An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

# An Introduction to Machine Learning

### Fabio A. González Ph.D.

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

August 26, 2010

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

# Content

**1** Patterns and Generalization
   Generalizing from patterns
   Overfitting/ Overlearning

**2** Learning Problems
   Supervised
   Non-supervised
   Active
   On-line

**3** Learning Techniques

**4** Main Questions
   How to State the Learning Problem?
   How to Solve the Learning Problem?
   How to Measure the Quality of a Solution?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization
Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization
Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# What is a pattern?

- Data regularities

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization
Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# What is a pattern?

- Data regularities
- Data relationships

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# What is a pattern?

- Data regularities
- Data relationships
- Redundancy

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# What is a pattern?

- Data regularities
- Data relationships
- Redundancy
- Generative model

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# 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?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# 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?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# 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?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# Inductive bias

- 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)

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# Inductive bias

- 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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# Inductive bias

- 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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization
Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization
Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# What is a good pattern?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization
Generalizing from
patterns
Overfitting/
Overlearning

Learning
Problems

Learning
Techniques

Main
Questions

# What is a good pattern?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization
Generalizing from
patterns
**Overfitting/
Overlearning**

Learning
Problems

Learning
Techniques

Main
Questions

# 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."**

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
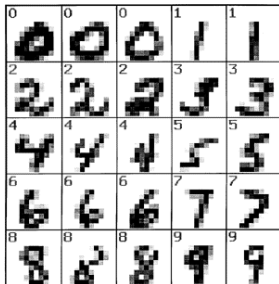Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

# Types

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions
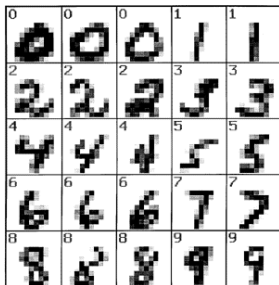
# Supervised learning

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions
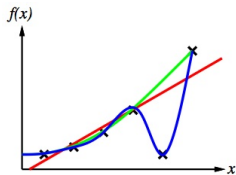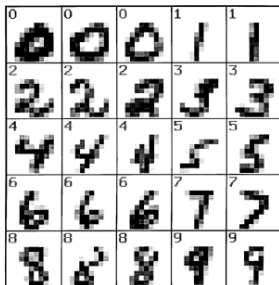
# Supervised learning



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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Supervised
Non-supervised
Active
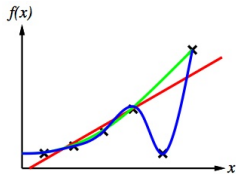On-line

Learning
Techniques

Main
Questions

# Supervised learning

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Supervised
Non-supervised
Active
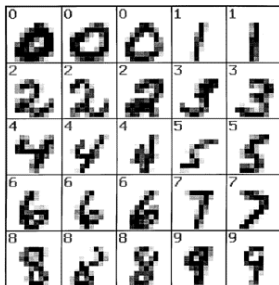On-line

Learning
Techniques

Main
Questions

# Supervised learning

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Non-supervised learning

- There are not labels for the training samples

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Non-supervised learning

- There are not labels for the training samples
- **Fundamental problem**: to find the subjacent structure of a training data set

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
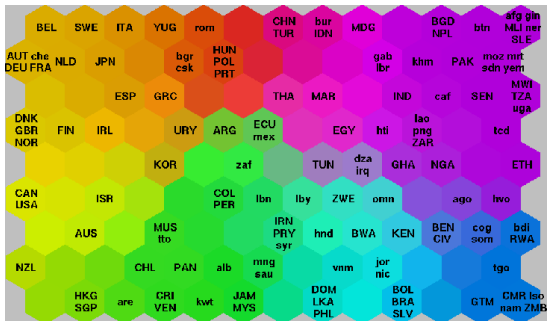Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# 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, data compression

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
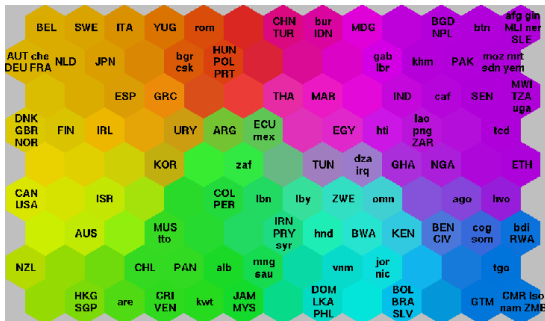Supervised
Non-supervised
Active
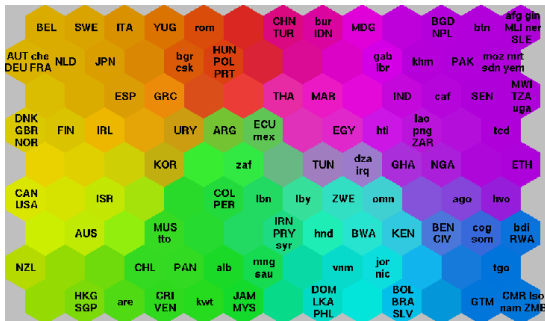On-line

Learning
Techniques

Main
Questions

# 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, data compression
- Some samples may have labels, in that case it is called semi-supervised learning

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Active/reinforcing learning

- Generally, it happens in the context of an agent acting in an environment

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Active/reinforcing learning

- Generally, it happens in the context of an agent acting in an environment
- The agent is not told whether it has make the right decision or not

