

Roma

20-23.03.2013

www.codemotionworld.com



Simone Scardapane

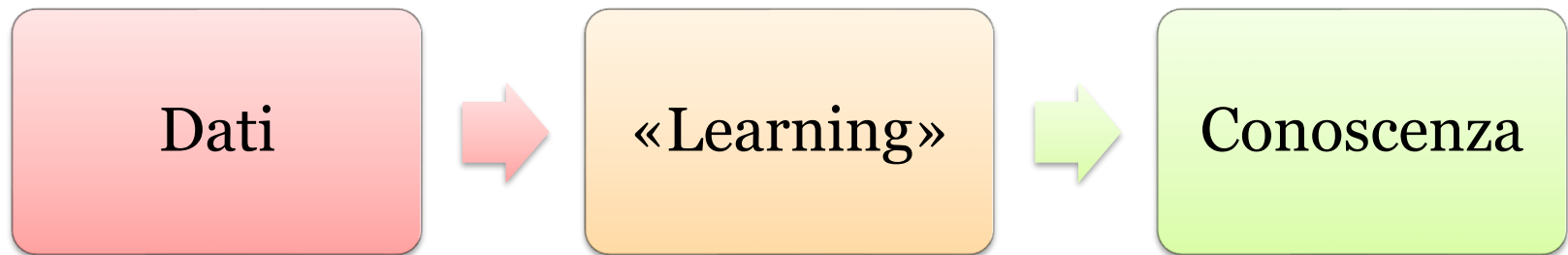
*How to make smarter programs.
A gentle introduction to Machine Learning*

simone.scardapane@uniroma1.it



1. Il Machine Learning (ieri ed oggi)
2. Un esempio pratico: *spam detection*
3. Cenni su Weka

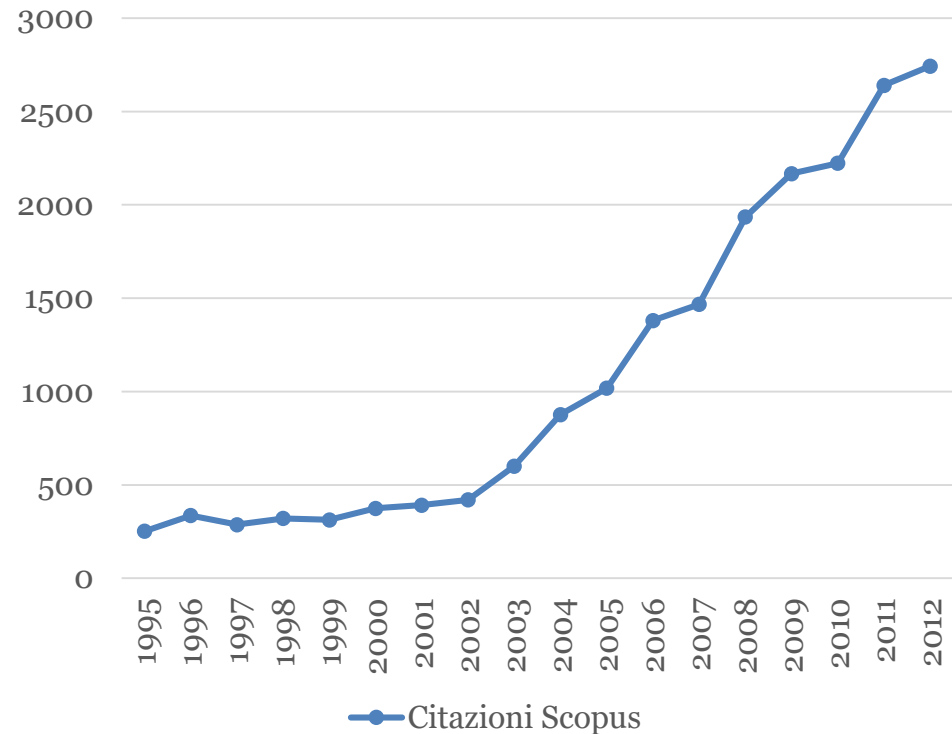
«Estrazione *automatica* di conoscenza a partire da un insieme di dati»



