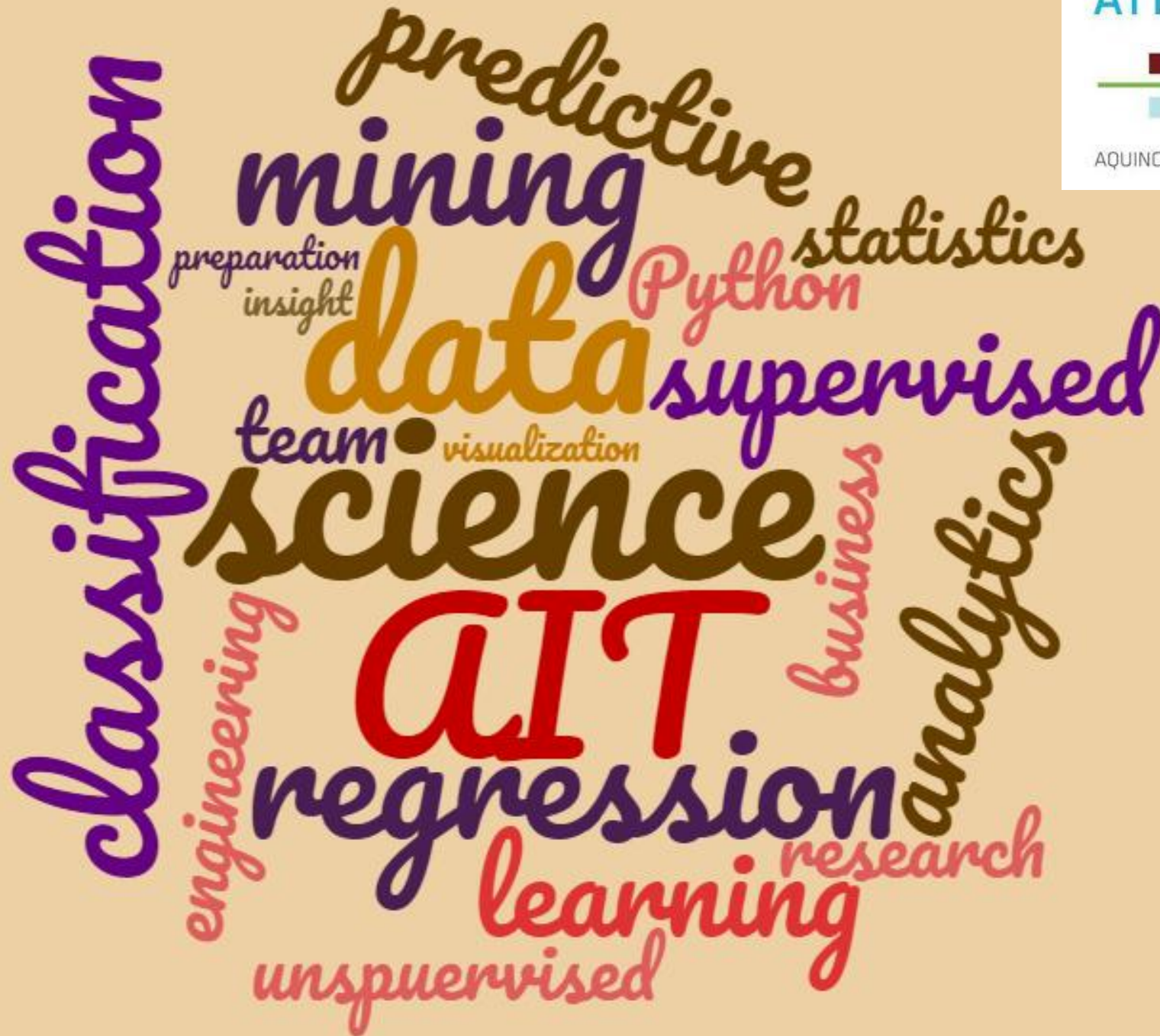


# Data Science

February 27, 2023.  
Decision tree



AIT-BUDAPEST



AQUINCUM INSTITUTE OF TECHNOLOGY

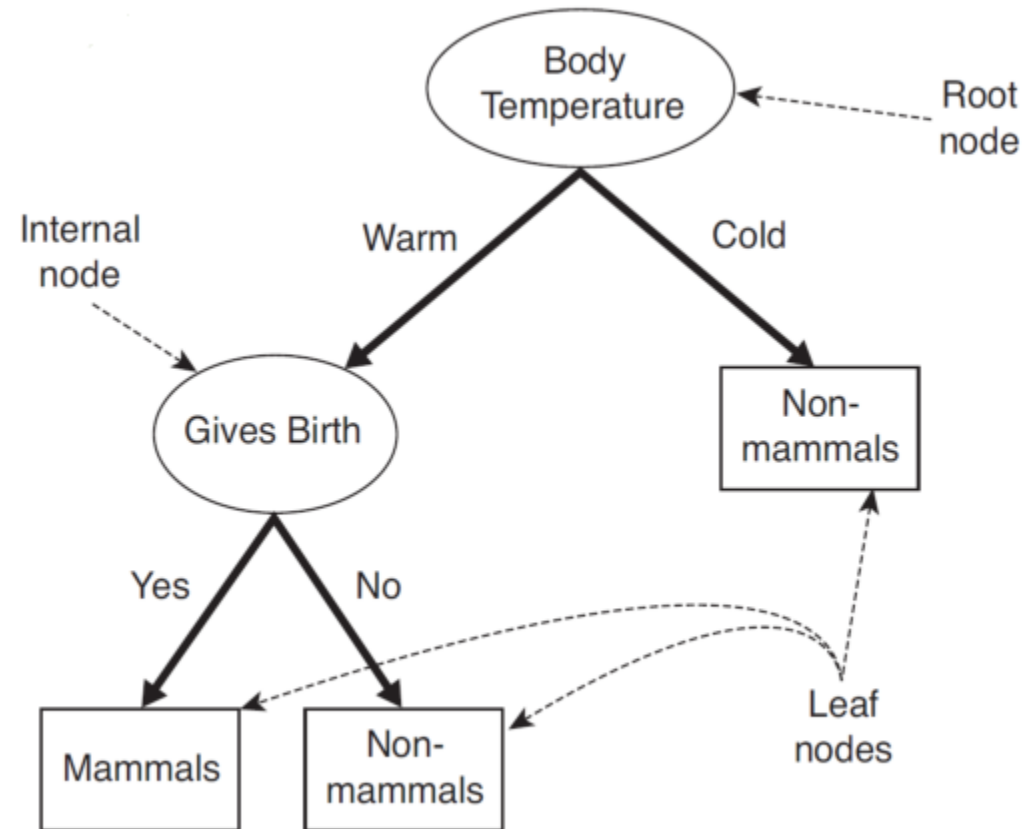
Roland Molontay

# Schedule of the semester

|                     | <i>Monday midnight</i> | <i>Tuesday class</i>  | <i>Friday class</i>   |
|---------------------|------------------------|-----------------------|-----------------------|
| <b>W1 (02/06)</b>   |                        |                       |                       |
| <b>W2 (02/13)</b>   |                        | HW1 out               |                       |
| <b>W3 (02/20)</b>   |                        |                       |                       |
| <b>W4 (02/27)</b>   | HW1 deadline + TEAMS   | HW2 out               |                       |
| <b>W5 (03/06)</b>   | PROJECT PLAN           |                       |                       |
| <b>W6 (03/13)</b>   | HW2 deadline           | HW3 out               |                       |
| <b>W7 (03/20)</b>   |                        |                       | MIDTERM               |
| <b>SPRING BREAK</b> |                        | SPRING BREAK          | SPRING BREAK          |
| <b>W8 (04/03)</b>   | HW3 deadline           |                       | GOOD FRIDAY           |
| <b>W9 (04/10)</b>   | MILESTONE 1            | HW4 out               |                       |
| <b>W10 (04/17)</b>  |                        |                       |                       |
| <b>W11 (04/24)</b>  | HW4 deadline           |                       |                       |
| <b>W12 (05/01)</b>  | MILESTONE 2            |                       |                       |
| <b>W13 (05/08)</b>  |                        |                       |                       |
| <b>W14 (05/15)</b>  |                        | FINAL                 | PROJECT presentations |
| <b>W15 (05/22)</b>  |                        | PROJECT presentations |                       |

