

An Introduction to Machine Learning

Fabio A. González Ph.D.

Depto. de Ing. de Sistemas e Industrial
Universidad Nacional de Colombia, Bogotá

February 9, 2018

1 Introduction

Example

How to State the Learning Problem?

How to Solve the Learning Problem?

2 Patterns and Generalization

Generalizing from patterns

Overfitting/ Overlearning

How to Measure the Quality of a Solution?

3 Learning Problems

Supervised

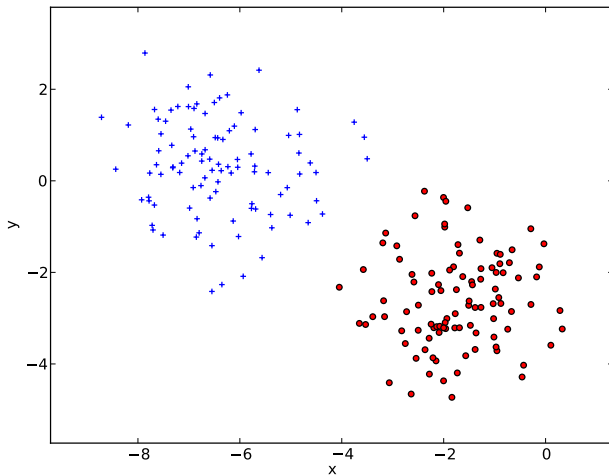
Non-supervised

Active

On-line

4 Learning Techniques

Two class classification problem



Introduction

Example

How to State the
Learning Problem?

How to Solve the
Learning Problem?

Patterns and Generalization

Learning Problems

Learning Techniques

How to solve it?

- We need to build a prediction function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ such that::

$$\text{Prediction}(x, y) = \begin{cases} C_1 & \text{si } f(x, y) \geq 0 \\ C_2 & \text{si } f(x, y) < 0 \end{cases}$$

- Training set: $D = \{((x_1, y_1), l_1), \dots, ((x_n, y_n), l_n)\}$

- Example:

$$D = \{((1, 2), -1), ((1, 3), -1), ((3, 1), 1), \dots\}$$

- Loss function:

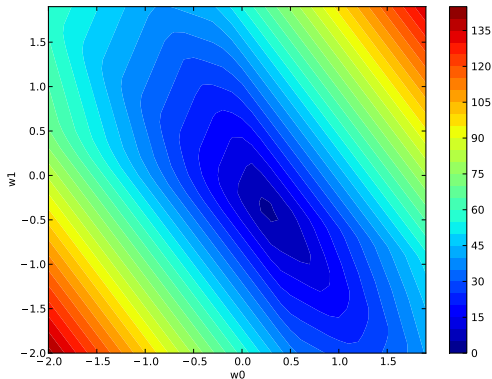
$$L(f, D) = \sum_{(x_i, y_i, l_i) \in D} \frac{|\text{sign}(f(x_i, y_i)) - l_i|}{2}$$

L_1 Error loss

$$f(x, y) = w_1 x + w_0 y$$

$$L(f, D) = \frac{1}{2} \sum_{(x_i, y_i, l_i) \in D} |f(x_i, y_i) - l_i|$$

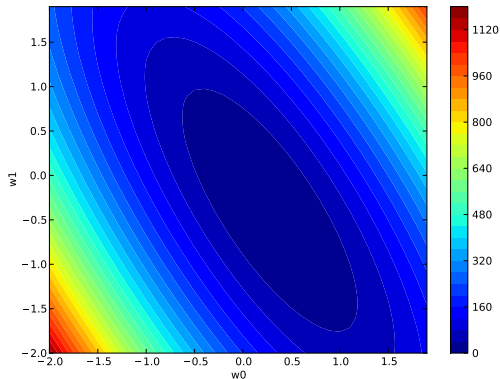
- Are there other alternative loss functions?