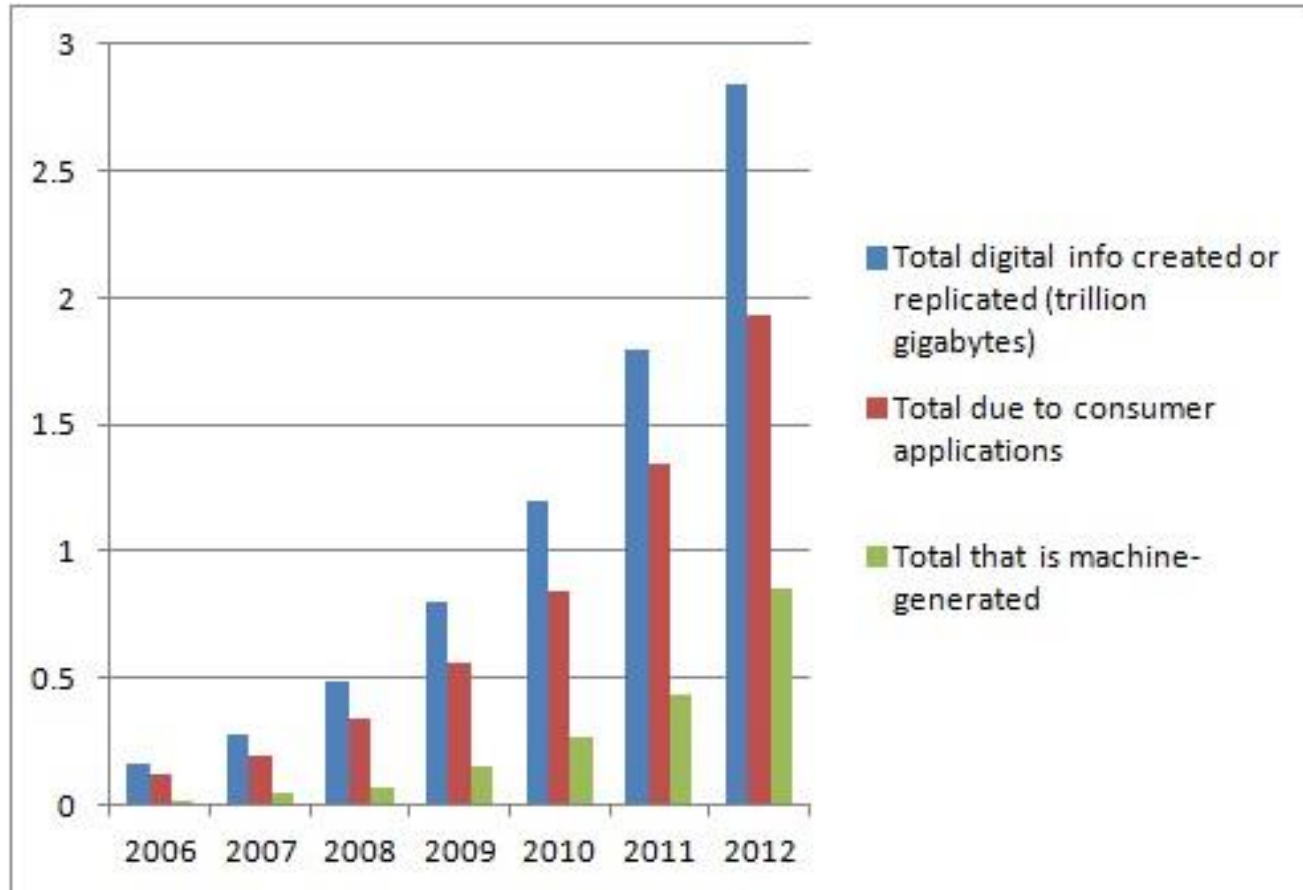
65 anni di ricerche da parte di:

- Ingegneri
- Informatici
- Statistici
- Matematici
- Fisici
- Neuroscienziati...

