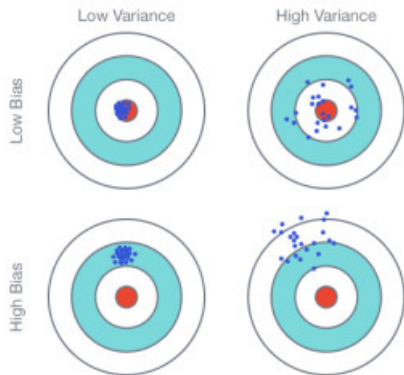# Apprendimento Automatico

## Practical Issues
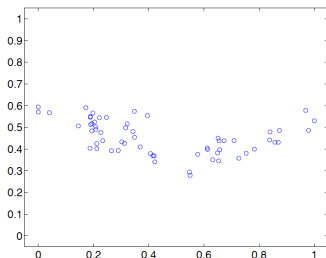
Fabio Aiolli

16 November 2016

# Bias and Variance

- The BIAS measures the *distortion* of an estimate
- The VARIANCE measures the *dispersion* of an estimate

# Some practical issues

**Underfitting/Overfitting and learning parameters**

- Suppose we have some data (60 points) that we want to fit a curve to
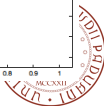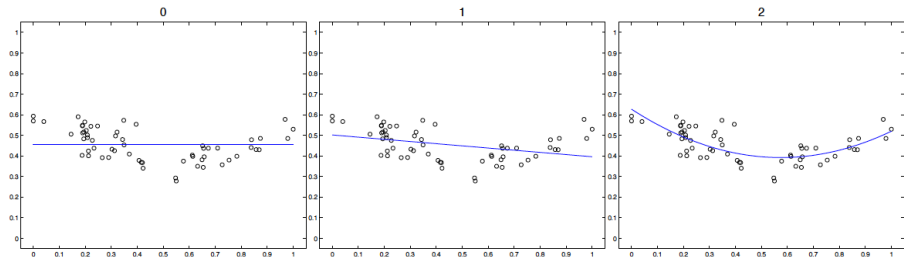


- Let fit a polynomial, of the form

$$y = w_0 + w_1 x + w_2 x^2 + ... + w_p x^p$$
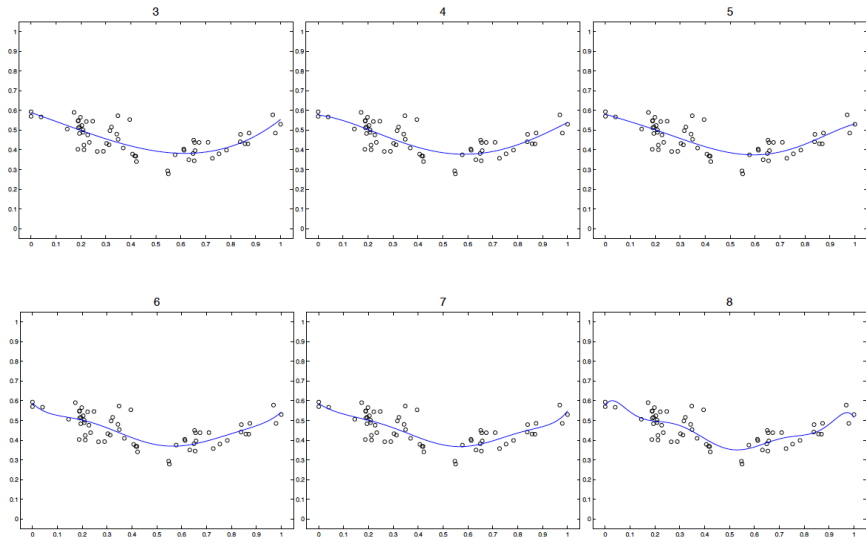
# Some practical issues

**Underfitting/Overfitting and learning parameters**

- How to choose $p$ ? (Hypothesis Space)
- For various $p$, we can find and plot the best polynomial, in terms of minimizing the Empirical Error (Mean Squared Error in this case)
- Here are the solutions for different values of $p$

