

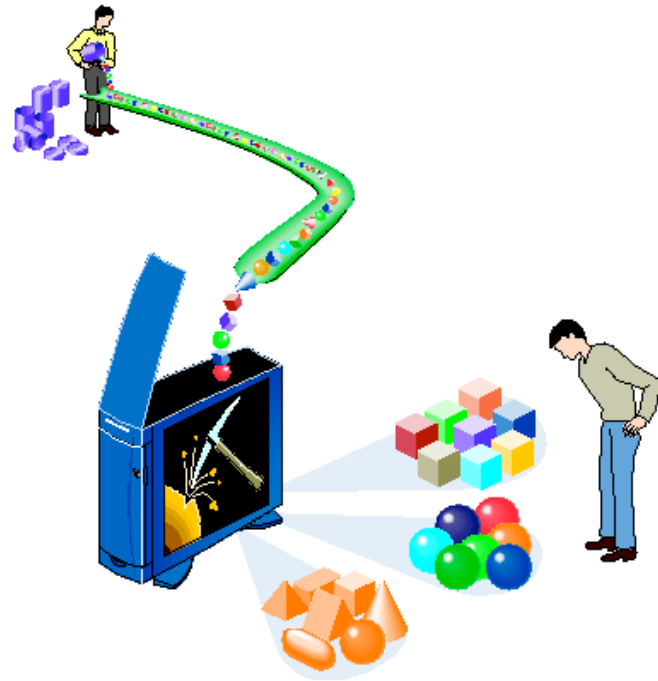
Supervised and Unsupervised Learning

Ciro Donalek

Ay/Bi 199 – April 2011

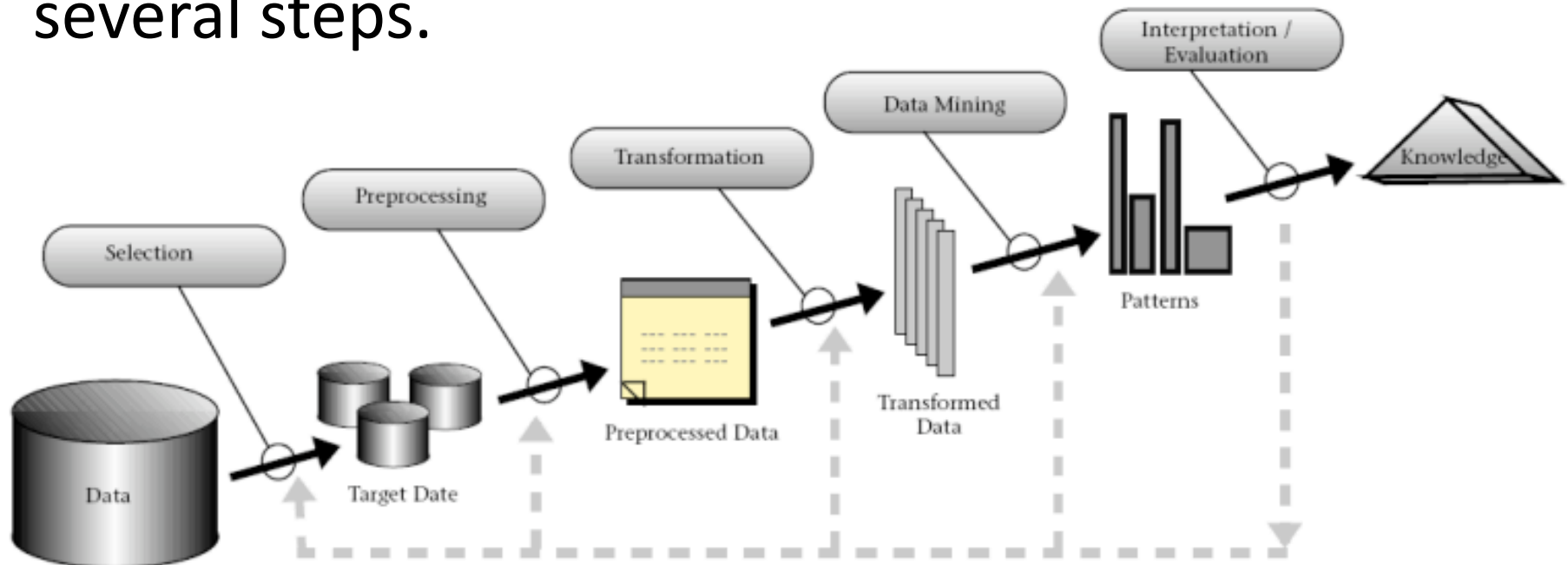
Summary

- KDD and Data Mining Tasks
- Finding the optimal approach
- Supervised Models
 - Neural Networks
 - Multi Layer Perceptron
 - Decision Trees
- Unsupervised Models
 - Kmeans
 - Self Organizing Maps
- Ensembles
- Links and References



Knowledge Discovery in Databases

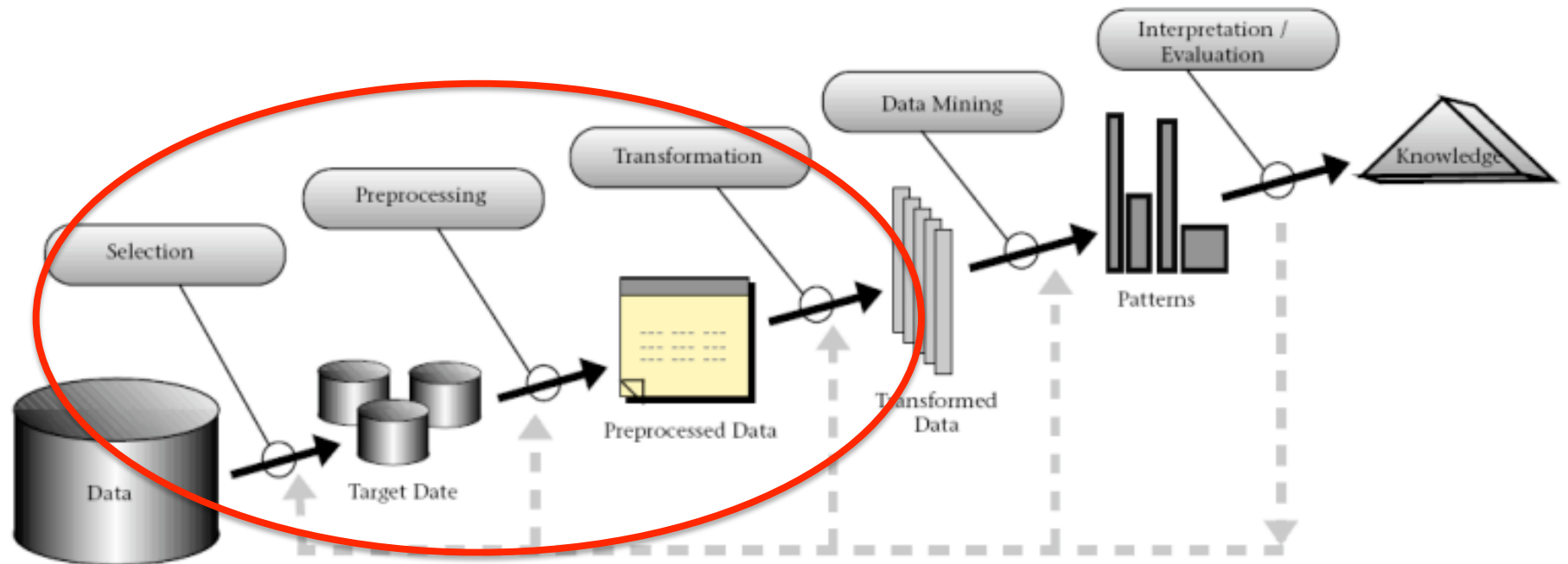
- KDD may be defined as: "*The non trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data*".
- KDD is an interactive and iterative process involving several steps.