Citazioni ML



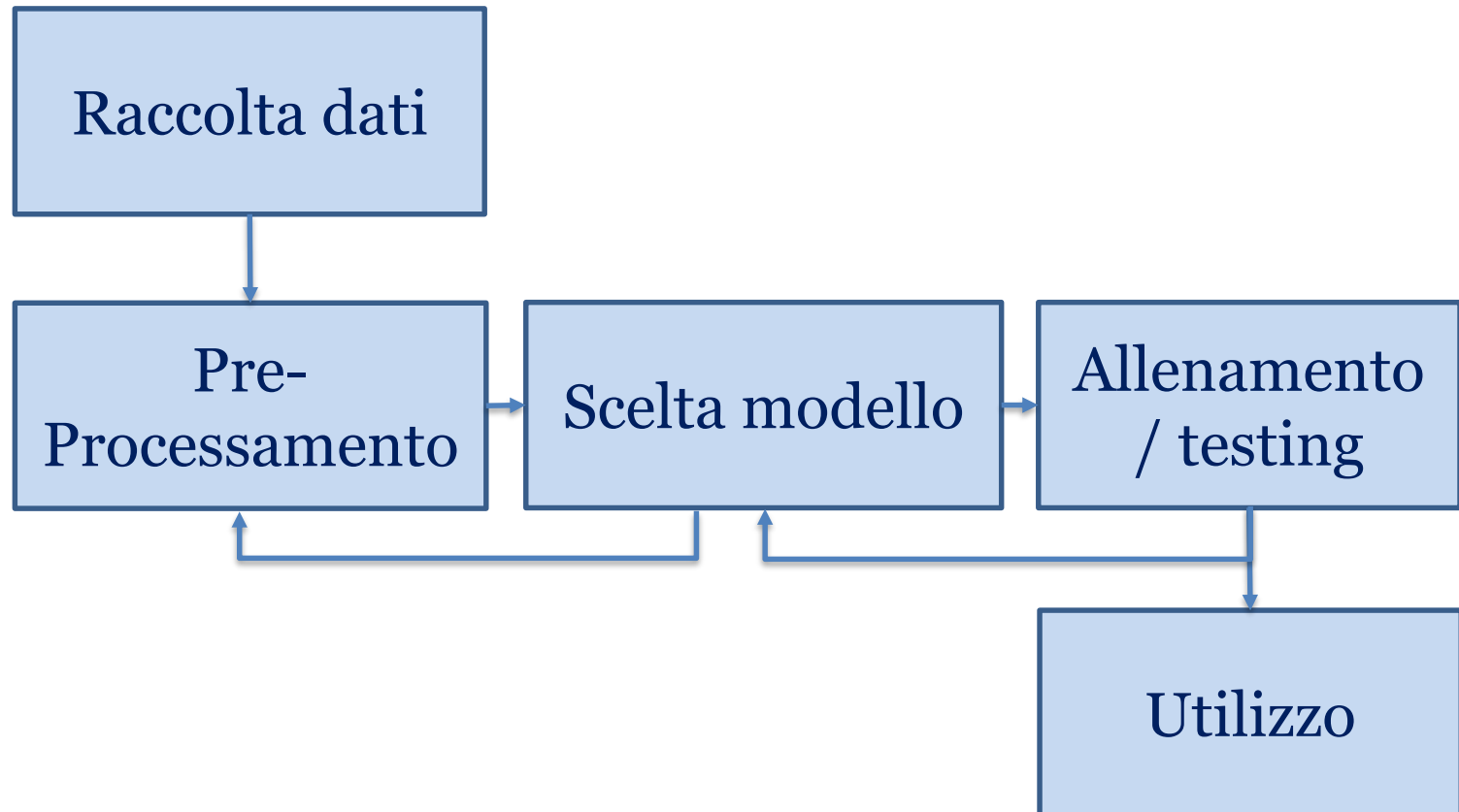




(Fonte: IDC)

Possibili operazioni su una libreria musicale:

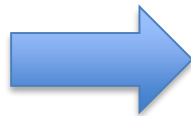
1. **Classificazione** del genere (o del *mood*)
2. Raggruppamento automatico (**clustering**)
3. **Tagging**
4. Ricerca per similitudini (**association rule**)
5. **Predizione** del prossimo ascolto
6. ...