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Active/reinforcing learning

- Generally, it happens in the context of an agent acting in an environment
- The agent is not told whether it has make the right decision or not
- The agent is punished or rewarded (not necessarily in an immediate way)

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Active/reinforcing learning

- Generally, it happens in the context of an agent acting in an environment

- The agent is not told whether it has make 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)

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# On-line learning

- Only one pass through the data

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# On-line learning

- Only one pass through the data
  - big data volume

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# On-line learning

- Only one pass through the data
    - big data volume
    - real time

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# On-line learning

- Only one pass through the data
    - big data volume
    - real time
- It may be supervised or unsupervised

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems
Supervised
Non-supervised
Active
On-line

Learning
Techniques

Main
Questions

# 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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

# 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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Two Class Classification Problem



• The idea is to buid a linear classifier function, $f : \mathbb{R}^2 \to \mathbb{R}$, such that:

$$f(x, y) = \begin{cases} < 0 & \text{if } (x, y) \in C_0 \\ > 0 & \text{if } (x, y) \in C_1 \end{cases}$$

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?
How to Solve the
Learning Problem?
How to Measure the
Quality of a
Solution?

# Loss Function

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Loss Function

- Training set: $S = \{((x_1, y_1), l_1), \ldots, ((x_n, y_n), l_n)\}$
- Loss function:

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Loss Function

- Training set: $S = \{((x_1, y_1), l_1), \ldots, ((x_n, y_n), l_n)\}$
- Loss function:

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

- Are there other alternative loss functions?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Square Error Loss

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# $L_1$ Error Loss

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Learning as Optimization

- General optimization problem:

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
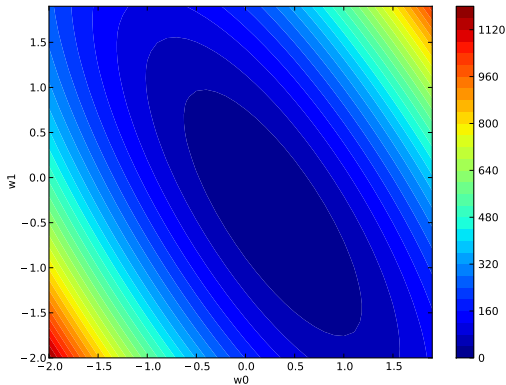Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Learning as Optimization

- General optimization problem:

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

- Two Class 2D Classification:

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

$$\min_{f \in H} L(f, S) = \min_{W \in \mathbb{R}^3} \frac{1}{2} \sum_{(x_i, y_i) \in S} (w_2 x_i + w_1 y_i + w_0 - l_i)^2$$

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions
How to State the
Learning Problem?
**How to Solve the
Learning Problem?**
How to Measure the
Quality of a
Solution?

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Gradient Descent

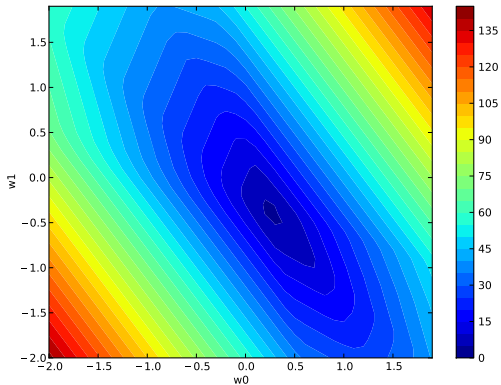Iterative optimization of the loss function:

```
initialize  W^0 = w_0, w_1, w_2
k ← 0
repeat
    k ← k + 1
    W^k ← W^{k-1} - η(k)∇L(f_{W^{k-1}}, S)
until |η(k)∇L(f_{W^{k-1}}, S)| < Θ
```

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Gradient Descent Iteration Example (1)

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Gradient Descent Iteration Example (2)

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?
How to Solve the
Learning Problem?
**How to Measure the
Quality of a
Solution?**

# Outline

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

**How to Measure the
Quality of a
Solution?**

# Training Error vs Generalization Error

- The loss function measures the error in the training set

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions
How to State the
Learning Problem?
How to Solve the
Learning Problem?
How to Measure the
Quality of a
Solution?

# 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?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
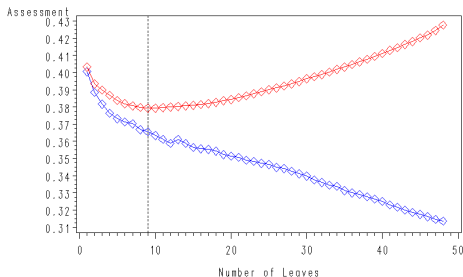Learning Problem?

How to Measure the
Quality of a
Solution?

# 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)



- Training
  Validation

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Generalization Error

- Generalization error:

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Generalization Error

- Generalization error:

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

- How to control the generalization error during training?

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

# Generalization Error

- Generalization error:

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

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?
How to Solve the
Learning Problem?
**How to Measure the
Quality of a
Solution?**

# Generalization Error

- Generalization error:

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

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?
How to Solve the
Learning Problem?
How to Measure the
Quality of a
Solution?

# Regularization

- Vapnik, 1995:

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

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

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

An
Introduction
to Machine
Learning

Fabio A.
González
Ph.D.

Patterns and
Generalization

Learning
Problems

Learning
Techniques

Main
Questions

How to State the
Learning Problem?

How to Solve the
Learning Problem?

How to Measure the
Quality of a
Solution?

Alpaydin, E. 2004 Introduction to Machine Learning
(Adaptive Computation and Machine Learning). The MIT
Press. (Cap 1,2)