



# Introduction to Machine Learning

## Decision Trees

---

Nikolay Manchev

6 December 2016

London Machine Learning Study Group

## **Next events**

<http://www.meetup.com/London-Machine-Learning-Study-Group>

## **Follow me**

<https://twitter.com/nikolaymanchev>

## **Slides and code**

Available at <https://github.com/nmanchev/MachineLearningStudyGroup>

**Flach talks about three types of Machine Learning models  
[Fla12]**

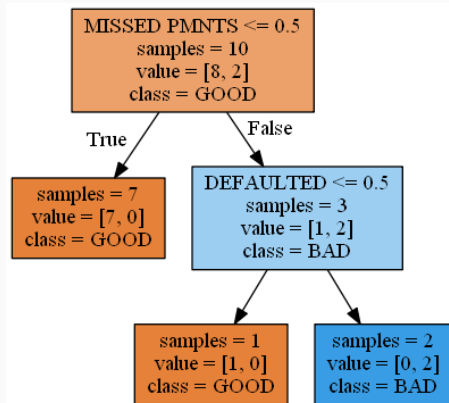
- Geometric models
- Logical models
- Statistical models

# Decision Trees

---

# Decision Tree Example

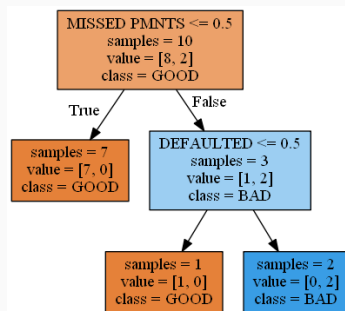
< 2 years at current job	missed pay- ments	default	rating
N	N	N	GOOD
Y	N	Y	GOOD
N	N	N	GOOD
N	N	N	GOOD
N	Y	Y	BAD
Y	N	N	GOOD
N	Y	N	GOOD
Y	Y	Y	BAD
Y	N	N	GOOD
Y	N	N	GOOD



*Credit rating example. Based on [Lew07]*

# Terminology

- **Root Node** – That is where the tree starts
- **Sub-tree** – A sub-section of the tree
- **Splitting** – The process of dividing a node into sub-nodes
- **Decision Node** – Node that splits into sub-nodes
- **Terminal Node** – Node that does not split



---

## Algorithm 1 Pseudocode for training a Decision tree

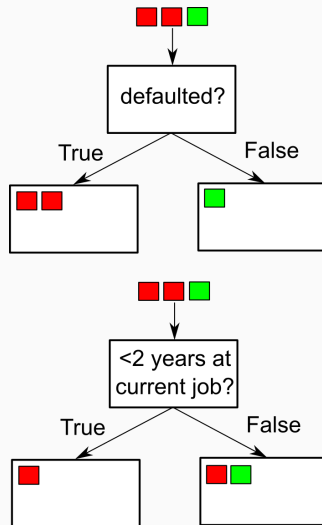
---

```
function TREE(examples, attributes, default_value)  
  if examples is empty then return a node with default_value  
  if all examples have the same class then return a node with that class  
  if attributes is empty then return MODE(examples)  
  attr  $\leftarrow$  SELECT_BEST(examples, attributes)  
  for each value  $v_i$  of attr do  
    examplesi  $\leftarrow$  {examples where attr =  $v_i$ }  
    subtree  $\leftarrow$  TREE(examplesi, attributes – attr, MODE(examplesi))  
    add a decision node  $v_i$  and subtree to tree  
  return tree
```

---

# Attribute Selection

< 2 years at current job	missed payments	default	rating
N	N	N	GOOD
Y	N	Y	GOOD
N	N	N	GOOD
N	N	N	GOOD
N	Y	Y	BAD
Y	N	N	GOOD
N	Y	N	GOOD
Y	Y	Y	BAD
Y	N	N	GOOD
Y	N	N	GOOD





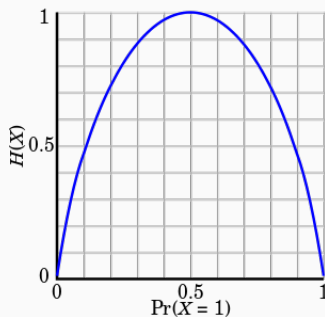
## Measure of uncertainty

- Information theory entropy  
For a discrete random variable  $X \in \{\chi_1, \dots, \chi_n\}$  and probability mass function  $P(X)$

$$H(X) = - \sum_{i=1}^n P(x_i) \log_2 P(x_i)$$

- Example – 10 samples in the credit rating data (8 GOOD, 2 BAD):

$$\begin{aligned} H(X) &= H\left(\frac{8}{10}, \frac{2}{10}\right) = \\ &= -\frac{8}{10} \log_2 \frac{8}{10} - \frac{2}{10} \log_2 \frac{2}{10} \approx 0.72 \end{aligned}$$



