



University
of Exeter

Predictive Modelling

Week 05-BEM2031

Term2: 2024/25

Today:

- What is random?
- What is predictive modelling?
- Supervised learning
 - Regression
 - Classification
- A decision tree step-by-step

Types of Analytics:

- Descriptive Analytics: WHAT happened (or is happening)?
- Diagnostic Analytics: WHY did it happen?
- Predictive Analytics: WHAT is likely to happen in the future?
- Prescriptive Analytics: WHAT can we do about it?

What do we mean by random?



[List of solar eclipses in the 21st century - Wikipedia](#)

Date	Time	Saros	Type	Magnitude	Duration	Location	Path width	Geography	
September 4, 2100	08:49:20	146	Total	−0.3384	1.0402	3:32 🌐 10.5°S 39.0°E	142	88	Congo, Rwanda, Uganda, Burundi, Tanzania, Mozambique, Madagascar Partial: Africa,



What do we mean by random?



Random in the context of prediction models means that some aspect of the model is determined by chance rather than by design.

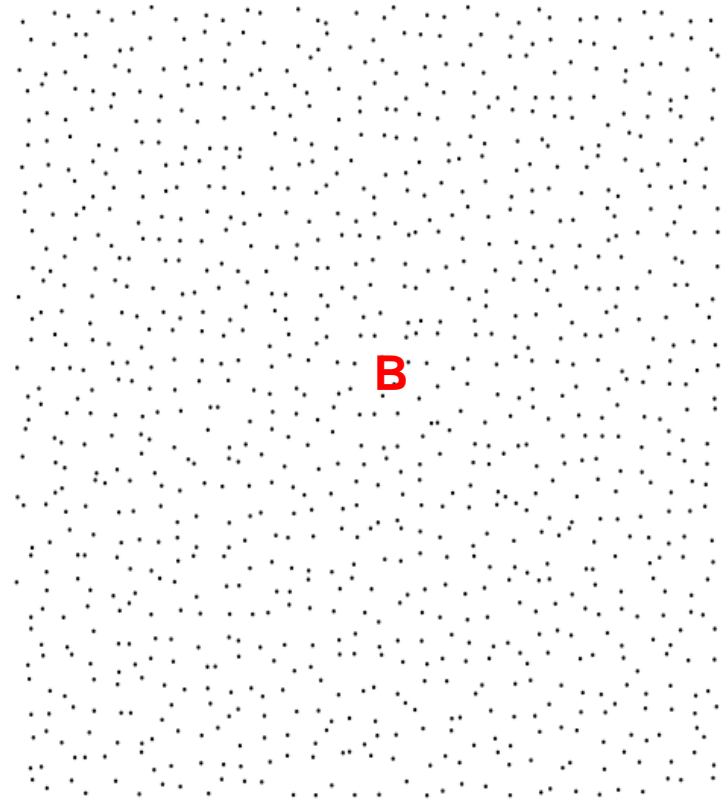
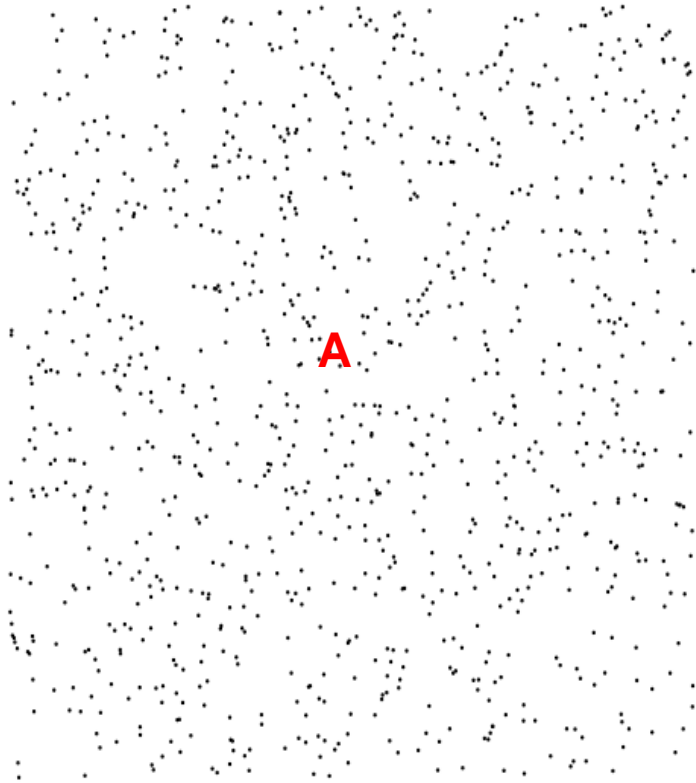


Randomness is also a source of uncertainty and error. Random phenomena are difficult to measure, predict, and control.

What do we mean by random?



University
of Exeter



Which of these could be modelled and which is completely random?

What do we mean by random?

MODULE MATERIALS

Lecture slides >

Workshop output files >

Workshop interactive code v

Week 1 practice

Week 2 practice

Week 3 practice

Week 4 practice

Week 5 practice

Random and regular fields of

The random field of points from Lecture 5. `d1` distribution:

R Code

Start Over

Run Code

```
1 N <- 1024
2 d1 <- tibble(x = runif(N), y = runif(N), type = 't1')
```

- Create these random fields of points
- Random coin throw, random dice

To create the evenly distributed bit (`d2`), you have to start with an even grid of points:

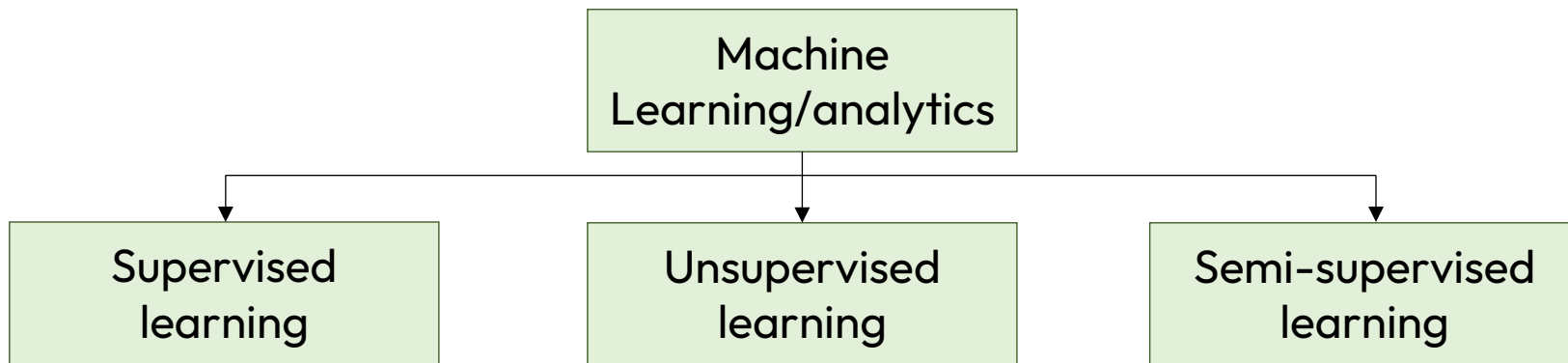
R Code

Start Over

Run Code

```
1 p <- seq(0,1, length.out = N / 32)
2
3 d2 <- expand.grid(p, p) %>%
4   data.frame %>%
5   rename(x = Var1, y = Var2) %>%
6   mutate(type = 't2')
```


Supervised vs Unsupervised methods



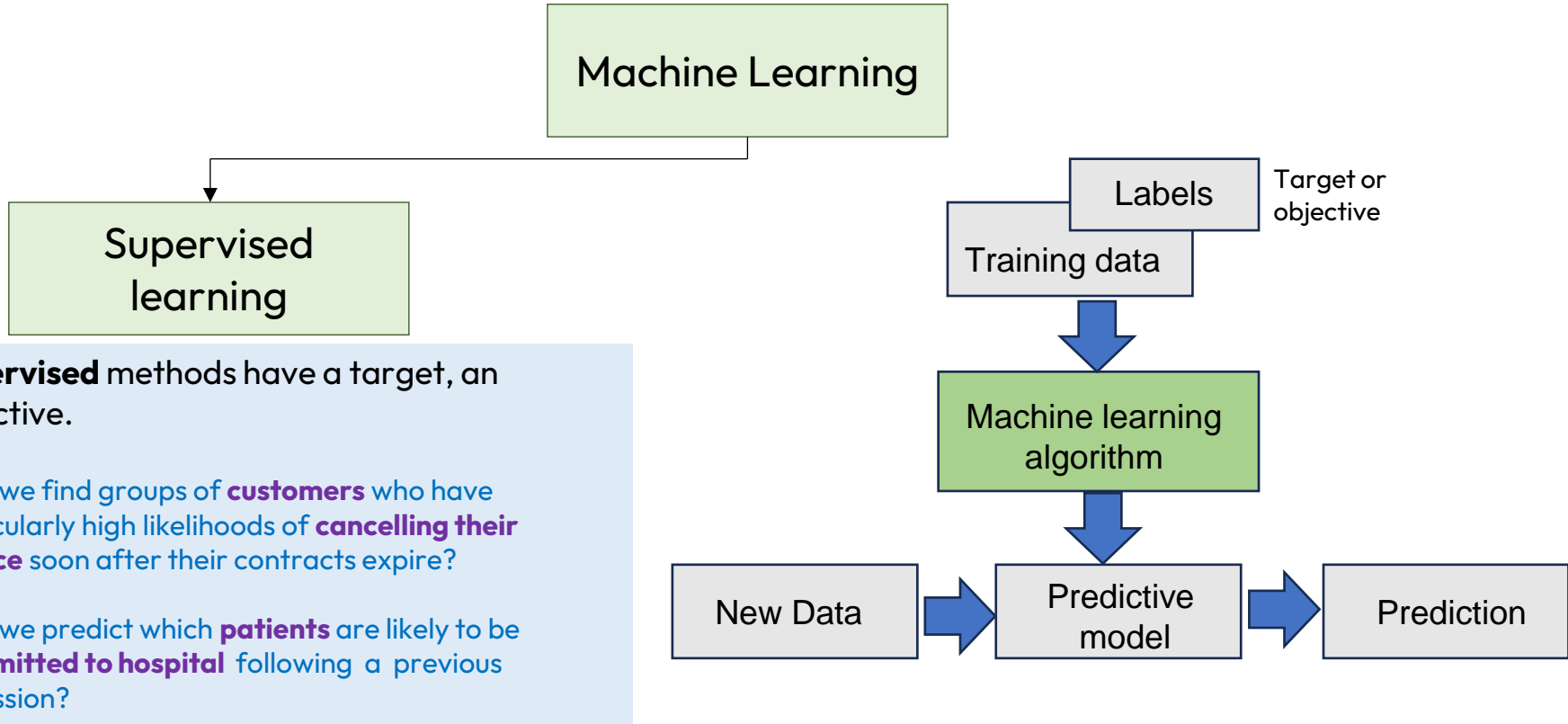
Supervised methods have a target, an objective.

