

Practical Issues

Machine Learning - A.Y. 2022/23 - Padova



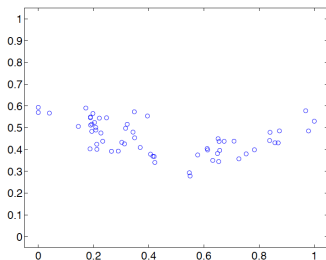
Fabio Aioli

Novembre 21st, 2022



Underfitting/Overfitting and learning parameters

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



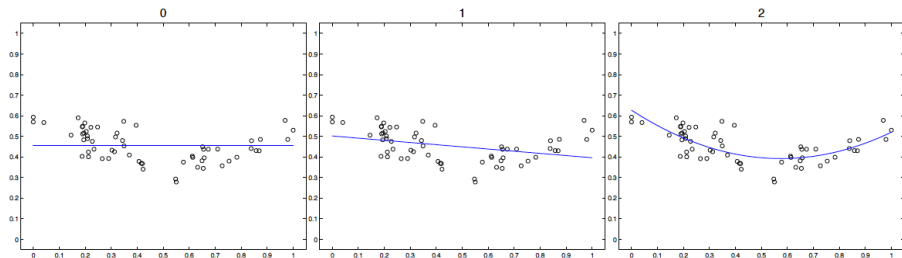
- Let fit a polynomial, of the form

$$y = w_0 + w_1x + w_2x^2 + \dots + w_px^p$$

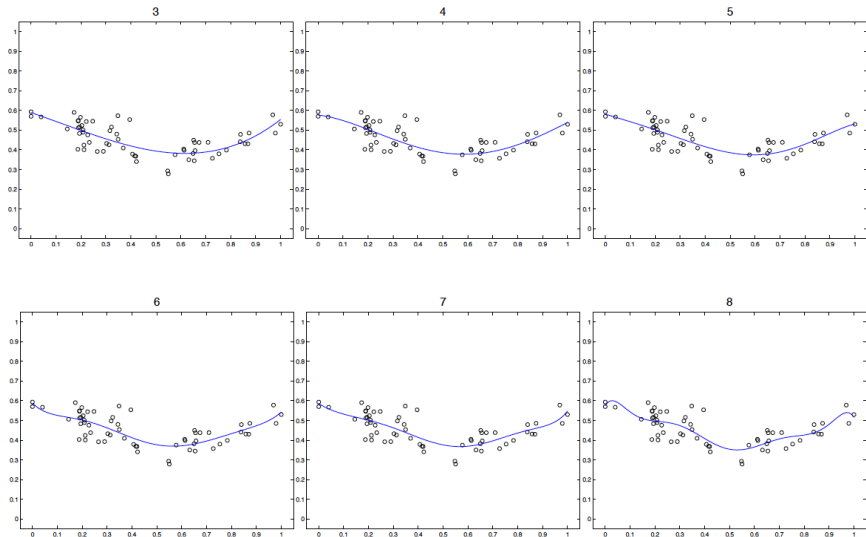


Underfitting/Overfitting and learning parameters

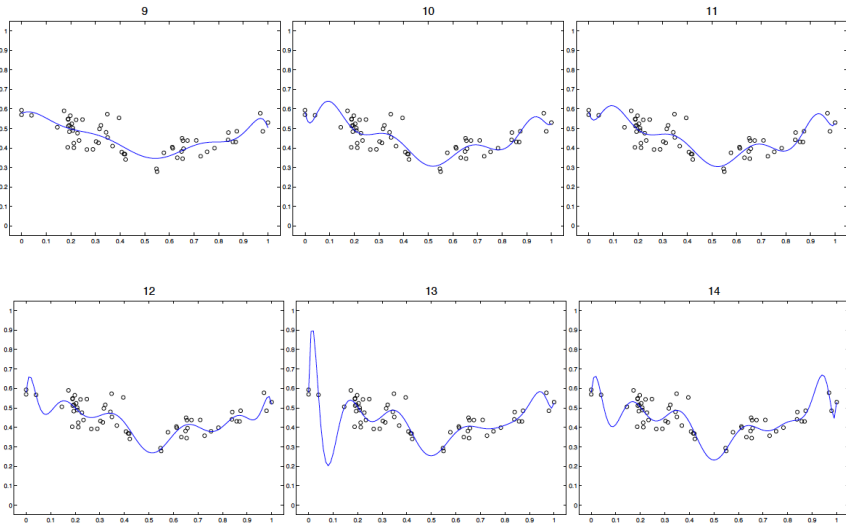
- How to choose p ? (Hypothesis Space)
- For various p , we can find and plot the best polynomial, in terms of minimizing the Empirical Error (Mean Squared Error in this case)
- Here are the solutions for different values of p