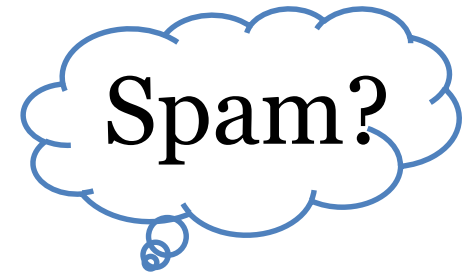
- **Online Learning**: l'algoritmo riceve dati in real-time e si adatta di conseguenza.
- **Active Learning**: durante la fase di learning, è possibile richiedere attivamente nuove informazioni.
- **Collaborative/Cooperative Learning...**

- Dati: insieme S di emails taggate come spam / non spam.
- Obiettivo: metodo automatico per individuare spam.
- Problemi:
 1. Come rappresentare l'email?
 2. Che modello utilizzare?
 3. Come allenarlo?

Email

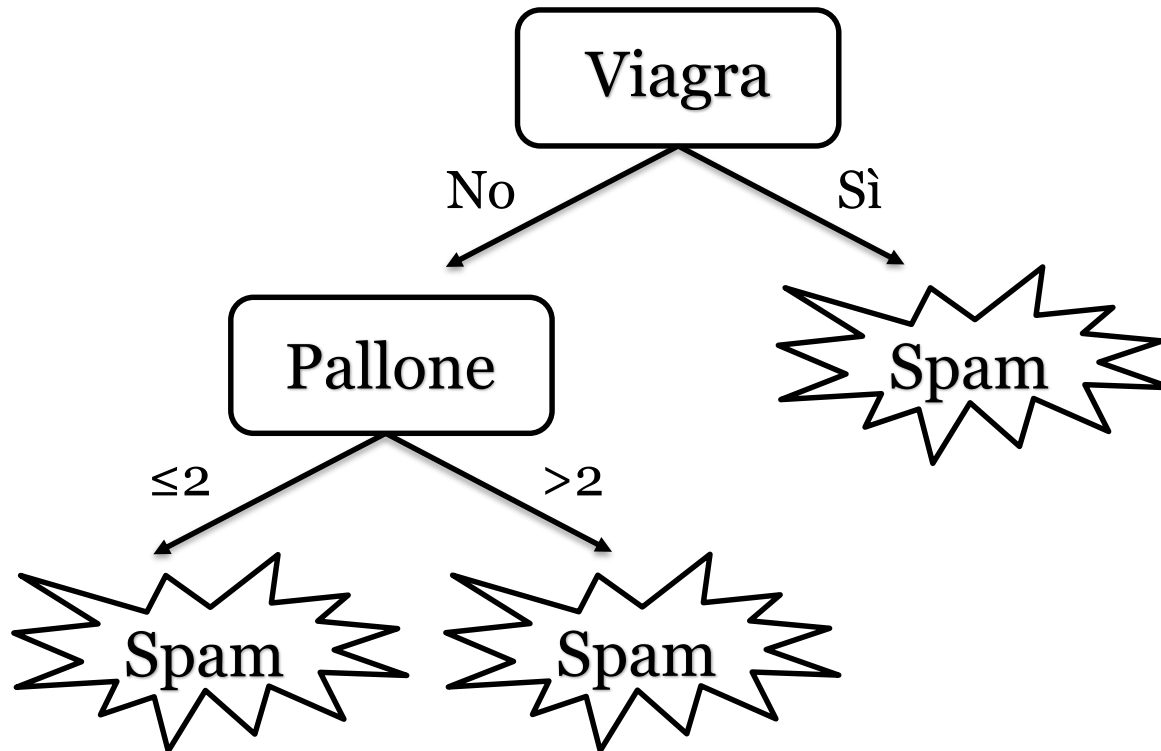


Parola	#
Viagra	2
Bambino	5
Macchina	0
Stereo	0
Cane	1
...	

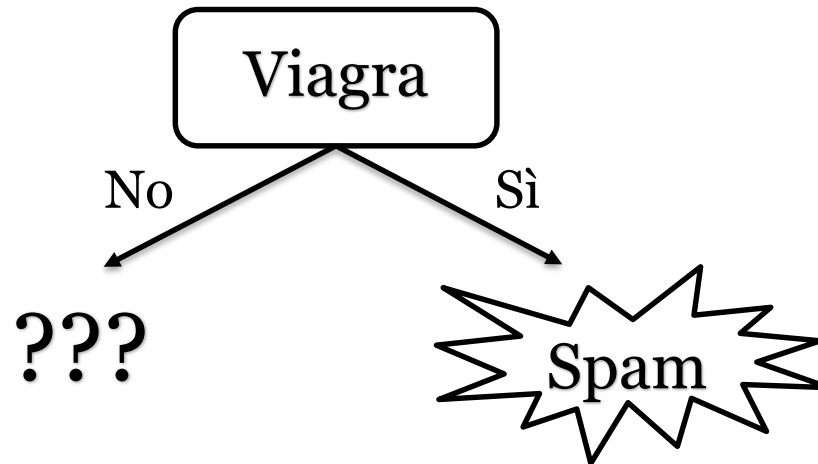


«Bag of words»