“Can we find groups of customers who have particularly high likelihoods of cancelling their service soon after their contracts expire?”

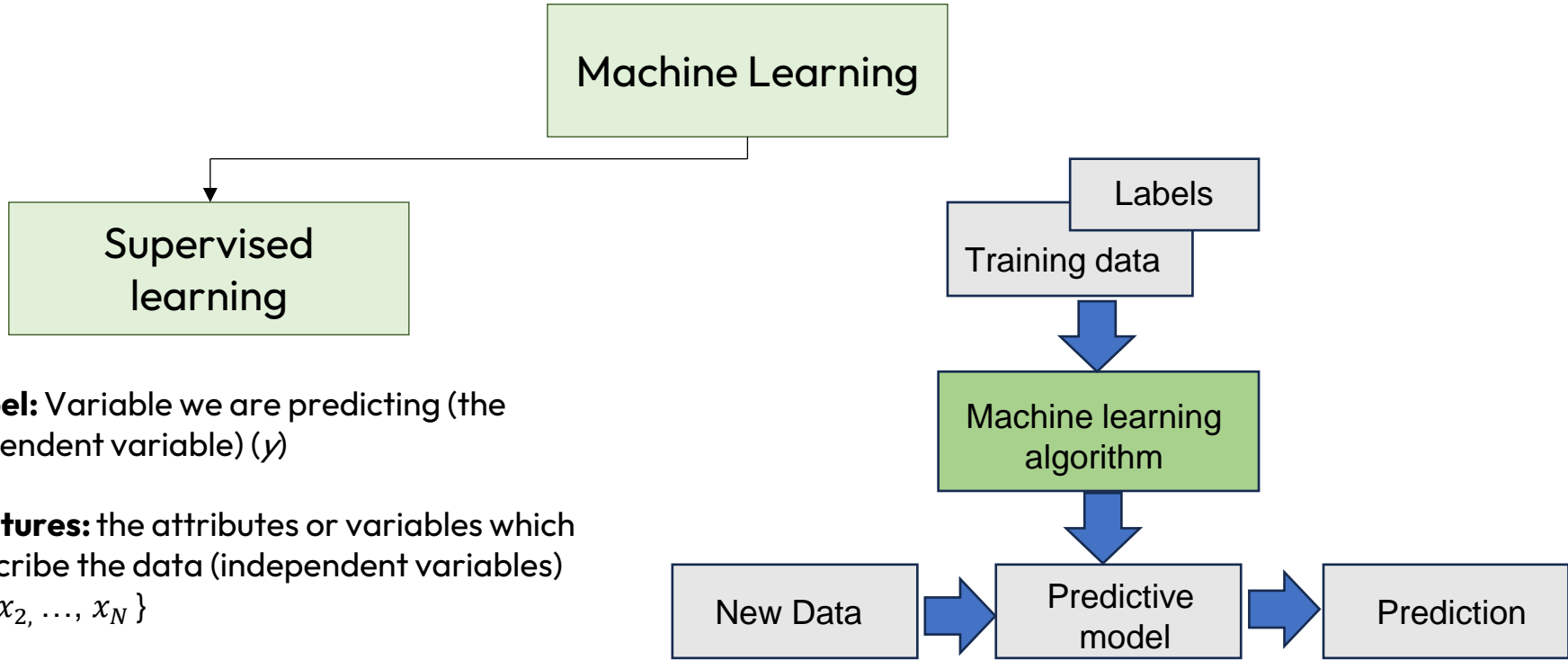
Unsupervised methods have no specific target.

“Do our customers naturally fall into different groups?”

Supervised vs Unsupervised methods



Supervised vs Unsupervised methods

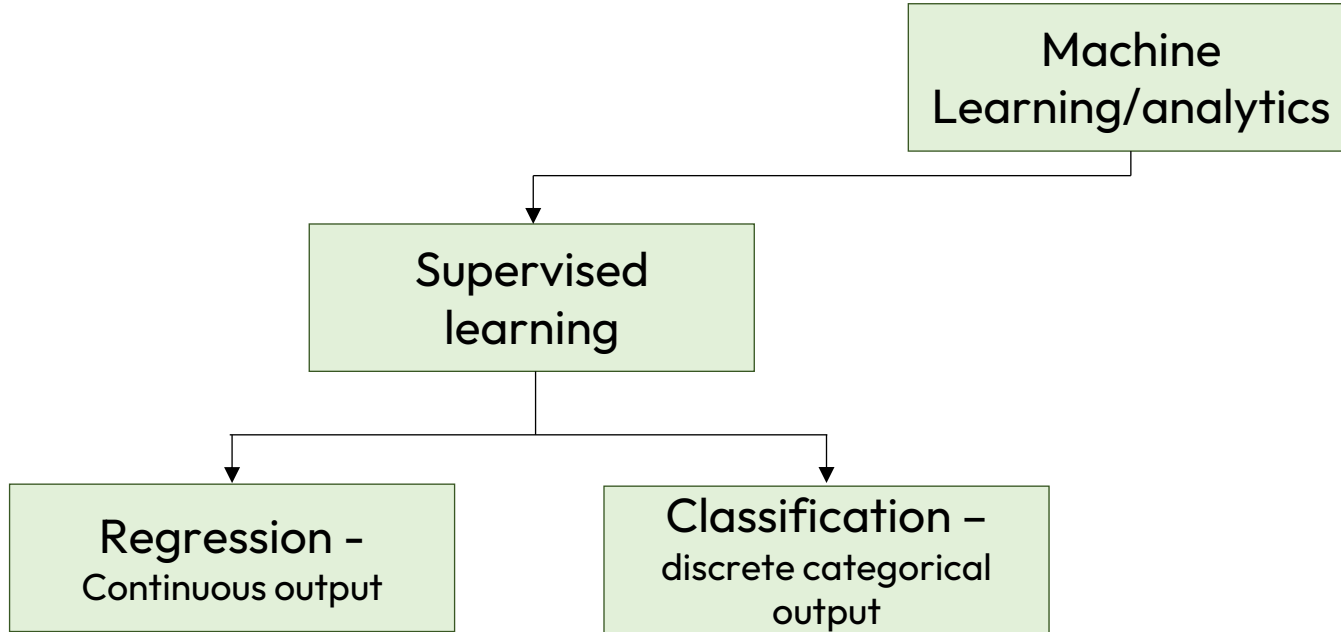


Label: Variable we are predicting (the dependent variable) (y)

Features: the attributes or variables which describe the data (independent variables) $\{x_1, x_2, \dots, x_N\}$

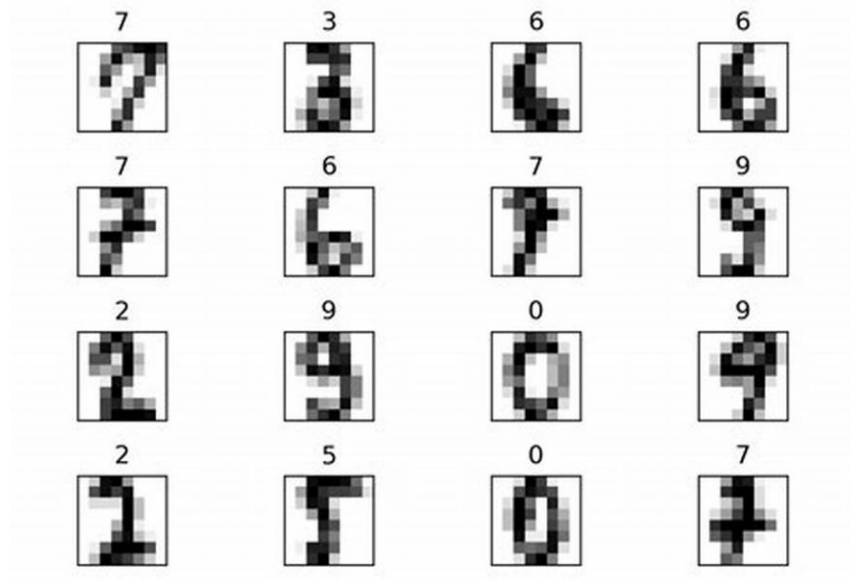
Unseen (new) data: test the performance of the model $y=f(x)$ on unlabelled data

