

Even inside the company, the chatbot's popularity has come as something of a shock.

OLLSCOIL TEICNEOLAÍOCHTA
BHAILE ÁTHA CLIATH
T U DUBLIN
TECHNOLOGICAL
UNIVERSITY DUBLIN

# Education



**Performance of ChatGPT on USMLE: Potential for AI-Assisted Medical Education Using Large Language Models**

Tiffany H. Kung, Morgan Cheatham, ChatGPT, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, Victor Tseng

**doi:** https://doi.org/10.1101/2022.12.19.22283643

https://www.medrxiv.org/content/10.1101/2022.12.19.22283643v2.full

🔔 **Follow this preprint**



https://www.cam.ac.uk/stories/ChatGPT-and-education



## Would Chat GPT3 Get a Wharton MBA?
A Prediction Based on Its Performance in the Operations Management Course
by Christian Terwiesch (terwiesch@wharton.upenn.edu)

https://mackinstitute.wharton.upenn.edu/wp-content/uploads/2023/01/Christian-Terwiesch-Chat-GTP.pdf

# GPTZero to detect AI



T- The study involves assessing the performance of REST APIs over a specific time frame. This time component is crucial for understanding how the APIs behave under different levels of user activity and stress, including any variations over an extended period. The time dimension is integral to the overall investigation.
Abstract:
The development and deployment of REST APIs have become increasingly popular in recent years, as they provide a simple and standardized way to access web services. With the increasing number of users and requests to these APIs, it is crucial to ensure that they can handle the load and maintain a high level of performance. This project aims to explore the process of performance testing and analysis of a REST API, also load testing will be conducted, stress testing and soak testing to evaluate
the API's response time and throughput under different levels of load.
Keywords: REST API, web services, performance testing, load testing, stress testing and soak testing.

**Detect Text**        ⬆ Upload File        2,331/15,000 Characters
(Get up to 100,000 here)

**Your Text is Most Likely AI/GPT generated**

11.63%
AI GPT*

For this part the PICOT (Problem, Intervention, Comparison, Outcome and Time) framework is used.

P- The research object is centered around investigating the performance of REST APIs. Specifically,

# But…

- ChatGPT can put together answers to questions but doesn't "know" anything

- No comprehensive understanding of the physical and social world, no ability to reason about relationships between concepts and entities

- Examples of failures https://github.com/giuven95/chatgpt-failures

A year ago…

This year…

This year…

A year ago…

# DeepSeek



- Nvidia 2022 export ban
- DS built at a fraction of the cost of industry-leading models like OpenAI - because it uses fewer advanced chips.
- It is reportedly as powerful as OpenAI's o1 model - released at the end of last year - in tasks including mathematics and coding.
- Like o1, R1 is a "reasoning" model. These models produce responses incrementally, simulating a process similar to how humans reason through problems or ideas. It uses less memory than its rivals, ultimately reducing the cost to perform tasks.

# Next…

- Applications will get more specific.

- Regulation

- EU AI Act

    https://artificialintelligenceact.eu/

# AI Systems are still brittle…

## AI camera operator repeatedly confuses bald head for soccer ball during live stream

https://www.theverge.com/tldr/2020/11/3/21547392/ai-camera-operator-football-bald-head-soccer-mistakes



## Tesla behind eight-vehicle crash was in 'full self-driving' mode, says driver

San Francisco crash is the latest in a series of accidents blamed on Tesla technology, which is facing regulatory scrutiny



https://www.theguardian.com/technology/2022/dec/22/tesla-crash-full-self-driving-mode-san-francisco

https://medium.com/thoughts-and-reflections/racial-bias-and-gender-bias-examples-in-ai-systems-7211e4c166a1

# Expert Systems (Rule-Based Systems)

- Rule:

**If**

– Income > Expenditure &

– Collateral > Loan

**Then**

– Risk = Low

Knowledge engineering

# Learning from historical data



Table of historical data

Each row

✓ description of instance

✓ decision (Yes/No)

Predict decision for new instance

# Learn what?

A. Learn expert decision-making
  – What if the expert gets it wrong sometimes?

B. Learn from outcomes
  – Outperform experts

# Supervised ML (Predictive Analytics)

**Training Step:**

- Learning a model from a set of historical data instances

# Supervised ML (Predictive Analytics)

**Prediction Step:**

• Using the model to make predictions

# Classification Task

Example: Credit scoring

- A training set with 10 examples (customers)

