

Bayesian Learning: An Introduction

João Gama

LIAAD-INESC TEC, University of Porto, Portugal

jgama@fep.up.pt

September 2020

- 1 Motivation
- 2 Introduction
- 3 Bayesian Network Classifiers
- 4 k-Dependence Bayesian Classifiers
- 5 Links and References

Motivation: Supervised Learning Task

- Given:
 - Historical data about clients credit operations
 - Characteristics of the clients
 - Characteristics of the credit application
 - The success of the operation constitutes the class
- Goal: Learn a model that accurately predicts the success (or not) of new credit applications

Supervision: Add Class Values

Client	amount	age	salary	account	Loan
Client 1	medium	junior	low	yes	no
Client 2	medium	junior	low	no	no
Client 3	low	junior	low	yes	yes
Client 4	high	medium	low	yes	yes
Client 5	high	senior	high	yes	yes
Client 6	high	senior	high	no	no

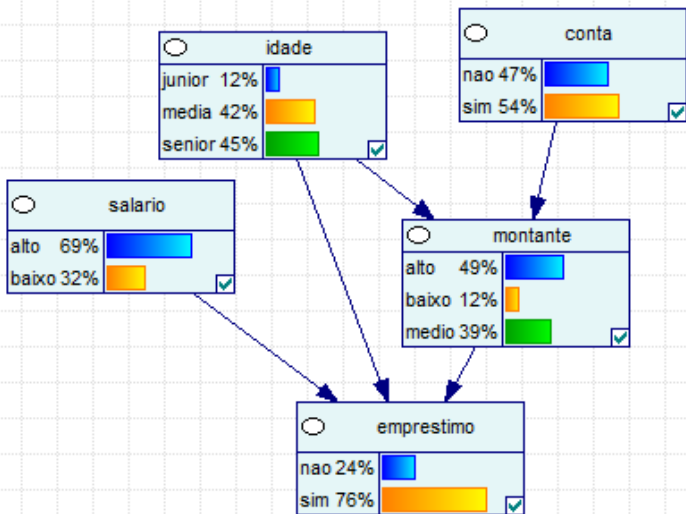
Learning Problems:

- Find a function:

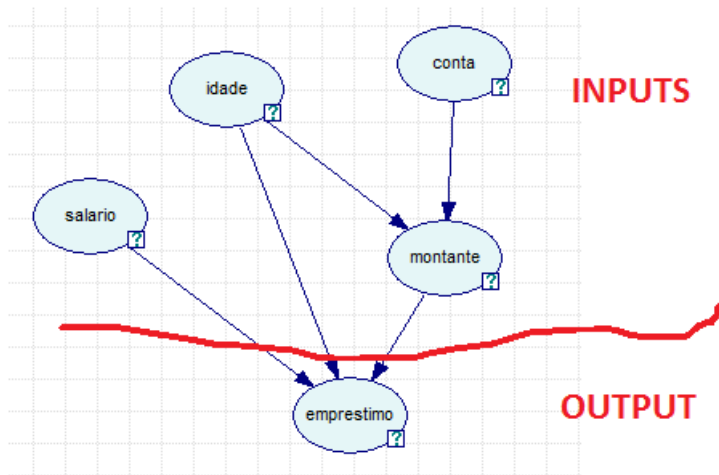
$$\text{Loan} = f(\text{amount}, \text{age}, \text{salary}, \text{account})$$
- Given the characteristics of a client, predict if the application will succeed or not.

... and Quantitative model

A set of contingency tables between dependent variables.

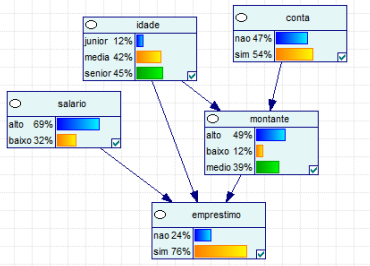


Causal Model: inputs and outputs.

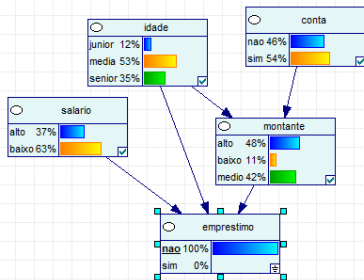


Propagating Evidence ...

A Bayesian net:



Diagnosis: Observing Outputs and propagating this evidence



Why Learn Bayesian Networks?