Types of supervised learning

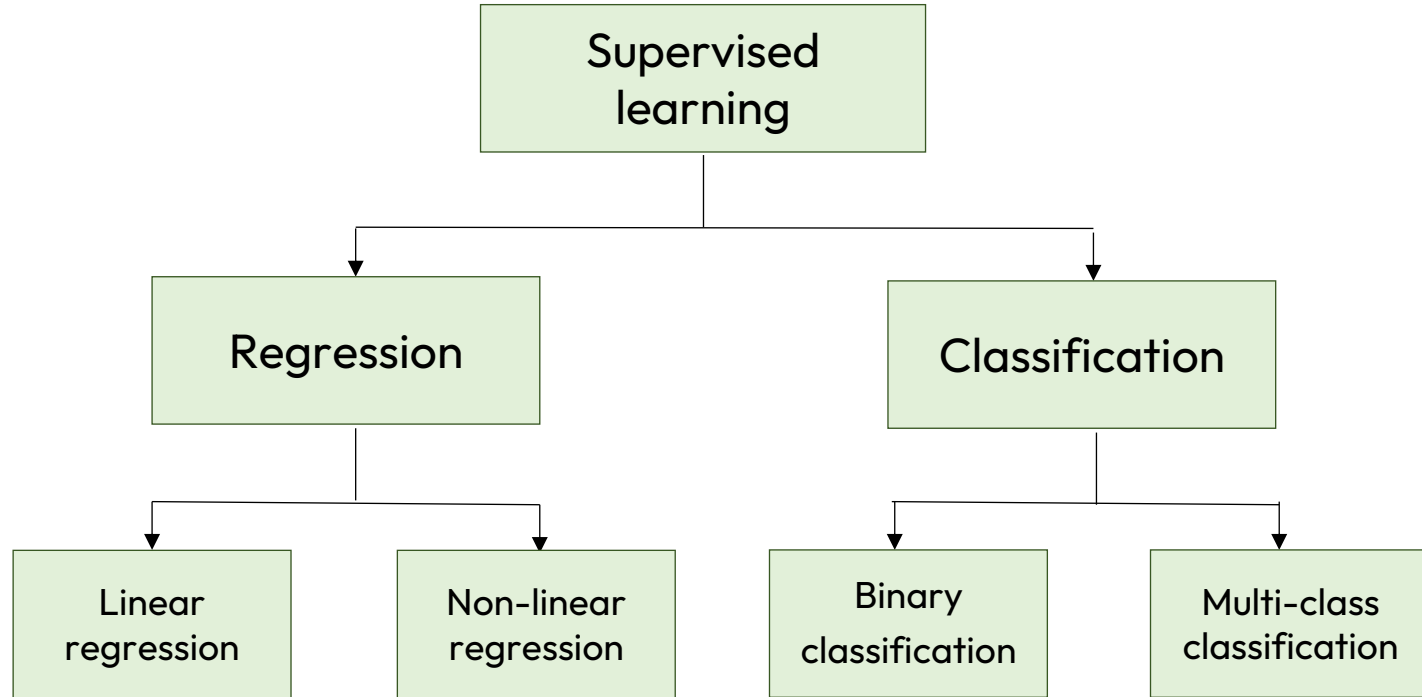


Regression or classification?

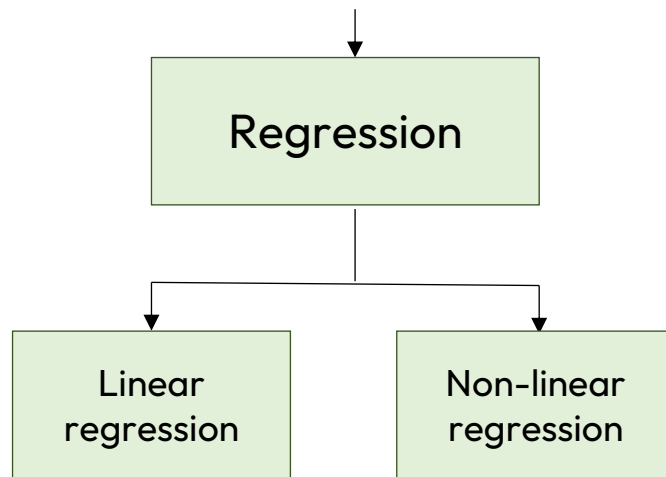
- Weather prediction
- Identification of cancer cells
- Identification of handwritten digits
- Oil price prediction
- Identification of fraudulent credit card transactions
- Monthly income prediction



Machine Learning



Regression modelling:



We'll talk
about this
now



We'll
cover this
in the
workshop

Linear Regression: House Price Prediction


Given a set of input features (which may influence the price of a house), the goal of the algorithm is to predict the price of a new house going to market

House Price Prediction

Square footage	Num of Rooms	Garden?	Parking Facility?	Num of floors	Price (\$) in 1000's
..	3	Yes	..	2	460
1700	..	No	No	1	320
..	5
..
..
..

 Labeled Example

Attributes/Features/Independent Variables (x_1, x_2, \dots, x_n)

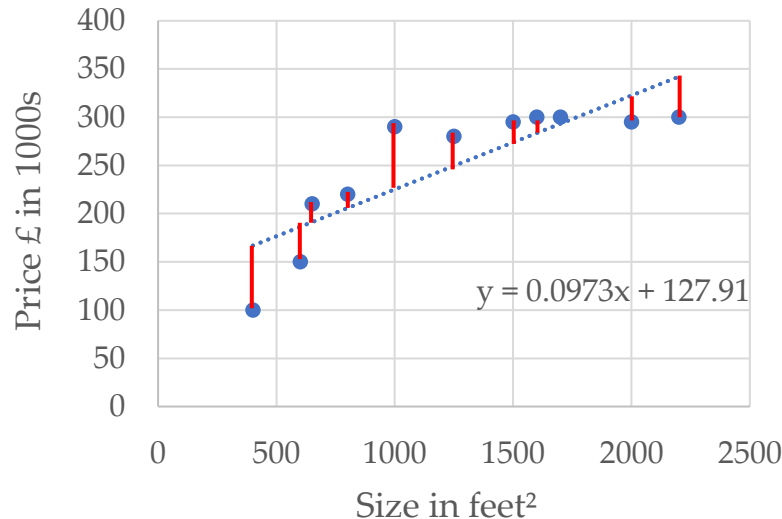


Target/Dependent Variable (y)

Simple Linear Regression

Consider the problem of **predicting house prices (y)**

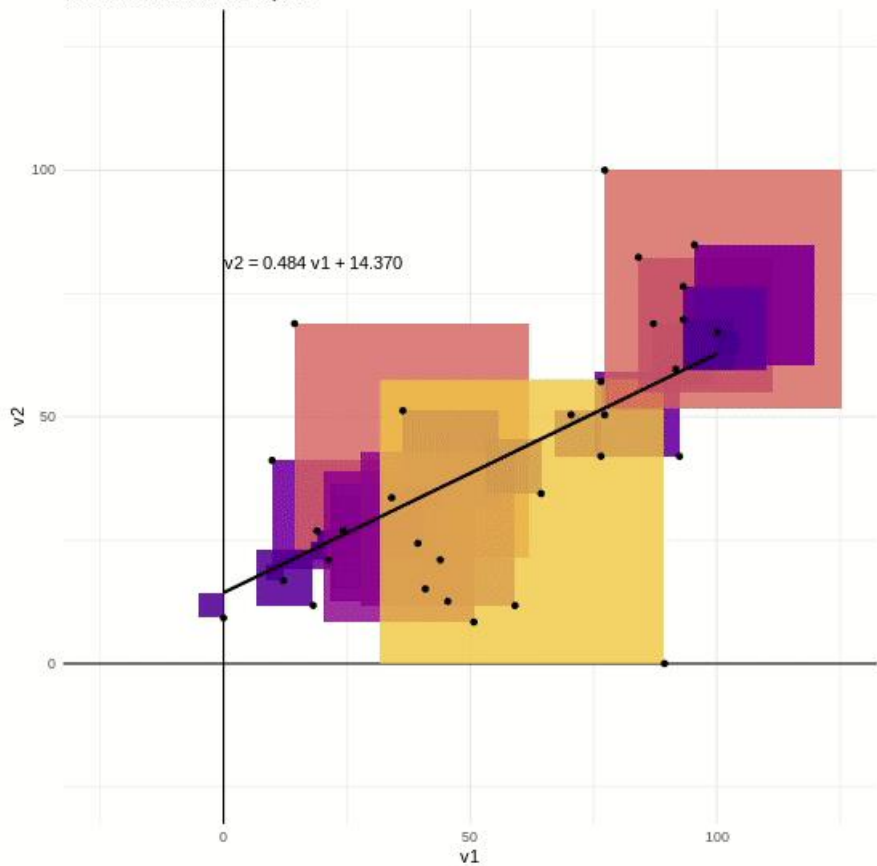
- **Feature Selection** – input variables that can be used to predict house prices
 - let's consider one input variable (size in sq.ft) → **univariate/simple regression**
- Simple linear regression finds a linear function (straight line) that predicts the target variable (y) as a function of the features or independent variables (x)
- $y = mx + c$ where m is the slope and c is the intercept



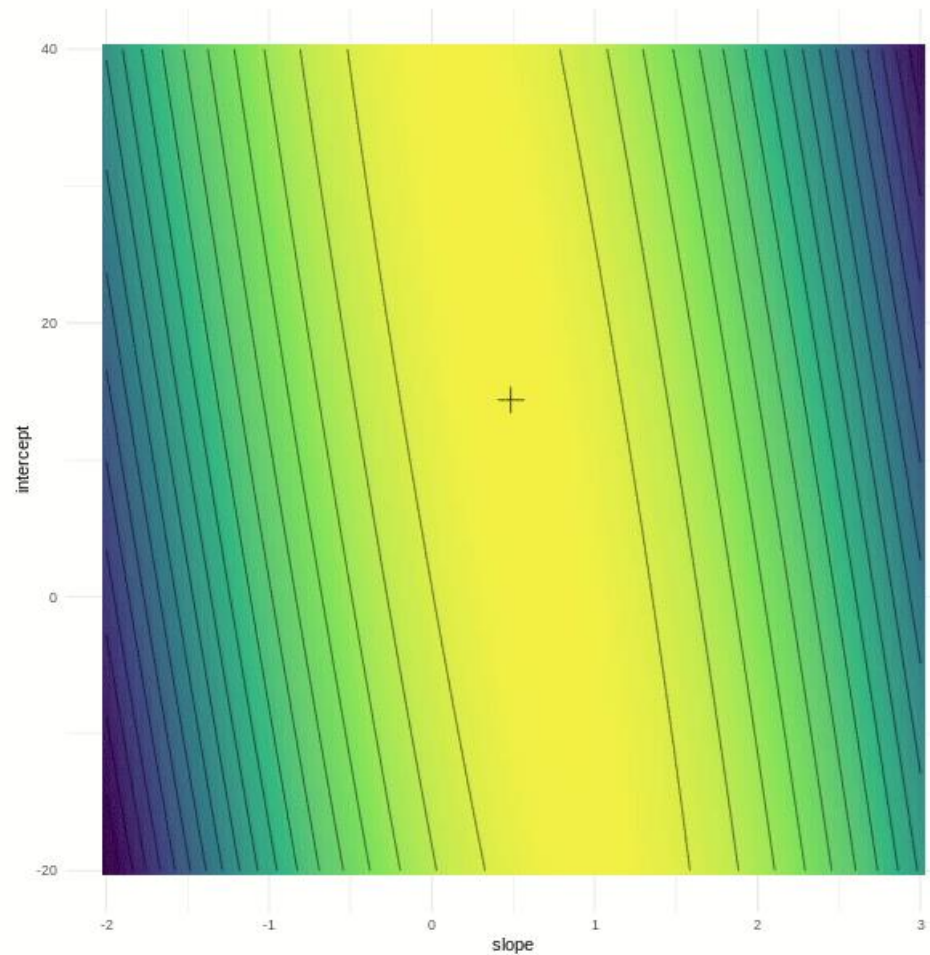
Features/independent variables (x)	Target/dependent variable (y)
Size in feet²	Price £ in 1000s
400	100
600	150
650	210
800	220
1000	290
1250	280
1500	295
1600	300
1700	300
2000	295
2200	300

