# Class: Machine Learning

# Decision Trees

**Instructor: Matteo Leonetti**

# Learning outcomes

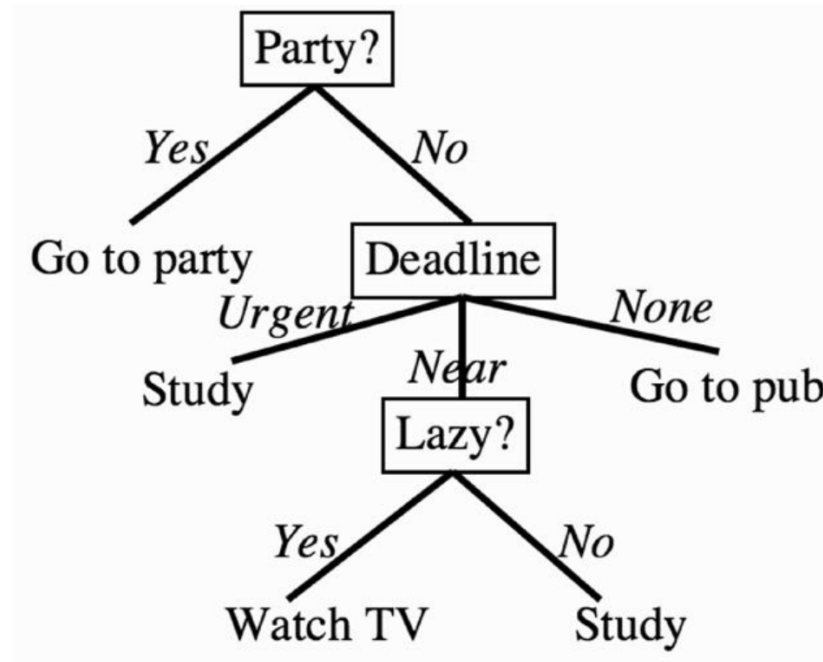- Define the entropy of a set
- Compute the entropy of a given set
- Define the information gain for a given feature
- Define the Gini Impurity of a set
- Implement the ID3 and CART algorithms

# Making Decisions

Nonmetric data

How to choose the variable for each split?

# History

1983 - Ross Quinlan (U. of Sidney)

*Learning efficient classification procedures and their application to chess end games*.

# Entropy and information

How much information do I receive, with a message X?