- Join Probability Distribution of a set of Variables.
- Conditional independences & graphical language capture structure of many real-world distributions
- Graph structure provides much insight into domain
- Learned model can be used for many tasks:
 - **Prediction:** Given the Inputs: which Outputs?
 - **Diagnosis:** Given the Outputs: which Inputs?
 - **Unsupervised:** Given Inputs and Outputs: Which structure?

Outline

- 1 Motivation
- 2 Introduction**
- 3 Bayesian Network Classifiers
- 4 k-Dependence Bayesian Classifiers
- 5 Links and References

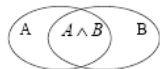
Basics on Statistics

Axioms:

- All probabilities between 0 and 1
 $0 \leq P(A) \leq 1$
- True proposition has probability 1, false has probability 0.
 $P(true) = 1$ and $P(false) = 0$
- The probability of a disjunction is:
 $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

Conditional Probability:

- $P(A|B)$ is the probability of A given B
 Assumes that B is all and only information known.
- Defined by: $P(A|B) = \frac{P(A \wedge B)}{P(B)}$



Bayes Theorem

$$P(h|D) = \frac{P(h)P(D|h)}{P(D)}$$

Combine the prior distribution and the likelihood of the observed data in order to derive the posterior distribution.

Proof:

$$P(h|D) = \frac{P(h \wedge D)}{P(D)}$$

$$P(D|h) = \frac{P(h \wedge D)}{P(h)}$$

$$P(h \wedge D) = P(D|h)P(h)$$

$$P(h|D) = \frac{P(h)P(D|h)}{P(D)}$$

Bayes Theorem

What is the most probable hypothesis h , given training data D ?

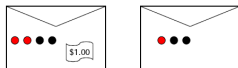
$$P(h|D) = \frac{P(h)P(D|h)}{P(D)}$$

Computing the probability of a hypothesis based on:

- its prior probability,
- the probability of observing the data given the hypothesis,
- The data itself

- $P(h)$ = prior probability of hypothesis h
- $P(D)$ = prior probability of training data D
- $P(h|D)$ = probability of h given D
- $P(D|h)$ = probability of D given h

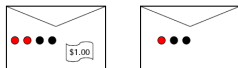
Illustrative Example



- The Win envelope has 1 dollar and four beads in it.
- The Lose envelope has three beads in it.

Someone draws an envelope at random and offers to sell it to you. Which one should we choose?

Illustrative Example



Before deciding, you are allowed to see one bead in one envelope.

- If it is black, Which one should we choose?
- And if it is red?



Illustrative Example

Prior Probabilities:

$$P(\textit{Win}) =$$

$$P(\textit{Lose}) =$$

$$P(\textit{red}) =$$

$$P(\textit{black}) =$$

$$P(\textit{black}|\textit{Win}) =$$

$$P(\textit{red}|\textit{Win}) =$$

After seeing the bead:

$$P(\textit{Win}|\textit{black}) =$$

$$P(\textit{Win}|\textit{red}) =$$

Illustrative Example

Prior Probabilities:

$$P(\textit{Win}) = 1/2$$

$$P(\textit{Lose}) = 1/2$$

$$P(\textit{red}) = 3/7$$

$$P(\textit{black}) = 4/7$$

$$P(\textit{black}|\textit{Win}) = 1/2$$

$$P(\textit{red}|\textit{Win}) = 1/2$$

After seeing the bead:

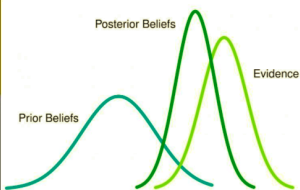
If bead = black:

$$P(\textit{Win}|\textit{black}) = \frac{P(\textit{Win})P(\textit{black}|\textit{Win})}{P(\textit{black})} = \frac{1/2 * 1/2}{4/7} = 0.4375$$