Decision Tree:



Come costruirlo?



Consideriamo l'algoritmo C4.5:

1. Scegliamo per il nodo l'attributo a che «divide» meglio i dati.
2. Suddividiamo l'insieme lungo i nodi.
3. Ci fermiamo quando i dati sono perfettamente divisi.

(Difficoltà: gestire dati continui, mancanti...)

Matematicamente, scegliamo l'attributo che riduce maggiormente l'*entropia*:

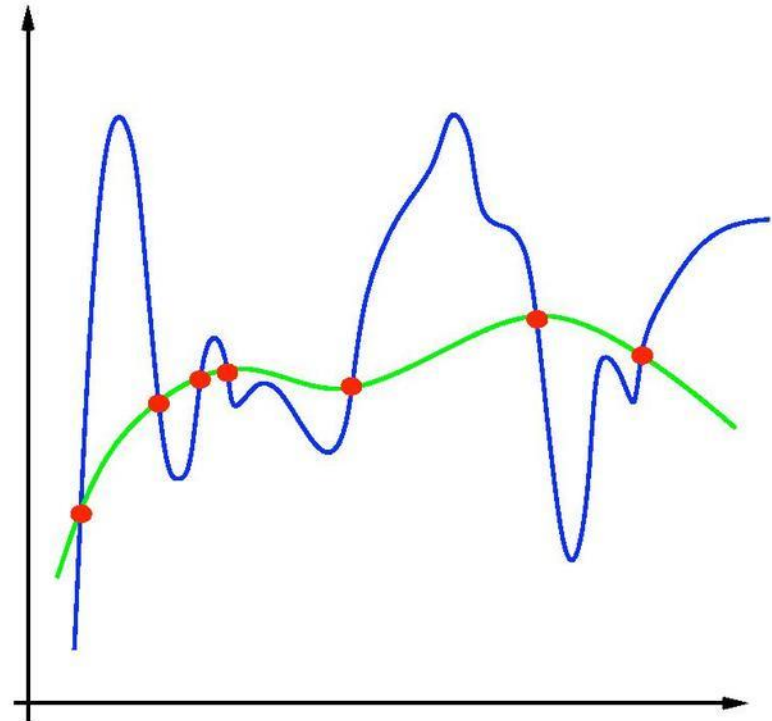
$$E(S) = -p_{spam} \log p_{spam} - p_{\neg spam} \log p_{\neg spam}$$

(In realtà ci basiamo su una misura correlata)

Il tema di *massimizzare* un criterio è cruciale nel ML.

Problema principale:
overfitting!

(Immagine con Copyright
Tomaso Poggio)



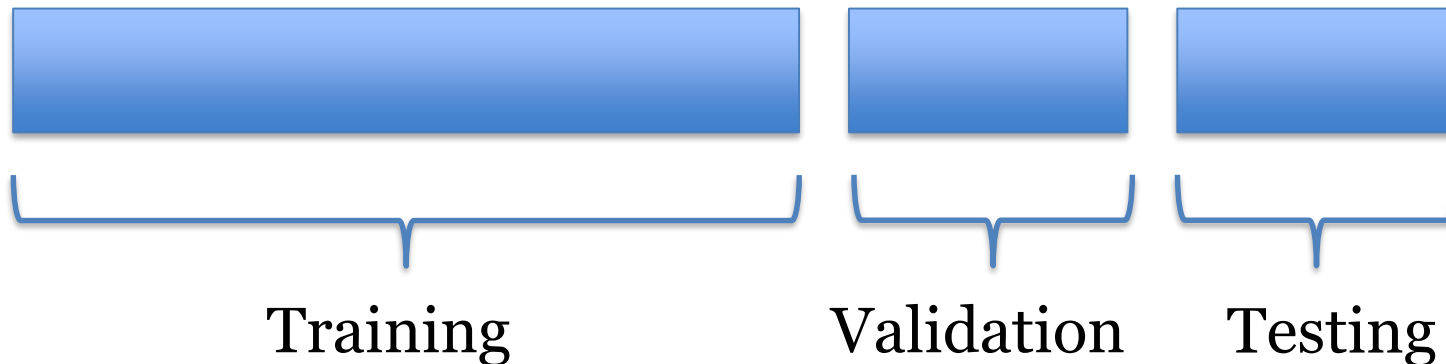
Possibile soluzione (per i decision trees):

