
Machine Learning Formulas*

García Prado, Sergio
sergio@garciparedes.me

Taboada Roderó, Ismael José
ismaeljose.taboada@alumnos.uva.es

May 28, 2017

Abstract

[TODO]

1 Introduction

[TODO]

2 Supervised Learning

[TODO]

Overfitting [TODO]

Error Measures [TODO]

- **Error Rate:** [TODO]
- **Resubstitution error:** [TODO]
- **Generalization Error:** [TODO]

Experimental Strategies [TODO]

- **Holdout:** [TODO]
- **Repeated Holdout:** [TODO]
- **Cross Validation:** [TODO]
- **Repeated Cross Validation:** [TODO]

2.1 Decision Trees

Decision trees are classifiers that works splitting the instances between the branches of a tree according to its attributes and assigning the feature's attribute in the leafs of the tree. An scheme of this structure is been included in the figure 1. Decision trees can represent logic functions. It tend to ovefit the data so the advanced heuristics try to avoid this problem.

*This document is being maintained in <https://github.com/garciparedes/machine-learning-formulas>

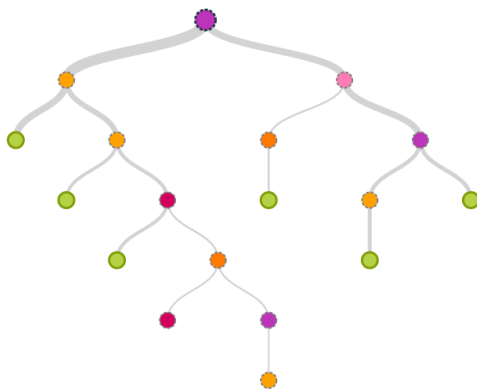


Figure 1: *General concept of Decision Tree. Figure from [?]*

2.1.1 Information Theory

Most popular decision trees use information theory heuristics. It's heuristics is based on the information generated by an attribute respect to the feature's attribute. In the equation (3) D is referred to the feature's attribute and A to the current attribute to be examined. Note that the cardinalities ($card()$) and proportions ($p()$) is been updated according to dataset instances to be examined.

$$I(D) = - \sum_{i \in \text{card}(D)} p(d_i) * \log_2(p(d_i)) \quad (1)$$

$$Rest(A) = \sum_{j \in card(D|A)} \frac{card(D|A_j)}{card(D)} * I(D|A_j) \quad (2)$$

$$Gain(D, A) = I(D) - Rest(A) \quad (3)$$

There is another heuristic who also uses the cardinality to give advantage to attributes with few categories. This heuristic called *Information Gain Ratio* is shown at the equation (5).

$$GainDivisor(D, A) = - \sum_{j \in card(D|A)} \frac{card(D|A_j)}{card(D)} * \log_2 \left(\frac{card(D|A_j)}{card(D)} \right) \quad (4)$$

$$GainRatio(D, A) = \frac{Gain(D, A)}{GainDivisor(D, A)} \quad (5)$$

The technique implemented by *C4.5* to deal with numeric attributes is the binary discretization by intervals distributing the values in two categories. The partition is selected picking the point where the attribute *Gain()* is maximized. One of the points where the feature changes its value.

2.1.2 ID3

ID3 algorithm builds a tree selecting the attribute who maximizes *Gain()* function respecto de the feature's attribute. It only works with nominal data and removes the attributes that have already been used. When there is no more attributes to select, it label the branch with the feature of the remainder instances.

The algorithm uses the bias concept referred as the a-priori intuitions used to learn concepts. The **inductive bias** used by *ID3* is divided in two parts: the *representational bias* who is related with de hypothesis space (logic functions), and the *preference bias* who is based in the idea that the important characteristics will be placed at the top of the tree (because of the *Gain()* heuristic and the Breadth-first search) that has as consequence the tendency of generate low level trees.

2.1.3 C4.5

C4.5 is an extension of ID3 tree. It includes some improvements like branch pruning, working with numerical attributes and unknown values. The pruning works in post-pruning mode, i.e. the tree is generated and after is pruned. The heuristic use to determine if branch is pruned is the improve of pesimistic error rate. There are two pruning operations:

- **Subtree Replacement:** Consist of remove the selected subtree and assign the most frecuent feature to this leaf.
- **Subtree raising:** Consis on replace a tree by one of it's childrens and redistribute the instances. The operation is expensive.

Unknown values are trated spliting the selection between branches and assigning probabilities to each of them. The feature where the instance is classified is the one who maximized the probability. During the train period attributes with unknown values are penalized.