Square error loss

$$f(x, y) = w_1 x + w_0 y$$

$$L(f, D) = \frac{1}{2} \sum_{(x_i, y_i, l_i) \in D} (f(x_i, y_i) - l_i)^2$$



Learning as optimization

- General optimization problem:

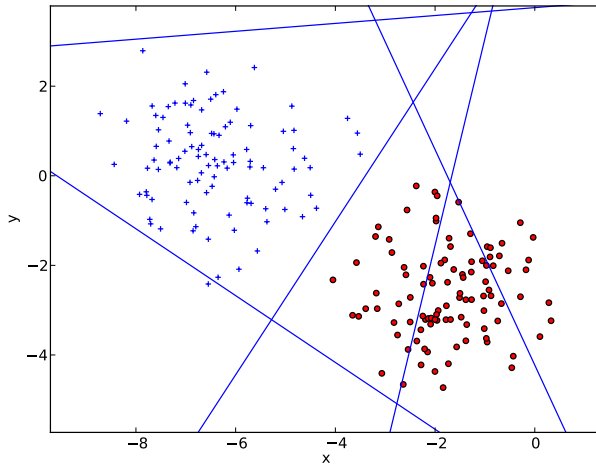
$$\min_{f \in H} L(f, D)$$

- Two Class 2D classification using linear functions:

$$H = \{f : f(x, y) = w_2x + w_1y + w_0, \forall w_0, w_1, w_2 \in \mathbb{R}\}$$

$$\min_{f \in H} L(f, D) = \min_{W \in \mathbb{R}^3} \frac{1}{2} \sum_{(x_i, y_i) \in D} (w_2x_i + w_1y_i + w_0 - l_i)^2$$

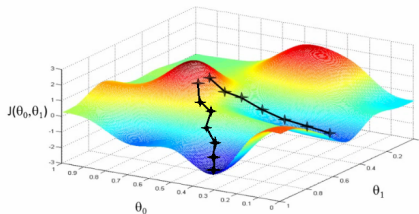
Hypothesis space



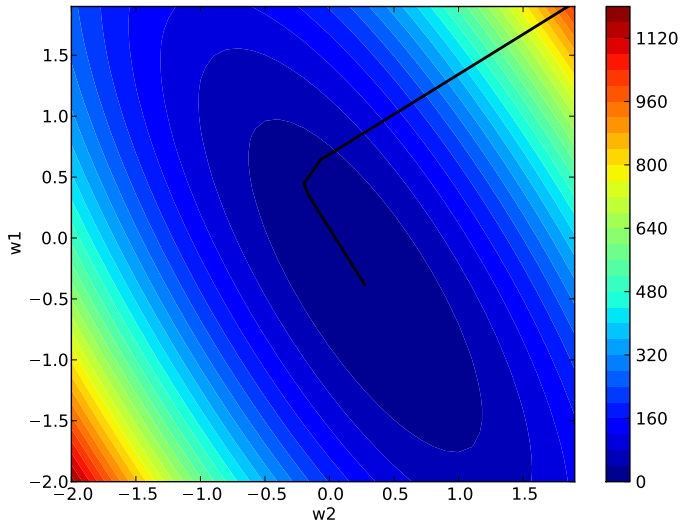
Gradient descent

Iterative optimization of the loss function:

```
initialize  $W^0 = w_0, w_1, w_2$   
 $k \leftarrow 0$   
repeat  
     $k \leftarrow k + 1$   
     $W^k \leftarrow$   
         $W^{k-1} - \eta(k) \nabla L(f_{W^{k-1}}, S)$   
until  $|\eta(k) \nabla L(f_{W^{k-1}}, S)| < \Theta$ 
```



Gradient descent iteration example (1)



Introduction

Example

How to State the
Learning Problem?