You got your data: what's next?

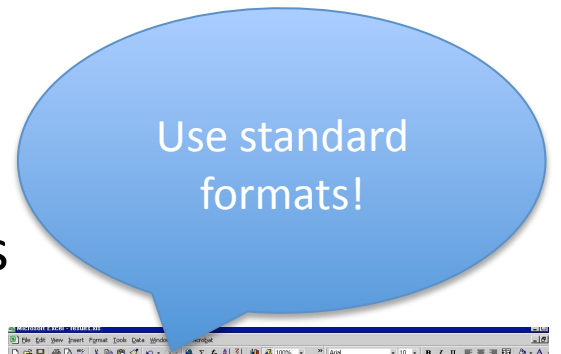


What kind of analysis do you need? Which model is more appropriate for it? ...



Clean your data!

- Data preprocessing transforms the raw data into a format that will be more easily and effectively processed for the purpose of the user.
- Some tasks
 - *sampling*: selects a representative subset from a large population of data;
 - *Noise treatment*
 - strategies to handle **missing data**: sometimes your raws will be incomplete, not all parameters are measured for all samples.
 - *normalization*
 - *feature extraction*: pulls out specified data that is significant in some particular context.

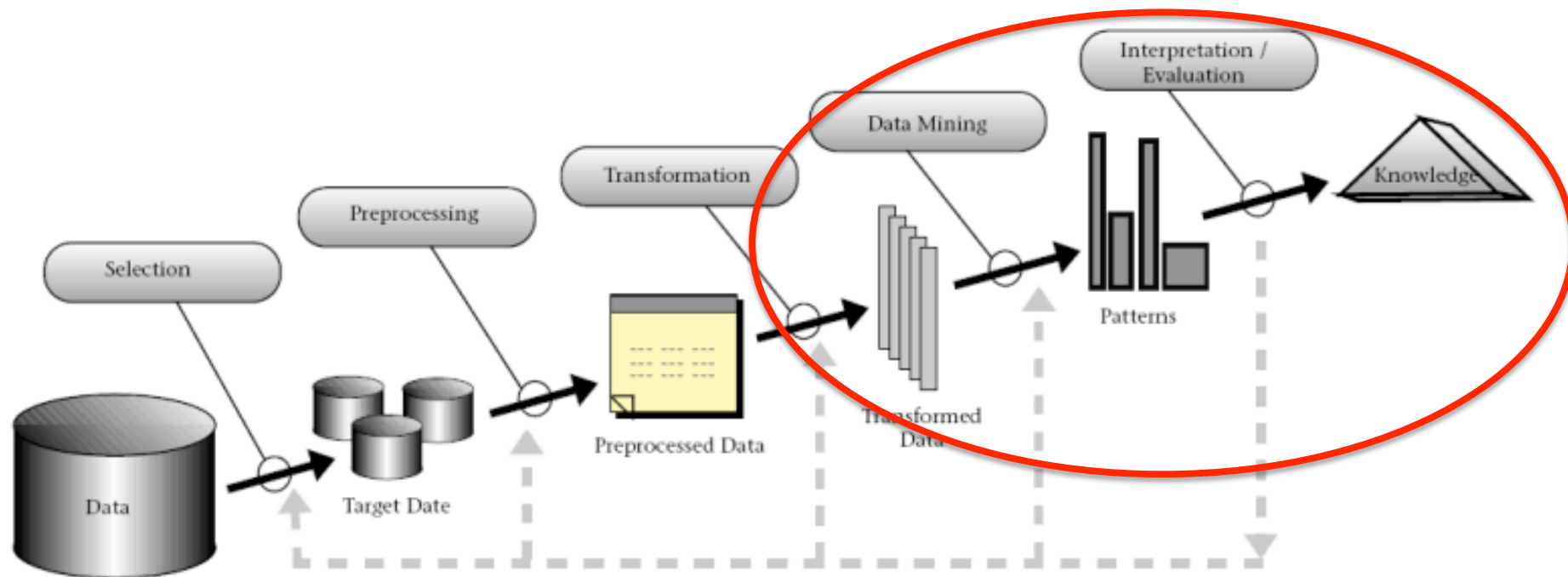


GO ID	GO term	Frequency	Genome frequency	P-value	Corrected P-value	Gene(s)
GO:0000281	DNA repair	1.00000	0.00332	6.5984E-213	2.1230E-211	MG1_1261809, MG1_1097010, MG1_1097076
GO:0000119	response to endogenous stim	1.00000	0.00387	4.3664E-204	1.6602E-202	MG1_1261809, MG1_1097010, MG1_1097076
GO:0000874	response to DNA damage str	1.00000	0.00387	4.3664E-204	1.6602E-202	MG1_1261809, MG1_1097010, MG1_1097076
GO:0000659	DNA metabolism	1.00000	0.00333	3.5066E-193	1.3232E-192	MG1_1261809, MG1_1097010, MG1_1097076
GO:0000660	response to stress	1.00000	0.02197	4.3608E-132	1.6570E-130	MG1_1261809, MG1_1097010, MG1_1097076
GO:0006139	nucleobase, nucleoside, nucleotide	1.00000	0.00332	4.1424E-97	1.5141E-95	MG1_1261809, MG1_1097010, MG1_1097076
GO:0000152	metabolism	1.00000	0.15150	8.7169E-85	2.5520E-83	MG1_1261809, MG1_1097010, MG1_1097076
GO:0007582	physiological processes	1.00000	0.28877	5.4946E-46	2.0379E-45	MG1_1261809, MG1_1097010, MG1_1097076
GO:0000310	DNA recombination	0.17349	0.00000	3.6564E-30	1.3000E-29	MG1_1261809, MG1_90895, MG1_1332098
GO:0000067	DNA replication and chromos	0.23077	0.00327	6.2807E-29	2.3711E-27	MG1_1332098, MG1_134881, MG1_134938
GO:0000360	DNA replication	0.21796	0.00287	1.58017E-28	5.3204E-27	MG1_1332098, MG1_134881, MG1_134938
GO:0000084	S phase of mitotic cell cycle	0.21796	0.00270	1.9079E-28	7.4899E-27	MG1_1332098, MG1_134881, MG1_134938
GO:0000289	nucleotide-excision repair	0.14103	0.00000	7.5910E-26	2.8499E-24	MG1_103667, MG1_103982, MG1_95400, MG1_103983
GO:0000279	mitotic cell cycle	0.23077	0.00427	1.0688E-24	4.4416E-23	MG1_1332098, MG1_134881, MG1_134938
GO:0007049	cell cycle	0.29487	0.01307	1.2336E-24	4.6084E-23	MG1_1332098, MG1_134881, MG1_134938
GO:0000083	cell proliferation	0.29487	0.01647	8.0132E-23	3.2725E-21	MG1_1332098, MG1_134881, MG1_134938
GO:0000084	base-excision repair	0.12056	0.00000	1.2953E-20	4.9227E-19	MG1_1261809, MG1_1097093, MG1_134111
GO:0000151	cell growth and/or maintenance	0.33333	0.00000	6.0640E-10	2.3119E-09	MG1_1332098, MG1_134881, MG1_1097076
GO:0000058	mismatch repair	0.05128	0.00000	1.4761E-09	5.5907E-09	MG1_134361, MG1_101816, MG1_101839
GO:0040005	maintenance of fidelity during meiosis	0.05128	0.00000	1.4761E-09	5.5907E-09	MG1_134361, MG1_101816, MG1_101839
GO:0007126	meiosis	0.06410	0.00000	3.7171E-09	9.7019E-09	MG1_1100512, MG1_102683, MG1_97090, MG1_102684
GO:0000080	nuclear division	0.07692	0.00257	5.1932E-08	1.9754E-06	MG1_1100512, MG1_102683, MG1_97090, MG1_102684
GO:0000079	M phase	0.07692	0.00287	6.5462E-08	2.4807E-06	MG1_1100512, MG1_102683, MG1_97090, MG1_102684
GO:0000081	chromosome-based repair	0.07692	0.00000	8.7921E-08	3.6661E-06	MG1_102779, MG1_102683, MG1_97090, MG1_102684