# Some practical issues

# Some practical issues

# Some practical issues

# Some practical issues

**Underfitting/Overfitting and learning parameters**

- Here is a summary of the Empirical Error ... and the Empirical Error over some new TEST data (100,000 extra points) from the same distribution, as a function of $p$:

# Some practical issues

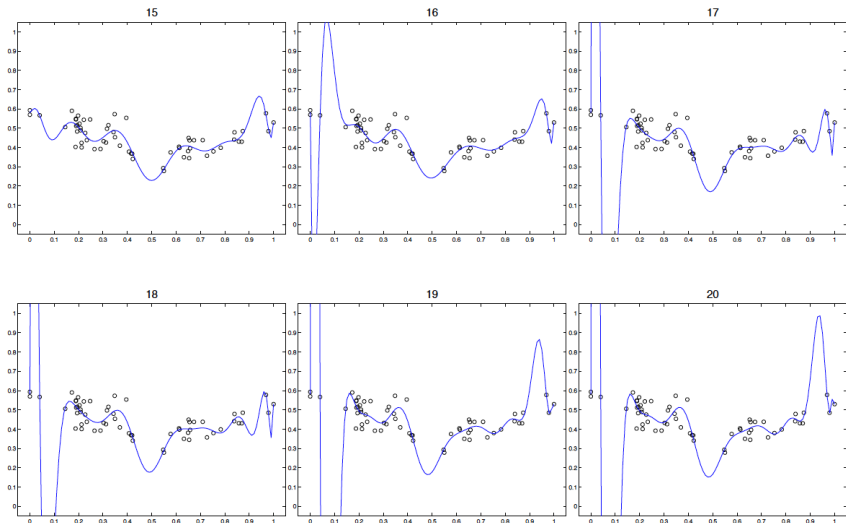**Underfitting/Overfitting and learning parameters**

- For very low $p$, the model is very simple, and so cannot capture the full complexities of the data (Underfitting! also called **bias**)

- For very high $p$, the model is complex, and so tends to overfit to spurious properties of the data (Overfitting! also called **variance**)

Unfortunately we cannot use the test set to pick up the right value of $p$!

> PRACTICAL PROBLEM: how can we use the training set to set $p$ ?

# Model Selection and Hold-out

We can hold out some of our original training data

**Hold-out procedure**

1. A small subset of $Tr$, called the validation set (or hold-out set), denoted $Va$, is identified

2. A classifier/regressor is learnt using examples in $Tr - Va$

3. Step 2 is performed with different values of the parameter(s) (in our example, $p$), and tested against the hold-out sample

In an operational setting, after parameter optimization, one typically re-trains the classifier on the entire training corpus, in order to boost effectiveness (debatable step!)

It is possible to show that the evaluation performed in Step 2 gives an unbiased estimate of the error performed by a classifier learnt with the same parameter(s) and with training set of cardinality $|Tr| - |Va| < |Tr|$

# K-fold Cross Validation

An alternative approach to model selection (and evaluation) is the K-fold cross-validation method

**K-fold CV procedure**

1. $K$ different classifiers/regressors $h_1, h_2, \ldots, h_k$ are built by partitioning the initial corpus $Tr$ into $k$ disjoint sets $Va_1, \ldots, Va_k$ and then iteratively applying the Hold-out approach on the $k$-pairs ($Tr_i = Tr - Va_i, Va_i$)

2. Final error is obtained by individually computing the errors of $h_1, \ldots, h_k$, and then averaging the individual results

The above procedure is repeated for different values of the parameter(s) and the setting (model) with smaller final error is selected

The special case $k = |Tr|$ of $k$-fold cross-validation is called **leave-one-out** cross-validation

# Back to our example

- Let's apply 5-fold CV



- Minimum error reached for $p = 3$, rather than the optimal $p = 12$
- Clearly, cross validation is no substitute for a large test set. However if we only have a limited training set, it is often the best option available.