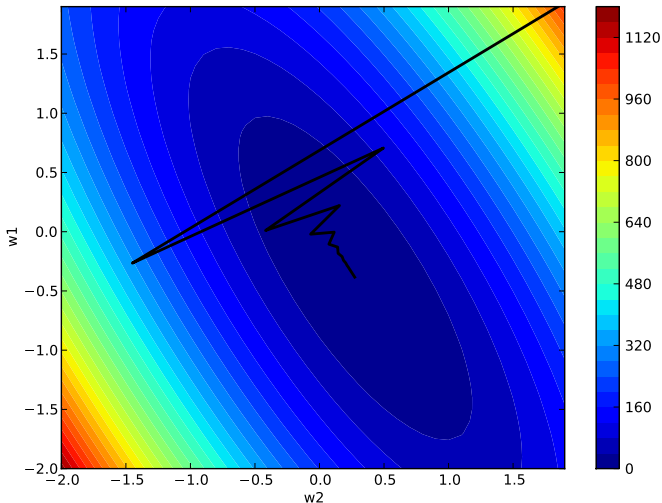
How to Solve the
Learning Problem?

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Gradient descent iteration example (2)



Introduction

Example

How to State the
Learning Problem?

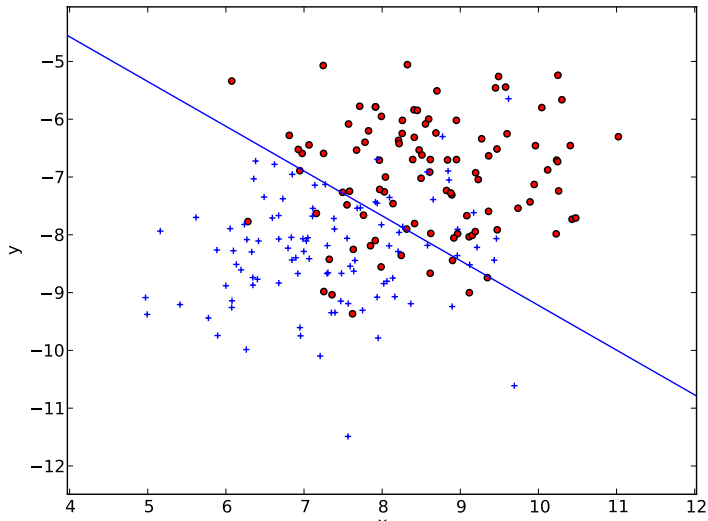
How to Solve the
Learning Problem?

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Non-separable data



What is a pattern?

- Data regularities
- Data relationships
- Redundancy
- Generative model

Learning a boolean function

x_1	x_2	f_1	f_2	...	f_{16}
0	0	0	0	...	1
0	1	0	0	...	1
1	0	0	0	...	1
1	1	0	1	...	1

- How many Boolean functions of n variables are?
- How many candidate functions are removed by a sample?
- Is it possible to generalize?

Inductive bias

Introduction

Patterns and Generalization

Generalizing from
patterns

Overfitting/
Overlearning

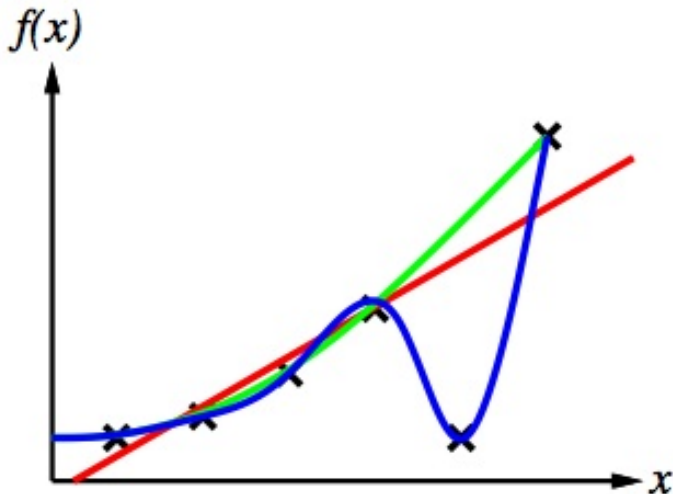
How to Measure the
Quality of a
Solution?