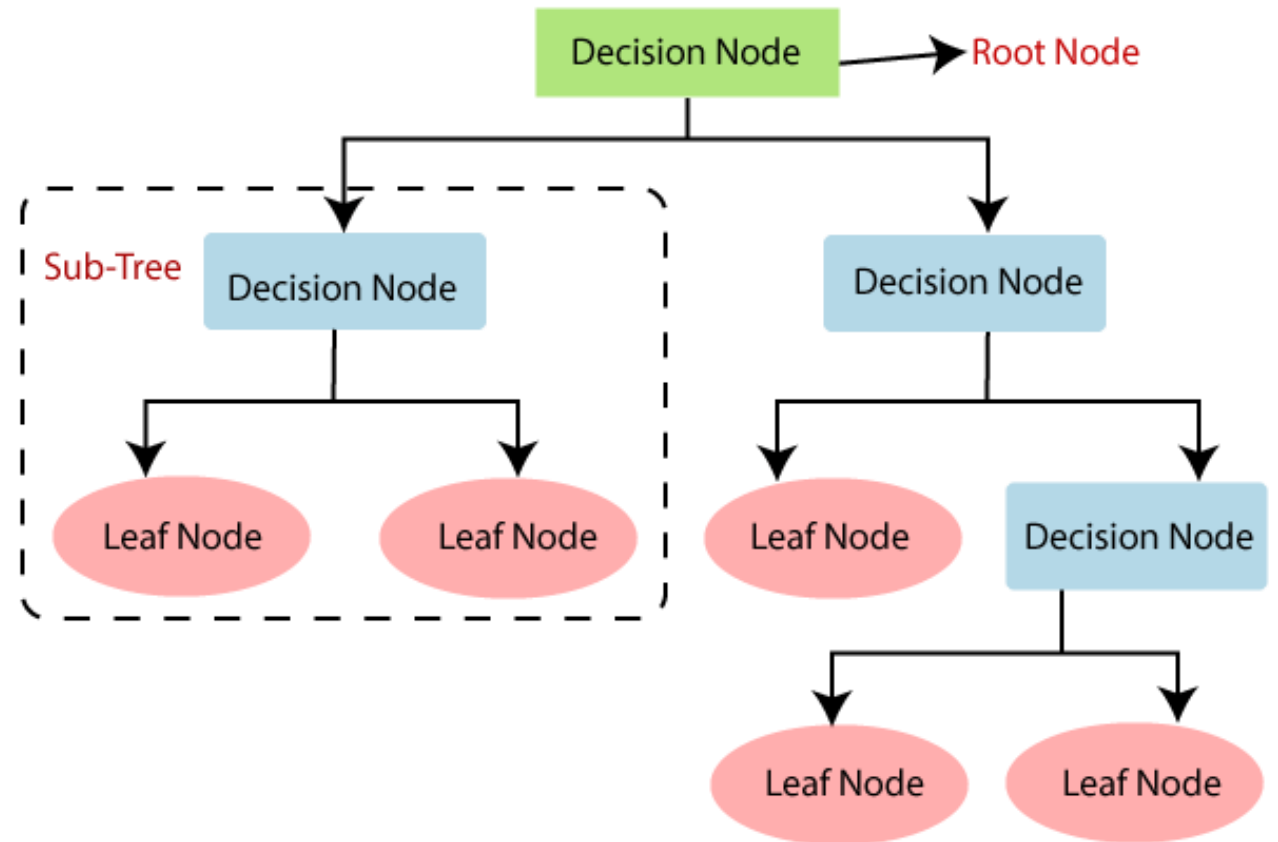
# Decision tree

- Rooted, directed (sometimes binary) tree
- Each inner node has a decision rule that assigns instances uniquely to child nodes of the actual node
- Each leaf node has a class label



# Prediction with decision tree

- Decision tree is one of the oldest predictive algorithm
  - We start at root node
  - At each interior node we evaluate the decision rule, branch to the child node picked by the decision rule
  - Once a leaf node is reached, predict the label assigned to that node
- It is mainly used for classification, but also suitable for regression problems

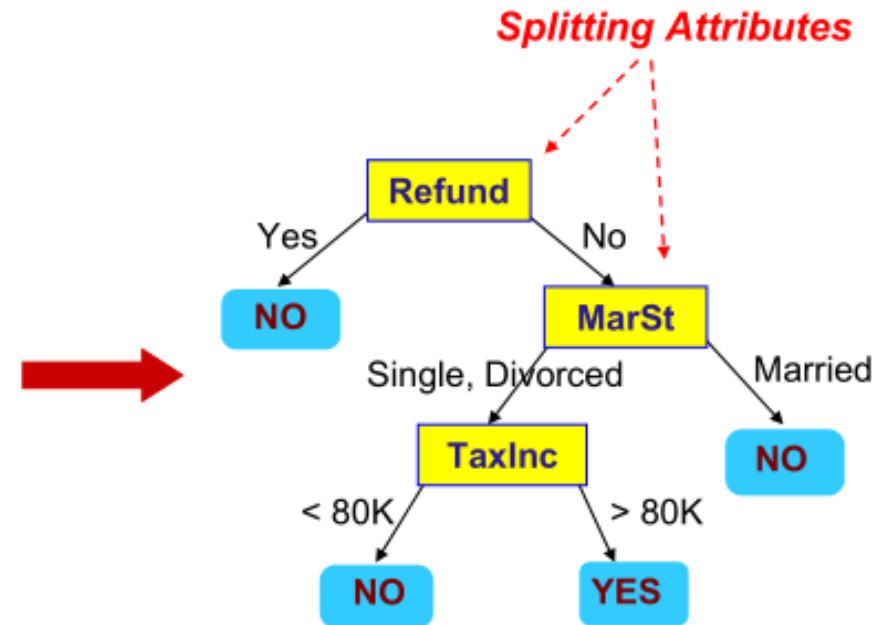


# Decision tree - learning phase

It works with  
continuous  
and categorical  
attributes

| <i>Tid</i> | <i>Refund</i> | <i>Marital Status</i> | <i>Taxable Income</i> | <i>Cheat</i> |
|------------|---------------|-----------------------|-----------------------|--------------|
| 1          | Yes           | Single                | 125K                  | No           |
| 2          | No            | Married               | 100K                  | No           |
| 3          | No            | Single                | 70K                   | No           |
| 4          | Yes           | Married               | 120K                  | No           |
| 5          | No            | Divorced              | 95K                   | Yes          |
| 6          | No            | Married               | 60K                   | No           |
| 7          | Yes           | Divorced              | 220K                  | No           |
| 8          | No            | Single                | 85K                   | Yes          |
| 9          | No            | Married               | 75K                   | No           |
| 10         | No            | Single                | 90K                   | Yes          |

