# Introduction

Introduction to machine learning

# Machine learning?
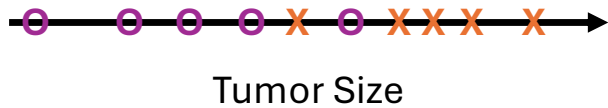
- Learning from data
  - Large datasets, from the growth of the internet, medical records, cameras & images are ubiquitous, …

- Applications we can't program by hand
  - Handwriting recognition, NLP, Computer Vision, …

- «Self-learning» algorithms
  - e.g. product or movie recommendations, spam filtering (with occasional/optional supervision input)

# Machine learning?

- Supervised learning
  - Classification, regression

- Unsupervised learning
  - Clustering, dimensionality reduction, density estimation

- Others: Reinforcement learning, sequence learning, semi-supervised learning, …

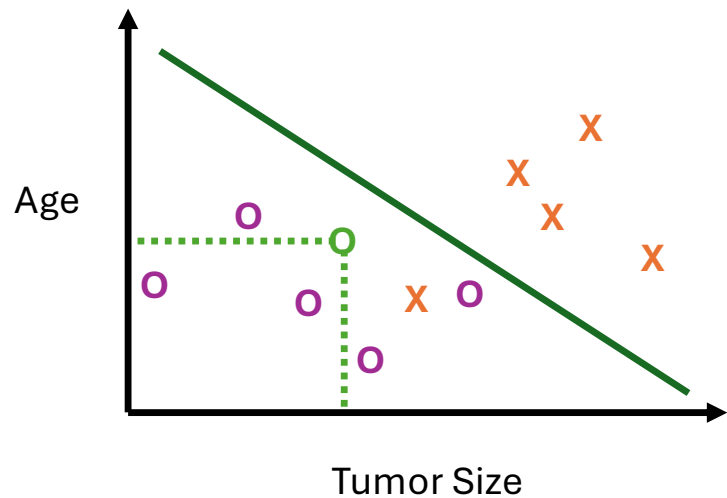# Supervised learning - Classification

Cancer data (malignant, benign)



Tumor Size

Discrete output

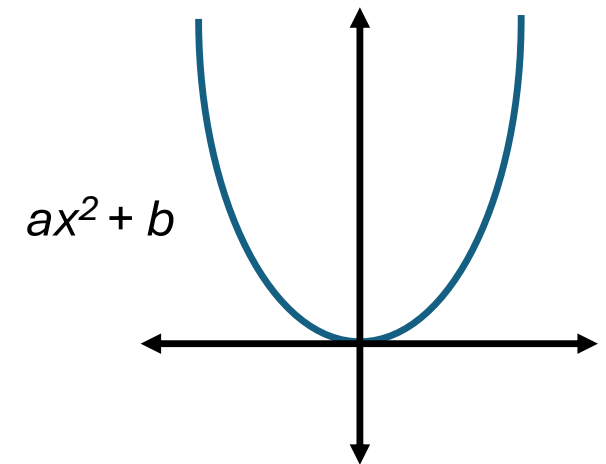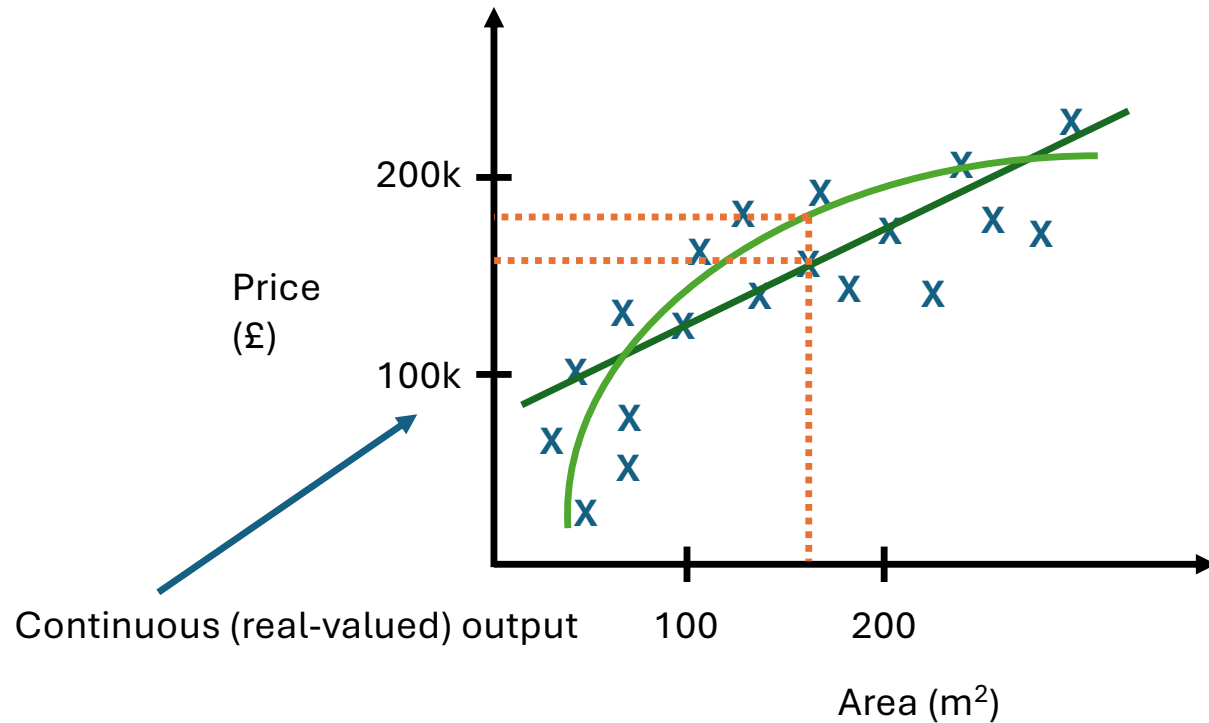**X = Malignant = 1**