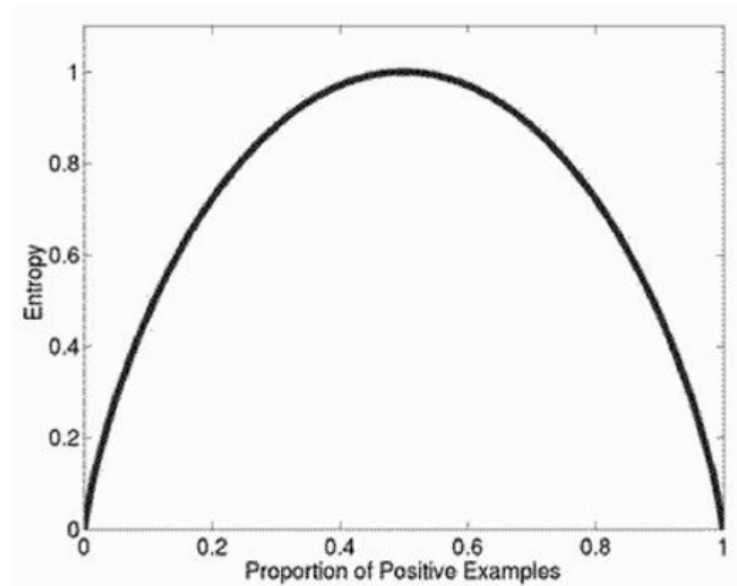X a random variable over possible messages
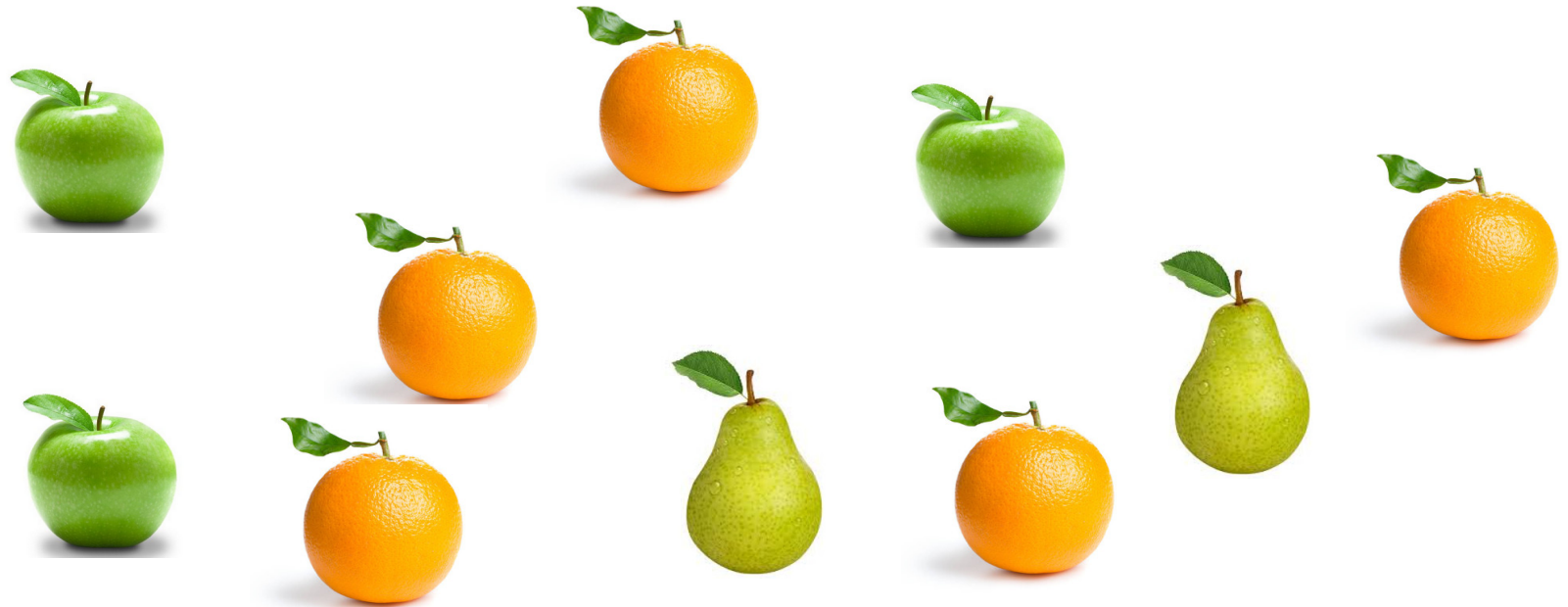
Information