Binary entropy plot [BD]

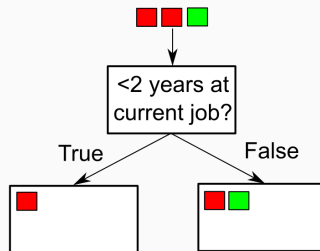
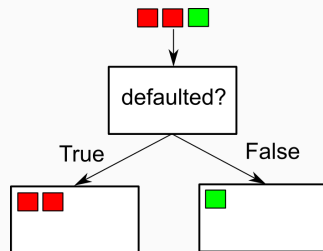
## Measure of uncertainty

- For the (defaulted? = True) terminal node:

$$H\left(\frac{2}{2}\right) = -\frac{2}{2}\log_2\frac{2}{2} = 0$$

- For the (<2 years? = False) terminal node:


$$H\left(\frac{1}{2}, \frac{1}{2}\right) = -\frac{1}{2}\log_2\frac{1}{2} - \frac{1}{2}\log_2\frac{1}{2} = \frac{1}{2} + \frac{1}{2} = 1$$



## UCI Machine Learning Repository –

[archive.ics.uci.edu/ml](http://archive.ics.uci.edu/ml)


- Great resource for Machine Learning data sets
- Over 330 freely available sets
- Auto MPG Data Set
  - Fuel consumption in MPG
  - Attributes: mpg, cylinders, displacement, horsepower, weight, acceleration etc.



**Machine Learning Repository**  
Center for Machine Learning and Intelligent Systems

### Auto MPG Data Set

Download: [Data Folder](#), [Data Set Description](#)



**Abstract:** Revised from CMU StatLib library, data concerns city-cycle fuel consumption

<b>Data Set Characteristics:</b>	Multivariate	<b>Number of Instances:</b>	398	<b>Area:</b>	N/A
<b>Attribute Characteristics:</b>	Categorical, Real	<b>Number of Attributes:</b>	8	<b>Date Donated</b>	1993-07-07
<b>Associated Tasks:</b>	Regression	<b>Missing Values?</b>	Yes	<b>Number of Web Hits:</b>	167833

## Entropy of a split

We can compute the weighted average over all sets in the split

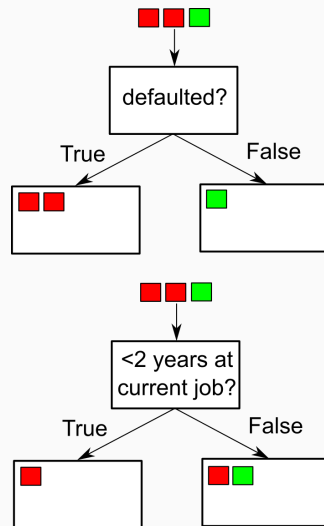
$$I(X, A) = \sum_{i=1}^m \frac{|X_i|}{|X|} \times H(X_i)$$

where  $m$  is the number of distinct values in  $A$ ,  $|X|$  is the size of  $X$ , and  $|X_i|$  is the size of  $X_i$

**Example:**

$$I(X, \text{defaulted}) = \frac{2}{3} \times 0 + \frac{1}{3} \times 0 = 0$$

$$I(X, <2 \text{ years?}) = \frac{1}{3} \times 0 + \frac{2}{3} \times 1 = \frac{2}{3}$$



# Selecting the best attribute

## Information Gain

Expected decrease of entropy after splitting on an attribute

$$IG(X, A) = H(X) - I(X, A)$$

**Example:**

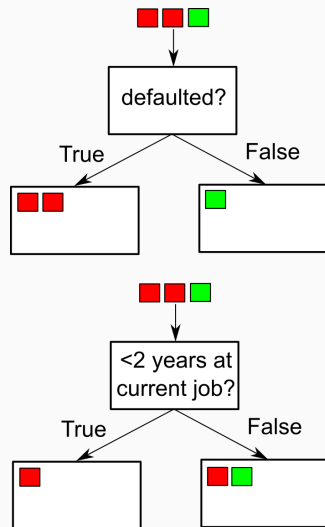
$$H(X) = H\left(\frac{2}{3}, \frac{1}{3}\right) = -\frac{2}{3}\log_2\frac{2}{3} - \frac{1}{3}\log_2\frac{1}{3} = 0.92$$