Missing Data

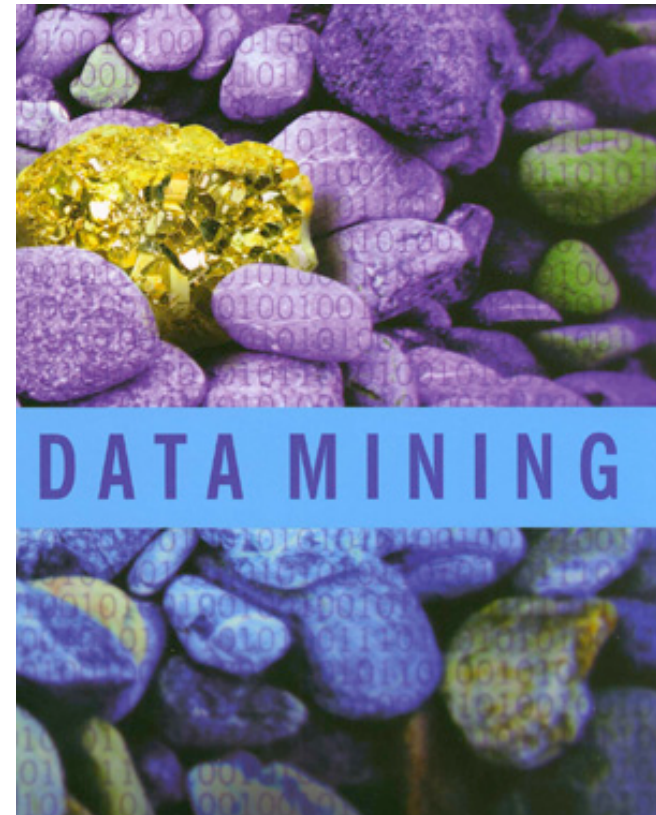
- Missing data are a part of almost all research, and we all have to decide how to deal with it.
- Complete Case Analysis: use only rows with all the values
- Available Case Analysis
- Substitution
 - Mean Value: replace the missing value with the mean value for that particular attribute
 - Regression Substitution: we can replace the missing value with historical value from similar cases
 - Matching Imputation: for each unit with a missing y , find a unit with similar values of x in the observed data and take its y value
 - Maximum Likelihood, EM, etc
- Some DM models can deal with missing data better than others.





Data Mining

- Data Mining is about automating the process of searching for patterns in the data.
- More in details, the most relevant DM tasks are:
 - association
 - sequence or path analysis
 - **clustering**
 - **classification**
 - **regression**
 - visualization



Finding Solution via Purposes

- What kind of analysis do you need?
- Regression
 - predict new values based on the past, inference
 - compute the new values for a dependent variable based on the values of one or more measured attributes
- Classification:
 - divide samples in classes
 - use a trained set of previously labeled data
- Clustering
 - partitioning of a data set into subsets (clusters) so that data in each subset ideally share some common characteristics
- Classification is in a some way similar to the clustering, but requires that the analyst know ahead of time how classes are defined.

Cluster Analysis

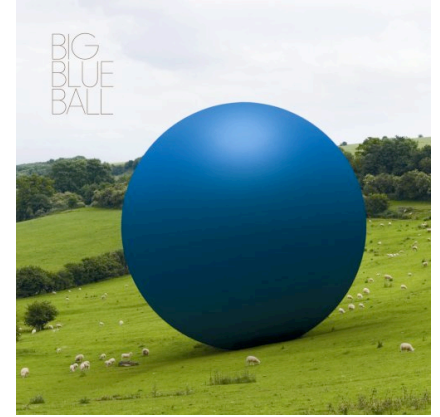


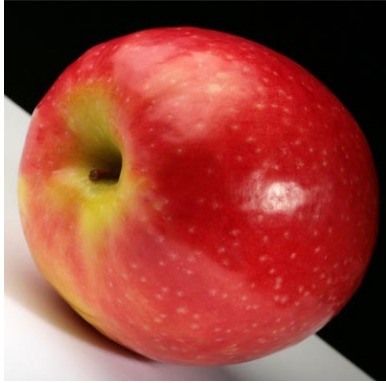
starshadow78/Flickr

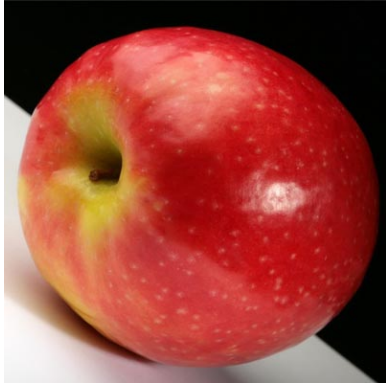
How many clusters do you expect?