Learning Problems

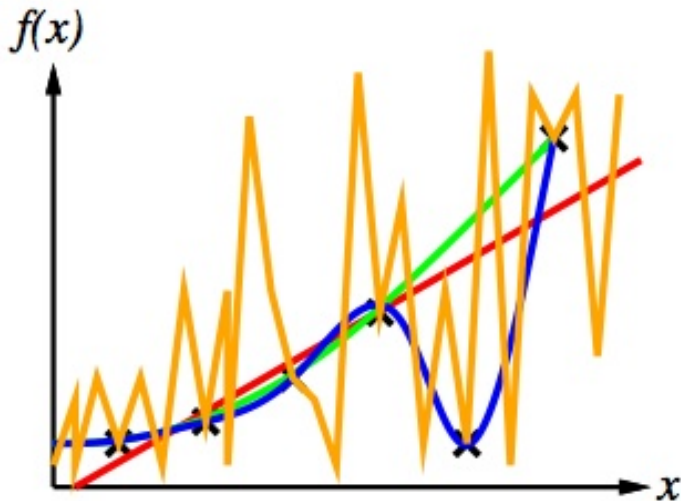
Learning Techniques

- In general, the learning problem is *ill-posed* (more than one possible solution for the same particular problem, solutions are sensitive to small changes on the problem)
- It is necessary to make additional assumptions about the kind of pattern that we want to learn
- **Hypothesis space**: set of valid patterns that can be learnt by the algorithm

What is a good pattern?



What is a good pattern?



Occam's razor

from Wikipedia:

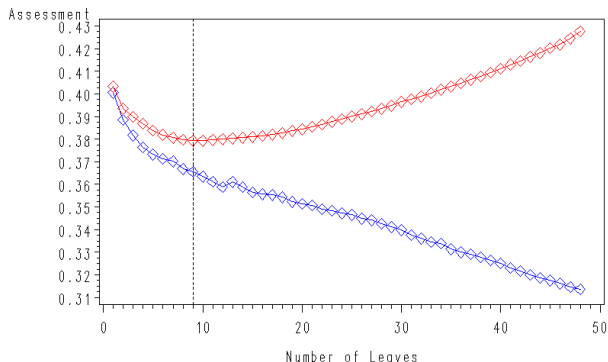
Occam's razor (also spelled Ockham's razor) is a principle attributed to the 14th-century English logician and Franciscan friar William of Ockham. The principle states that the explanation of any phenomenon should make as few assumptions as possible, eliminating, or "shaving off", those that make no difference in the observable predictions of the explanatory hypothesis or theory. The principle is often expressed in Latin as the *lex parsimoniae* (law of succinctness or parsimony).

"All things being equal, the simplest solution tends to be the best one."

Training error vs generalization error

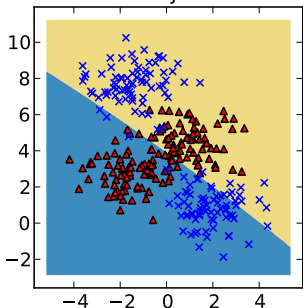
- The loss function measures the error in the training set
- Is this a good measure of the quality of the solution?

Average Square Error (Gini index)

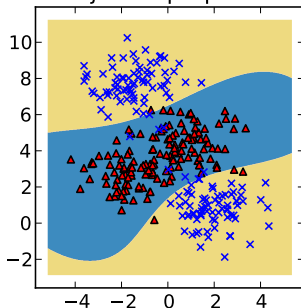


Over-fitting and under-fitting

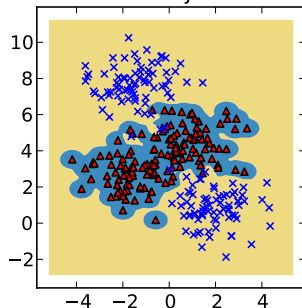
Subajuste



Ajuste apropiado



Sobreajuste



Generalization error

- Generalization error:

$$E[(L(f_w, S))]$$

- How to control the generalization error during training?
 - Cross validation
 - Regularization