$$IG(X, \text{defaulted}) = H\left(\frac{2}{3}, \frac{1}{3}\right) -$$

$$I(X, \text{defaulted}) = 0.92 - 0 = 0$$

$$I(X, <2 \text{ years?}) = H\left(\frac{2}{3}, \frac{1}{3}\right) -$$

$$I(X, <2 \text{ years?}) = 0.92 - \frac{2}{3} = 0.26$$



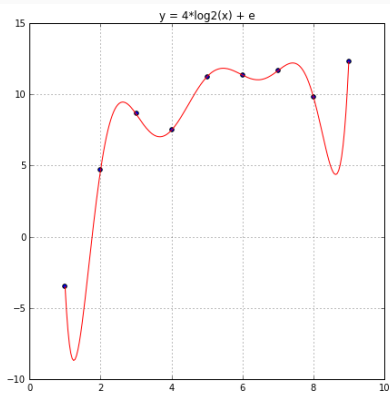
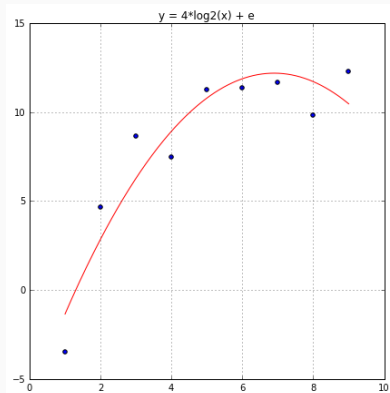
## Pros

- Interpretability (easy to understand)
- Mixed data type – numerical and categorical variables in the same model
- Less data preparation

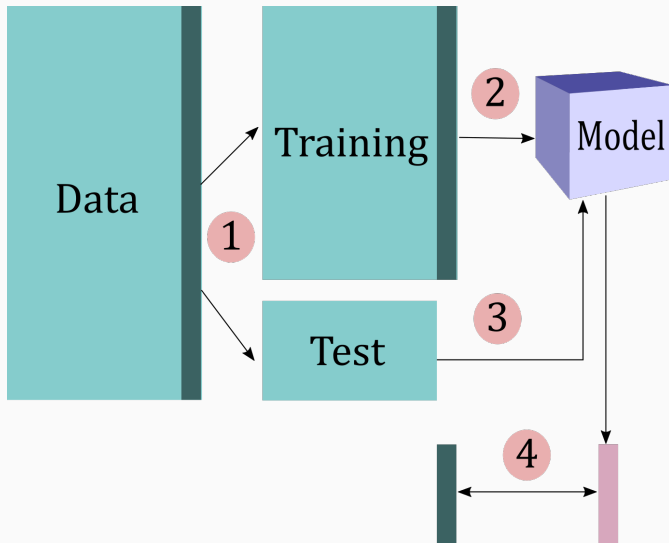
## Cons

- They tend to get complex and overfit

# Overfitting



# Holdout method





- An attribute with high cardinality is typically considered a good candidate for a split
  - Worst case would be a unique identifier – splitting on it produces pure nodes (one example per node), but such tree would be useless for prediction
  - One way to avoid this is by using different metrics (e.g. Gain Ratio, which captures the ratio of information gain to the intrinsic information)

## Reducing complexity

- Other selection criteria
  - Gain Ratio
  - Chi-Square
- Restriction on splits
  - Binary and multi-way splits
  - Minimum samples to split
- Early stopping
- Pruning
  - Reduced Error Pruning

## Pros

- Interpretability (easy to understand)
- Mixed data type – numerical and categorical variables in the same model
- Less data preparation

## Cons

- They tend to get complex and overfit
- Instability
- Inadequate for predicting continuous values




# Decision Trees vs. Linear Regression

## Decision Trees

- Can solve both classification and regression problems
- If the relationship is non-linear it will outperform Linear Regression
- Can build models that are easy to explain

## Linear Regression

- In a linear relationship Linear Regression will likely outperform Decision Trees
- Cannot easily handle categorical variables

-  Brona and A. Damato, *Binary entropy plot*.
-  Peter Flach, *Machine learning: The art and science of algorithms that make sense of data*, Cambridge University Press, New York, NY, USA, 2012.
-  Michael S. Lewicki, *Artificial intelligence: Learning and decision trees*, Carnegie Mellon, 2007.