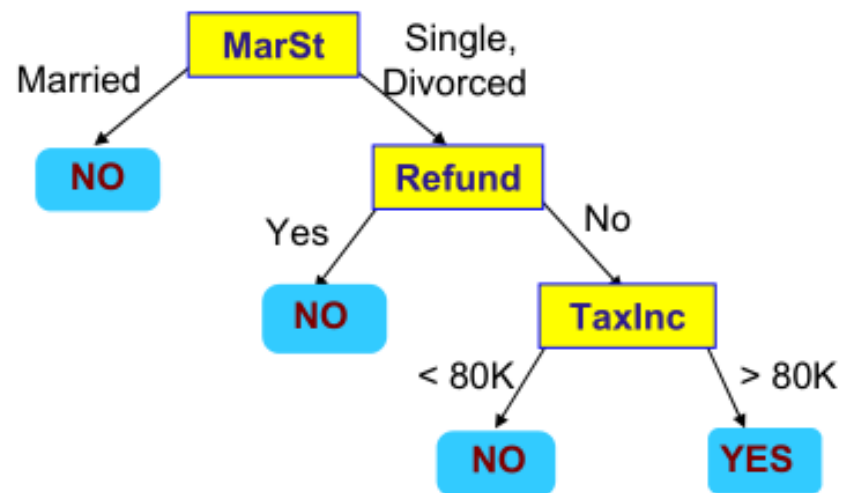
Training Data



Model: Decision Tree

# Another tree for the same data

| <i>Tid</i> | <i>Refund</i> | <i>Marital Status</i> | <i>Taxable Income</i> | <i>Cheat</i> |
|------------|---------------|-----------------------|-----------------------|--------------|
| 1          | Yes           | Single                | 125K                  | No           |
| 2          | No            | Married               | 100K                  | No           |
| 3          | No            | Single                | 70K                   | No           |
| 4          | Yes           | Married               | 120K                  | No           |
| 5          | No            | Divorced              | 95K                   | Yes          |
| 6          | No            | Married               | 60K                   | No           |
| 7          | Yes           | Divorced              | 220K                  | No           |
| 8          | No            | Single                | 85K                   | Yes          |
| 9          | No            | Married               | 75K                   | No           |
| 10         | No            | Single                | 90K                   | Yes          |

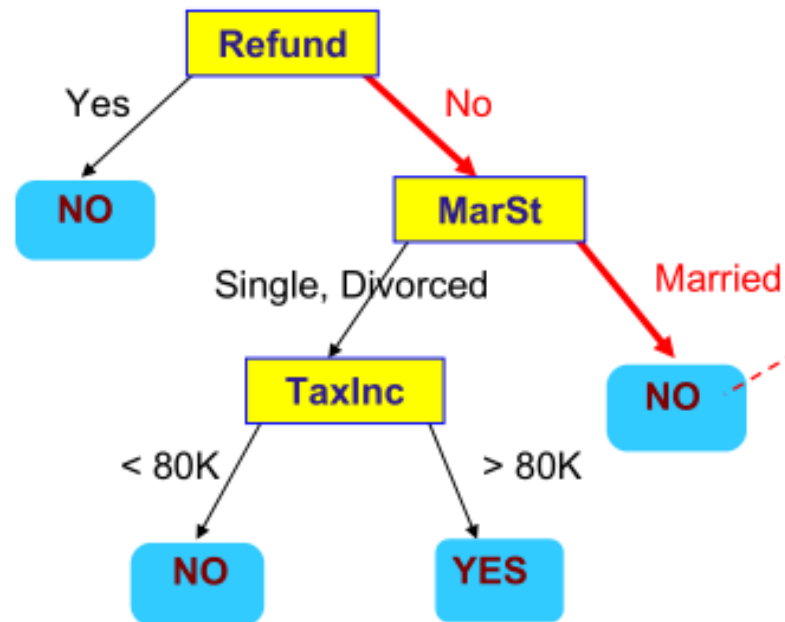


There could be more than one tree that fits the same data!

# Applying the model to new data

Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No     | Married        | 80K            | ?     |



Assign Cheat to "No"

# Problems

## Ex 1

How many logical (Boolean,  $f : \{0,1\}^N \rightarrow \{0,1\}$ ) functions can be generated on  $N$  binary attributes? What are the possible functions for  $N = 2$ ?

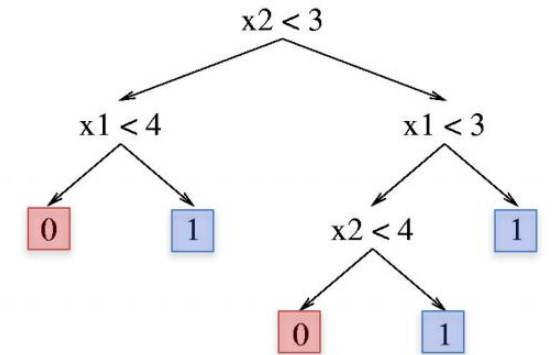
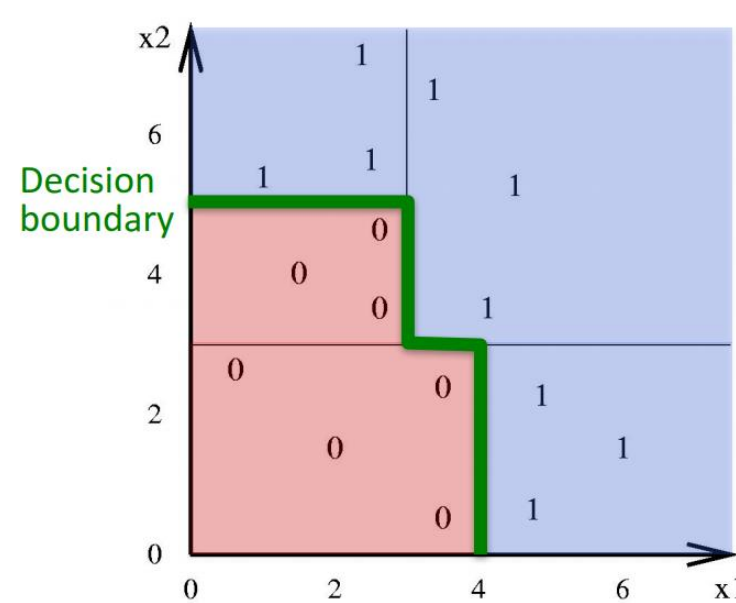
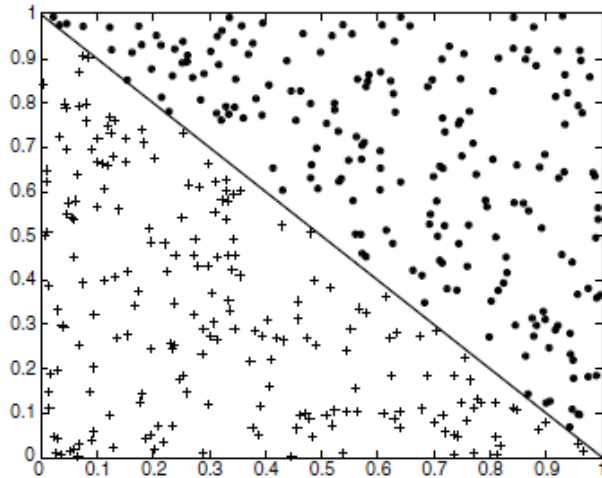
## Ex 2

Can decision trees learn logical (Boolean) functions? How to represent the following functions with a decision tree:  $A \text{ OR } B$ ,  $A \text{ AND } B$ ,  $A \text{ XOR } B$ , where  $A$  and  $B$  are logical variables.



