
Machine Learning Formulas*

García Prado, Sergio
sergio@garciparedes.me

Taboada Rodero, 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>

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]

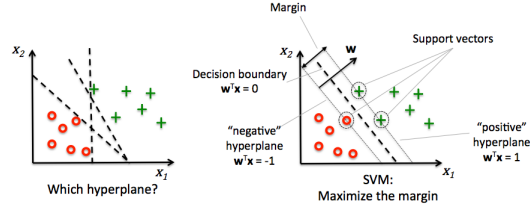


Figure 2: *Support Vector Machine*

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

Maximum margin intuition

$$\begin{aligned} w_0 + w^T x_{pos} &= 1 \\ w_0 + w^T x_{neg} &= -1 \end{aligned} \Rightarrow w^T (x_{pos} - x_{neg}) = 2 \quad (6)$$

Length of weight vector

$$\|w\| = \sqrt{\sum_{j=1}^m w_j^2} \quad (7)$$

Maximum margin intuition normalized

$$\frac{w^T (x_{pos} - x_{neg})}{\|w\|} = \frac{2}{\|w\|} \quad (8)$$

$$y^{(i)}(w_0 + w^T x^{(i)}) \geq 1 \quad \forall i \quad (9)$$

Objective function of SVM

$$\frac{1}{2} \|w\|^2 \quad (10)$$

Nonlinearly separable case using slack variables ξ

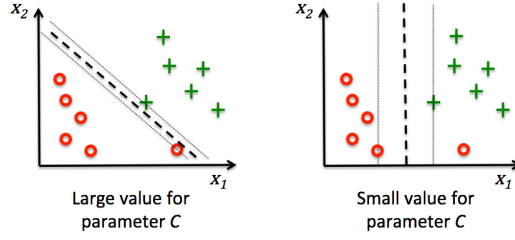


Figure 3: Effects of ξ parameter over same dataset

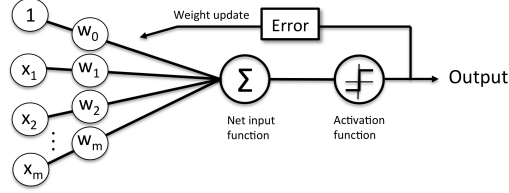


Figure 4: General concept of perceptron

Linea constraints

$$\begin{aligned} \mathbf{w}^T \mathbf{x}^{(i)} &\geq 1 - \xi^{(i)} \text{ if } y^{(i)} = 1 \\ \mathbf{w}^T \mathbf{x}^{(i)} &\leq -1 + \xi^{(i)} \text{ if } y^{(i)} = -1 \end{aligned} \quad (11)$$

Objective function of SVM

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \left(\sum_i \xi^{(i)} \right) \quad (12)$$

2.6 Neural Networks

2.6.1 Perceptron model

Perceptron's structures

$$\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}, \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \quad (13)$$

Ouput equation

$$z = w_1 x_1 + \dots + w_m x_m = \mathbf{w}^T \mathbf{x} \quad (14)$$

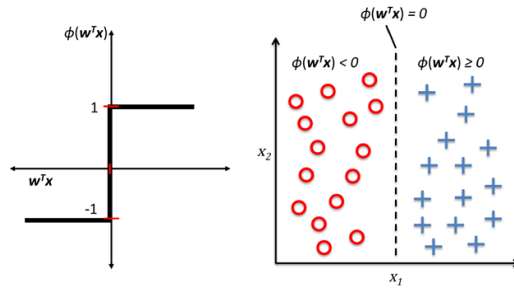


Figure 5: Activation function

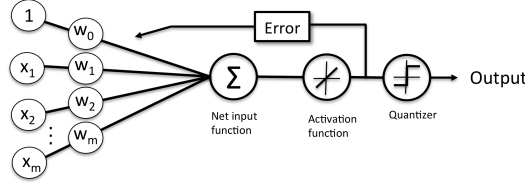


Figure 6: General concept of adaline perceptron

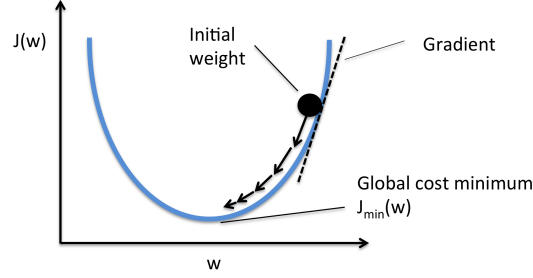


Figure 7: Gradient descent

Activation function

$$\phi(z) = \begin{cases} 1 & \text{if } z \geq \theta \\ -1 & \text{otherwise} \end{cases} \quad (15)$$

Update of weight vector

$$w_j := w_j + \Delta w_j \quad (16)$$

$$\Delta w_j = \eta(y^{(i)} - \hat{y}^{(i)})x_j^{(i)} \quad (17)$$

2.6.2 Adaptive linear neurons (ADALINE)

Cost function

$$J(w) = \frac{1}{2} \sum_i (y^{(i)} - \phi(z^{(i)}))^2 \quad (18)$$

Update of weight vector

$$w_j := w_j + \Delta w_j \quad (19)$$

$$\Delta w_j = -\eta \nabla J(w) = \eta \sum_i (y^{(i)} - \phi(z^{(i)}))x_j^{(i)} \quad (20)$$