## Back to our example

Why cross-validation selected a simpler model than optimal ?

- Notice that with 60 points, 5-fold cross validation effectively tries to pick the polynomial that makes the best bias-variance tradeoff for 48 ($60 * \frac{4}{5}$) points

- With 10-fold cross validation, it would instead try to pick the best polynomial for 54 ($60 * \frac{9}{10}$) points

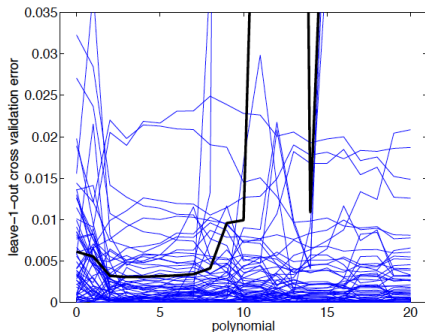> Thus, cross validation biases towards simpler models

- **leave-one-out** cross-validation reduces this tendency to the minimum possible by doing 60-fold cross validation

# Back to our example

- So let's try **leave-one-out** cross-validation



- We still get $p = 3$!

- Cross validation is a good technique, but it doesn't work miracles: there is only so much information in a small dataset.

# Analysis of Cross Validation

- What happens varying $k$?
- For highest $k$'s we have larger training sets, hence less bias! Smaller validation sets, hence more variance!
- For lower $k$'s we have smaller training sets, hence more bias! Larger validation sets, hence less variance!

# In general...

Almost all learning algorithms have (hyper)parameters!

- Support Vector Machines: $C$, type of kernel (polynomial, RBF, etc.), kernel parameter (degree of polynomial, width of RFB, etc.)
- Neural Networks: nonlinear/linear neurons, number of hidden units, $\eta$, other learning parameters we have not discussed (e.g., momentum $\mu$)

Hold-out or Cross-Validation can be used to select the "optimal" values for the (hyper)parameters (i.e., select the "optimal" model).

**After** model selection, the training set is used to evaluate the goodness of the selected model

# Evaluation for unbalanced data - Beyond accuracy

Classification accuracy:

- Very common in ML,
- Proportion of correct decisions,
- Not appropriate when the number of positive examples is much lower than the number of negative examples (or viceversa)
- Precision, Recall and $F_1$ are better in these cases!

## Contingency table

|  | Relevant | Not Relevant |
|---|---|---|
| Retrieved | True Positive (TP) | False Positive (FP) |
| Not Retrieved | False Negative (FN) | True Negative (TN) |

$$\pi = \frac{TP}{TP + FP} \qquad \rho = \frac{TP}{TP + FN}$$

why not using the accuracy $\alpha = \frac{TP + TN}{TP + TN + FP + FN}$ ?

# Effectiveness for Unbalanced Classification

If relevance is assumed to be binary-valued, effectiveness is typically measured as a combination of

- Precision is the "degree of soundness" of the system:
  $$P(\text{RELEVANT}|\text{RETURNED})$$
- Recall is the "degree of completeness" of the system:
  $$P(\text{RETURNED}|\text{RELEVANT})$$

# F Measure

How can one trade-off between precision and recall?
F-measure (weighted harmonic mean of the precision and the recall)

$$F_\beta = \frac{(1 + \beta^2)\pi\rho}{\beta^2\pi + \rho}$$

$\beta < 1$ emphasizes precision

$$\beta = 1 \ \Rightarrow \ F_1 = 2\frac{\pi\rho}{\pi + \rho}$$

# Many other algorithms/approaches and error measures

Please, remember that we have only presented **some** of the proposed learning algorithms/approaches, as well as possibile learning tasks and related error functions (loss functions)!

Just to name a few popular learning approaches

- Probabilistic approaches
- Decision trees
- Boosting
- Ensembles/Committees
- Genetic algorithms
- Prototype methods

Unsupervised Approaches!! Other model selection criteria: BIC, MDL, Bootstrap,...