Entropy

$$I(x) = -\log_2 P(x)$$

$$H = E[I] = \sum_i -p_i \log_2 p_i$$

$$0 \log_2 0 = 0$$

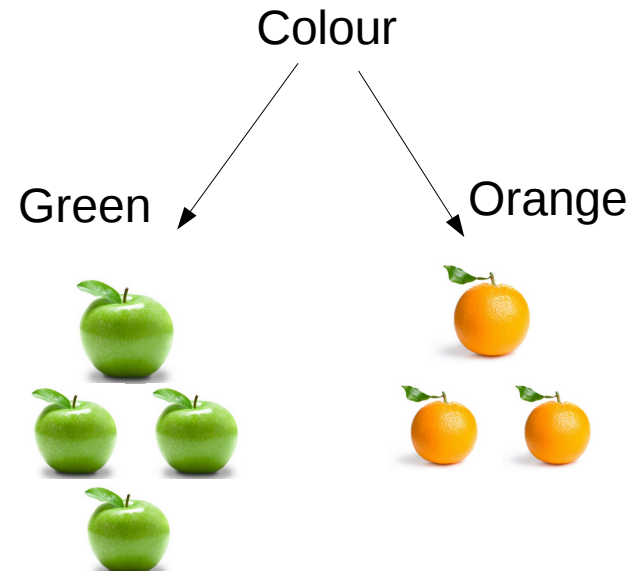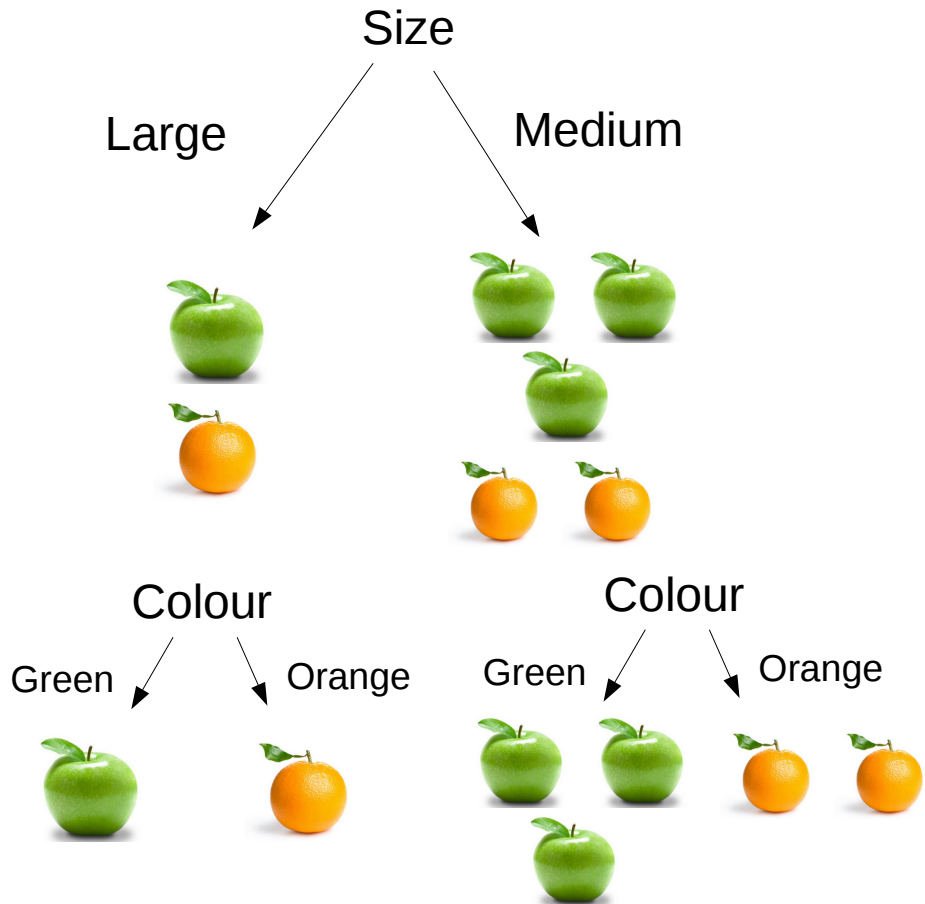# Fruits