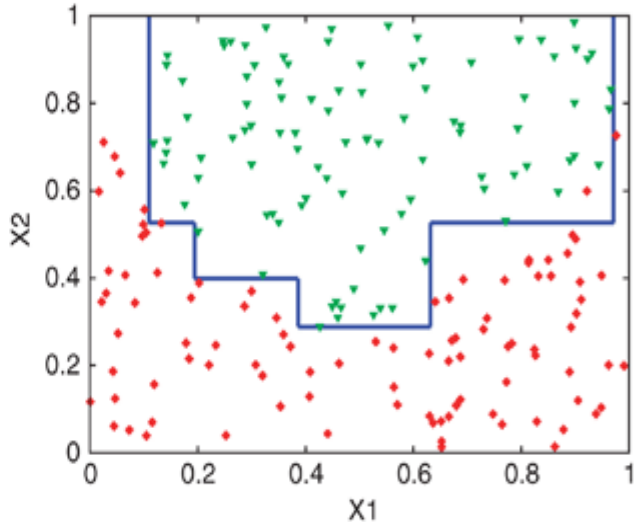
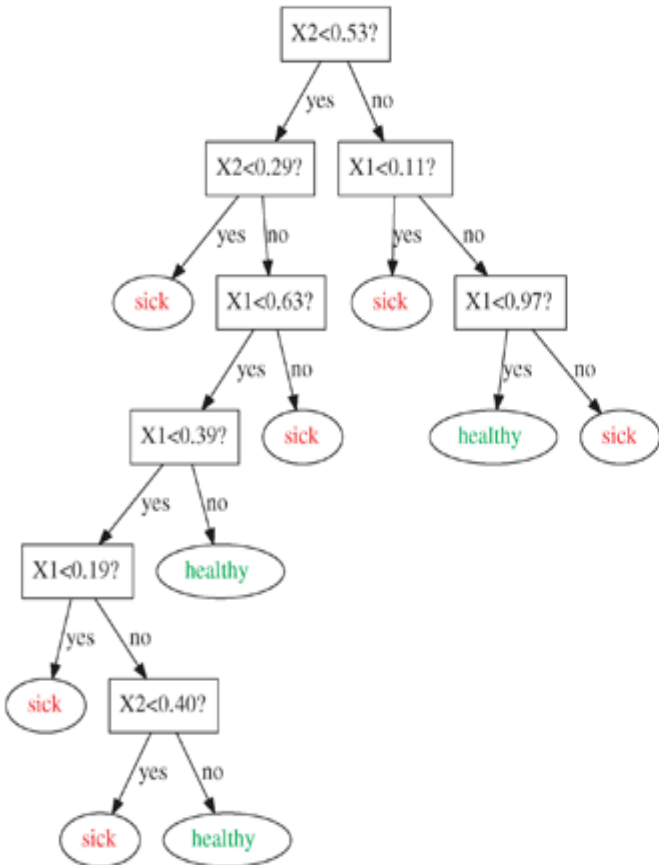
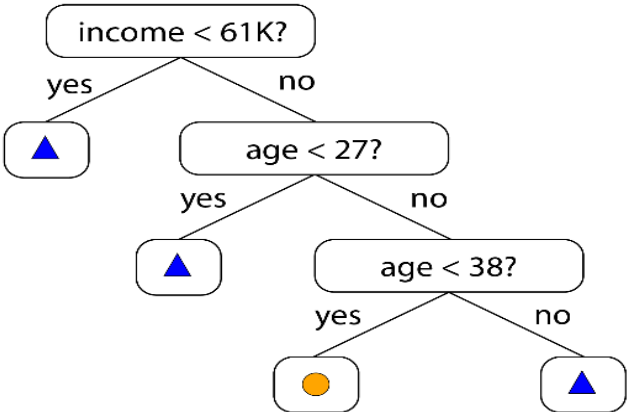
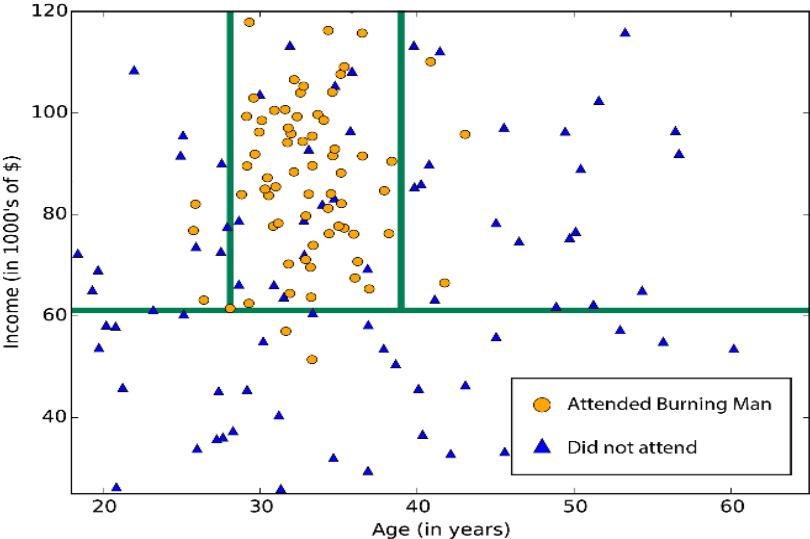
# Partitioning the feature space

The boundaries are always parallel to the axes  $\rightarrow$  the decision regions are always unions of (hyper-)rectangles



How could a decision tree represent a data set on the left?

## Partitioning the feature space - examples



# Algorithm for building decision trees

- We consider a general algorithm
  - Special versions of the algorithm are implemented in various programs
- We do not aim to find the best tree (it is underdefined what is best) but a good enough tree
  - There are exponentially many possible trees
- With a greedy method, deciding locally, quickly
- There are various algorithms (and their variants)
  - Hunt algorithm
  - CART
  - ID3, C4.5 (J48)
  - SLIQ, SPRINT

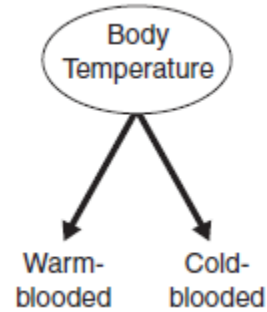
# Sketch of Hunt algorithm

- We start with one node (the root), all the records are there
  - Its label is the majority label
- Further steps: we choose a node that is worth splitting on
- End: until there is no node worth splitting on
  - In each node all the records have the same label
  - There are no good splits, e.g. in all nodes, all the records agree in each attribute (aside from the label)

# Hunt algorithm - questions

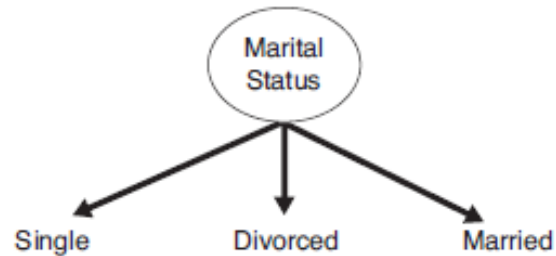
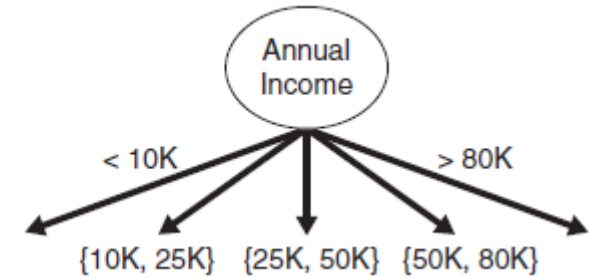
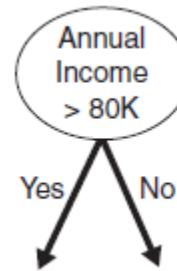
- When does it terminate?
  - When no more nodes left that are worth splitting on
- When is it not worth splitting on?
  - In each node all the records have the same label
  - There are no good splits, e.g. in all nodes, all the records agree in each attribute (aside from the label)
  - If we want to avoid to have a too deep tree (more details later)
- Which node to split on if there are more possibilities?
  - E.g. traversing nodes according to BFS (breadth-first search), DFS (depth-first search)
- How to split?
  - Splitting a node should increase the homogeneity (with respect to the label) of resultant sub-nodes
- How to decide on the label?
  - Using majority voting