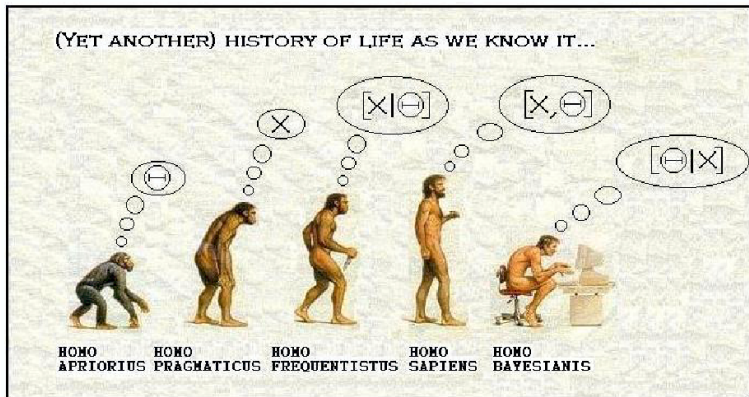
$$\text{If bead = red: } P(\textit{Win}|\textit{red}) = \frac{P(\textit{Win})P(\textit{red}|\textit{Win})}{P(\textit{red})} = \frac{1/2 * 1/2}{3/7} = 0.583$$

Bayes Theorem ...

<p>The Posterior</p>	<p>The Evidence</p> <p>The probability of getting this evidence if this hypothesis were true</p>	<p>The Prior</p> <p>The probability of H being true, before gathering evidence</p>
<p>$P(H E)$</p> <p>The probability that the hypothesis (H) is true given the evidence (E)</p>	$= \frac{P(E H)P(H)}{P(E)}$ <p>The marginal probability of the evidence (Prob of E over all possibilities)</p>	



Bayesians ...



Outline

- 1 Motivation
- 2 Introduction
- 3 Bayesian Network Classifiers**
- 4 k-Dependence Bayesian Classifiers
- 5 Links and References

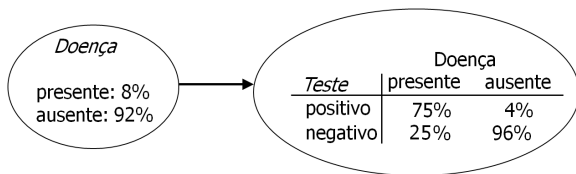
An Illustrative Problem:

A patient takes a lab test and the result comes back positive.

- The test returns a correct positive result in only 75% of the cases in which the disease is actually present,
- and a correct negative result in only 96% of the cases in which the disease is not present.
- Furthermore, 8% of the entire population have this cancer.

How to represent that information?

Representation:



It is useful to represent this information in a graph.

- The graphical information is qualitative
- The nodes represent variables.
- Arcs specify the (in)dependence between variables. Direct arcs represent influence between variables.

The direction of the arc tell us that the value of the variable *disease* influences the value of the variable *test*.

Naive Bayes Classifier

Assume target function $f : X \rightarrow Y$, where each instance x described by attributes $\langle x_1, x_2 \dots x_n \rangle$.

Most probable value of $f(x)$ is:

$$Y_{MAP} = \operatorname{argmax}_{y_j \in Y} P(y_j | x_1, x_2 \dots x_n)$$

$$Y_{MAP} = \operatorname{argmax}_{y_j \in Y} \frac{P(x_1, x_2 \dots x_n | y_j) P(y_j)}{P(x_1, x_2 \dots x_n)}$$

$$= \operatorname{argmax}_{y_j \in Y} P(x_1, x_2 \dots x_n | y_j) P(y_j)$$

Naive Bayes Classifier

Naive Bayes assumption: Attributes are independent given the class.

$P(x_1, x_2 \dots x_n | y_j) = \prod_i P(x_i | y_j)$ which gives Naive Bayes classifier:

$$Y_{NB} = \underset{y_j \in V}{\operatorname{argmax}} P(y_j) \prod_i P(x_i | y_j)$$

Naive Bayes

- Assume a decision problem with p variables.
- Each variable assume k values.
- The joint probability requires to estimate k^p probabilities.
- Assuming that variables are conditionally independent given the class, only requires to estimate $k \times p$ probabilities.

Naive Bayes Formulae

- Naive Bayes can be expressed in additive form:

$$P(y_i|\vec{x}) \propto \ln(P(y_i)) + \sum \ln(P(x_j|y_i))$$

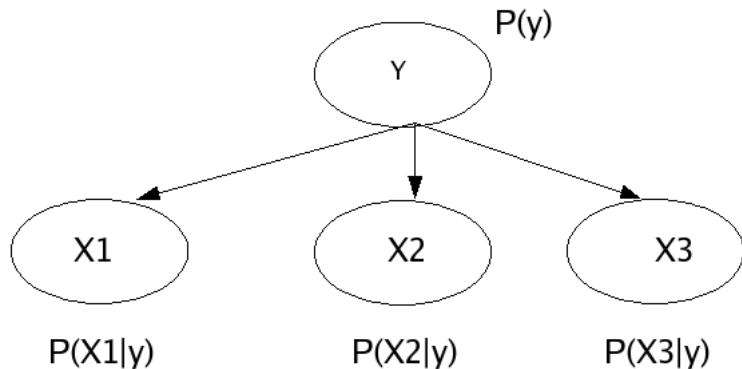
- Points out the contribution of each attribute to the decision.
- Two class problems:

$$\ln \frac{P(y_+|\vec{x})}{P(y_-|\vec{x})} \propto \ln \frac{P(y_+)}{P(y_-)} + \sum \ln \frac{P(x_j|y_+)}{P(x_j|y_-)}$$

$$\ln \frac{P(y_+|\vec{x})}{P(y_-|\vec{x})} \propto \ln \frac{P(y_+)}{1 - P(y_+)} + \sum \ln \frac{P(x_j|y_+)}{1 - P(x_j|y_+)}$$

The sign of each term indicates the class the attribute contributes to.

Naive Bayes as a Bayesian Net



$$p(Y|X_1, X_2, X_3) = P(Y)P(X_1|Y)P(X_2|Y)P(X_3|Y)$$

Naive Bayes Algorithm

Naive_Bayes_Learn(*examples*)

Initialize all counts to zero

For each example $\{X, y_j\}$

$Nr = Nr + 1$

$Count(y_j) = Count(y_j) + 1$

For each attribute value $x_i \in X$

$Count(x_i, y_j) = Count(x_i, y_j) + 1$

Classify_New_Instance(x)

$$y_{NB} = \operatorname{argmax}_{y_j \in Y} \hat{P}(y_j) \prod_{x_i \in X} \hat{P}(x_i | y_j)$$

Naive Bayes: Example

Weather	Temperature	Humidity	Wind	Play
Rainy	71	91	Yes	No
Sunny	69	70	No	Yes
Sunny	80	90	Yes	No
Overcast	83	86	No	Yes
Rainy	70	96	No	Yes
Rainy	65	70	Yes	No
Overcast	64	65	Yes	Yes
Overcast	72	90	Yes	Yes
Sunny	75	70	Yes	Yes
Rainy	68	80	No	Yes
Overcast	81	75	No	Yes
Sunny	85	85	No	No
Sunny	72	95	No	No
Rainy	75	80	No	Yes

Discretization: Basic Methods

- **Equal width discretization.** Divides the range of observed values for a feature into k equally sized bins. Pos: simplicity, Neg: the presence of outliers.
- **Equal frequency discretization.** Divides the range of observed values into k bins, where (considering n instances) each bin contains n/k values.
- **k-means.** An iterative method that begins with an equal-width discretization, iteratively adjust the boundaries to minimize a squared-error function and only stops when it can not change any value to improve the previous criteria.

1 3 6 7 8 9.5 10 11



Naive Bayes: Analysis

- 1 Conditional independence assumption is often violated

$$P(x_1, x_2 \dots x_n | y_j) = \prod_i P(x_i | y_j)$$

- ...but it works surprisingly well anyway. Note don't need estimated posteriors $\hat{P}(y_j|x)$ to be correct; need only that
$$\operatorname{argmax}_{y_j \in Y} \hat{P}(y_j) \prod_i \hat{P}(x_i | y_j) = \operatorname{argmax}_{y_j \in Y} P(y_j) P(x_1 \dots, x_n | y_j)$$
- see [Domingos & Pazzani, 1996] for analysis
- Naive Bayes posteriors often unrealistically close to 1 or 0

Naive Bayes: Analysis

- Robust to the presence of irrelevant attributes.

Suppose a two class problem, where X_i is an irrelevant attribute: $p(x_i|y_1) = p(x_i|y_2)$.

$$p(Y|x_1, \dots, x_i, \dots, x_n) \propto$$

$$p(Y)p(x_i|c) \prod_{l=1}^{i-1} p(x_l|Y) \prod_{l=i+1}^n p(x_l|Y) \text{ and}$$

$$p(Y|x_1, \dots, x_i, \dots, x_n) \propto p(Y|x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

- Redundant variables must be taken into account.