$$H = E[I] = \sum_{i \in classes} -p_i \log_2 p_i$$



$$H = -\overbrace{\frac{3}{10} \log_2 \frac{3}{10}}^{Apples} - \overbrace{\frac{5}{10} \log_2 \frac{5}{10}}^{Oranges} - \overbrace{\frac{2}{10} \log_2 \frac{2}{10}}^{Pears} = 1.485$$
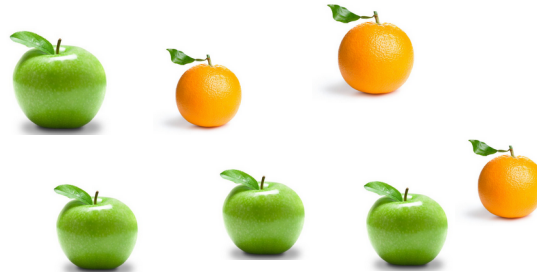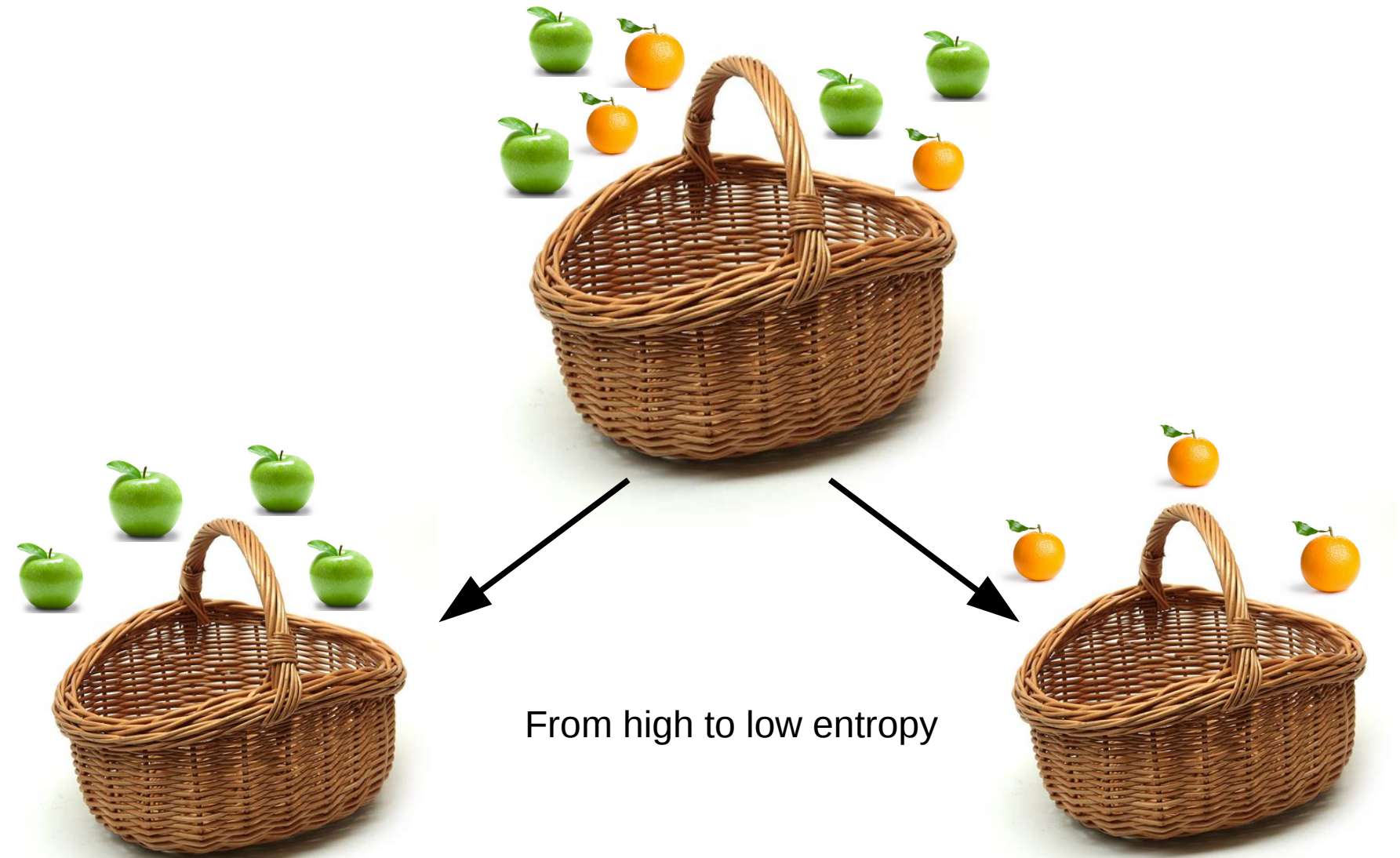
# Apples and Oranges