2.2 Rule Based Systems

[TODO]

2.2.1 1R

[TODO]

2.2.2 PRISM

[TODO]

2.2.3 IREP

[TODO]

2.2.4 RIPPER

[TODO]

2.2.5 PART

[TODO]

2.3 Instance Based Learning

[TODO]

2.3.1 K - Nearest Neighbors

[TODO]

2.3.2 IB3

[TODO]

2.4 Bayes Learning

[TODO]

2.4.1 Naive Bayes

[TODO]

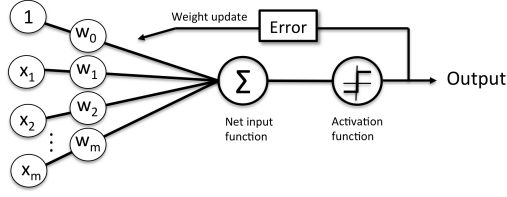


Figure 2: *General concept of perceptron*

2.4.2 K2

[TODO]

2.4.3 TAN

[TODO]

2.5 Linear Classifiers

[TODO]

2.5.1 Linear Regression

[TODO]

2.5.2 Logistic Regression

[TODO]

2.5.3 Support Vector Machines

[TODO]

2.6 Neural Networks

2.6.1 Perceptron model

Perceptron's structures

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}, x = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \quad (6)$$

Output equation

$$z = w_1x_1 + \cdots + w_mx_m = \mathbf{w}^T \mathbf{x} \quad (7)$$

Activation function

$$\phi(z) = \begin{cases} 1 & \text{if } z \geq \theta \\ -1 & \text{otherwise} \end{cases} \quad (8)$$

Update of weight vector

$$w_j := w_j + \Delta w_j \quad (9)$$

$$\Delta w_j = \eta(y^{(i)} - \hat{y}^{(i)})x_j^{(i)} \quad (10)$$

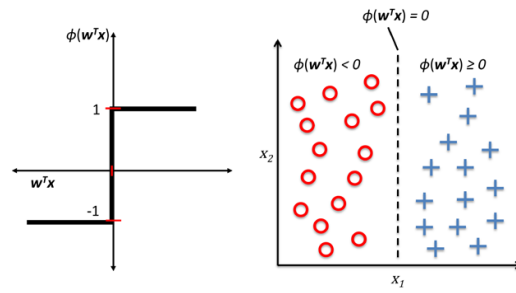


Figure 3: *Activation function*

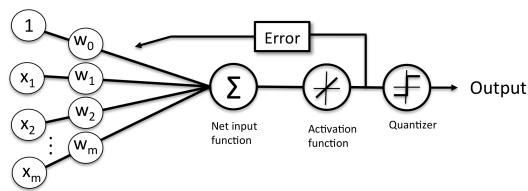


Figure 4: *General concept of adaline perceptron*

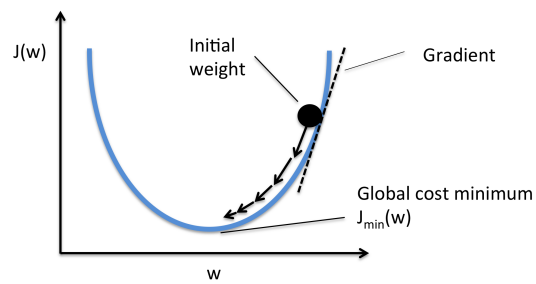


Figure 5: *Gradient descent*

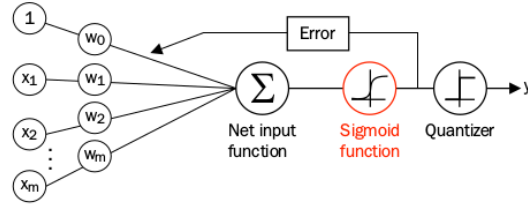


Figure 6: General concept of logistic regression model

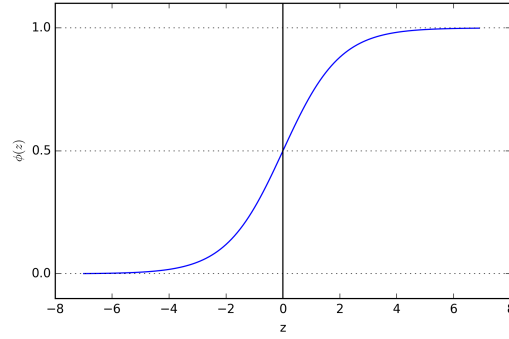


Figure 7: Sigmoid function

2.6.2 Adaptive linear neurons (ADALINE)

Cost function

$$J(w) = \frac{1}{2} \sum_i (y^{(i)} - \phi(z^{(i)}))^2 \quad (11)$$

Update of weight vector

$$w_j := w_j + \Delta w_j \quad (12)$$

$$\Delta w_j = -\eta \nabla J(w) = \eta \sum_i (y^{(i)} - \phi(z^{(i)})) x_j^{(i)} \quad (13)$$