Estimating OLS Regression

Nelder-Mead minimizer step: 75



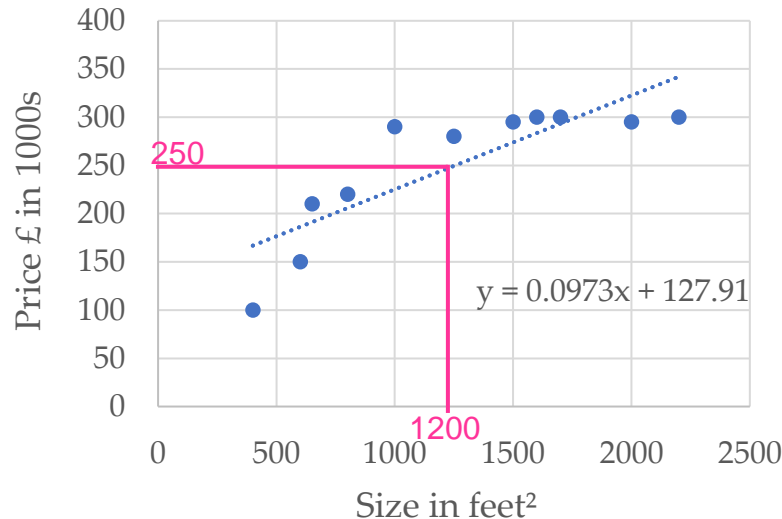
Data source: <https://xkcd.com/2048/>
Created by @jessemfagan



Simple Linear Regression

Consider the problem of **predicting house prices (y)**

- **Feature Selection** – input variables that can be used to predict house prices
 - let's consider one input variable (size in sq.ft) → **univariate/simple regression**
- Simple linear regression finds a linear function (straight line) that predicts the target variable (y) as a function of the features or independent variables (x)
- $y = mx + c$ where m is the slope and c is the intercept



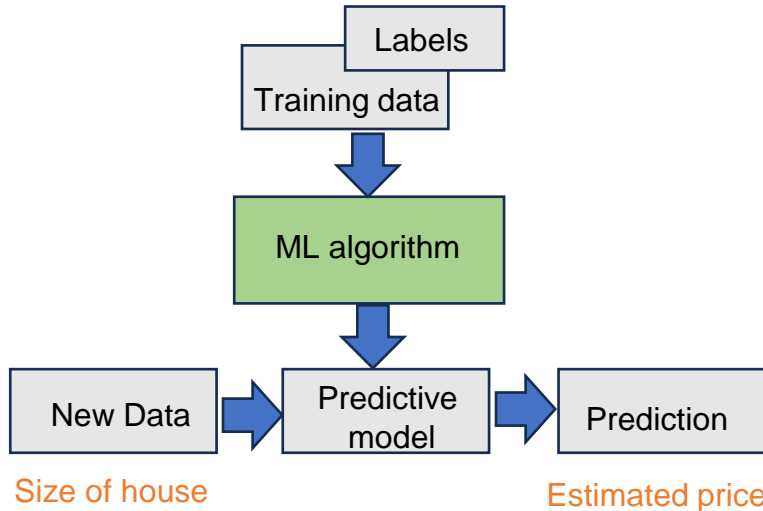
Features/independent variables (x) Target/dependent variable (y)

Size in feet²	Price £ in 1000s
400	100
600	150
650	210
800	220
1000	290
1250	280
1500	295
1600	300
1700	300
2000	295
2200	300

Simple Linear Regression

Consider the problem of **predicting house prices (y)**

- **Feature Selection** – input variables that can be used to predict house prices
 - let's consider **one input variable** (size in sq.ft) → **univariate/simple regression**
- Simple linear regression finds a linear function (straight line) that predicts the target variable (y) as a function of the features or independent variables (x)
- $y = mx + c$ where m is the slope and c is the intercept

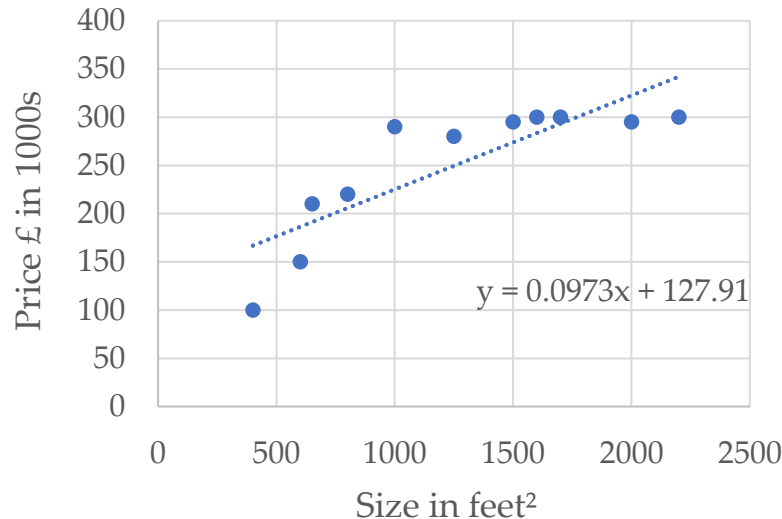


Features/independent variables (y)	Target/dependent variable (x)
Size in feet ²	Price £ in 1000s
400	100
600	150
650	210
800	220
1000	290
1250	280
1500	295
1600	300
1700	300
2000	295
2200	300

Simple Linear Regression

Consider the problem of **predicting house prices (y)**

- **Feature Selection** – input variables that can be used to predict house prices
 - let's consider **one input variable** (size in sq.ft) → **univariate/simple regression**
- Simple linear regression finds a linear function (straight line) that predicts the target variable (y) as a function of the features or independent variables (x)
- $y = mx + c$ where m is the slope and c is the intercept



Features/independent variables (x)	Target/dependent variable (y)
Size in feet²	Price £ in 1000s
400	100
600	150
650	210
800	220
1000	290
1250	280
1500	295
1600	300
1700	300
2000	295
2200	300

Multiple Linear Regression



Consider the problem of **predicting house prices (y)**

- **Feature Selection** – input variables that can be used to predict house prices
 - let's consider **multiple input variable** → **multiple linear regression**
- Multiple linear regression models a linear function that predicts a target variable as a function of the independent variables: $y = a_0 + a_1x_1 + a_2x_2 + a_3x_3 + \dots$

Square footage	Num of Rooms	Garden?	Parking Facility?	Num of floors	Price (\$) in 1000's
..	3	Yes	..	2	460
1700	..	No	No	1	320
..	5
..
..
..

m training
examples

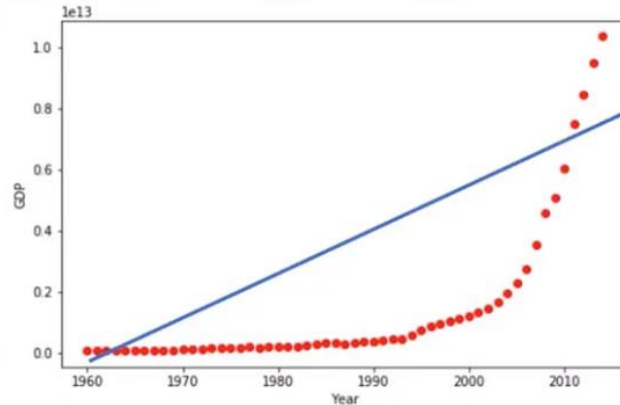
Attributes/Features/Independent Variables (x_1, x_2, \dots, x_n)

Target/Dependent Variable (y)

Problems faced: Underfitting

Should we use linear regression?

	Year	Value
0	1960	5.918412e+10
1	1961	4.955705e+10
2	1962	4.668518e+10
3	1963	5.009730e+10
4	1964	5.906225e+10
5	1965	6.970915e+10
6	1966	7.587943e+10
7	1967	7.205703e+10
8	1968	6.999350e+10
9	1969	7.871882e+10
...



These data points correspond to China's gross domestic product (GDP) from 1960–2014.

Model is not complex enough to capture the underlying patterns in the data.

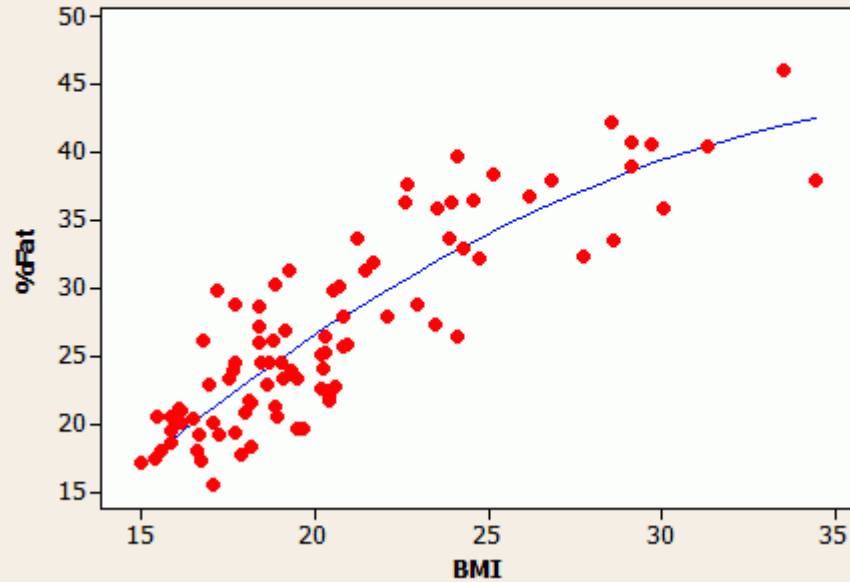
Leads to **bias**:

The amount of error introduced by approximating real-world phenomena in a simplified model