# Entropy of the set



$$H = -p_O \log_2(p_O) - p_A \log_2(p_A) = -\frac{3}{7}\log_2\left(\frac{3}{7}\right) - \frac{4}{7}\log_2\left(\frac{4}{7}\right) = 0.985$$
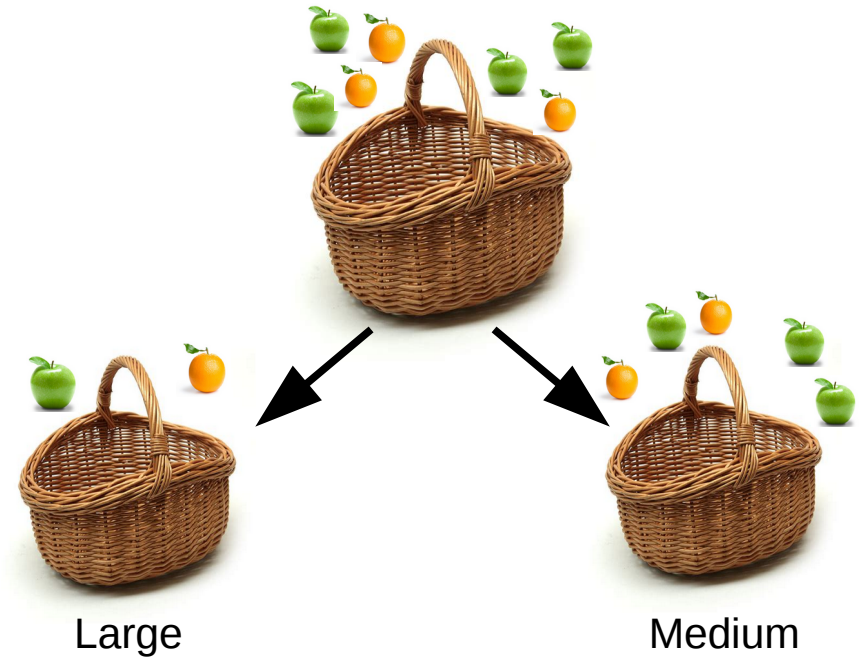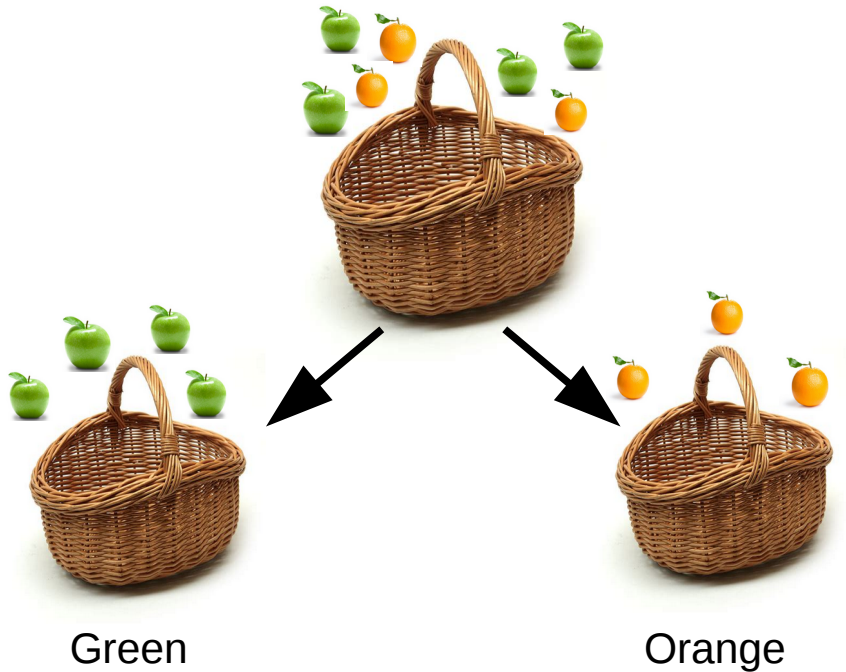
# Entropy of the set