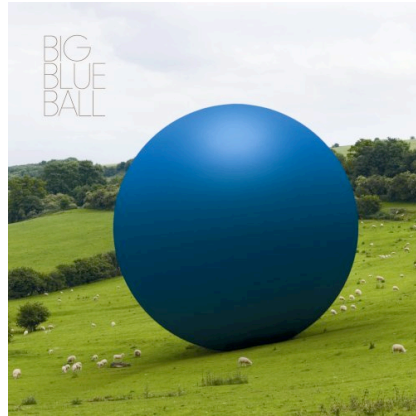
starshadow78/Flickr



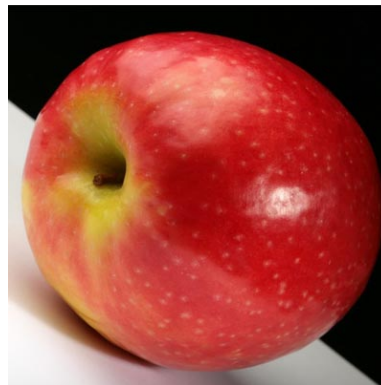
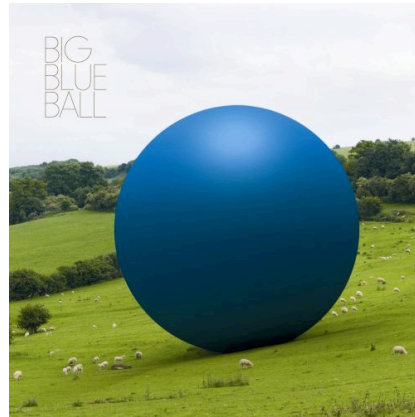




starshadow78/Flickr



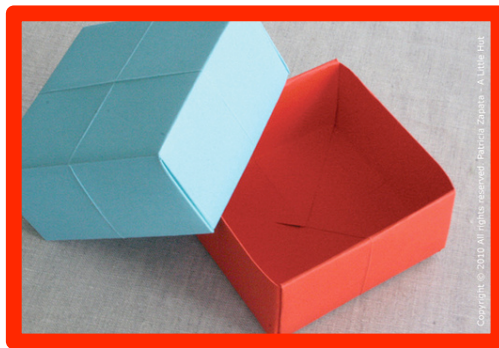
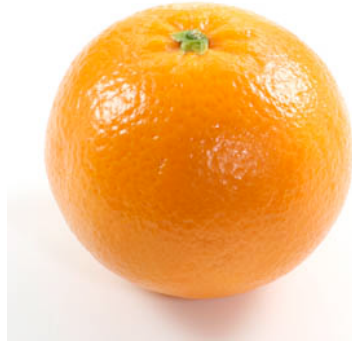
Search for Outliers



Classification

- Data mining technique used to predict group membership for data instances. There are two ways to assign a new value to a given class.
- **Crispy classification**
 - given an input, the classifier returns its label
- **Probabilistic classification**
 - given an input, the classifier returns its probabilities to belong to each class
 - useful when some mistakes can be more costly than others
 - winner take all and other rules
 - assign the object to the class with the highest probability (WTA)
 - ...but only if its probability is greater than 40% (WTA with thresholds)

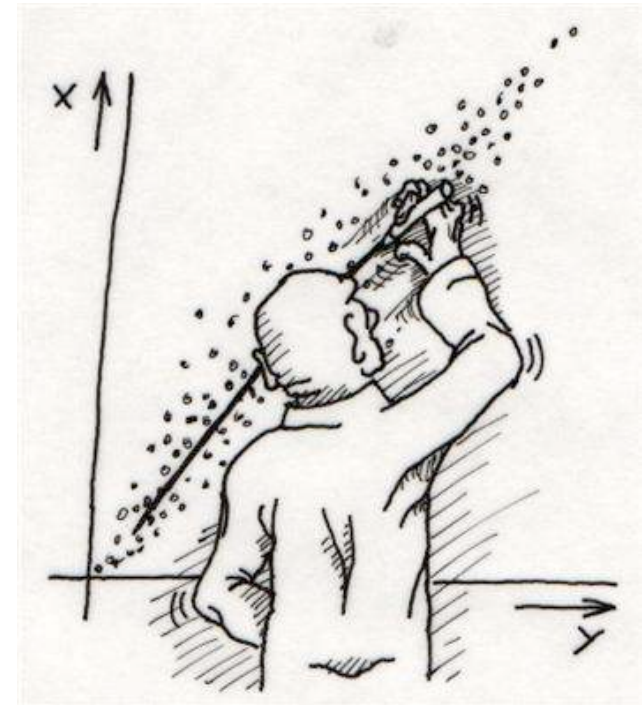




starshadow78/Flickr

Regression / Forecasting

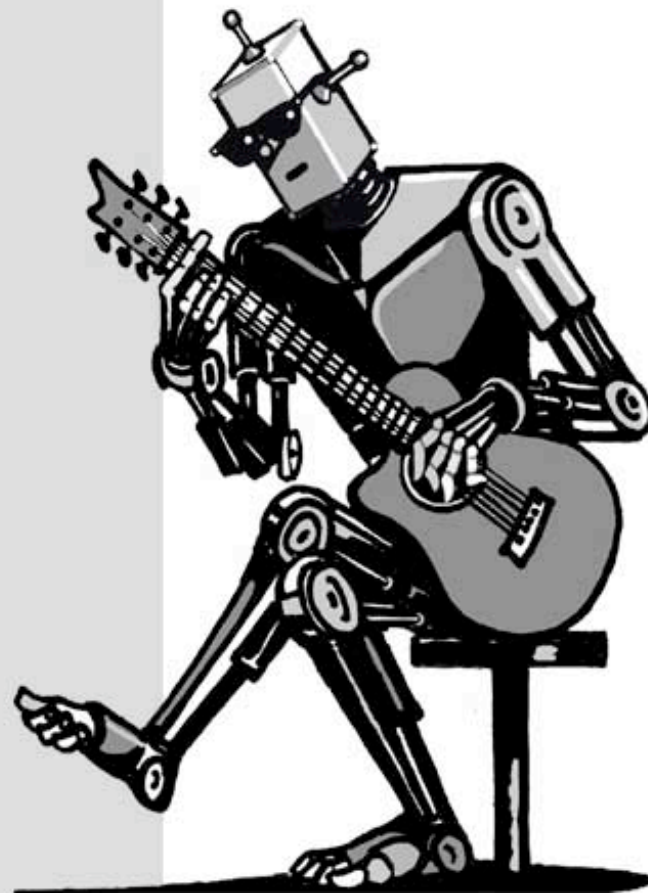
- Data table statistical correlation
 - mapping without any prior assumption on the functional form of the data distribution;
 - machine learning algorithms well suited for this.
- Curve fitting
 - find a well defined and known function underlying your data;
 - theory / expertise can help.



Machine Learning

- To learn: *to get knowledge of by study, experience, or being taught.*
- Types of Learning
 - Supervised
 - Unsupervised

J. SCHMIDHUBER, 2006



COGNITIVE ROBOTICS

Unsupervised Learning

- The model is not provided with the correct results during the training.
- Can be used to cluster the input data in classes on the basis of their statistical properties only.
- Cluster significance and labeling.
- The labeling can be carried out even if the labels are only available for a small number of objects representative of the desired classes.

Supervised Learning

- Training data includes both the input and the desired results.
- For some examples the correct results (targets) are known and are given in input to the model during the learning process.
- The construction of a proper training, validation and test set (Bok) is crucial.
- These methods are usually fast and accurate.
- Have to be able to **generalize**: give the correct results when new data are given in input without knowing a priori the target.