Non-linear regression

Using BMI to Predict Fat Percentage

$$\% \text{Fat} = -23.19 + 3.286 \text{ BMI} - 0.03999 \text{ BMI}^2$$



S	3.53399
R-Sq	76.1%
R-Sq(adj)	75.5%

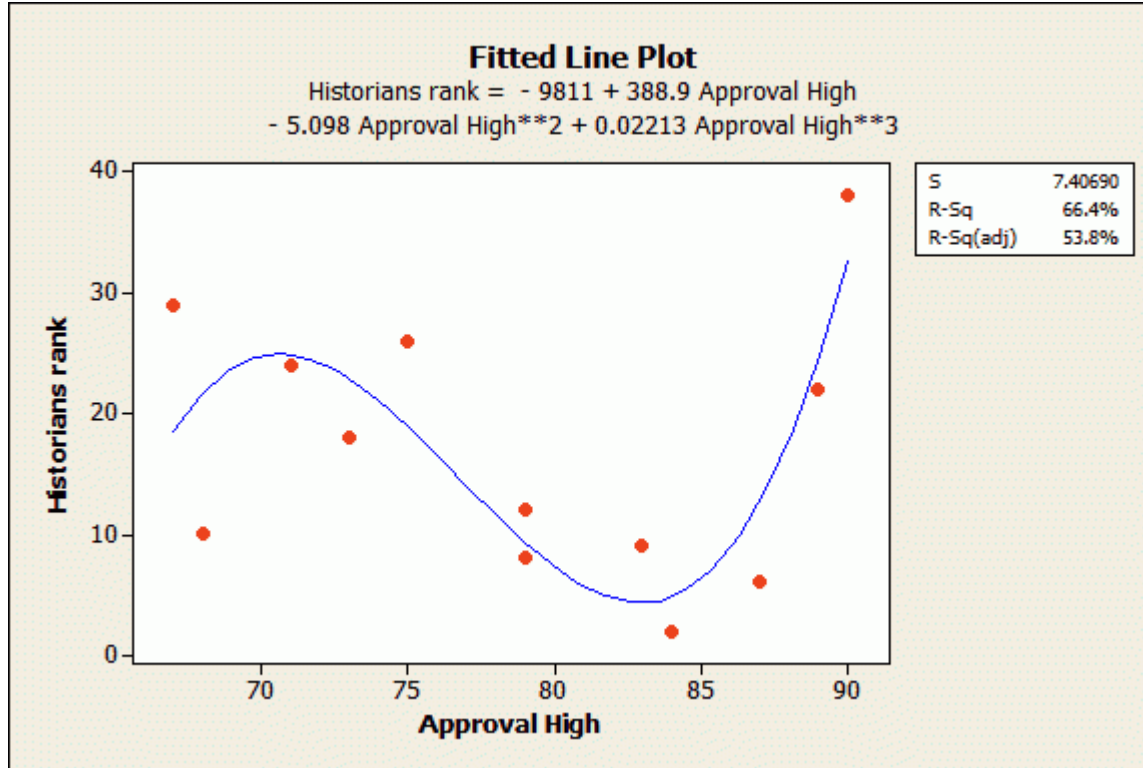
If a regression equation doesn't follow the rules for a linear model, then it must be a nonlinear model.

The regression example models the relationship between body mass index (BMI) and body fat percent.

It is a linear model that uses a quadratic (squared) term to model the curved relationship.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2$$

Problems faced: Overfitting



If you try to estimate too many parameters, you will overfit!

The size of your dataset restricts the number of terms you can safely add to your model

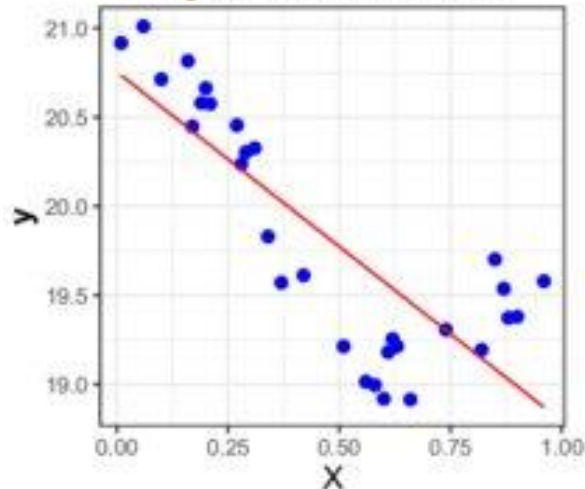
If your study calls for a complex model, you must collect a relatively large sample size.

Problems faced: Overfitting

Polynomial fit degree 1

Training error: 0.4

Generalization error: 0.42

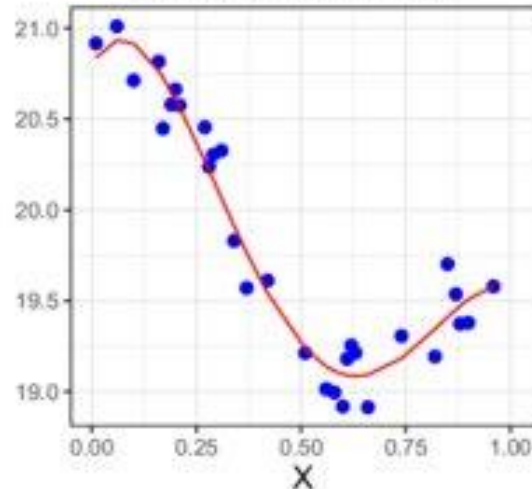


Underfit

Polynomial fit degree 4

Training error: 0.14

Generalization error: 0.17

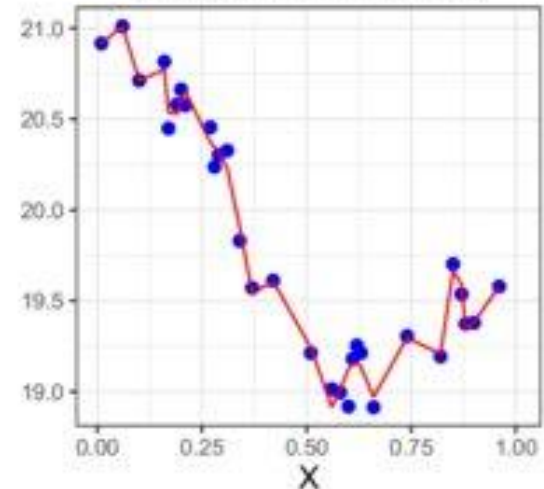


Good fit

Polynomial fit degree 20

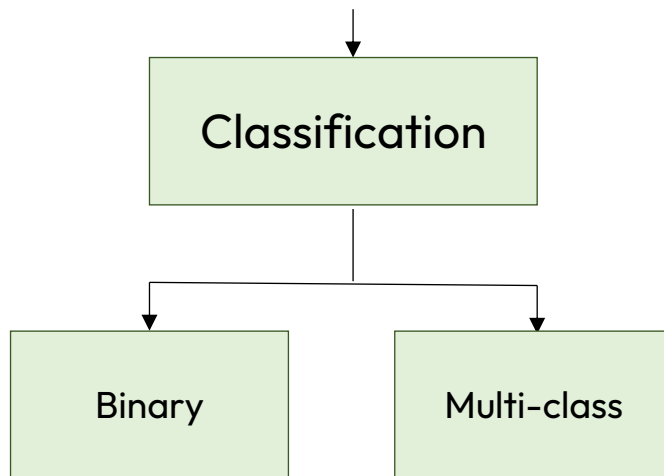
Training error: 0.07

Generalization error: 2000



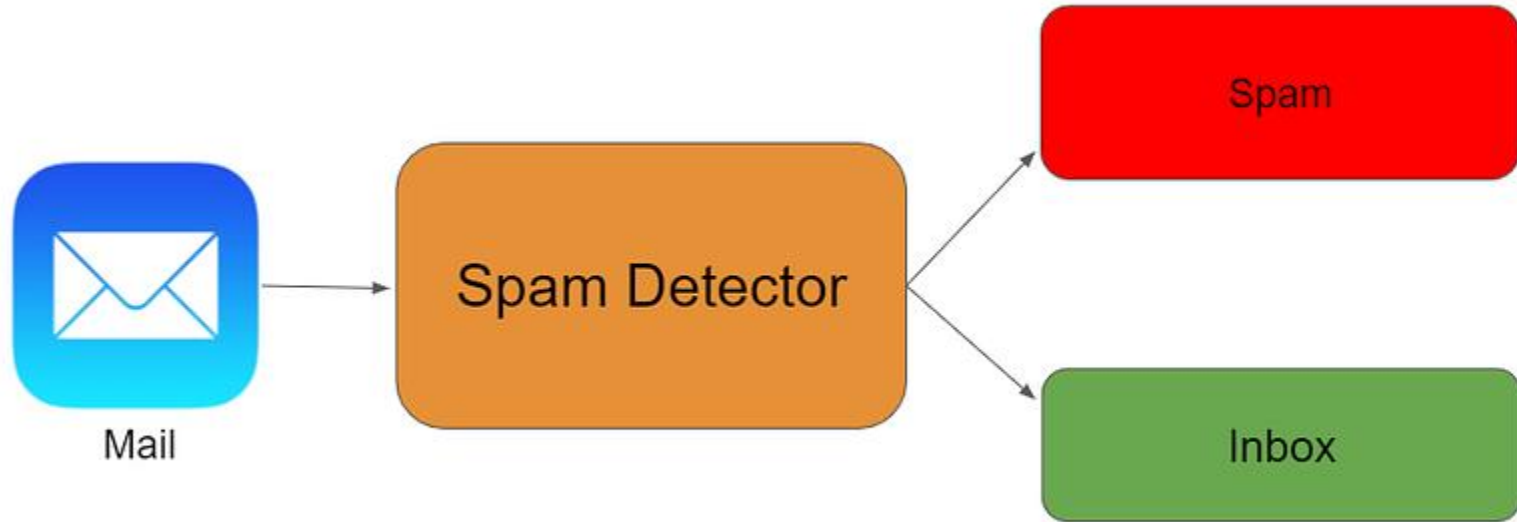
Overfit

Classification modelling:

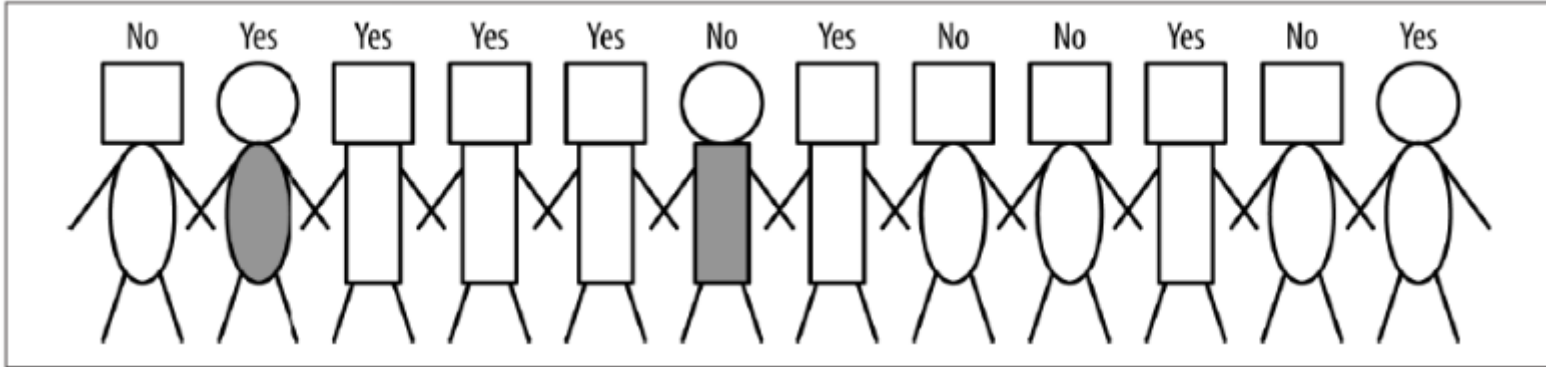


Binary Classification: Email spam prediction

Trained on a large number of spam and non-spam emails, the algorithm's goal is to predict whether or not an email is spam



Classification



Attributes:

Head shape: square, circle

Body shape: rectangle, oval

Body colour: grey, white

Target:

Write-off: yes, no

Attributes rarely split a group perfectly

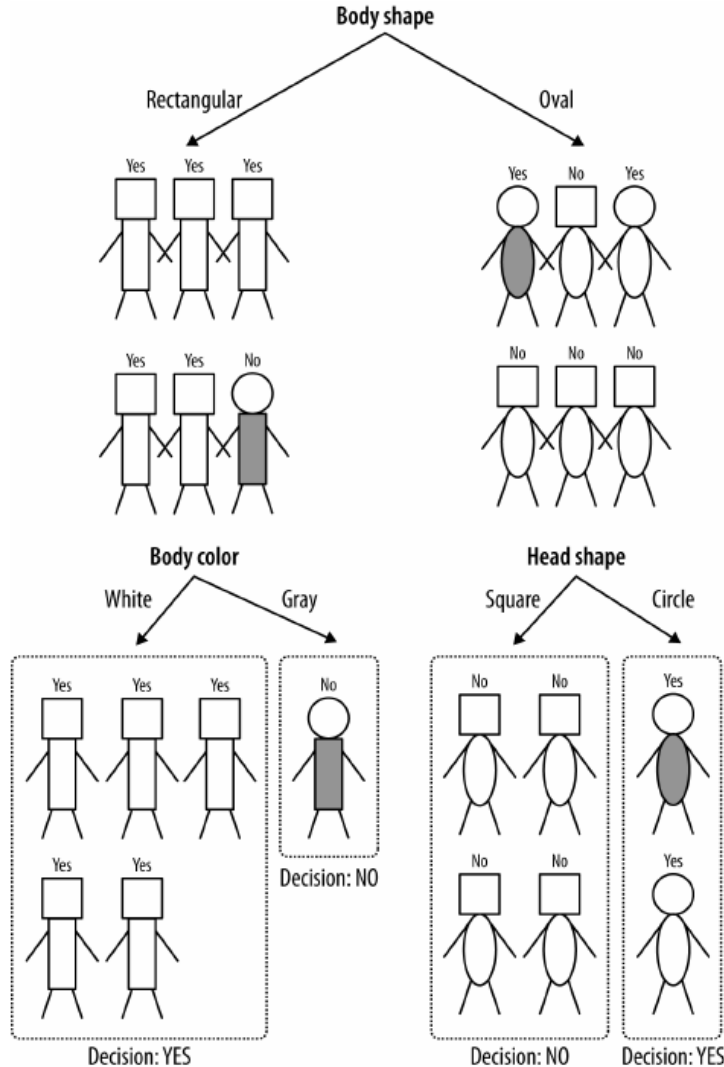
Not all attributes are binary

How do we segment for numeric values?

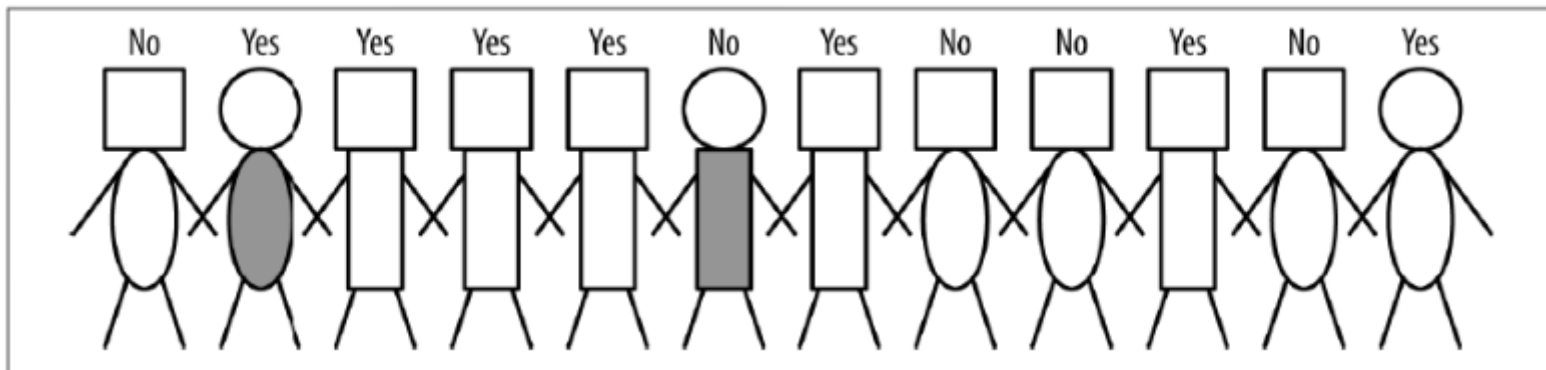
A decision tree is a model with a number of branching options that lead to a decision at the end.

Each point on the tree is called a **node**. The **depth** of the of the tree is maximum number of steps to reach a decision.

A **leaf** of the tree is where the decision is made (when there's no more splitting).



Classification – decision trees



Entropy:

The amount of uncertainty or randomness in a system

→ How *impure* a node (how mixed the training data assigned to that node is)

Information gain:

The reduction in entropy or uncertainty after a dataset is split based on a feature.

→ Helps the algorithm decide which feature to split on at each step. Features with the highest information gain are selected because they reduce uncertainty the most.

Features that result in a higher information gain are considered more important – as they provide more information

Classification - entropy



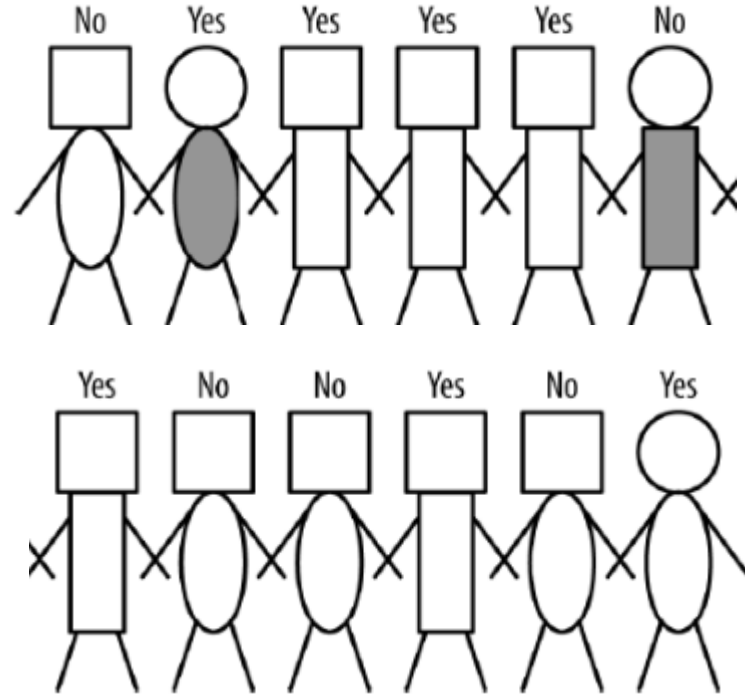
An **information gain** is how much an attribute improves (or decreases) **entropy** (uncertainty) of the model prediction.

$$\text{entropy} = -p_1 \log(p_1) - p_2 \log(p_2) - \dots$$

$$p_{\text{yes}} = 7/12$$

$$p_{\text{no}} = 5/12$$

$$\begin{aligned} \text{entropy}(S) &= -\left[\left(\frac{7}{12}\right) \times \log_2\left(\frac{7}{12}\right) + \left(\frac{5}{12}\right) \times \log_2\left(\frac{5}{12}\right)\right] \\ &= 0.98 \end{aligned}$$

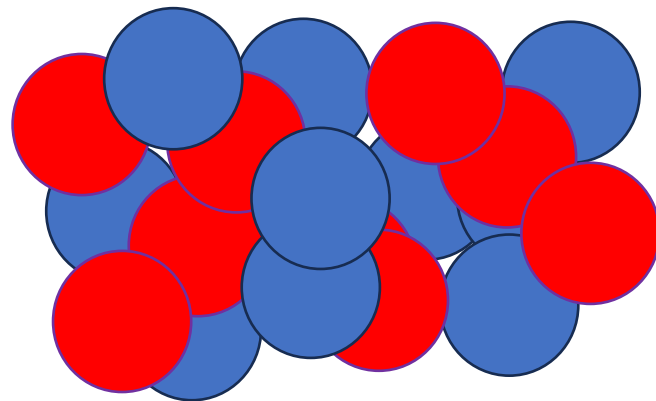


Classification - entropy

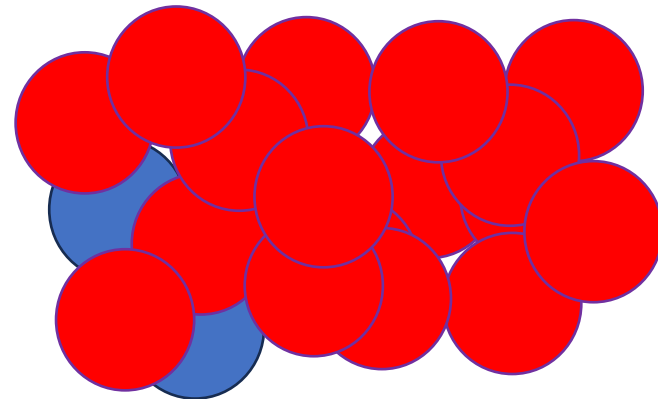


Pick a ball at random and guess the colour. The chances of being right or wrong depends on the mix of colours.