**O = Benign = -1**

(We could also have more than two output classes – this would be called *multi-class classification*.)
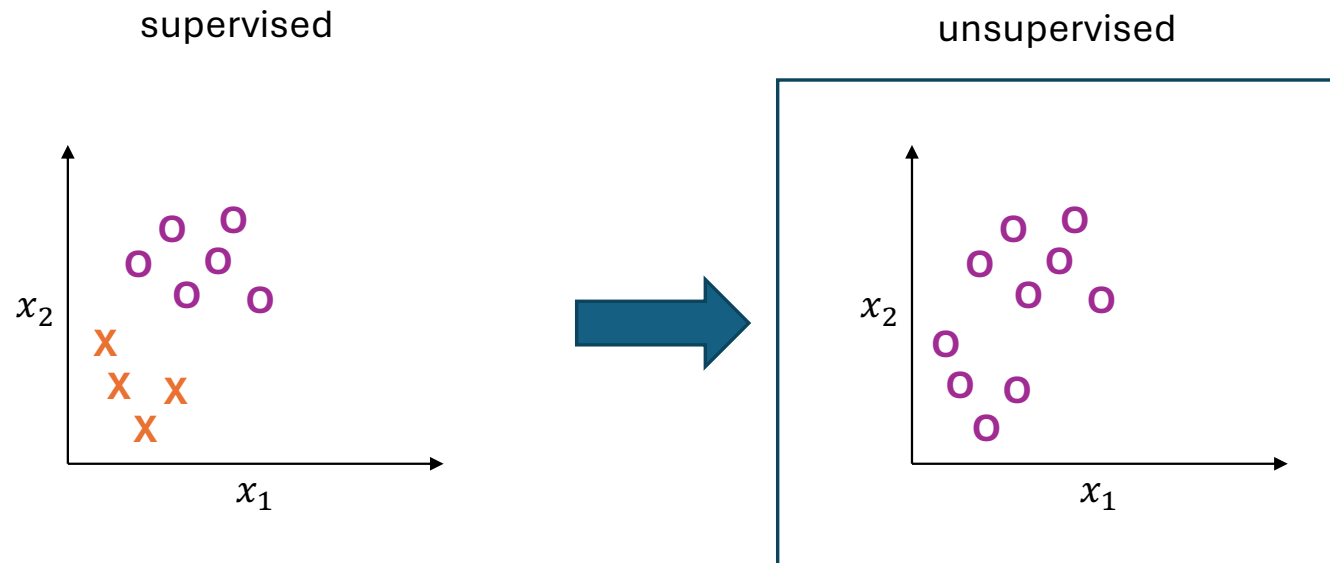