- Si tiene da parte un insieme di dati.
- Si eliminano i rami non necessari (*pruning*) in base a quei dati (*error-based pruning*).

Più generalmente si usano tecniche di *cross-validation*.

Possiamo tenere da parte un secondo insieme per *testare* l'accuratezza dell'algoritmo.

Dividiamo quindi i nostri dati in tre parti:





Tool di *data mining* sviluppato dalla Waikato University in Java:

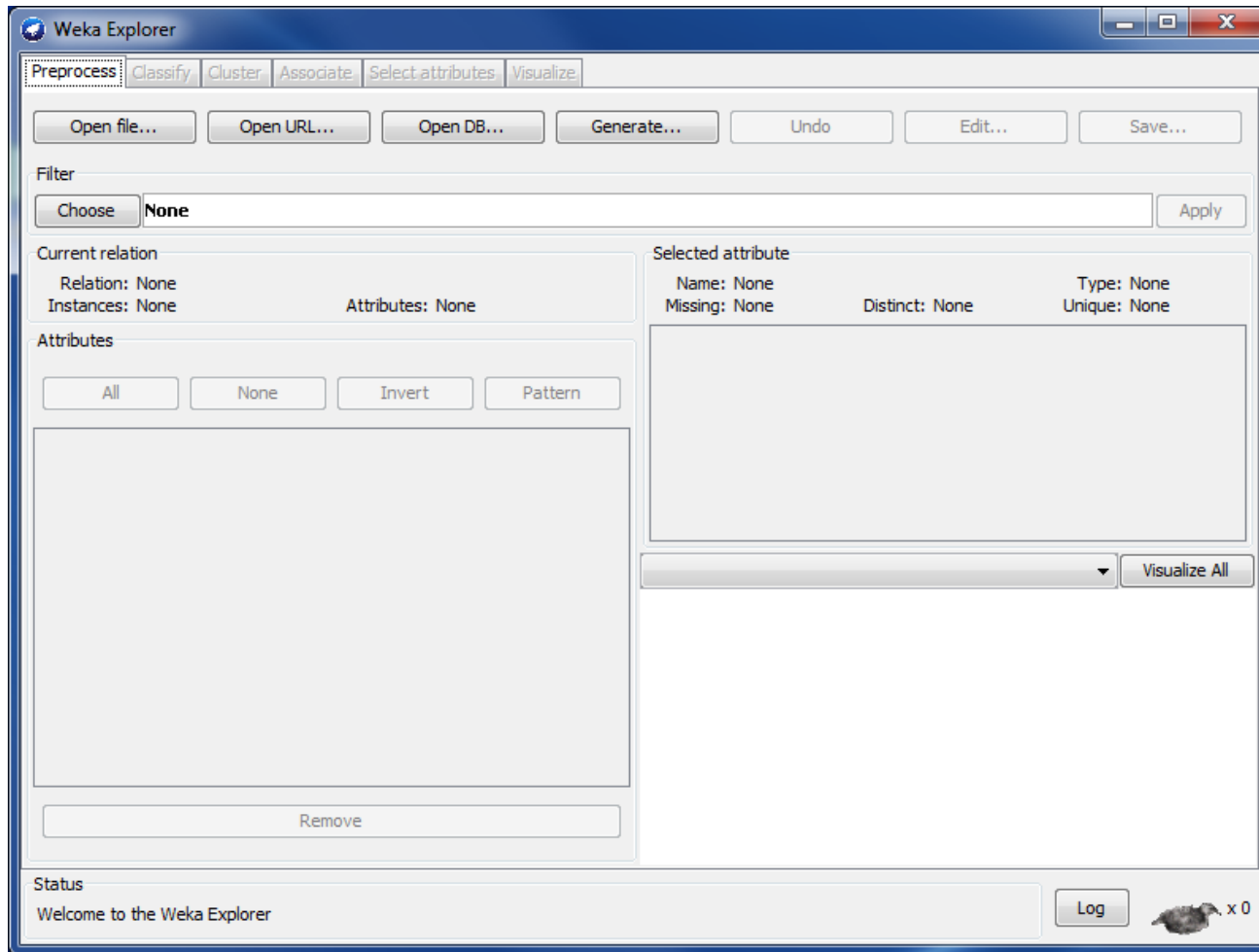
1. Ampio numero di funzioni.
2. Estremamente portabile.
3. Interfaccia di facile utilizzo.