- Vapnik, 1995:

$$R(\alpha) = \int \frac{1}{2} |y - f(\mathbf{x}, \alpha)| dP(\mathbf{x}, y)$$

$$R_{emp}(\alpha) = \frac{1}{2l} \sum_{i=1}^l |y_i - f(\mathbf{x}_i, \alpha)|.$$

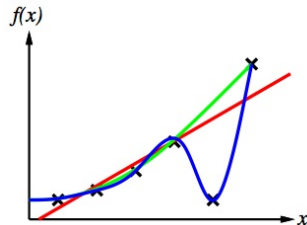
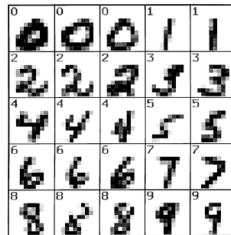
$$R(\alpha) \leq R_{emp}(\alpha) + \sqrt{\left(\frac{h(\log(2l/h) + 1) - \log(\eta/4)}{l} \right)}$$

Types

- Supervised learning
- Non-supervised learning
- Semi-supervised learning
- Active/reinforcement learning
- On-line learning

- **Fundamental problem:** to find a function that relates a set of inputs with a set of outputs
- Typical problems:
 - Classification
 - Regression

Supervised learning



Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

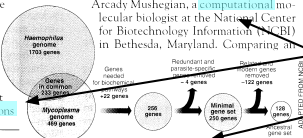
Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers** game, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

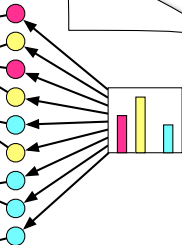


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

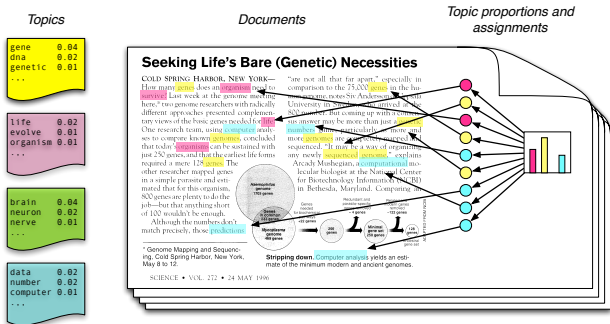
Stripping down. **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Non-supervised learning



- There are not labels for the training samples
- **Fundamental problem:** to find the subjacent structure of a training data set
- Typical problems: clustering, probability density estimation, dimensionality reduction, latent topic analysis, data compression
- Some samples may have labels, in that case it is called semi-supervised learning

Active/reinforcement learning

- Generally, it happens in the context of an agent acting in an environment
- The agent is not told whether it has made the right decision or not
- The agent is punished or rewarded (not necessarily in an immediate way)
- **Fundamental problem:** to define a policy that allows to maximize the positive stimulus (reward)



On-line learning

- Only one pass through the data
 - big data volume
 - real time
- It may be supervised or unsupervised
- **Fundamental problem:** to extract the maximum information from data with minimum number of passes

Representative techniques

- Computational
 - Decision trees
 - Nearest-neighbor classification
 - Graph-based clustering
 - Association rules
- Statistical
 - Multivariate regression
 - Linear discriminant analysis
 - Bayesian decision theory
 - Bayesian networks
 - K-means
- Computational-Statistical
 - SVM
 - AdaBoost
- Bio-inspired
 - Neural networks
 - Genetic algorithms
 - Artificial immune systems



Alpaydin, E. 2010 Introduction to Machine Learning (Adaptive Computation and Machine Learning). The MIT Press. (Chap 1,2)

Introduction

Patterns and
Generalization

Learning
Problems

Learning
Techniques