From high to low entropy

# Entropy of the set



Green

Orange

Large

Medium

$$H_{colour} = \overbrace{\frac{4}{7}}^{\text{fraction in Green}} \overbrace{(0)}^{\text{entropy of Green}} + \overbrace{\frac{3}{7}}^{\text{fraction in Orange}} \overbrace{(0)}^{\text{entropy of Orange}} = 0$$

$$H_{size} = \overbrace{\frac{2}{7}}^{\text{fraction in Large}} \overbrace{\left(-\frac{1}{2}\log_2\frac{1}{2} - \frac{1}{2}\log_2\frac{1}{2}\right)}^{\text{entropy of Large}} + \overbrace{\frac{5}{7}}^{\text{fraction in Medium}} \overbrace{\left(-\frac{2}{5}\log_2\frac{2}{5} - \frac{3}{5}\log_2\frac{3}{5}\right)}^{\text{entropy of Medium}} = 0.98$$

# Entropy of the set

$H = 0.985$

$H_{colour} = 0$

$H_{size} = 0.98$

$G(\text{Colour}) = H - H_{colour} = 0.985$

$G(\text{Size}) = H - H_{size} = 0.005$

# Information gain

Set of elements

elements in S with feature F = f

$$G(S,F) = H(S) - \sum_{f \in values(F)} \frac{|S_f|}{|S|} H(S_f)$$

Feature

compare with:

fraction in Large     entropy of Large     fraction in Medium     entropy of Medium

$$H_{size} = \frac{2}{7} \left( -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{5}{7} \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) = 0.98$$