Generalization

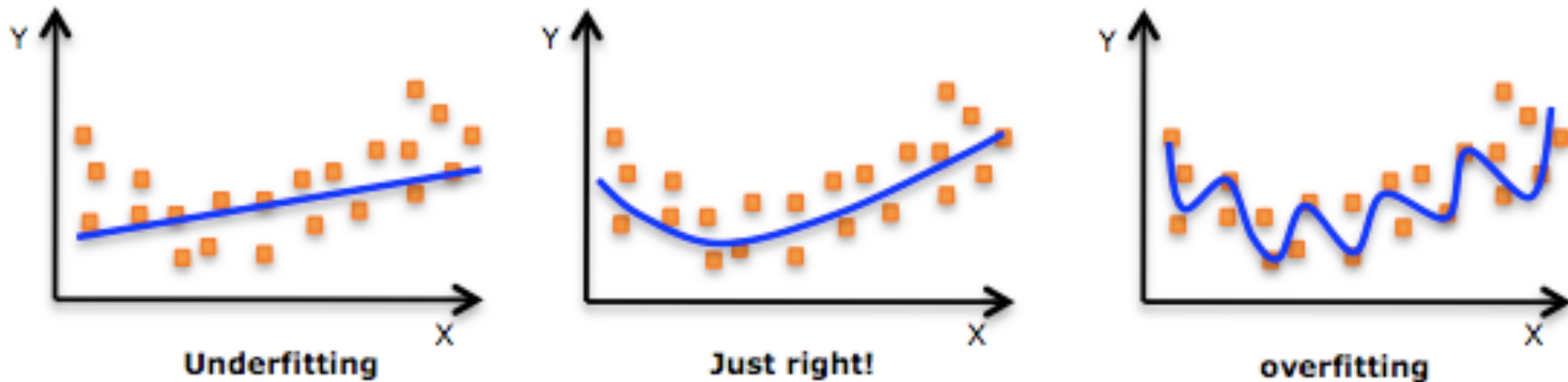
- Refers to the ability to produce reasonable outputs for inputs not encountered during the training.



In other words: NO PANIC when "never seen before" data are given in input!

A common problem: OVERFITTING

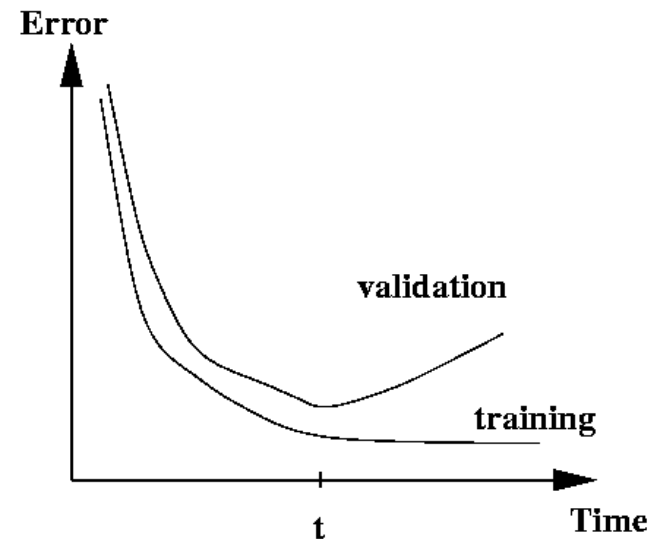
- Learn the “data” and not the underlying function
- Performs well on the data used during the training and poorly with new data.



Use proper training sets, early stopping.

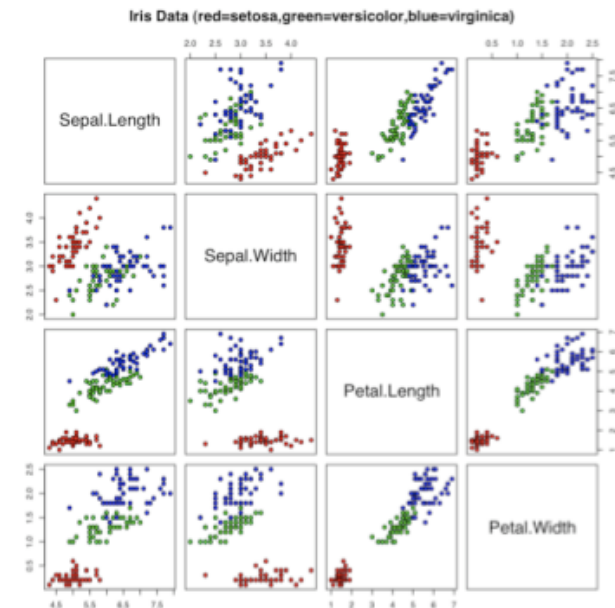
Datasets

- **Training set:** a set of examples used for learning, where the target value is known.
- **Validation set:** a set of examples used to tune the architecture of a classifier and estimate the error.
- **Test set:** used only to assess the performances of a classifier. It is **never used** during the training process so that the error on the test set provides an unbiased estimate of the generalization error.



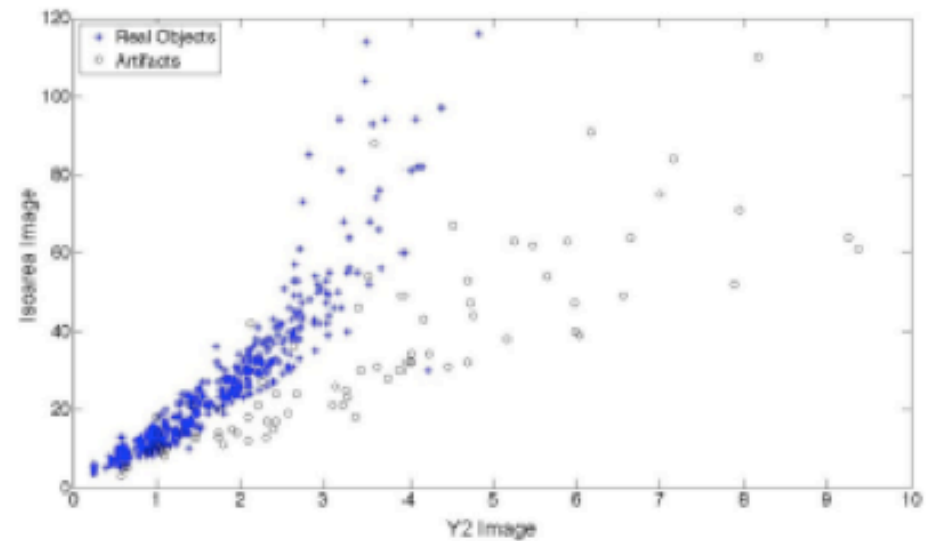
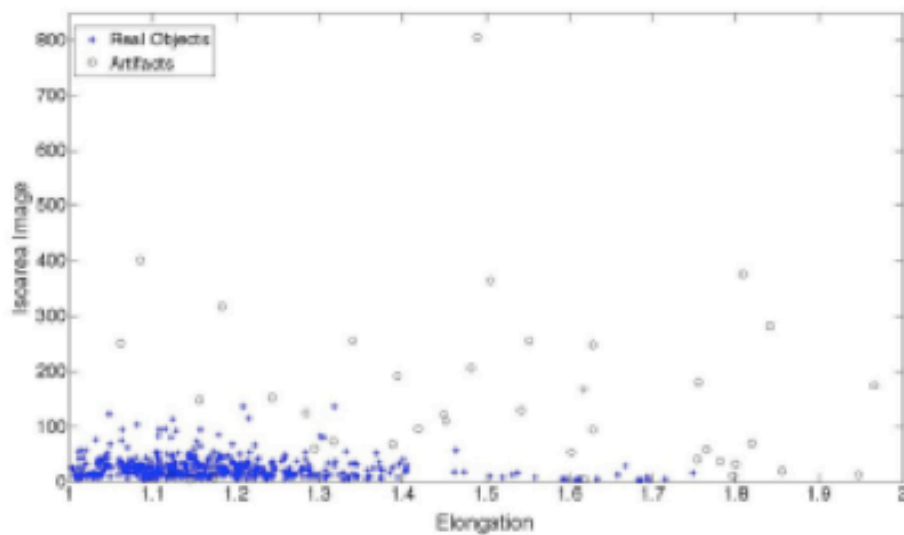
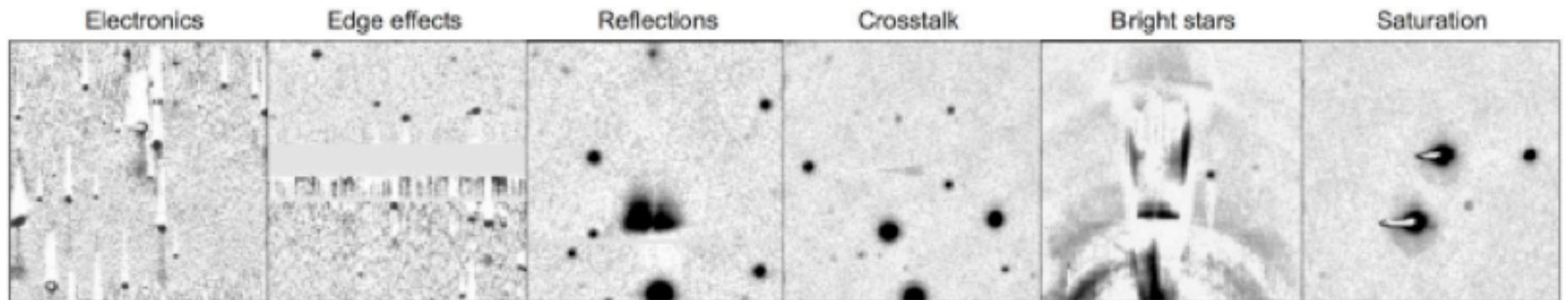
IRIS dataset