# Possible splits

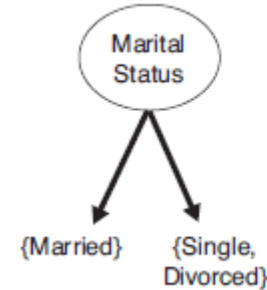


Binary

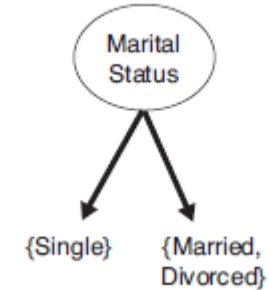
Numerical



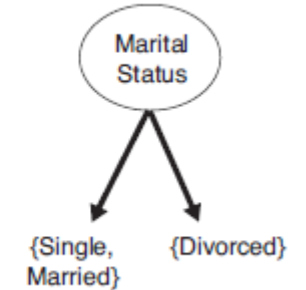
Nominal



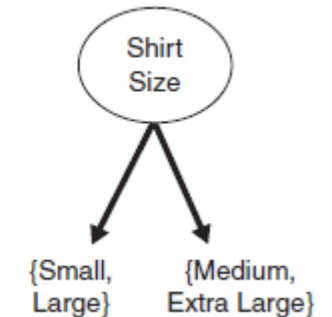
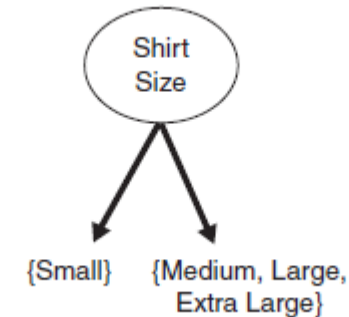
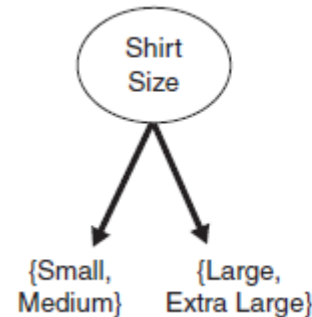
OR



OR

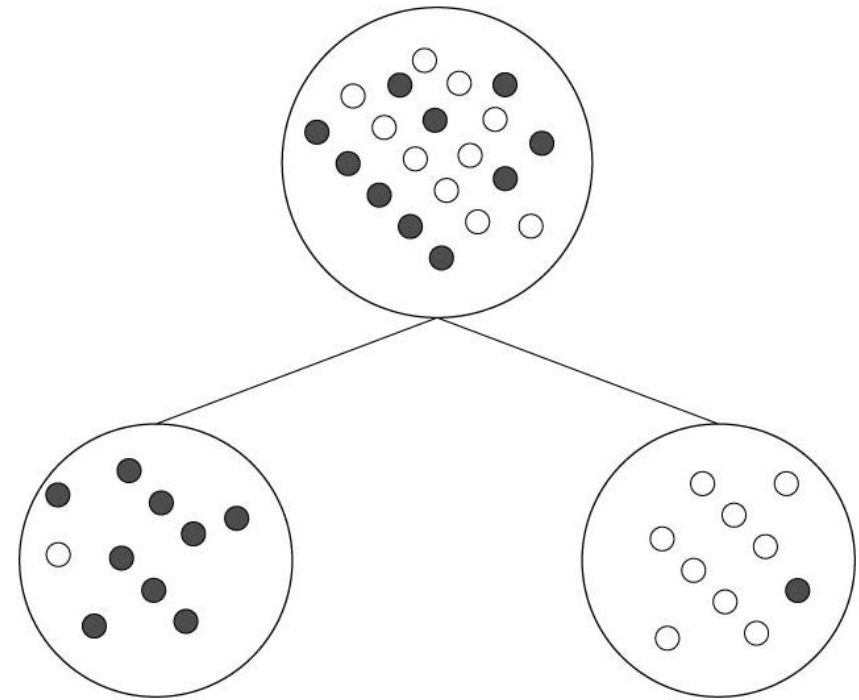


Ordinal



# What defines a good split?

- It increases the homogeneity („purity”) of the target variable in the emerging child nodes compared to the parent node
- A good split creates child nodes of similar size (or at least does not create very small nodes)



# Example

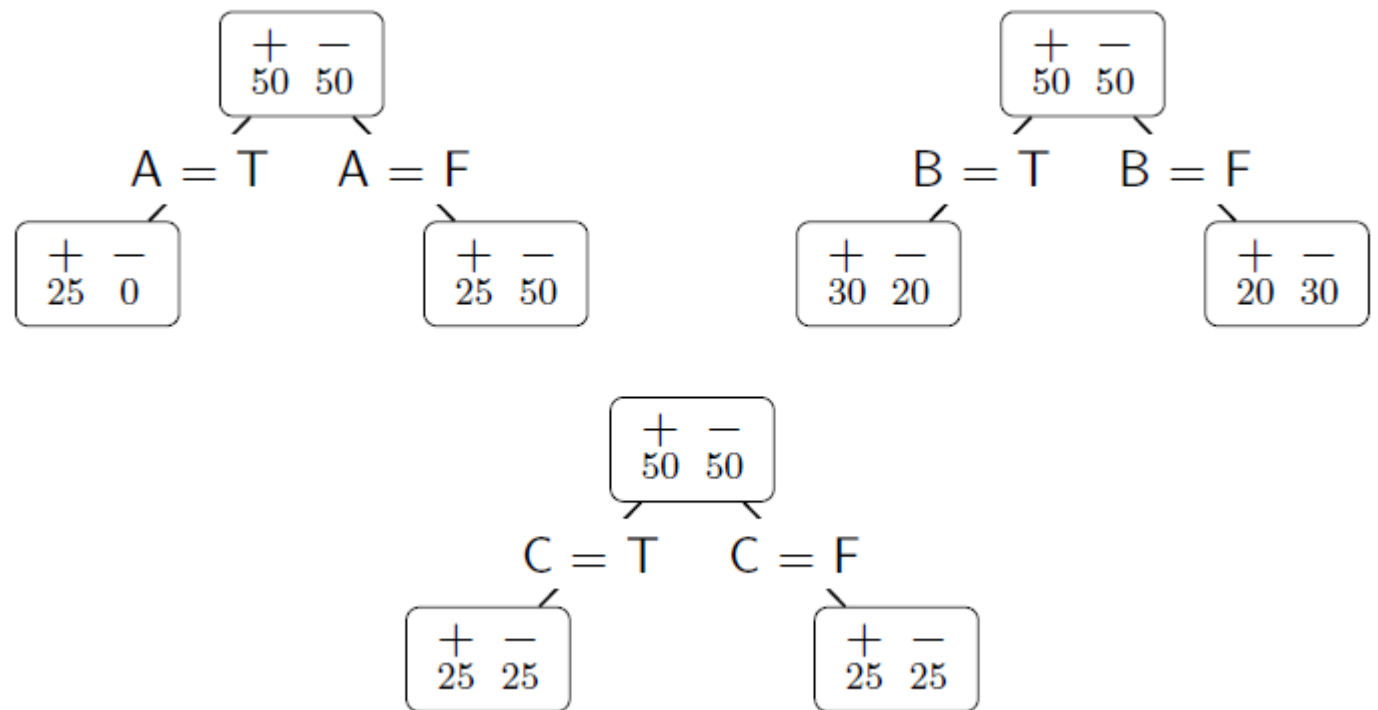
- The following table summarizes a data set with three binary attributes (A, B, C) and two class labels (+, -).
- Which is the best split? What attribute would you split on?

| A | B | C | Number of instances |          |
|---|---|---|---------------------|----------|
|   |   |   | class: +            | class: - |
| T | T | T | 5                   | 0        |
| F | T | T | 0                   | 20       |
| T | F | T | 20                  | 0        |
| F | F | T | 0                   | 5        |
| T | T | F | 0                   | 0        |
| F | T | F | 25                  | 0        |
| T | F | F | 0                   | 0        |
| F | F | F | 0                   | 25       |



# Possible splits

| A | B | C | Number of instances |          |
|---|---|---|---------------------|----------|
|   |   |   | class: +            | class: - |
| T | T | T | 5                   | 0        |
| F | T | T | 0                   | 20       |
| T | F | T | 20                  | 0        |
| F | F | T | 0                   | 5        |
| T | T | F | 0                   | 0        |
| F | T | F | 25                  | 0        |
| T | F | F | 0                   | 0        |
| F | F | F | 0                   | 25       |