Age

Tumor Size

Often, we have more than two input features. Here, that additionally could be tumor clump thickness, uniformity of cell size, uniformity of cell shape, etc.
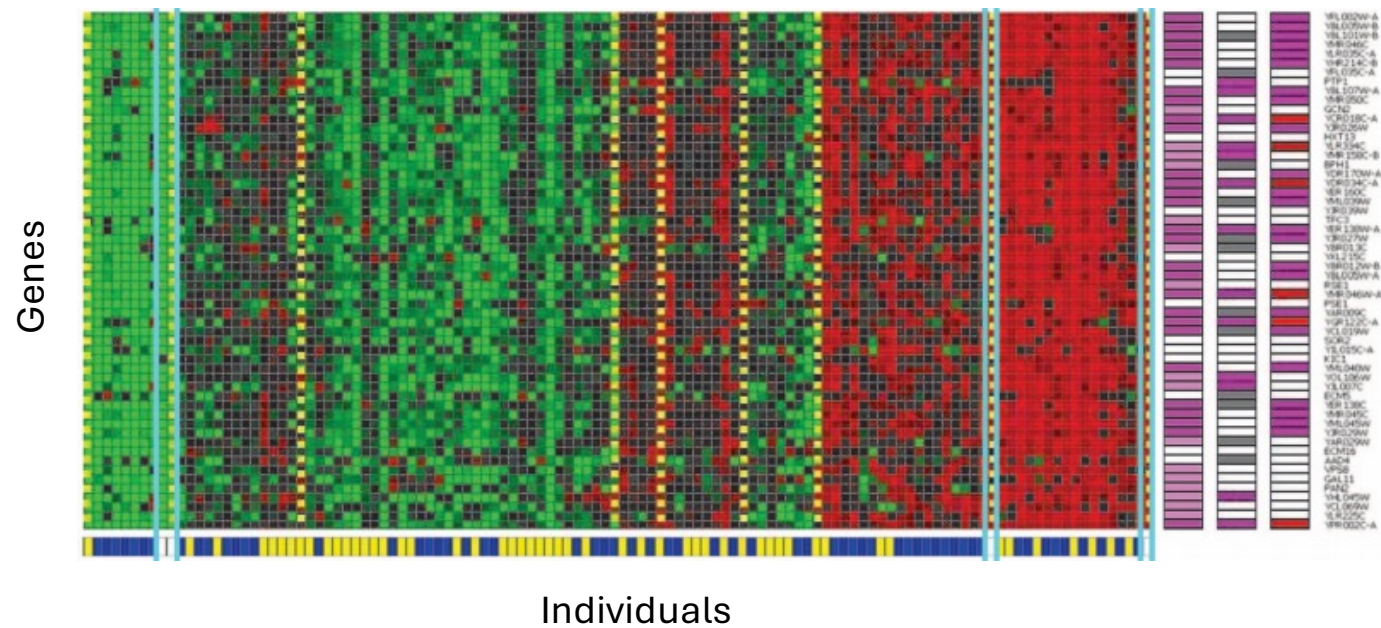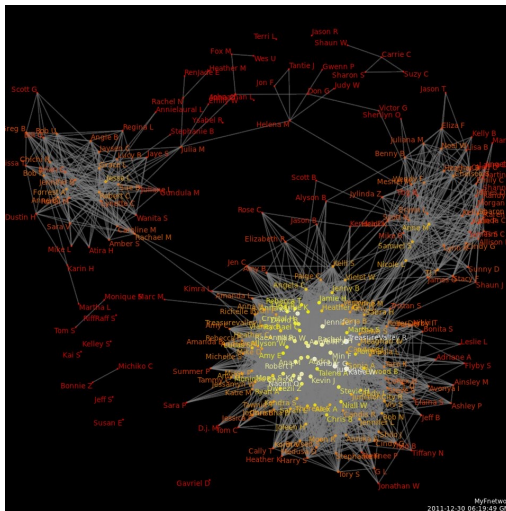
# Supervised learning - Regression

# Unsupervised learning

supervised

unsupervised

# Unsupervised learning



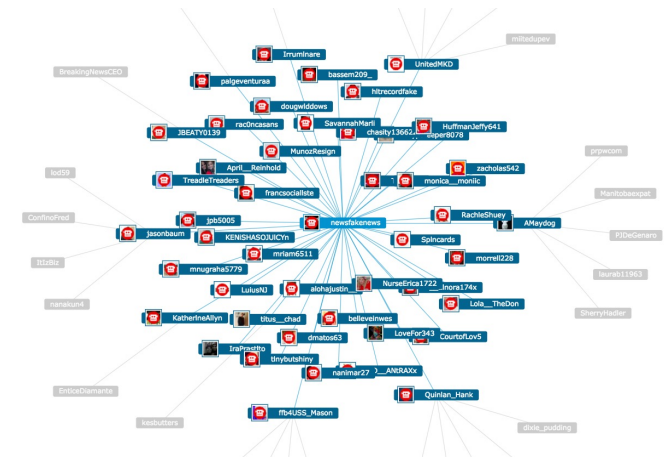Genes (vertical axis label)

Individuals (horizontal axis label)