- IRIS
 - consists of 3 classes, 50 instances each
 - 4 numerical attributes (sepal and petal length and width in cm)
 - each class refers to a type of Iris plant (Setosa, Versicolor, Verginica)
 - the first class is linearly separable from the other two while the 2nd and the 3rd are not linearly separable



Artifacts Dataset

- PQ Artifacts
 - 2 main classes and 4 numerical attributes
 - classes are: true objects, artifacts



Data Selection

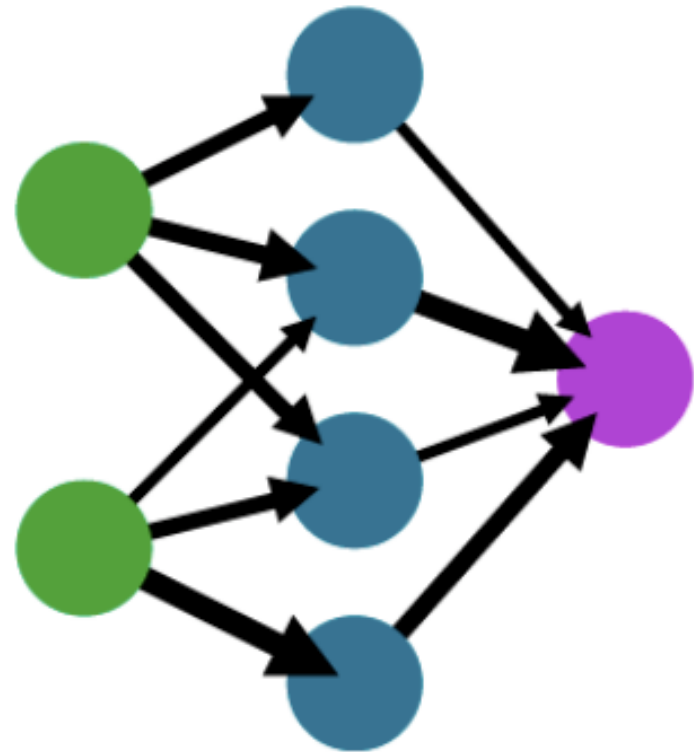
- “**Garbage in, garbage out**”: training, validation and test data must be representative of the underlying model
- All eventualities must be covered
- Unbalanced datasets
 - since the network minimizes the overall error, the proportion of types of data in the set is critical;
 - inclusion of a loss matrix (Bishop,1995);
 - often, the best approach is to ensure even representation of different cases, then to interpret the network's decisions accordingly.

Artificial Neural Network

An Artificial Neural Network is an information processing paradigm that is inspired by the way biological nervous systems process information:

“a large number of highly interconnected simple processing elements (neurons) working together to solve specific problems”

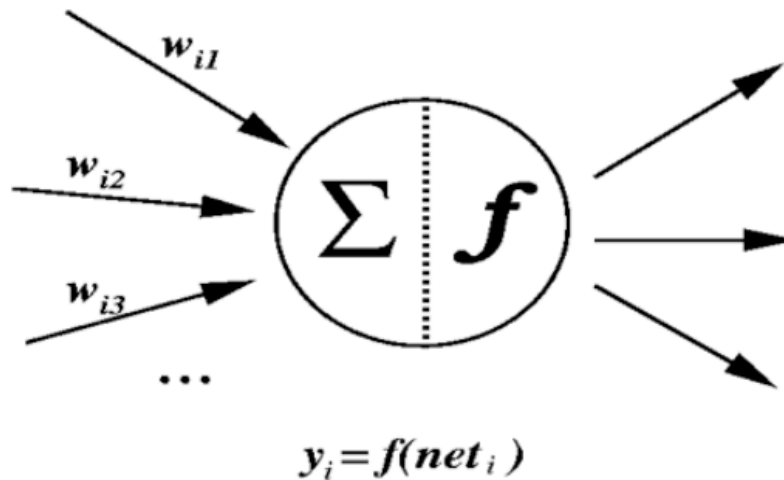
A simple neural network
input layer hidden layer output layer