- Each example has one of two class labels = {High-risk, Low-risk}

| Example | Income | Savings | Married | Gender | Age | Class |
|---------|--------|---------|---------|--------|-----|-------|
| 1 | 35,000 | 2,000 | Y | M | 32 | High‗risk |
| 2 | 51,000 | 18,000 | N | M | 34 | High‗risk |
| 3 | 70,000 | 42,000 | Y | F | 41 | Low‗risk |
| 4 | 26,500 | 4,500 | N | M | 22 | High‗risk |
| 5 | 32,000 | 11,000 | N | F | 25 | High‗risk |
| 6 | 53,000 | 37,000 | N | F | 39 | Low‗risk |
| 7 | 88,000 | 46,000 | Y | M | 48 | Low‗risk |
| 8 | 55,000 | 5,700 | N | M | 55 | High‗risk |
| 9 | 90,000 | 35,000 | Y | F | 61 | Low‗risk |
| 10 | 43,000 | 24,000 | Y | M | 33 | High‗risk |

# Classification task

- Manually classify customers into two categories (low-risk and high-risk) based on savings and income data.

**If**

    income > α & savings > β

**Then**

    Risk - low

# Classification Task

| Example | Income | Savings | Married | Gender | Age | Class |
|---------|--------|---------|---------|--------|-----|-------|
| 1 | 35,000 | 2,000 | Y | M | 32 | High₋risk |
| 2 | 51,000 | 18,000 | N | M | 34 | High₋risk |
| 3 | 70,000 | 42,000 | Y | F | 41 | Low₋risk |
| 4 | 26,500 | 4,500 | N | M | 22 | High₋risk |
| 5 | 32,000 | 11,000 | N | F | 25 | High₋risk |
| 6 | 53,000 | 37,000 | N | F | 39 | Low₋risk |
| 7 | 88,000 | 46,000 | Y | M | 48 | Low₋risk |
| 8 | 55,000 | 5,700 | N | M | 55 | High₋risk |
| 9 | 90,000 | 35,000 | Y | F | 61 | Low₋risk |
| 10 | 43,000 | 24,000 | Y | M | 33 | High₋risk |

- Q. To which class does this new customer belong?

| Example | Income | Savings | Married | Gender | Age | Class |
|---------|--------|---------|---------|--------|-----|-------|
| X | 66,000 | 13,000 | Y | M | 44 | ??? |

- Q. Can we train an algorithm to learn to automatically classify new customers as either low-risk or high-risk?

# Supervised Learning

- Supervised Machine Learning algorithms automate the process of learning a model that captures the relationship between the **descriptive features** and the **target feature** in a **training dataset.**

# Supervised Learning

- ML algorithms search through all possible patterns that exist between a set of descriptive features and a target feature to find the best model that is **consistent** with the training data (i.e. agrees with all the training instances).

- Useful predictive models must be able to **generalise** well, i.e. make predictions for queries that are not present in the training data.

# What can go wrong?

- **Underfitting** occurs when the prediction model is too simplistic to represent the underlying relationship between the descriptive and target features.

- **Overfitting** occurs when the model is so complex that it fits the data too closely and becomes sensitive to noise (e.g. mislabelled feature values).



Dataset      Underfitting      Overfitting      Just Right

# Supervised Learning Algorithms

- Typically, training data does not have enough information to choose a single best model, so additional assumptions are needed to drive the model selection, known as the **inductive bias.**

- **Restrictive bias** constrains the set of models that will be considered (e.g. linear regression considers models that produce predictions based on a linear combination of descriptive features)

- **Preference bias** guides the algorithm to prefer certain models over others (e.g., the decision tree prefers shallower trees)

- Different algorithms have a different inductive bias,

- It is important to identify the machine algorithm that will fit the predictive task most appropriately.

# Supervised ML (aka Predictive Analytics)

# Supervised ML (aka Predictive Analytics)

# Supervised vs. Unsupervised Learning

- **Supervised Learning:**
  - An algorithm that learns a function from examples of its inputs and outputs. It uses manually labelled example data (i.e. a training set ) to predict the correct answer for new unseen query inputs.

    e.g. Classification, regression algorithms

- **Unsupervised Learning:**
  - An algorithm that finds structure in data where no manually labelled examples are available as inputs - i.e. there is no training set. These algorithms are more focused on data exploration and knowledge discovery.
  - e.g. Clustering, topic modelling algorithms

# Supervised Learning

- **Classification:**
  - Examples are represented by a set of features, which help decide the target class to which a new query input belongs (i.e. the output is a class label or target feature).