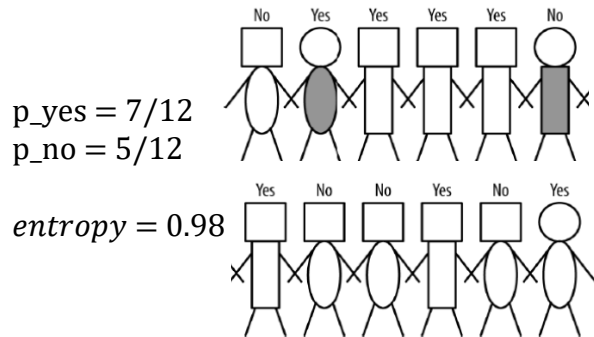
If your bag has an equal number of red and blue balls, your uncertainty is highest – this is a state of **high entropy**.



If the bag has mostly red balls, and not many blue balls, you'd probably guess red, and you would be right most of the time. This is a state of **low entropy**.



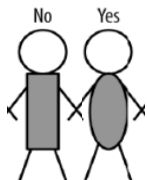
Entropy quantifies the uncertainty.



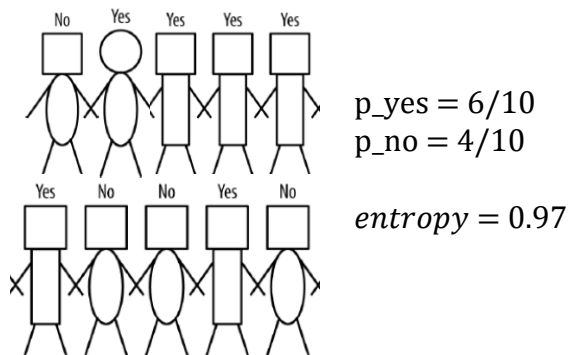
Body colour

Grey

$p_{\text{yes}} = 1/2$
 $p_{\text{no}} = 1/2$
 $\text{entropy} = 1.0$

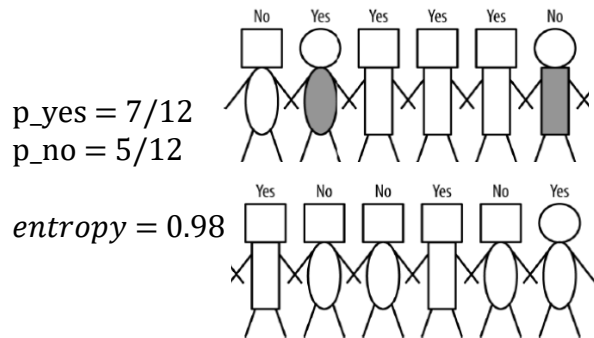


White



$$IG(\text{parent}, \text{children}) = \text{entropy}(\text{parent}) - p(c_1) \times \text{entropy}(c_1) + p(c_2) \times \text{entropy}(c_2) + \dots$$

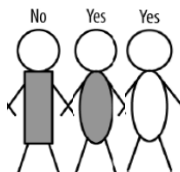
$$IG = 0.98 - (0.17 \times 1.0 + 0.83 \times 0.97) = 0.005$$



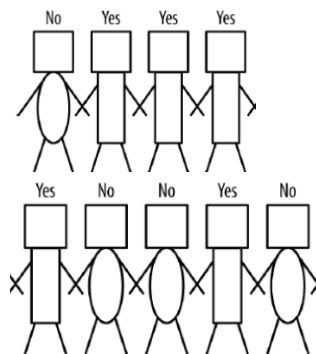
Head Shape

Circle

$p_{\text{yes}} = 2/3$
 $p_{\text{no}} = 1/3$
 $\text{entropy} = 0.92$



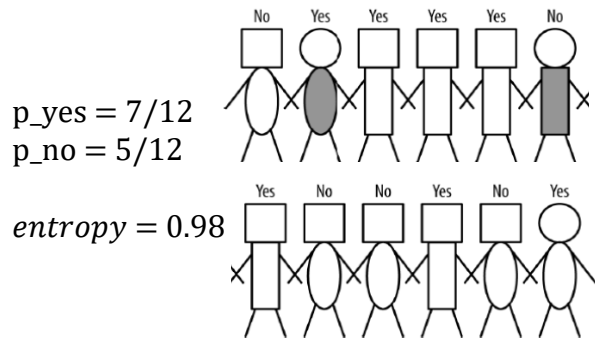
Square



$p_{\text{yes}} = 5/9$
 $p_{\text{no}} = 4/9$
 $\text{entropy} = 0.99$

$$IG = \text{entropy}(\text{base}) - p(c_1) \times \text{entropy}(c_1) + p(c_2) \times \text{entropy}(c_2) + \dots$$

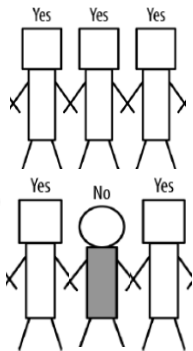
$$IG = 0.98 - (0.25 \times 0.92 + 0.75 \times 0.99) = 0.0075$$



Body Shape

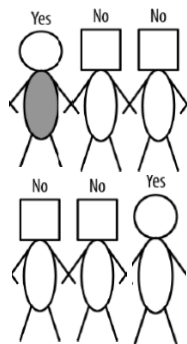
Rectangle

$p_{\text{yes}} = 5/6$
 $p_{\text{no}} = 1/6$
 $\text{entropy} = 0.650$



Oval

$p_{\text{yes}} = 2/6$
 $p_{\text{no}} = 4/6$
 $\text{entropy} = 0.918$



$$IG = \text{entropy}(\text{base}) - p(c_1) \times \text{entropy}(c_1) + p(c_2) \times \text{entropy}(c_2) + \dots$$

$$IG = 0.98 - (0.5 \times 0.650 + 0.5 \times 0.918) = 0.196$$

Body Shape(rectangle)

Body Shape(oval)

$$p_{\text{yes}} = 5/6$$

$$p_{\text{no}} = 1/6$$

$$p_{\text{yes}} = 2/6$$

$$p_{\text{no}} = 4/6$$

$$\text{entropy} = 0.650$$

$$\text{entropy} = 0.918$$

Body colour

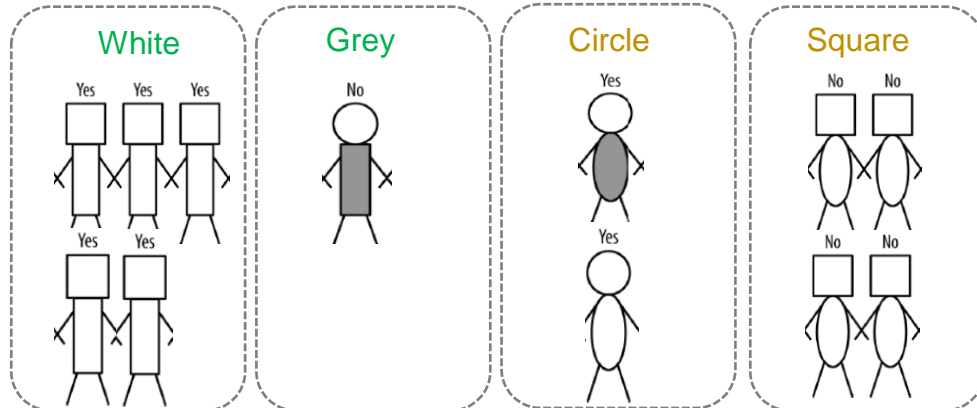
Head Shape

White

Grey

Circle

Square



Decision: YES

Decision: NO

Decision: YES

Decision: NO

$$IG = \text{entropy}(\text{base}) -$$

$$p(c_1) \times \text{entropy}(c_1) + p(c_2)$$

$$\times \text{entropy}(c_2) + \dots$$

$$IG = 0.650 - (0.17 \times 0 + 0.83 \times 0)$$

$$= 0.650$$

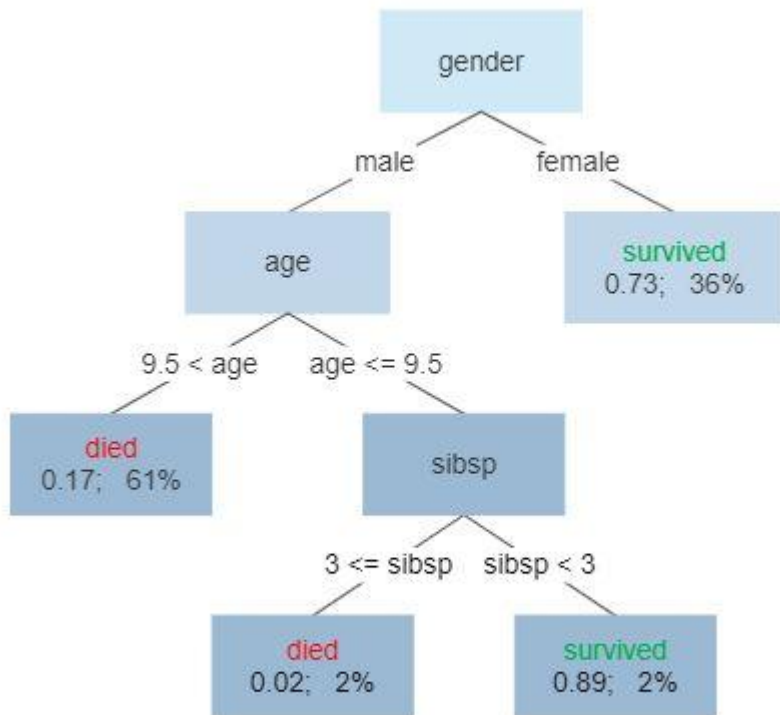
$$IG = 0.919 - (0.33 \times 0 + 0.67 \times 0)$$