A simple artificial neuron

- The basic computational element is often called a node or unit. It receives input from some other units, or from an external source.
- Each input has an associated weight w , which can be modified so as to model synaptic learning.
- The unit computes some function of the weighted sum of its inputs:

$$y_i = f\left(\sum_j w_{ij} y_j\right)$$



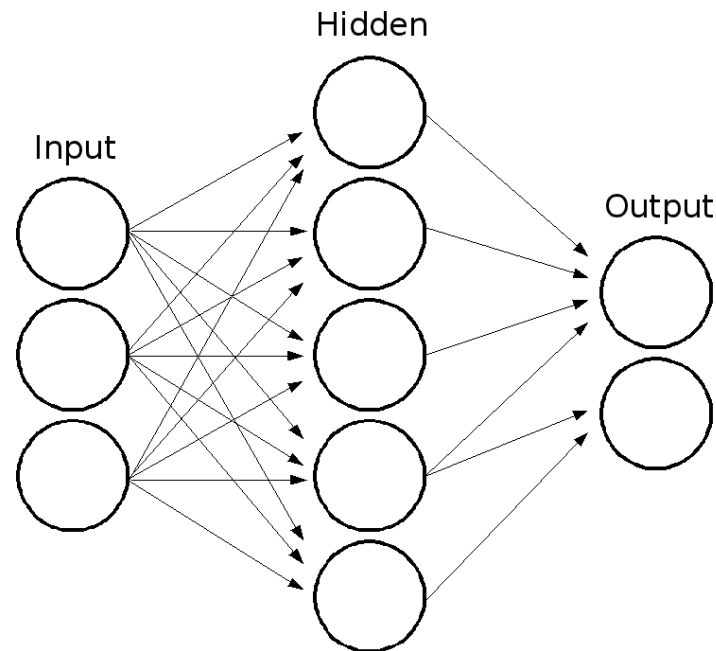
Neural Networks

A Neural Network is usually structured into an input layer of neurons, one or more hidden layers and one output layer.

Neurons belonging to adjacent layers are usually fully connected and the various types and architectures are identified both by the different topologies adopted for the connections as well by the choice of the activation function. The values of the functions associated with the connections are called “weights”.

The whole game of using NNs is in the fact that, in order for the network to yield appropriate outputs for given inputs, the weight must be set to suitable values.

The way this is obtained allows a further distinction among modes of operations.



Neural Networks: types

Feedforward: Single Layer Perceptron, MLP, ADALINE (Adaptive Linear Neuron), RBF

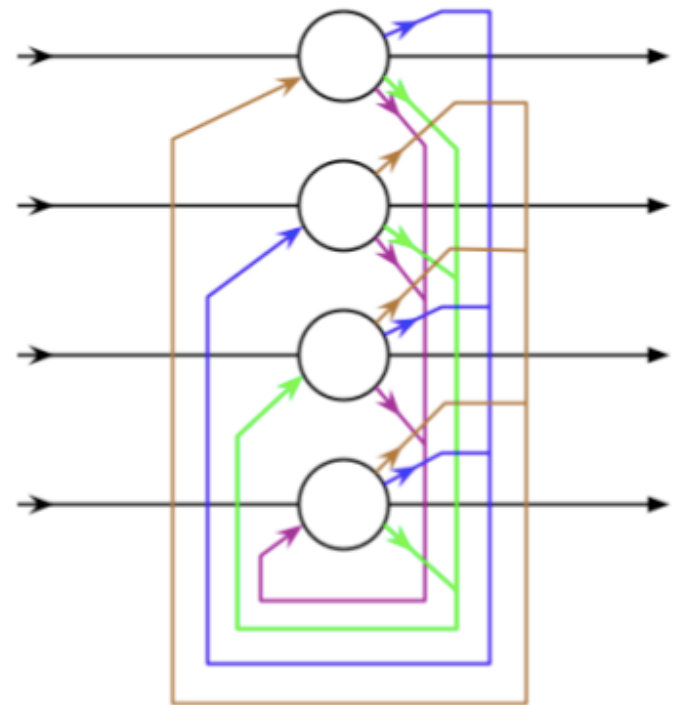
Self-Organized: SOM (Kohonen Maps)

Recurrent: Simple Recurrent Network, Hopfield Network.

Stochastic: Boltzmann machines, RBM.

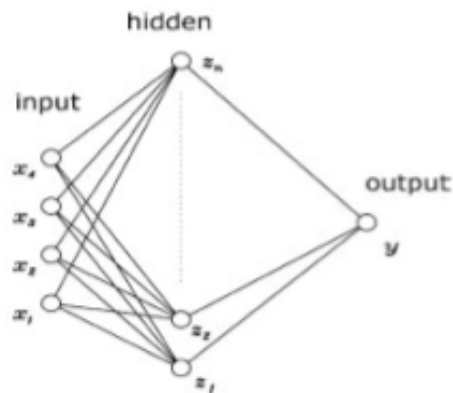
Modular: Committee of Machines, ASNN (Associative Neural Networks), Ensembles.

Others: Instantaneously Trained, Spiking (SNN), Dynamic, Cascades, NeuroFuzzy, PPS, GTM.

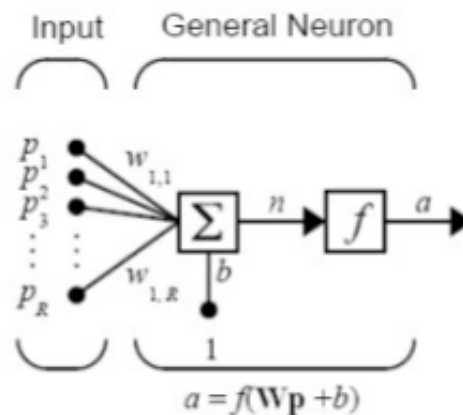


Multi Layer Perceptron

- The MLP is one of the most used supervised model: it consists of multiple layers of computational units, usually interconnected in a feed-forward way.
- Each neuron in one layer has direct connections to all the neurons of the subsequent layer.

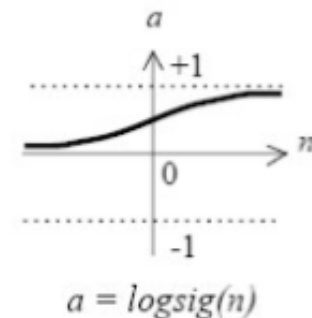


The architecture of a two layer MLP.



Where

R = number of elements in input vector



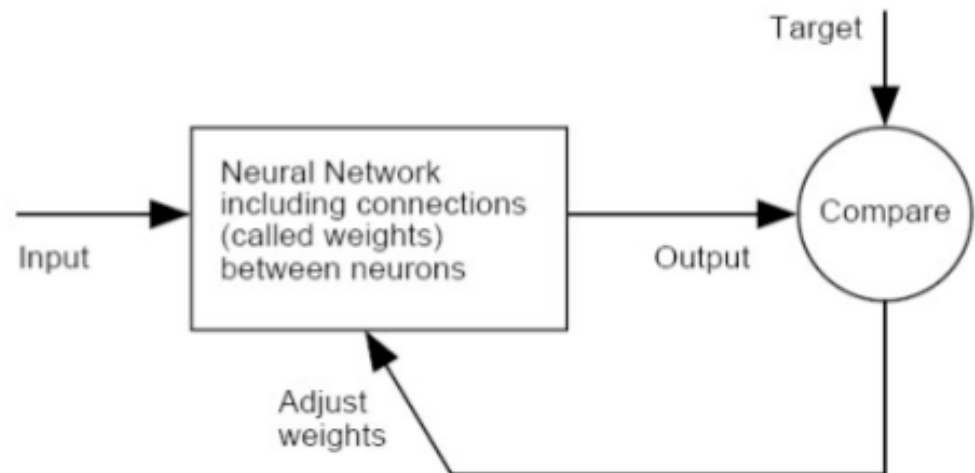
$$a = \text{logsig}(n)$$



Learning Process

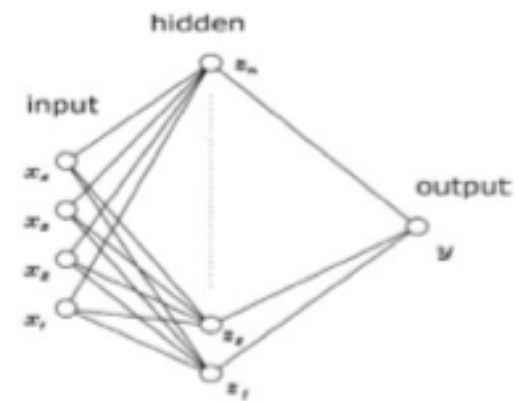
- Back Propagation
 - the output values are compared with the target to compute the value of some predefined error function
 - the error is then feedback through the network
 - using this information, the algorithm adjusts the weights of each connection in order to reduce the value of the error function

After repeating this process for a sufficiently large number of training cycles, the network will usually converge.