- **Regression:**
  - Examples are characterized by a set of features which help decide the value of a continuous output variable (i.e. the output is a number).



Labelled Examples     Input    Output

# Classification Tasks

- **Binary classification:** Assign a new query input to one of two possible target class labels.



Input             2 classes

- **Multiclass classification:** Assign a new query input to one of  M > 2 different target class labels.



Input             4 classes

# Representing Data

- Commonly, we use a tabular structure to represent a dataset, often referred to as the **analytics base table (ABT)**.

- Each row represents a different example and comprises a set of **descriptive features**.

- For prediction, each row also has a **target class label** or **target feature** for classification and regression, respectively - i.e. the "correct answer".

| | Descriptive Features | | | | | | | | | Target Class / Feature |
|---|---|---|---|---|---|---|---|---|---|---|
| Examples | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ...... |
| | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ...... |
| | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ...... |
| | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ....... | ...... |

# Representing Data

The descriptive features used to represent examples can be distinguished by the type and number of values they can take.

- **Binary:** Takes only two values - a boolean True/False decision
  e.g. married={True, False}, test_result={Pass, Fail}
- **Categorical (Nominal):** A feature that takes values from a finite set of values, with no intrinsic ordering to the values
  e.g. blood_group={A,B,AB,O}, nationality={French, Irish, Italian}
- **Ordinal:** A categorical variable with a clear ordering of the variables.
  e.g. grade={A, B, C, D, E, F}, dosage={Low, Medium, High}
- **Interval:** Values that allow ordering and subtraction but do not allow other arithmetic operations
  e.g date, time
- **Continuous:** Numeric measurements, with or without a fixed range for the values.
  e.g. temperature, price, age, weight, height, latitude, longitude etc.

# Typical Classification Task

- The training set with N=10 examples (customers). Each is described by D=5 features: 3 continuous, 2 categorical

- Each example has one of two class labels = {High-risk, Low-risk}

| Example | Income | Savings | Married | Gender | Age | Class |
|---------|--------|---------|---------|--------|-----|-------|
| 1 | 35,000 | 2,000 | Y | M | 32 | High_risk |
| 2 | 51,000 | 18,000 | N | M | 34 | High_risk |
| 3 | 70,000 | 42,000 | Y | F | 41 | Low_risk |
| 4 | 26,500 | 4,500 | N | M | 22 | High_risk |
| 5 | 32,000 | 11,000 | N | F | 25 | High_risk |
| 6 | 53,000 | 37,000 | N | F | 39 | Low_risk |
| 7 | 88,000 | 46,000 | Y | M | 48 | Low_risk |
| 8 | 55,000 | 5,700 | N | M | 55 | High_risk |
| 9 | 90,000 | 35,000 | Y | F | 61 | Low_risk |
| 10 | 43,000 | 24,000 | Y | M | 33 | High_risk |

Q. To which class does this new customer belong?

| Example | Income | Savings | Married | Gender | Age | Class |
|---------|--------|---------|---------|--------|-----|-------|
| X | 66,000 | 13,000 | Y | M | 44 | ??? |

# Algorithms

- Many different learning algorithms exist for prediction (e.g. k -nearest neighbour, decision tree, neural network, support vector machine).

- Due to processing, memory, and storage constraints, problem dimensions will often determine which algorithm will be practically applicable.

1. Number of input examples N.

    → Sometimes millions of input examples.

2. Number of features (dimensions) D representing each input example.

    → Often 10-1000, but sometimes far higher.

3. For classification, the number of target classes M.

    → Often small (binary), but sometimes far higher.

# Machine Learning - Overview

- Supervised Learning
- Classification: KNNs, Decision Trees, Naive Bayes
  - Neural Networks
  - Linear regression, Logistic Regression
- Dimensionality Reduction
  - Feature Selection, PCA
- The ML Process
  - Data Preprocessing, Missing Values, Scaling
  - Model Selection, Hyperparameters
- Evaluation
- Working with Text
- Unsupervised Learning

# Questions?