Suppose that $X_i = X_{i-1}$ then $p(x_{i-1}|Y) = p(x_i|Y)$

$$p(Y|x_1, \dots, x_{i-1}, x_i, \dots, x_n) \propto$$

$$p(Y)p(x_{i-1}|Y)p(x_i|Y) \prod_{l=1}^{i-2} p(x_l|Y) \prod_{l=i+1}^n p(x_l|Y)$$

and

$$p(Y|x_1, \dots, x_{i-1}, x_i, \dots, x_n) \propto$$

$$p(Y)p(x_i|Y)^2 \prod_{l=1}^{i-2} p(x_l|Y) \prod_{l=i+1}^n p(x_l|Y)$$

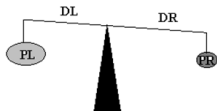
Naive Bayes: Summary

- The variability of a dataset is summarized in contingency tables.
 - Requires a single scan over the dataset.
 - The algorithm is Incremental (incorporation of new examples) and decremental (forgetting old examples).
- The dimension of the decision model is independent of the number of examples.
 - Low variance: stable with respect to small perturbations of the training set.
 - High bias: the number of possible states is finite.

Successful Stories

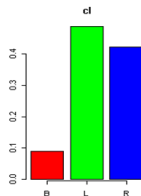
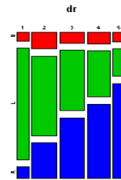
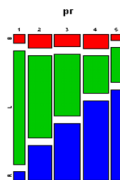
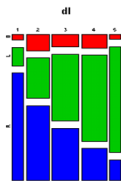
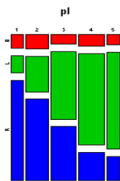
- KDDCup 1998: Winner - Boosting naive Bayes;
- Coil 1999: Winner - Simple naive Bayes;
- The most used classifier in Text Mining;

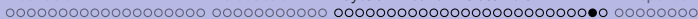
The Balance-scale Problem



A balança pode estar em três estados:
 .Equilibrada (B)
 .Inclinada para a direita (R)
 .Inclinada para a esquerda (L)

	pl	dl	pr	dr	cl
1	4	1	3	2	R
2	3	3	4	3	R
3	5	1	4	2	R
4	3	4	5	4	R
5	5	4	2	3	L
6	3	4	3	2	L
7	5	3	2	4	L
8	4	4	4	1	L





Knime: Naive Bayes

Naive B

The node creates a Bayesian model based on the data. It calculates the probability distribution for numerical attributes and the Bayes predictor to predict the class.

Dialog Options

Classification Column
The class value column.

Skip missing values (incl. class c...
The node ignores missing values. If not ticked, the node treats them as a separate class and considers them during the learning process.

Maximum number of unique nominal values
All nominal columns with more than this number of unique values will be skipped during learning. If not set, the 'Skip missing values' option is used.

Ports

Data Input

Class counts for class

Class:	Iris-setosa	Iris-versicolor	Iris-virginica
Count:	50	50	50

Gaussian distribution for petallength per class value

	Iris-setosa	Iris-versicolor	Iris-virginica
Count:	50	50	50
Mean:	1.464	4.26	5.552
Std. Deviation:	0.17351	0.46991	0.55189
Rate:	50/150	50/150	50/150

Gaussian distribution for petalwidth per class value

	Iris-setosa	Iris-versicolor	Iris-virginica
Count:	50	50	50
Mean:	0.244	1.326	2.026
Std. Deviation:	0.10721	0.19775	0.27465
Rate:	50/150	50/150	50/150

Conditional Entropy

conditional entropy quantifies the remaining entropy (i.e. uncertainty) of a random variable Y given that the value of a second random variable X is known.

Uncertainty about X after knowing Y :

$$H(X|Y) = - \sum_y P(y) \sum_x P(x|y) \log P(x|Y)$$

Mutual Information

Reduction in the uncertainty of X when Y is known:

$$I(X, Y) = H(X|Y) - H(X)$$

$$I(X, Y) = \sum_i \sum_j P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(x_i)P(y_j)}$$

- $I(X, Y) = 0$ if and only if X and Y are independent random variables;
- $I(X, Y)$ increases with the increase of the degree of dependence between X and Y ;
- nonnegative: $I(X, Y) \geq 0$;
- symmetric: $I(X, Y) = I(Y, X)$.

k-Dependency Bayesian Networks

- Compute $I(X_i, C)$ and $I(X_i, X_j|C)$ for all pairs of variables
- Iterate
 - Choose the variable X_{max} not yet in the model that maximizes $I(X_i|C)$
 - Choose the k parents of X_{max} : those with greater $I(X_j, X_{max}|C)$