Features standarization

$$x'_j = \frac{x_j - \mu_j}{\sigma_j} \quad (14)$$

2.6.3 Logistic regression model

logit function

$$\text{logit}(p) = \log \frac{p}{(1-p)} \quad (15)$$

logit function according to class $y = 1$

$$\text{logit}(p(y=1|\mathbf{x})) = w_0 x_0 + w_1 x_1 + \dots + w_m x_m = \sum_{i=0}^m w_i x_i = \mathbf{w}^T \mathbf{x} \quad (16)$$

Sigmoid function (activation function)

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad (17)$$

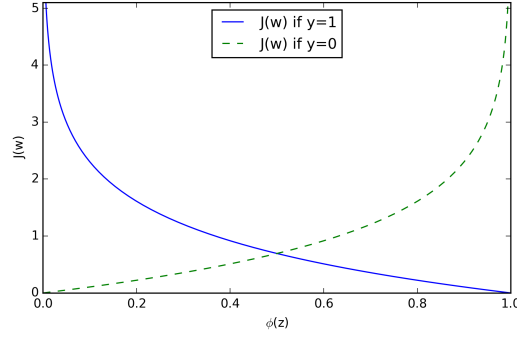


Figure 8: Cost of classification of single-sample instance with logistic regression

Quantizer (unit step function)

$$\hat{y} = \begin{cases} 1 & \text{if } \phi(z) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{y} = \begin{cases} 1 & \text{if } z \geq 0.0 \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

Likelihood L function

$$L(w) = P(y|x; w) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}; w) = \prod_{i=1}^n (\phi(z^{(i)}))^{y^{(i)}} (1 - \phi(z^{(i)}))^{1-y^{(i)}} \quad (19)$$

Log-likelihood function

$$l(w) = \log L(w) = \sum_{i=1}^n [y^{(i)} \log(\phi(z^{(i)})) + (1 - y^{(i)}) \log(1 - \phi(z^{(i)}))] \quad (20)$$

Cost function

$$J(w) = \sum_{i=1}^n [-y^{(i)} \log(\phi(z^{(i)})) - (1 - y^{(i)}) \log(1 - \phi(z^{(i)}))] \quad (21)$$

$$J(\phi(z), y; w) = -y \log(\phi(z)) - (1 - y) \log(1 - \phi(z))$$

Cost of classification

$$J(\phi(z), y; w) = \begin{cases} -\log(\phi(z)) & \text{if } y = 1 \\ -\log(1 - \phi(z)) & \text{if } y = 0 \end{cases} \quad (22)$$

Result cost function

$$J(w) = - \sum_i y^{(i)} \log(\phi(z^{(i)})) + (1 - y^{(i)}) \log(1 - \phi(z^{(i)})) \quad (23)$$

Update of weight vector

$$w := w + \Delta w, \Delta w = -\mu \nabla J(w)$$

$$w_j := w_j + \mu \sum_{i=1}^n (y^{(i)} - \phi(z^{(i)})) x_j^{(i)} \quad (24)$$

L2 regularization

$$\frac{\lambda}{2} \|w\|^2 = \frac{\lambda}{2} \sum_{j=1}^m w_j^2 \quad (25)$$

C inverse regularization parameter

$$C = \frac{1}{\lambda} \tag{26}$$

Cost function with regularization

$$J(w) = C \left[\sum_{i=1}^n (-y^{(i)} \log(\phi(z^{(i)})) - (1 - y^{(i)}) \log(1 - \phi(z^{(i)}))) \right] + \frac{1}{2} \|w\|^2 \tag{27}$$

2.6.4 Single Layer Neural Networks

Simple Perceptron [TODO]

ADALINE [TODO]

2.6.5 Multi Layer Perceptron

[TODO]

2.6.6 Radial Basis Functions

[TODO]

2.6.7 Convolutional Neural Networks

[TODO]

2.6.8 Recurrent Neural Networks

[TODO]

3 Unsupervised Learning

[TODO]