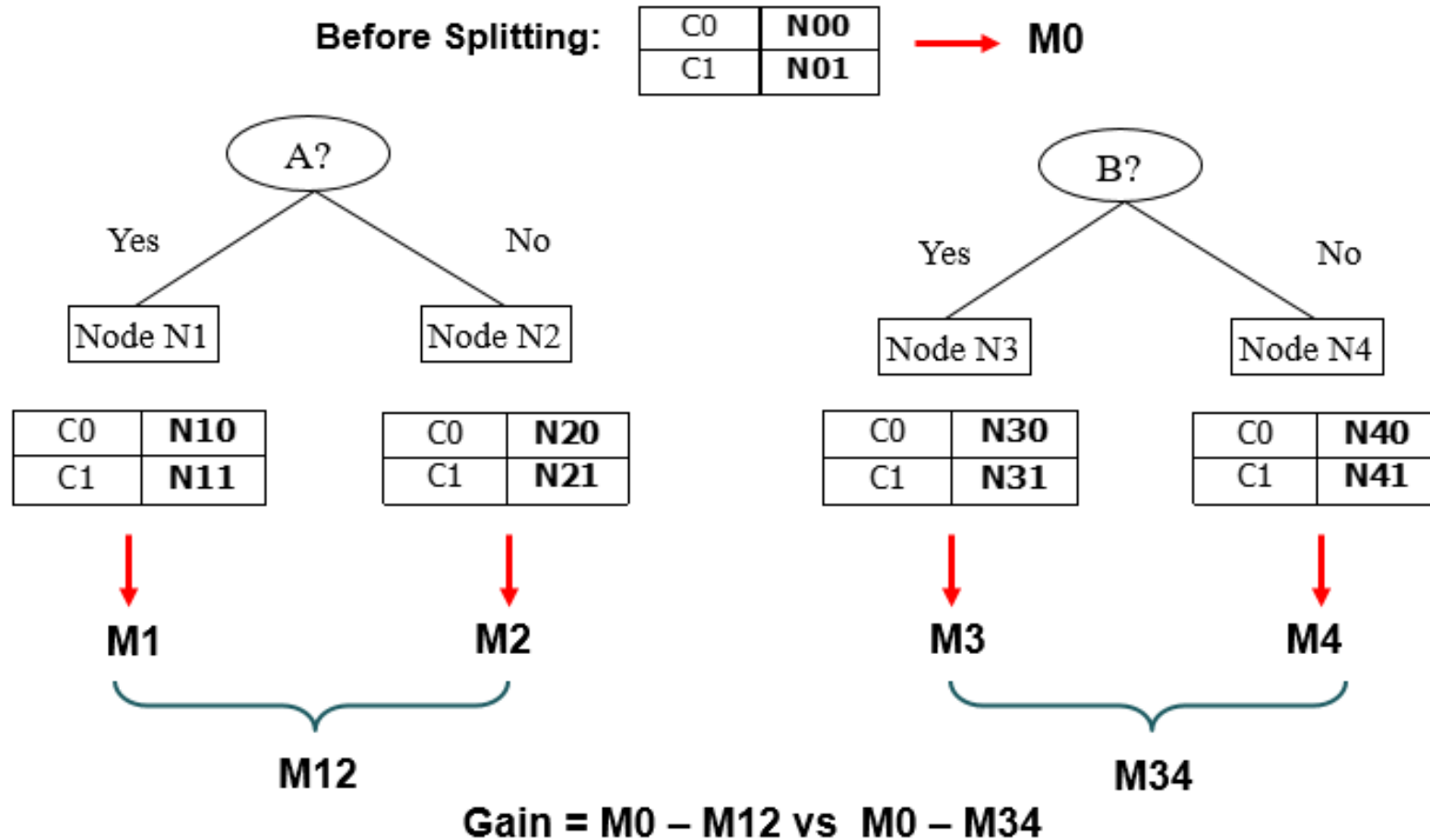


A is the best split!

# How to measure the goodness of a split?

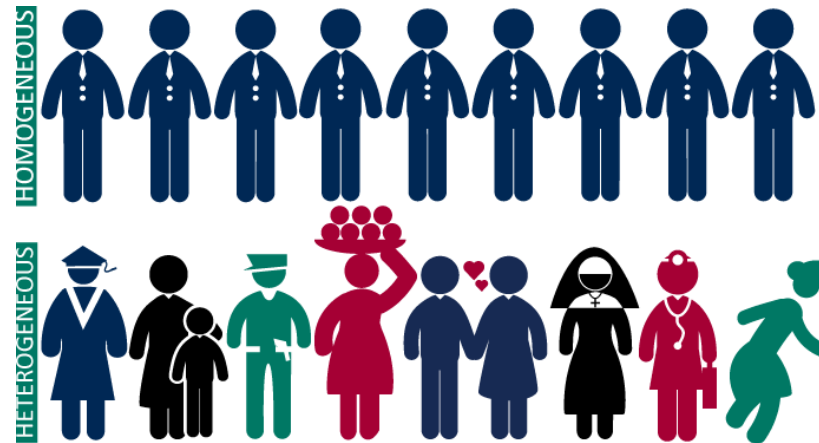
- We consider three possible metrics
- The main principle is the same for all:
  - We define a metric for a set of records that measures the degree of homogeneity (purity) of the target variable within that set
    - We consider three metrics: Gini coefficient, entropy, misclassification error
  - The goodness of a split is measured by the difference between the impurity of the parent node and the emergent child nodes
    - How much do we gain if we split the node? What is the degree of homogeneity increase?
    - We can also consider the size of the emerging child nodes, punishing the too small splits

# Which split is better?



# Measuring inhomogeneity

- Let  $t$  be the node of a decision tree (i.e. a set of records), we aim to measure its homogeneity with respect to a target variable that has  $c$  possible values
- Let  $p(i|t)$  denote the relative frequency of records with label  $i$  in node  $t$



# Misclassification error

- Misclassification error or classification error:

$$\text{Classification error}(t) = 1 - \max_i [p(i|t)]$$

- Its maximal value is:  $1 - 1/c$ , when the records are distributed equally in all classes
- Its minimal value is 0, when all the records have the same label

|    |          |
|----|----------|
| C1 | <b>0</b> |
| C2 | <b>6</b> |

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Error} = 1 - \max(0, 1) = 1 - 1 = 0$$

|    |          |
|----|----------|
| C1 | <b>1</b> |
| C2 | <b>5</b> |

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Error} = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

|    |          |
|----|----------|
| C1 | <b>2</b> |
| C2 | <b>4</b> |

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Error} = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

# Gini coefficient

- Gini coefficient

$$\text{Gini}(t) = 1 - \sum_{i=1}^c p(i|t)^2$$

|    |          |
|----|----------|
| C1 | <b>0</b> |
| C2 | <b>6</b> |

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Gini} = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

- Its value ranges between 0 and  $1 - 1/c$
- Usually Gini is default setting for splitting criterion

