

## EDAN20

## Language Technology

<http://cs.lth.se/edan20/>

## Chapter 4: Topics in Information Theory and Machine Learning

Pierre Nugues

Lund University

[Pierre.Nugues@cs.lth.se](mailto:Pierre.Nugues@cs.lth.se)

[http://cs.lth.se/pierre\\_nugues/](http://cs.lth.se/pierre_nugues/)

September 3, 2015



# Entropy

Information theory models a text as a sequence of symbols.

Let  $x_1, x_2, \dots, x_N$  be a discrete set of  $N$  symbols representing the characters.

The information content of a symbol is defined as

$$I(x_i) = -\log_2 p(x_i) = \log_2 \frac{1}{p(x_i)},$$

Entropy, the average information content, is defined as:

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x),$$

By convention:  $0 \log_2 0 = 0$ .



# Entropy of a Text

The entropy of the text is

$$\begin{aligned} H(X) &= - \sum_{x \in X} p(x) \log_2 p(x). \\ &= -p(A) \log_2 p(A) - p(B) \log_2 p(B) - \dots \\ &\quad - p(Z) \log_2 p(Z) - p(\grave{A}) \log_2 p(\grave{A}) - \dots \\ &\quad - p(\ddot{Y}) \log_2 p(\ddot{Y}) - p(\text{blanks}) \log_2 p(\text{blanks}). \end{aligned}$$

Entropy of Gustave Flaubert's *Salammbô* in French is  $H(X) = 4.39$ .



# Cross-Entropy

The cross entropy of  $m$  on  $p$  is defined as:

$$H(p, m) = - \sum_{x \in X} p(x) \log_2 m(x).$$

We have the inequality  $H(p) \leq H(p, m)$ .

	Entropy	Cross entropy	Difference
<i>Salammbô</i> , chapters 1-14, training set	4.39481	4.39481	0.0
<i>Salammbô</i> , chapter 15, test set	4.34937	4.36074	0.01137
<i>Notre Dame de Paris</i> , test set	4.43696	4.45507	0.01811
<i>Nineteen Eighty-Four</i> , test set	4.35922	4.82012	0.46090



# Entropy, Decision Trees, and Classification

Decision trees are useful devices to classify objects into a set of classes. Entropy can help us learn automatically decision trees from a set of data. The algorithm is one of the simplest machine-learning techniques to obtain a classifier.

There are many other machine-learning algorithms, which can be classified along two lines: supervised and unsupervised

Supervised algorithms need a training set.

Their performance is measured against a test set.

We can also use  $N$ -fold cross validation, where the test set is selected randomly from the training set  $N$  times, usually 10.

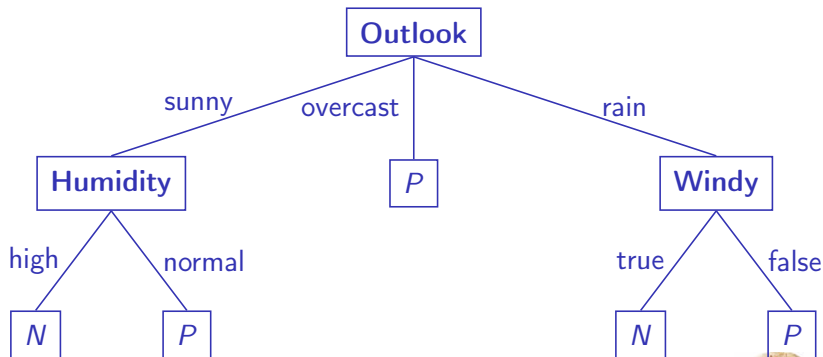


# Objects, Classes, and Attributes. After Quinlan (1986)

Object	Attributes				Class
	Outlook	Temperature	Humidity	Windy	
1	Sunny	Hot	High	False	<i>N</i>
2	Sunny	Hot	High	True	<i>N</i>
3	Overcast	Hot	High	False	<i>P</i>
4	Rain	Mild	High	False	<i>P</i>
5	Rain	Cool	Normal	False	<i>P</i>
6	Rain	Cool	Normal	True	<i>N</i>
7	Overcast	Cool	Normal	True	<i>P</i>
8	Sunny	Mild	High	False	<i>N</i>
9	Sunny	Cool	Normal	False	<i>P</i>
10	Rain	Mild	Normal	False	<i>P</i>
11	Sunny	Mild	Normal	True	<i>P</i>
12	Overcast	Mild	High	True	<i>P</i>
13	Overcast	Hot	Normal	False	<i>P</i>
14	Rain	Mild	High	True	<i>N</i>



# Classifying Objects with Decision Trees. After Quinlan (1986)



# Decision Trees and Classification

Each object is defined by an attribute vector (or feature vector)

$\{A_1, A_2, \dots, A_v\}$

Each object belongs to a class  $\{C_1, C_2, \dots, C_n\}$

The attributes of the examples are:

$\{Outlook, Temperature, Humidity, Windy\}$  and the classes are:  $\{N, P\}$ .

The nodes of the tree are the attributes.

Each attribute has a set of possible values. The values of Outlook are

$\{Sunny, Rain, Overcast\}$

The branches correspond to the values of each attribute

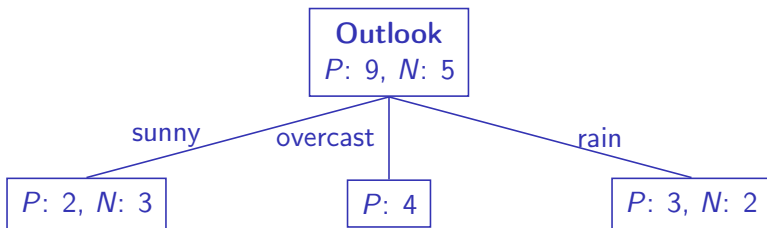
The optimal tree corresponds to a minimal number of tests.