# Unsupervised learning


Social network analysis


Market / customer segmentation


Identifying fake news

Sources:
https://en.wikipedia.org/wiki/Social_network_analysis#/media/File:Kencf0618FacebookNetwork.jpg
https://towardsdatascience.com/clustering-algorithms-for-customer-segmentation
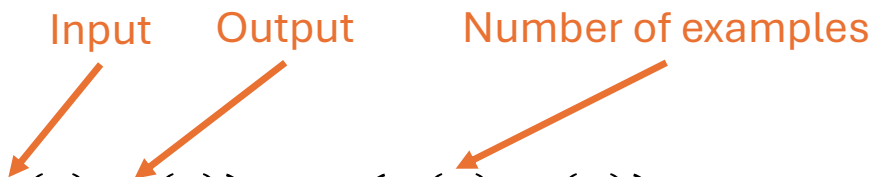https://medium.com/hackernoon/the-fake-news-arms-race-448675592803

# Machine learning – A magic box?

- Data
- Space of possible solutions
- Characterise objective
- Find algorithm
- Run
- Validate result

# Supervised Learning

# Data

- Dataset: $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

- $x^{(i)} \in \mathbb{R}^d$         $y^{(i)} \in \{+1, -1\}$     Binary Classification

- $\varphi(x)$: feature representation $\in \mathbb{R}^d$

13

# Hypotheses

- A hypothesis: $y = h(x; \theta)$

- $h \in \mathcal{H}$ (hypothesis class)

# Loss function

- $L(g, a)$   Guess    Actual

  $g \in \{+1, -1\}$

  $a \in \{+1, -1\}$

- How bad was it that we predicted $g$ when $a$ is the true answer

# Evaluating hypotheses

- Ideally: Small loss on **new** data

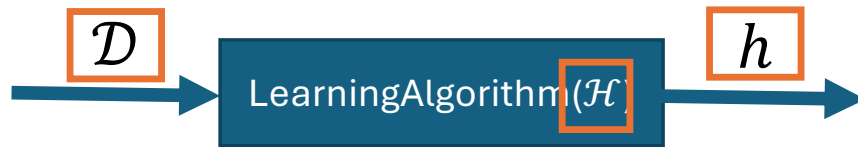$$\mathcal{E}(h) = \frac{1}{n'}\sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$$

test error

- What we can do (for now): Small loss on **training** data

$$\mathcal{E}_n(h) = \frac{1}{n}\sum_{i=1}^{n} L(\underbrace{h(x^{(i)})}_{g}, \underbrace{y^{(i)}}_{a})$$

training error

# Learning algorithms



$$\mathcal{D} \rightarrow \boxed{\text{LearningAlgorithm}(\mathcal{H})} \rightarrow h$$

- How to come up with learning algorithms:
  - Be a clever (or not so clever) human
  - Use optimisation methods

# Linear Classifiers

# Linear Classifiers

- Linear classifiers: A choice of $\mathcal{H}$

$$\overset{\mathbb{R}^d \quad \mathbb{R}^d \quad \mathbb{R}}{h(x\ ;\theta,\theta_0)} = sign(\underbrace{\overset{dot}{\theta^T x} + \theta_0}_{\mathbb{R}}) = \begin{cases} +1 \text{ if } \theta^T x + \theta_0 > 0 \\ -1 \text{ otherwise} \end{cases}$$
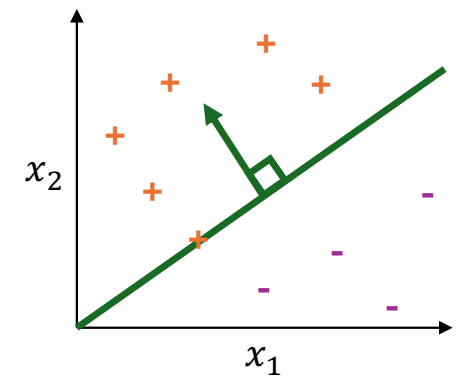
$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \qquad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$\theta : d \times 1$
$x : d \times 1$

$(1 \times d).(d \times 1)$

$\theta_1.x_1 + \theta_2.x_2 + \theta_0 = 0$  Implicit representation

$y = ax + b$  Parametric representation

slope



20

# The random linear classifier algorithm

```
random_linear_classifier(D, k):
for j=1 to k
```

$\quad \theta^{(j)}$ `= random(`$\mathbb{R}^d$`);` $\theta_0^{(j)}$ `= random(`$\mathbb{R}$`)`

`j* = ` $\displaystyle\operatorname*{argmin}_{j \in \{1..k\}} \mathcal{E}_n(\theta^{(j)}, \theta_0^{(j)})$

`return(` $\theta^{(j^*)}, \theta_0^{(j^*)}$ `)`