$$= 0.918$$

Decision trees



Survival of passengers on the Titanic



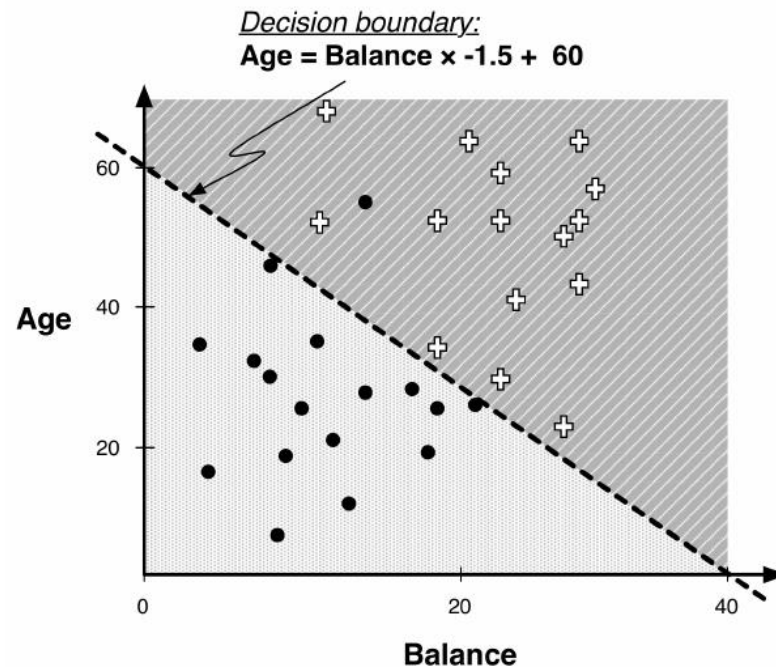
- Widely used for regression and classification problems
 - Root at top, leaves at bottom
-
- Titanic survival model predicts survival for:
 - females
 - males younger than 9.5 years with less than 2.5 siblings
 - The figures under the leaves show the probability of survival and the percentage of observations in the leaf.

Fitting classification models

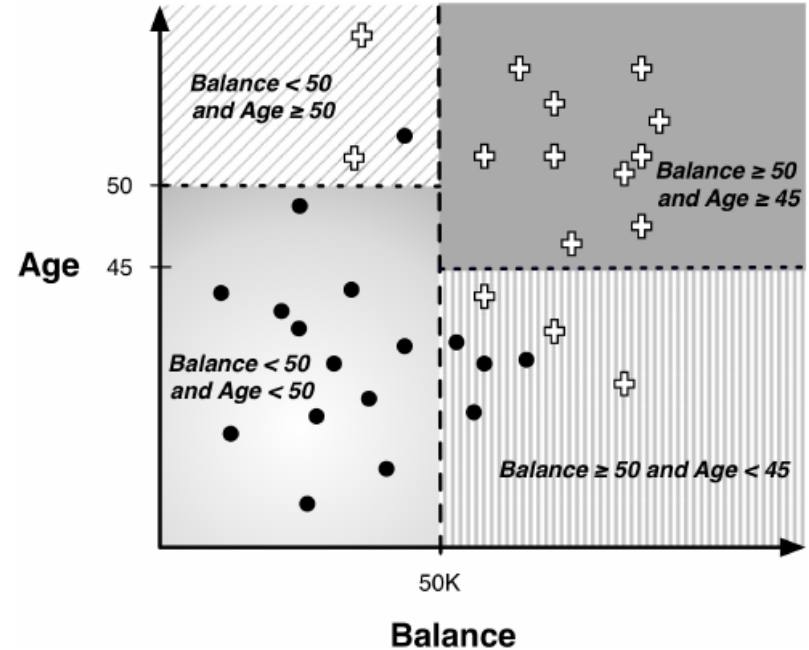
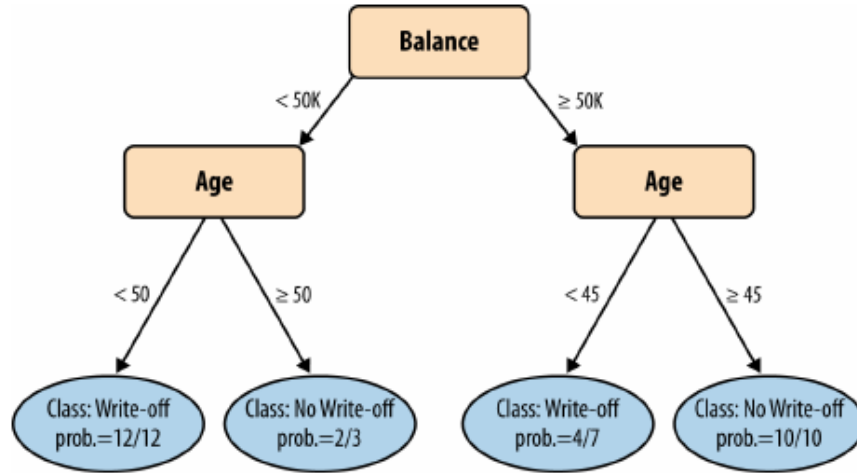


A poorly chosen decision boundary can lead to underfitting (oversimplifying) or overfitting (too closely fitting the data)

- **How do we know how best to draw the boundary?**



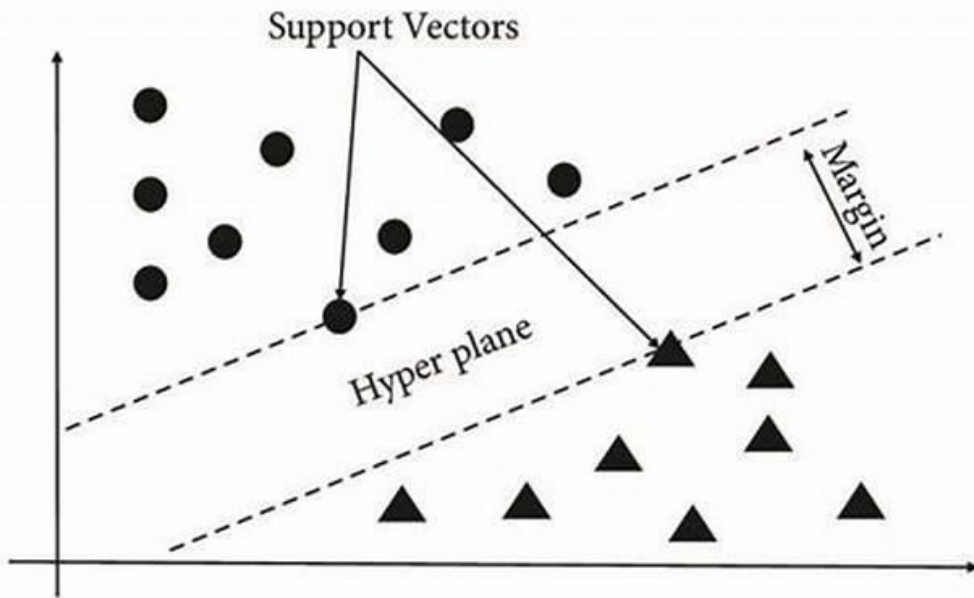
Decision trees: hyperplanes



A decision tree can be plotted with each decision segmenting a space into boxes.

The decision boundary is called a **hyperplane**.

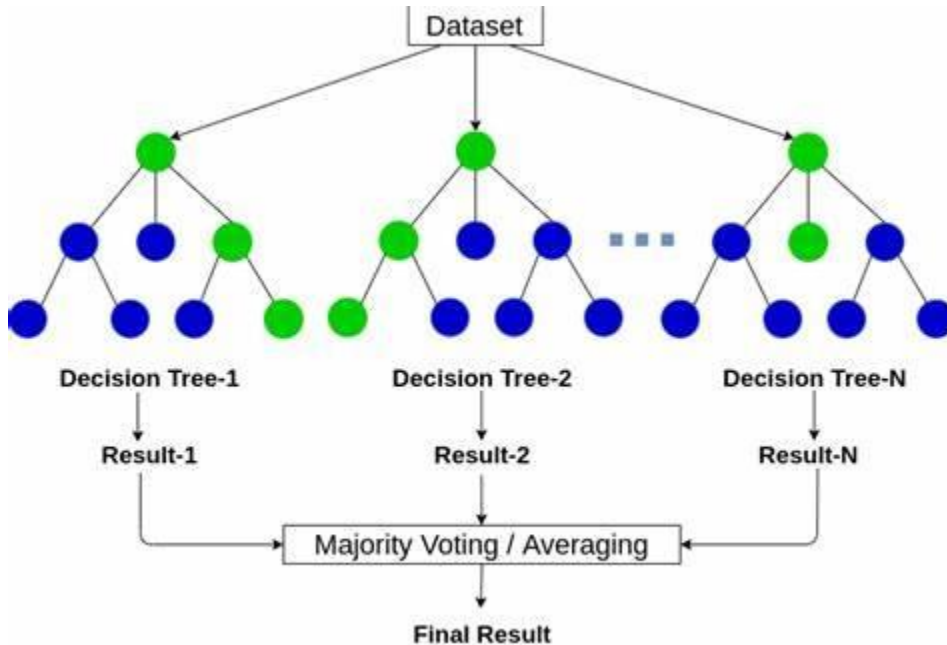
Support Vector Machines



For SVM, part of the objective function (the goal) uses not only the accuracy of the prediction, but also maximises the width of the margin between categories.

- **Improves Generalisation:**
- A larger margin reduces the risk of overfitting and improves the model's ability to generalize to new data.
- **Robustness:**
- By focusing on the support vectors, SVM ignores other points, making it less sensitive to noise or outliers.

Random Forests



Widely used for regression and classification

Consists of a 'forest' of decision trees:

All fit on *random bootstrap samples* of the data (each tree is trained on a slightly different dataset).

At each split in a tree, Random Forest considers only a *random subset of features* rather than all features.

Each decision tree is *trained independently* on its respective bootstrapped dataset and feature subset.

- The results are averaged for regression
- Majority vote for classification

Next Week: Reading Week – For week 7:



- Read Data Science for Business, chapters 5, 7, 8



- Read [AUC-ROC](#): a really good article



- Watch StatQuest: [ROC and AUC, Clearly Explained! – YouTube](#)



- Watch StatQuest: [Bias and Variance](#)



- Watch StatQuest: [Cross validation](#)



- Watch StatQuest: [Sensitivity and Specificity](#)



Any questions?

?