Features standarization

$$x'_j = \frac{x_j - \mu_j}{\sigma_j} \quad (21)$$

2.6.3 Logistic regression model

logit function

$$\text{logit}(p) = \log \frac{p}{(1-p)} \quad (22)$$

logit function according to class $y = 1$

$$\text{logit}(p(y=1|\mathbf{x})) = w_0x_0 + w_1x_1 + \dots + w_mx_m = \sum_{i=0}^m w_ix_i = \mathbf{w}^T \mathbf{x} \quad (23)$$

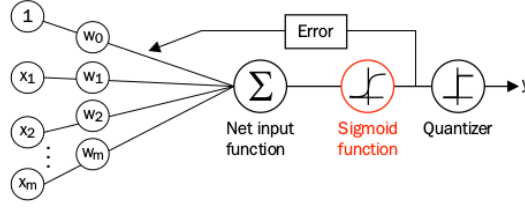


Figure 8: General concept of logistic regression model

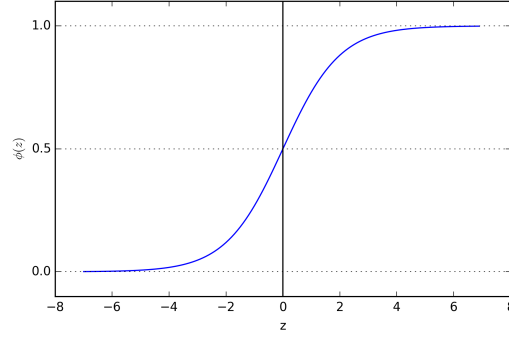


Figure 9: Sigmoid function

Sigmoid function (activation function)

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad (24)$$

Quantizer (unit step function)

$$\hat{y} = \begin{cases} 1 & \text{if } \phi(z) \geq 0.5 \\ 0 & \text{otherwise} \end{cases} \quad (25)$$

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

Likelihood L function

$$L(w) = P(y|x; w) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}; w) = \prod_{i=1}^n (\phi(z^{(i)}))^{y^{(i)}} (1 - \phi(z^{(i)}))^{1-y^{(i)}} \quad (26)$$

Log-likelihood function

$$l(w) = \log L(w) = \sum_{i=1}^n [y^{(i)} \log(\phi(z^{(i)})) + (1 - y^{(i)}) \log(1 - \phi(z^{(i)}))] \quad (27)$$

Cost function

$$J(w) = \sum_{i=1}^n [-y^{(i)} \log(\phi(z^{(i)})) - (1 - y^{(i)}) \log(1 - \phi(z^{(i)}))] \quad (28)$$

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

Cost of classification

$$J(\phi(z), y; w) = \begin{cases} -\log(\phi(z)) & \text{if } y = 1 \\ -\log(1 - \phi(z)) & \text{if } y = 0 \end{cases} \quad (29)$$

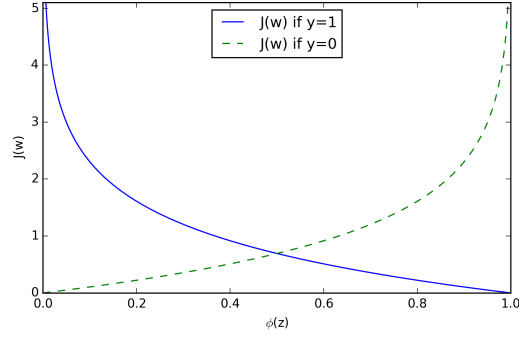


Figure 10: Cost of classification of single-sample instance with logistic regression

Result cost function

$$J(w) = - \sum_i y^{(i)} \log(\phi(z^{(i)})) + (1 - y^{(i)}) \log(1 - \phi(z^{(i)})) \quad (30)$$

Update of weight vector

$$\begin{aligned} w &:= w + \Delta w, \Delta w = -\mu \nabla J(w) \\ w_j &:= w_j + \mu \sum_{i=1}^n (y^{(i)} - \phi(z^{(i)})) x_j^{(i)} \end{aligned} \quad (31)$$

L2 regularization

$$\frac{\lambda}{2} \|w\|^2 = \frac{\lambda}{2} \sum_{j=1}^m w_j^2 \quad (32)$$

C inverse regularization parameter

$$C = \frac{1}{\lambda} \quad (33)$$

Cost function with regularization

$$J(w) = C \left[\sum_{i=1}^n (-y^{(i)} \log(\phi(z^{(i)})) - (1 - y^{(i)}) \log(1 - \phi(z^{(i)}))) \right] + \frac{1}{2} \|w\|^2 \quad (34)$$

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]

References

- [CCAG17] Teodoro Calonge Cano and Carlos Javier Alonso González. Técnicas de Aprendizaje Automático, 2016/17.
- [ML13] Big ML. A New Way to Visualize Decision Trees. <https://blog.bigml.com/2013/04/19/a-new-way-to-visualize-decision-trees/>, 2013.
- [Ols16] Dr. Randal S. Olson. Python Machine Learning, 2016.