# The ID3 algorithm
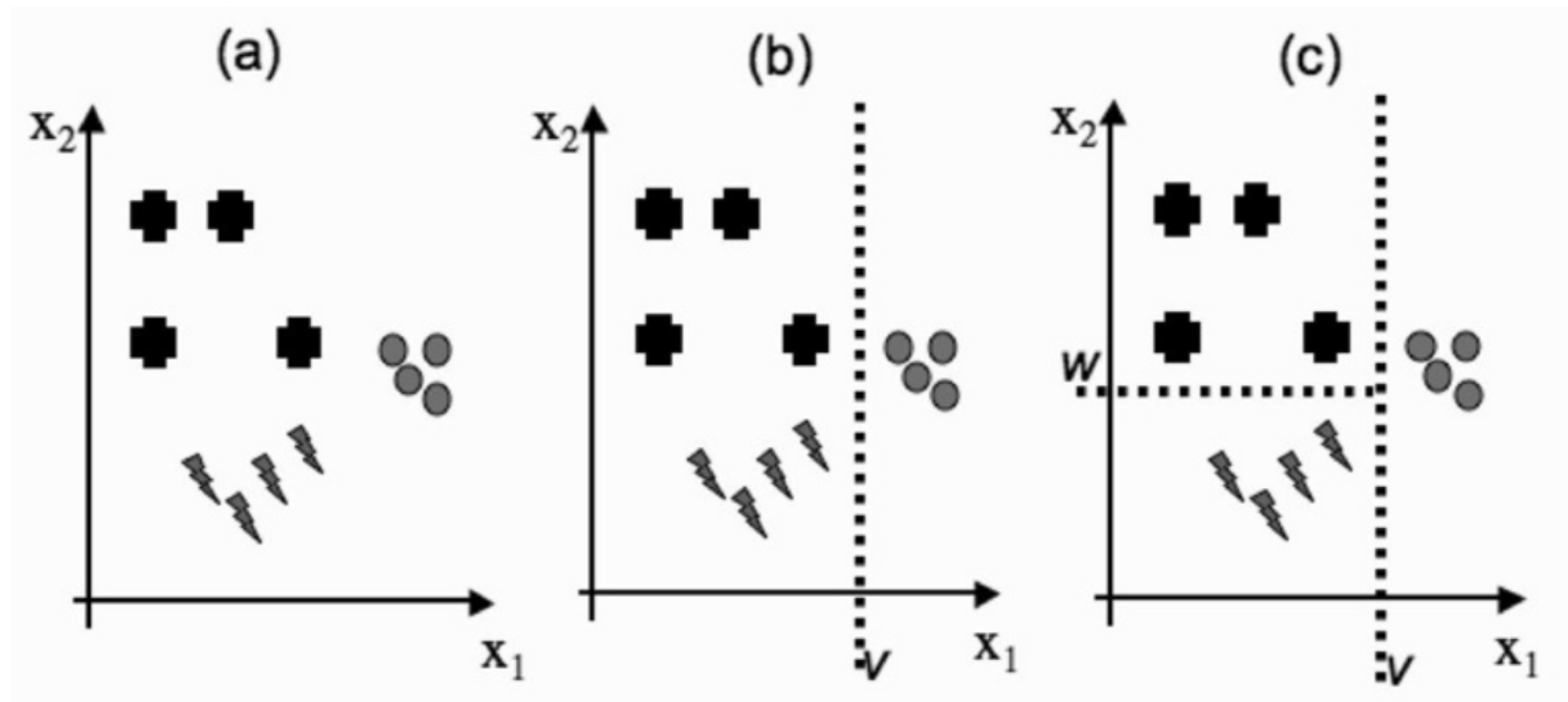
- If all examples have the same label:

    - return a leaf with that label

- Else if there are no features left to test:

    - return a leaf with the most common label

- Else:

    - choose the feature $\hat{F}$ that maximises the information gain of $S$ to be the next node using **Equation (12.2)**
    - add a branch from the node for each possible value $f$ in $\hat{F}$
    - for each branch:

        * calculate $S_f$ by removing $\hat{F}$ from the set of features
        * recursively call the algorithm with $S_f$, to compute the gain relative to the current set of examples

# Visualizing splits

# Characteristics

Greedy with respect to G → potential local minimum

Deals with noisy data (by assigning the label to most common class)