|    |          |
|----|----------|
| C1 | <b>1</b> |
| C2 | <b>5</b> |

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

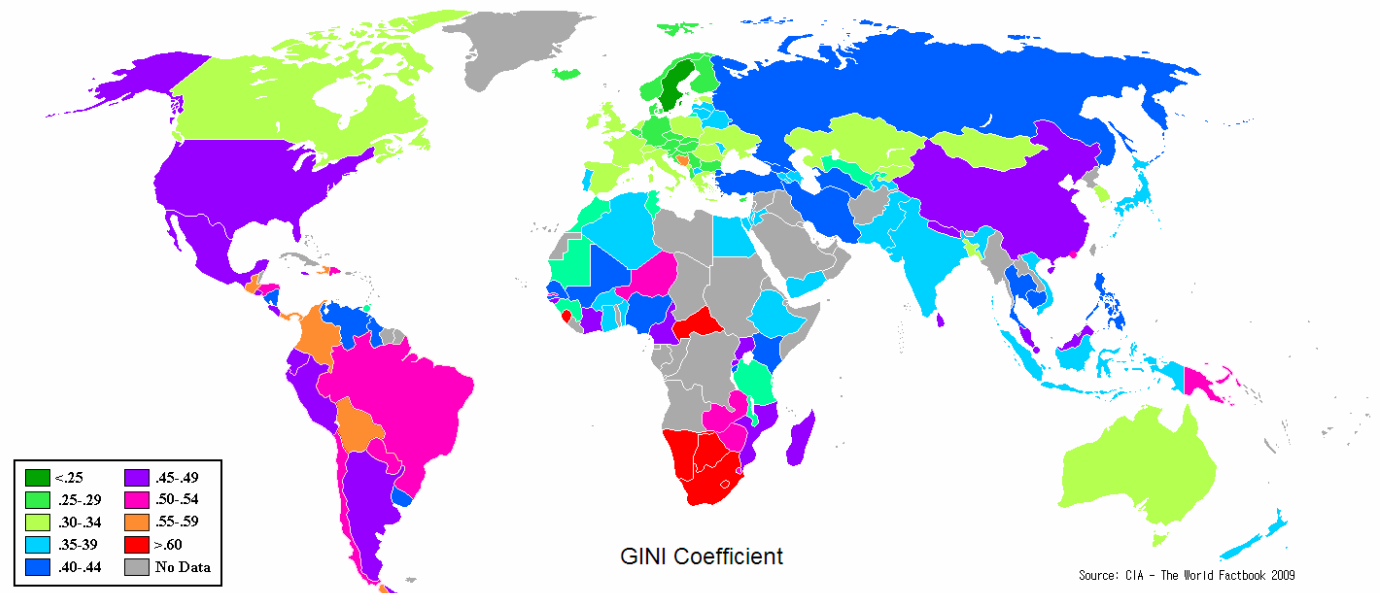
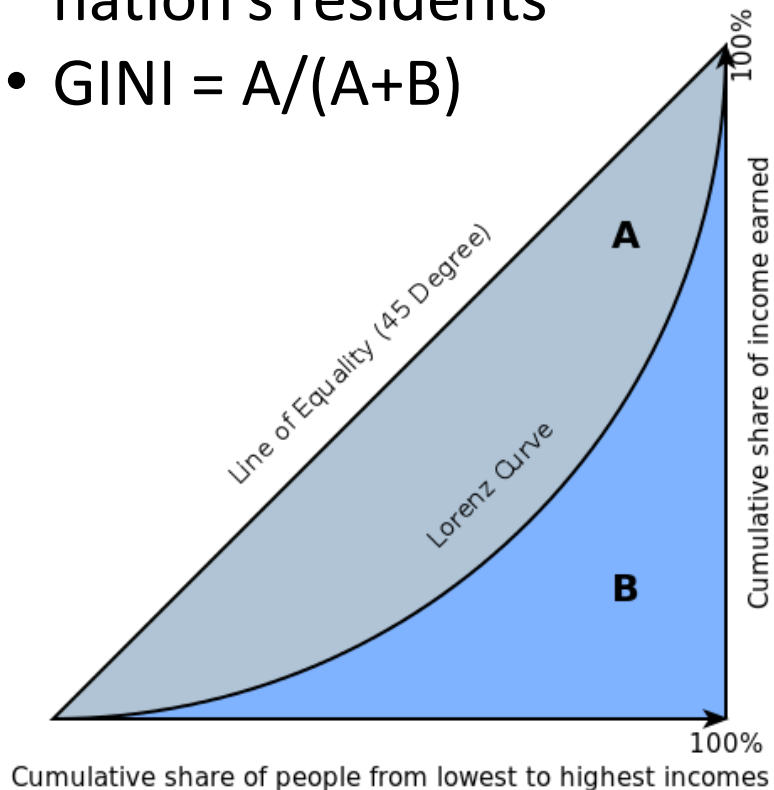
|    |          |
|----|----------|
| C1 | <b>2</b> |
| C2 | <b>4</b> |

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Gini} = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

# Outlook – Gini coefficient in economics

- A similar concept but not the same, do not mix them
- It is used to measure the inequality in income or wealth distribution of a nation's residents
- $GINI = A/(A+B)$



# Entropy

$$\text{Entropy}(t) = - \sum_{i=1}^c p(i|t) \log_2 p(i|t)$$

|    |          |
|----|----------|
| C1 | <b>0</b> |
| C2 | <b>6</b> |

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Entropy} = -0 \log 0 - 1 \log 1 = -0 - 0 = 0$$

- Its value ranges between 0 and  $\log_2 c$

|    |          |
|----|----------|
| C1 | <b>1</b> |
| C2 | <b>5</b> |

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Entropy} = - (1/6) \log_2 (1/6) - (5/6) \log_2 (5/6) = 0.65$$

|    |          |
|----|----------|
| C1 | <b>2</b> |
| C2 | <b>4</b> |

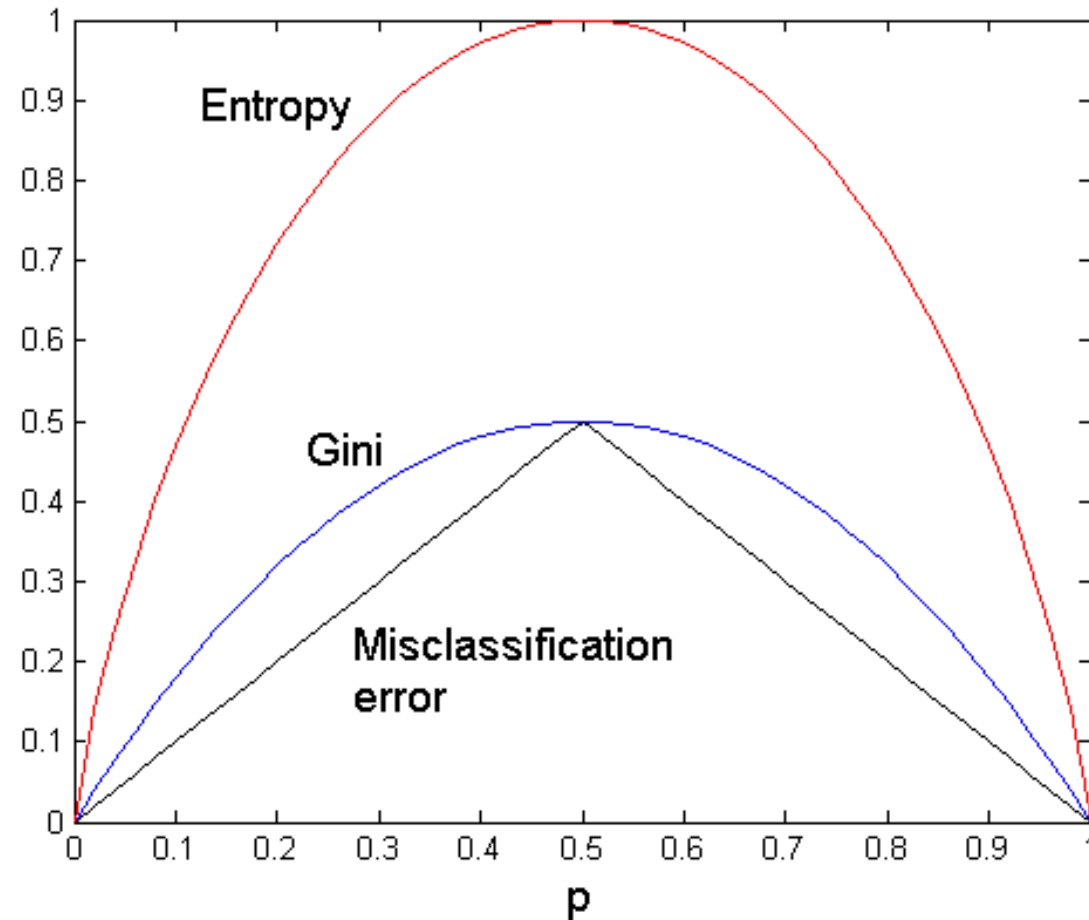
$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Entropy} = - (2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$$



# Comparing homogeneity metrics

For binary target variable ( $c=2$ )



# The goodness of a split – the gain

- How much do we gain by the split?
  - Gain ( $\Delta$ )

$$\Delta = I(\text{parent}) - \sum_{i=1}^k \frac{n_i}{n} I(\text{child}_i)$$

- $I()$ : one of the three inhomogeneity metrics
- $n_i$ : number of records in child node  $i$ ,  $n$ : number of records in the parent node
  - We weight the inhomogeneity of the child nodes by their relative size