Hidden Units

- The best number of hidden units depend on:
 - number of inputs and outputs
 - number of training case
 - the amount of noise in the targets
 - the complexity of the function to be learned
 - the activation function



The architecture of a two layer MLP.

- Too few hidden units => high training and generalization error, due to underfitting and high statistical bias.
- Too many hidden units => low training error but high generalization error, due to overfitting and high variance.
- Rules of thumb don't usually work.

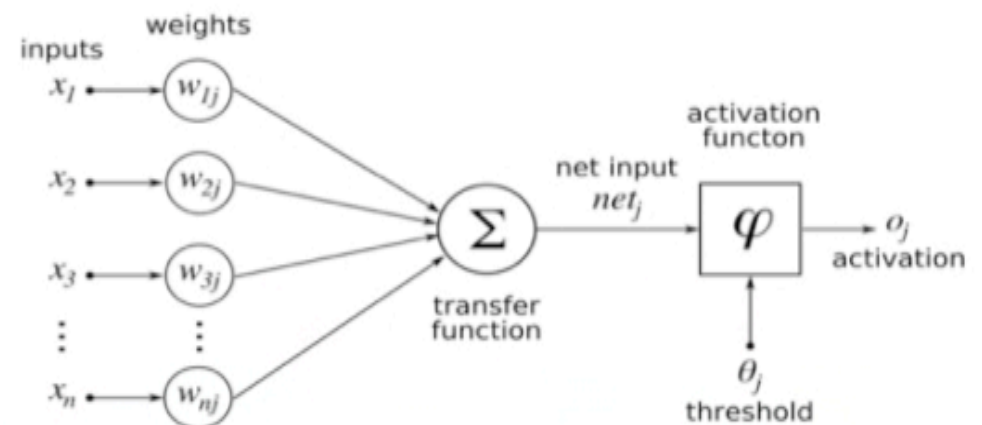
Activation and Error Functions

- **Error Functions**

- measure of the discrepancy between the network output values and the target;
- sum of the squared errors (SSE), cross entropy (CE), etc.

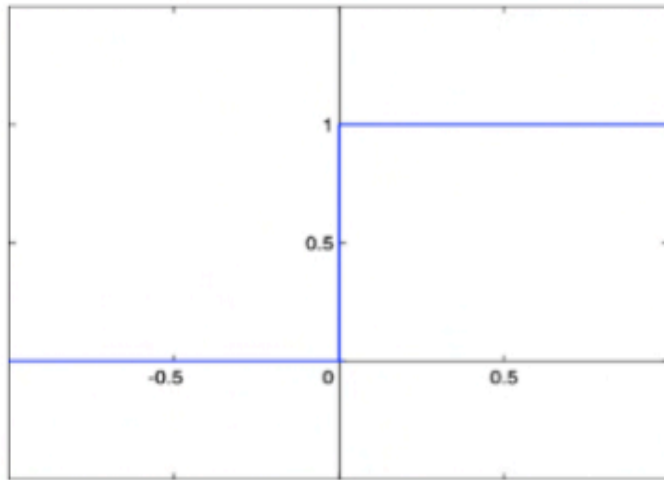
- **Activation Functions**

- used by most units to transform their inputs;
- needed to introduce non linearity into the network
- linear, logistic, tanh, softmax...



Using a Multilayer Perceptron with a softmax activation function and cross-entropy error, the network outputs can be interpreted as the conditional probabilities $p(C_1 | \mathbf{x})$ and $p(C_2 | \mathbf{x})$ where \mathbf{x} is the input vector, C_1 the first class, C_2 the second class.

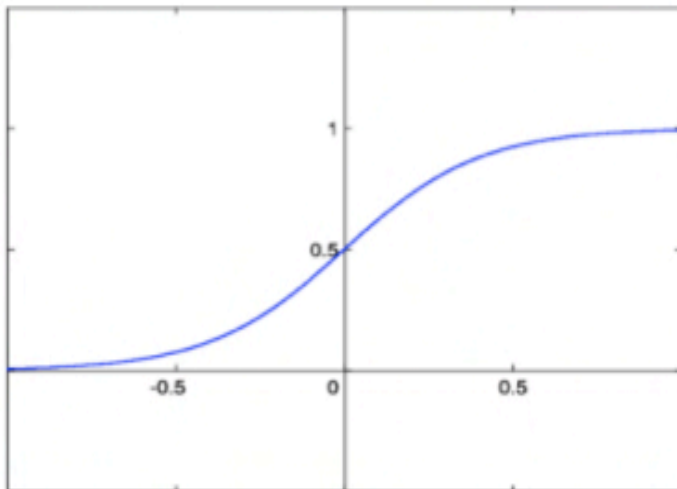
Activation Functions



Step function

The output is a certain value A_1 , if the input sum is above a certain threshold and A_0 if the input sum is below a certain threshold.

When we want to classify an input pattern into one of two groups, we can use a binary classifier with a step activation function.



Sigmoid function

Has the property of being similar to the step function, but with the addition of a region of uncertainty.

Sigmoid functions in this respect are very similar to the input-output relationships of biological neurons.

Results: confusion matrix

In the confusion matrix the network prediction Y are compared with the target T: the rows represent the true classes and the columns the predicted classes.

Training set
Classification rate: 97.35%

Galaxy	1009	34
Star	19	938
	Galaxy	Star

Test set
Classification rate: 91.975%

Galaxy	1641	65
Star	256	2038
	Galaxy	Star

Results: completeness and contamination

The performances of the classifiers are rated based on the following three criteria. Supposing we have 2 classes A and B:

- ✓ **completeness**: the percentage of objects of class A correctly classified as such;
- ✓ **contamination**: the percentage of objects of class A incorrectly classified as objects belonging to the class B;
- ✓ **classification rate**: the overall percentage of objects correctly classified.

Exercise: compute completeness and contamination for the previous confusion matrix (test set)

Decision Trees

- Is another classification method.
- A decision tree is a set of simple rules, such as "if the sepal length is less than 5.45, classify the specimen as setosa."
- Decision trees are also nonparametric because they do not require any assumptions about the distribution of the variables in each class.