Always uses all the features → prone to overfitting

Pruning

Continuous variables ⟶ C4.5

Missing attributes
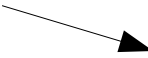
Colour

Green          Orange

Size

Large

P    A

P    P

O

Medium

A    A

A

O

O

$$G(S) = \overbrace{\frac{4}{10}}^{apples} (\overbrace{\frac{3}{10} + \frac{3}{10}}^{not\ apples}) + \overbrace{\frac{3}{10}}^{pears} (\overbrace{\frac{4}{10} + \frac{3}{10}}^{not\ pears}) + \overbrace{\frac{3}{10}}^{oranges} (\overbrace{\frac{4}{10} + \frac{3}{10}}^{not\ oranges}) = \sum_{i}^{C} p_i(1 - p_i)$$
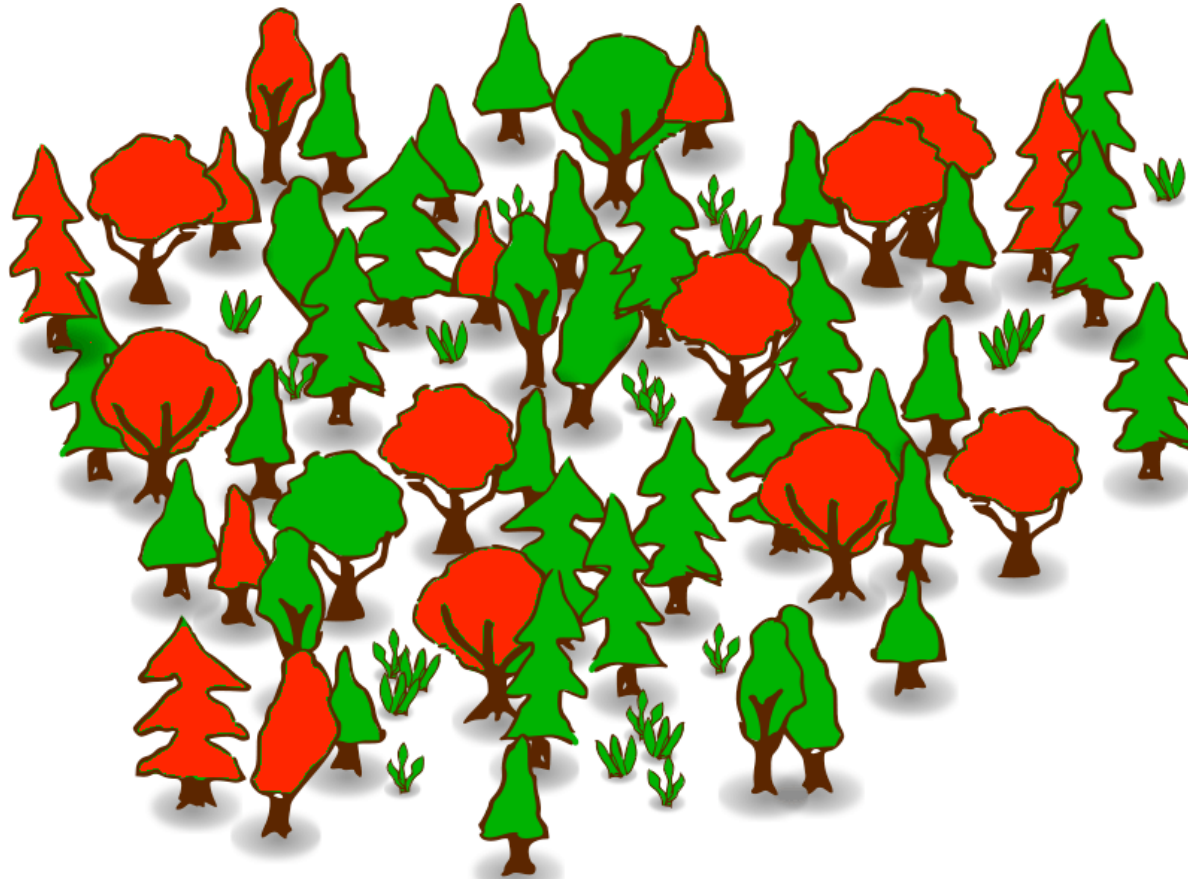
Gini split:

# of classes

$$G(S)=\sum_i^C p_i(1-p_i)=\sum_i^C (p_i-p_i^2)=\sum_i^C p_i-\sum_i^C p_i^2=1-\sum_i^C p_i^2$$

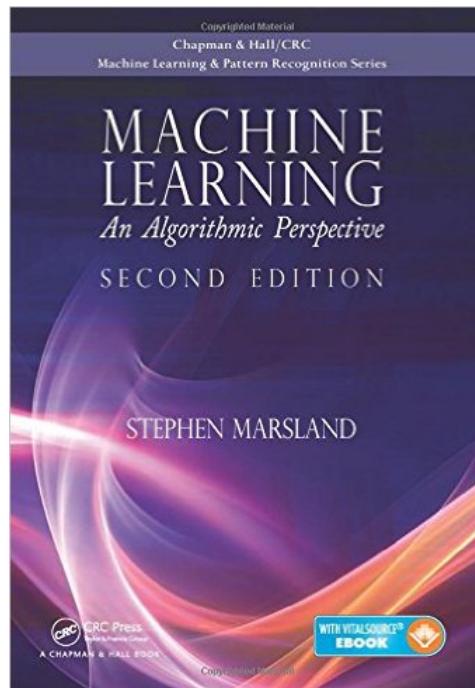$$G(S,F)=G(S)-\sum_{f\in values(F)}\frac{|S_f|}{|S|}G(S_f)$$

# Random forests

# Conclusion

# Learning outcomes

- Define the entropy of a set
- Compute the entropy of a given set
- Define the information gain for a given feature
- Define the Gini Impurity of a set
- Implement the ID3 and CART algorithms

# Chapter 12