# Problem

- The following table summarizes a data set with three binary attributes (A, B, C) and two class labels (+, -).
- Using misclassification error as inhomogeneity measure, calculate the gains for splitting on each attribute. Which attribute gives the best split?

| A | B | C | Number of instances |          |
|---|---|---|---------------------|----------|
|   |   |   | class: +            | class: - |
| T | T | T | 5                   | 0        |
| F | T | T | 0                   | 20       |
| T | F | T | 20                  | 0        |
| F | F | T | 0                   | 5        |
| T | T | F | 0                   | 0        |
| F | T | F | 25                  | 0        |
| T | F | F | 0                   | 0        |
| F | F | F | 0                   | 25       |

# Splitting by a continuous attribute

- Where to split?
- Sort the records according to the attribute values
- Calculate the inhomogeneity metric one by one increasing the value of the cut point
  - Which cut-point corresponds to the highest gain?

| Class             | No            |   | No    |   | No    |   | Yes   |   | Yes   |   | Yes   |   | No           |   | No    |   | No    |   | No    |   |       |   |
|-------------------|---------------|---|-------|---|-------|---|-------|---|-------|---|-------|---|--------------|---|-------|---|-------|---|-------|---|-------|---|
| Sorted Values →   | Annual Income |   |       |   |       |   |       |   |       |   |       |   |              |   |       |   |       |   |       |   |       |   |
|                   | 60            |   | 70    |   | 75    |   | 85    |   | 90    |   | 95    |   | 100          |   | 120   |   | 125   |   | 220   |   |       |   |
|                   | 55            |   | 65    |   | 72    |   | 80    |   | 87    |   | 92    |   | 97           |   | 110   |   | 122   |   | 172   |   | 230   |   |
| Split Positions → | <=            | > | <=    | > | <=    | > | <=    | > | <=    | > | <=    | > | <=           | > | <=    | > | <=    | > | <=    | > | <=    | > |
|                   |               |   |       |   |       |   |       |   |       |   |       |   |              |   |       |   |       |   |       |   |       |   |
| Yes               | 0             | 3 | 0     | 3 | 0     | 3 | 0     | 3 | 1     | 2 | 2     | 1 | 3            | 0 | 3     | 0 | 3     | 0 | 3     | 0 | 3     | 0 |
| No                | 0             | 7 | 1     | 6 | 2     | 5 | 3     | 4 | 3     | 4 | 3     | 4 | 3            | 4 | 4     | 3 | 5     | 2 | 6     | 1 | 7     | 0 |
| Gini              | 0.420         |   | 0.400 |   | 0.375 |   | 0.343 |   | 0.417 |   | 0.400 |   | <u>0.300</u> |   | 0.343 |   | 0.375 |   | 0.400 |   | 0.420 |   |

# Should we split on ID?

- $\Delta$  gain would be the highest if we split on ID
- The method thus prefers to split the parent node to many small nodes
- It is usually not favorable
  - E.g. if we split on ID, that is useless
  - If the emergent child nodes are too small the model is more likely to have a worse generalization ability
- Possible solutions
  - Only allow for binary splits
  - Filter out those attributes that are not reasonable to split on
  - Instead of  $\Delta$  gain use an other method to measure the goodness of split: gain ratio

# Building a decision tree

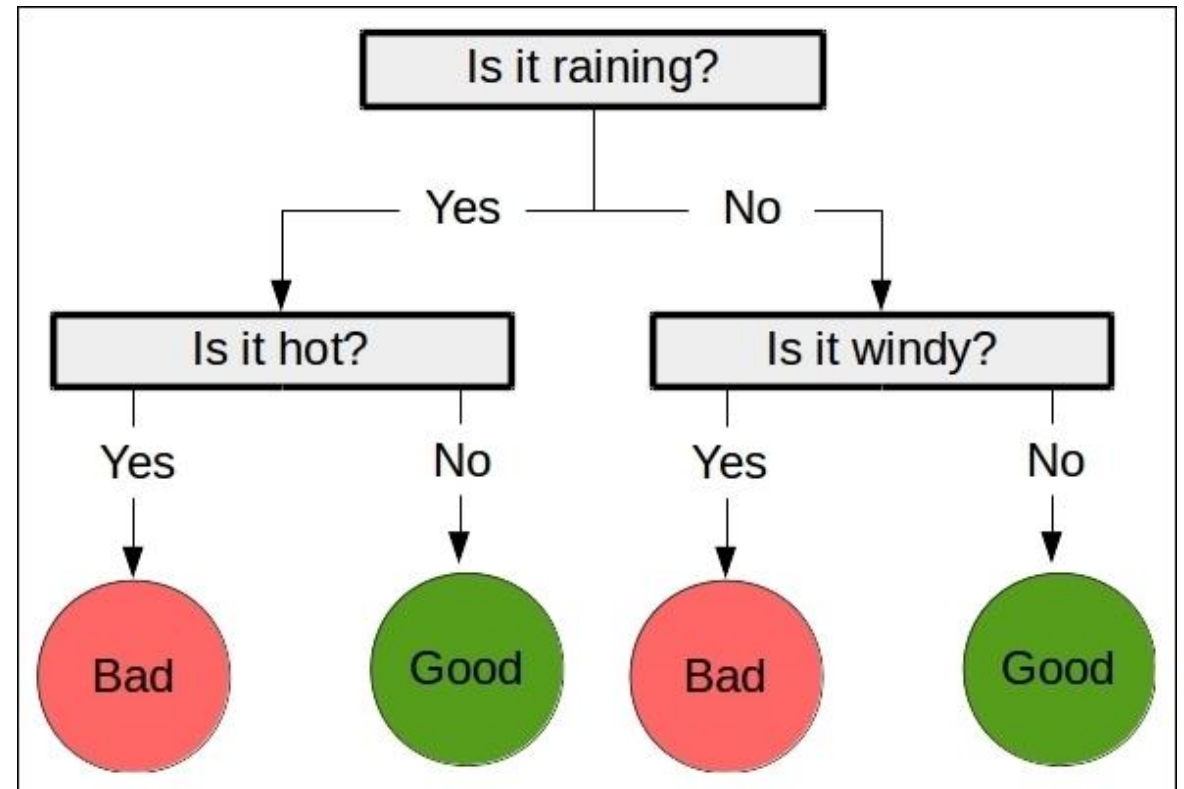
- Relatively fast, easy to interpret
- After building the decision tree, it is fast to predict new labels for unseen data points
- There are not many hyperparameter to set (but there are some)
  - What is the measure of inhomogeneity?
  - Do we allow for multi-split or just for binary?
  - How to traverse nodes?

# Termination criteria

- If every node is homogeneous
  - Or every node is quasi-homogeneous (almost homogeneous)
- All the records agree in each attribute (aside from the label)
- Global termination criteria:
  - Bound for the number of leaf nodes
  - Bound for the number of levels (depth of the tree)
- We do not search for the „best” tree
  - It is not practical since it is just the „best” on the training data set

# Advantages of decision tree

- Easy to interpret
- It works with a similar approach as human decision process
- Easy to visualize
- Insensitive to irrelevant attributes
- It can be used for a wide variety of problems (for classification/regression, with numerical/categorical attributes)





# Acknowledgement

- András Benczúr, Róbert Pálovics, SZTAKI-AIT, DM1-2
- Krisztián Buza, MTA-BME, VISZJV68
- Bálint Daróczy, SZTAKI-BME, VISZAMA01
- Judit Csimá, BME, VISZM185
- Gábor Horváth, Péter Antal, BME, VIMMD294, VIMIA313
- Lukács András, ELTE, MM1C1AB6E
- Tim Kraska, Brown University, CS195
- Dan Potter, Carsten Binnig, Eli Upfal, Brown University, CS1951A
- Erik Sudderth, Brown University, CS142
- Joe Blitzstein, Hanspeter Pfister, Verena Kaynig-Fittkau, Harvard University, CS109
- Rajan Patel, Stanford University, STAT202
- Andrew Ng, John Duchi, Stanford University, CS229