Factorization:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | Parents(X_i))$$

where $Parents(X_i)$ denotes immediate predecessors of X_i in graph

Weka: TAN

Program Applications Tools Visualization Windows Help

Explorer

Preprocessor Classify Cluster Associate Select attributes Visualize

Classifier

Choose **BayesNet** -D -Q weka.classifiers.bayes.net.search.local.TAN -- -S BAYES -E weka.classifiers.bayes.net.estimate.SimpleEstimator -- -A 0.5

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)

13:56:53 - bayes.BayesNet

13:57:30 - bayes.BayesNet

Classifier output

Kappa statistic 1

Mean absolute error 0.1703

Root mean squared error 0.1921

Relative absolute error 36.6716 %

Root relative squared error 40.063 %

Total Number of Instances

=== Detailed Accuracy By Class ===

TP Rate	FP Rate	Precision	Recall
1	0	1	1
1	0	1	1

=== Confusion Matrix ===

a b <-- classified as

9 0 | a = yes

0 5 | b = no

Weka Classifier Graph Visualizer: 13:57:30 - bayes.BayesNet

```

graph TD
    play((play)) --> temperature((temperature))
    play --> outlook((outlook))
    play --> humidity((humidity))
    temperature --> outlook
    temperature --> humidity
    outlook --> windy((windy))
    humidity --> windy
  
```

Probability Distribution Table For outlook

play	temperature	sunny	overcast	rainy
yes	hot	0,143		0,143
yes	mild	0,273	0,273	0,455
yes	cool	0,333	0,333	0,333
no	hot	0,714	0,143	0,143
no	mild	0,429	0,143	0,429
no	cool	0,2	0,2	0,6

Outline

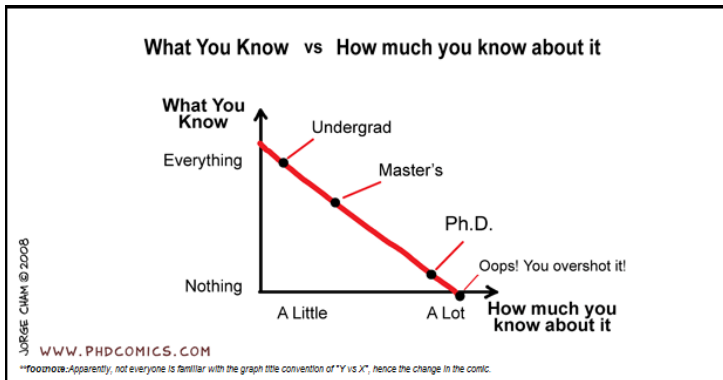
- 1 Motivation
- 2 Introduction
- 3 Bayesian Network Classifiers
- 4 k-Dependence Bayesian Classifiers
- 5 Links and References

Software Available

- R (package e1071)
- Weka (naive Bayes, TAN models, k-dependence Bayesian classifiers)
- Genie: Bayesian Networks, Influence Diagrams, Dynamic Bayesian Networks (<http://genie.sis.pitt.edu/>)
- Hugin (<http://www.hugin.com/>)
- Elvira (<http://www.ia.uned.es/elvira>)
- Kevin Murphy's MATLAB toolbox - supports dynamic BNs, decision networks, many exact and approximate inference algorithms, parameter estimation, and structure learning (<http://www.ai.mit.edu/murphyk/Software/BNT/bnt.html>)
- Free Windows software for creation, assessment and evaluation of belief networks. (<http://www.research.microsoft.com/dtas/msbn/default.htm>)
- Open source package for the technical computing language R, developed by Aalborg University (<http://www.math.auc.dk/novo/deal>)
- <http://www.snn.ru.nl/nijmegen>

Bibliography

- Gama, Carvalho, Faceli, Lorena, Oliveira *Extração de Conhecimento de Dados*, Silabo, 2012 (Cap 5)
- Tom Mitchell, *Machine Learning*, (chapter 6), McGraw-Hill, 1997
- R. Duda, P. Hart, D. Stork; *Pattern Classification*, J. Willey & Sons, 2000
- J.Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan Kaufmann, 1988
- P. Domingos, M.Pazzani; *On the Optimality of the Simple Bayes Classifier under zero-one loss*, Machine Learning, 29
- R. Neapolitan, *Learning Bayesian Networks*, Prentice Hall, 2004



Would you like to learn more? Wait for SAD ...