Usiamo il dataset *SpamBase* dal repository UCI:
<http://archive.ics.uci.edu/ml/datasets/Spambase>

4601 email rappresentate da 48 frequenze di parole (più qualche informazione aggiuntiva).

I dati sono salvati in formato ARFF (file di testo):

1. Header con descrizione degli attributi.
2. Elenco delle email.



Apriamo il file:

Preprocess **Classify** Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate...

Current relation
Relation: Spambase
Instances: 4601
Attributes: 58

Attributes

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> word_freq_make
2	<input type="checkbox"/> word_freq_address
3	<input type="checkbox"/> word_freq_all
4	<input type="checkbox"/> word_freq_3d
5	<input type="checkbox"/> word_freq_our
6	<input type="checkbox"/> word_freq_over
7	<input type="checkbox"/> word_freq_remove
8	<input type="checkbox"/> word_freq_internet
9	<input type="checkbox"/> word_freq_order
10	<input type="checkbox"/> word_freq_mail
11	<input type="checkbox"/> word_freq_receive
12	<input type="checkbox"/> word_freq_will
13	<input type="checkbox"/> word_freq_people

Remove

Selected attribute
Name: word_freq_make
Missing: 0 (0%)
Distinct: 142
Type: Numeric
Unique: 31 (1%)

Statistic	Value
Minimum	0
Maximum	4.54
Mean	0.105
StdDev	0.305

Class: is_spam (Nom) Visualize All

0 2.27 4.54

Scegliamo il modello:

Classifier

J48 -R -N 3 -Q 1 -M 2

weka.classifiers.trees.J48

About

Class for generating a pruned or unpruned C4.

binarySplits

confidenceFactor

debug

minNumObj

numFolds

Risultati:

Time taken to build model: 0.42 seconds

=== Evaluation on test split ===

=== Summary ===

Correctly Classified Instances	836	90.8696 %
Incorrectly Classified Instances	84	9.1304 %
Kappa statistic	0.8105	
Mean absolute error	0.1235	
Root mean squared error	0.2763	
Relative absolute error	25.7409 %	
Root relative squared error	56.2152 %	
Total Number of Instances	920	


```
import weka.core.converters.ConverterUtils.DataSource;
...
DataSource source = new DataSource("/some/where/data.arff");
Instances data = source.getDataSet();
// setting class attribute if the data format does not provide this information
// For example, the XRFF format saves the class attribute information as well
if (data.classIndex() == -1)
    data.setClassIndex(data.numAttributes() - 1);
```

```
import weka.classifiers.trees.J48;
...
String[] options = new String[1];
options[0] = "-U"; // unpruned tree
J48 tree = new J48(); // new instance of tree
tree.setOptions(options); // set the options
tree.buildClassifier(data); // build classifier
```

<http://weka.wikispaces.com/Use+WEKA+in+your+Java+code>

Programming Collective Intelligence, di **Toby Segaran**. Publisher: O'Reilly Media (2007).

Data Mining: Practical Machine Learning Tools and Techniques, di **Witten, Frank et Hall**. Publisher: Morgan Kaufmann (2011).

Introduction to Machine Learning, di **Alpaydin**. Publisher: the MIT Press (2009).

Roma

20-23.03.2013

www.codemotionworld.com

Fine!

Simone Scardapane

simone.scardapane@uniroma1.it

Grazie per l'attenzione! 😊