# ID3 (Quinlan, 1986)

Each attribute scatters the set into as many subsets as there are values for this attribute.



At each decision point, the “best” attribute has the maximal separation power, the maximal information gain



# ID3 (Quinlan, 1986)

The entropy of a two-class set  $p$  and  $n$  is:

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}.$$

The weighted average of all the nodes below an attribute is:

$$\sum_{i=1}^v \frac{p_i + n_i}{p+n} I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right).$$

The information gain is defined as  $I_{\text{before}} - I_{\text{after}}$



# Example

$$I_{\text{before}}(p, n) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940.$$

Outlook has three values: sunny, overcast, and rain.

$$I(p_1, n_1) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971.$$

$$I(p_2, n_2) = 0.$$

$$I(p_3, n_3) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971.$$

The gain is  $0.940 - 0.694 = 0.246$ , the highest possible among the attributes



# Other Supervised Machine-Learning Algorithms

Linear classifiers:

- 1 Perceptron
- 2 Logistic regression
- 3 Support vector machines



# The Weka Toolkit

Weka: A powerful collection of machine-learning algorithms

<http://www.cs.waikato.ac.nz/ml/weka/>.

**Weka Explorer**

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose None Apply

Current relation: Relation: weather.symbolic Attributes: 5 Instances: 14 Sum of weights: 14

Attributes: All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> outlook
2	<input type="checkbox"/> temperature
3	<input type="checkbox"/> humidity
4	<input type="checkbox"/> windy
5	<input type="checkbox"/> play

Remove

Selected attribute: Name: outlook Missing: 0 (0%) Type: Nominal Distinct: 3 Unique: 0 (0%)

No.	Label	Count	Weight
1	sunny	5	5.0
2	overcast	4	4.0
3	rainy	5	5.0

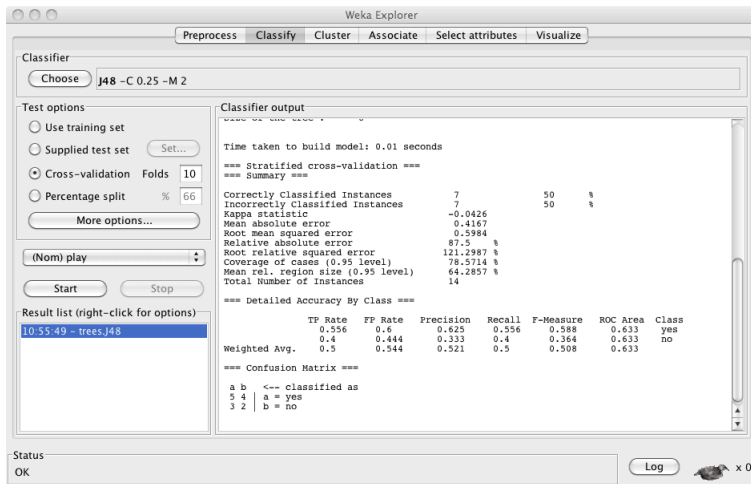
Class: play (Nom) Visualize All

Status: OK Log x 0



# The Weka Toolkit

## Running ID3



**Weka Explorer**

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose **J48 -C 0.25 -M 2**

Test options:

- ☐ Use training set
- ☐ Supplied test set (Set...)
- ☒ Cross-validation Folds **10**
- ☐ Percentage split % **66**

More options...

(Nom) play

Start Stop

Result list (right-click for options):

- 10:55:49 - trees.J48

**Classifier output**

Time taken to build model: 0.01 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	7	50	%
Incorrectly Classified Instances	7	50	%
Kappa statistic	-0.0426		
Mean absolute error	0.4167		
Root mean squared error	0.5984		
Relative absolute error	87.5	%	
Root relative squared error	121.2987	%	
Coverage of cases (0.95 level)	78.5714	%	
Mean rel. region size (0.95 level)	64.2857	%	
Total Number of Instances	14		

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
Weighted Avg.	0.556	0.444	0.333	0.4	0.364	0.633	yes
	0.5	0.544	0.521	0.5	0.508	0.633	no

=== Confusion Matrix ===

a \ b	<-- classified as	
	a = yes	b = no
5 4		
3 2		

Status: OK

Log x 0

# ARFF: The Weka Data Format

Storing Quinlan's data set in Weka's attribute-relation file format (ARFF)

<http://weka.wikispaces.com/ARFF>:

```
@relation weather.symbolic
```

```
@attribute outlook {sunny, overcast, rainy}
```

```
@attribute temperature {hot, mild, cool}
```

```
@attribute humidity {high, normal}
```

```
@attribute windy {TRUE, FALSE}
```

```
@attribute play {yes, no}
```

```
@data
```

```
sunny,hot,high,FALSE,no
```

```
sunny,hot,high,TRUE,no
```

```
overcast,hot,high,FALSE,yes
```

```
rainy,mild,high,FALSE,yes
```

```
rainy,cool,normal,FALSE,yes
```

```
rainy,cool,normal,TRUE,no
```

```
overcast,cool,normal,TRUE,yes
```