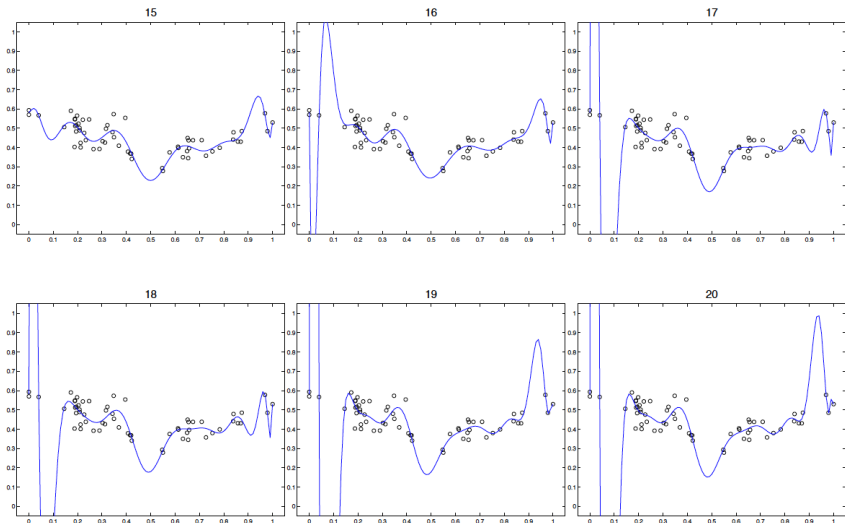
Some practical issues



Some practical issues



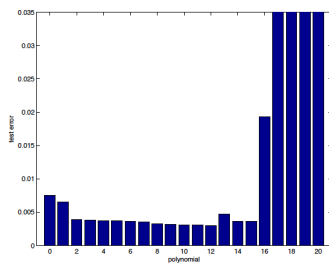
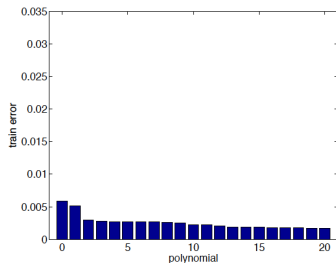
Some practical issues



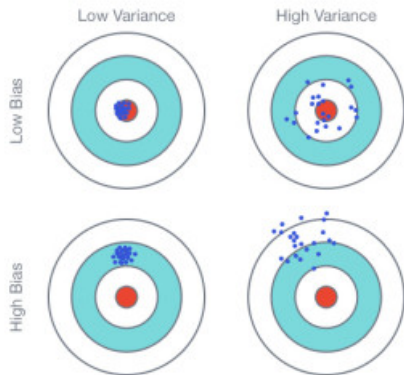


Underfitting/Overfitting and learning parameters

- Here is a summary of the training error ... and the error on some new TEST data (100,000 extra points) from the same distribution, as a function of p :



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



Bias-Variance Decomposition for Regression



Let $y = f(\mathbf{x}) + \epsilon$ the target function, where ϵ has zero mean, variance σ^2 .

We want to find a function $\hat{f}(\mathbf{x})$ that approximates the true function $f(\mathbf{x})$ as well as possible, by means of some learning algorithm.

Given any pair (\mathbf{x}, y) , the following holds:

$$\mathbb{E}[(y - \hat{f}(\mathbf{x}))^2] = \left(\text{Bias}[\hat{f}(\mathbf{x})]\right)^2 + \text{Var}[\hat{f}(\mathbf{x})] + \sigma^2$$

where,

$$\text{Bias}[\hat{f}(\mathbf{x})] = \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x})$$

$$\text{Var}[\hat{f}(\mathbf{x})] = \mathbb{E}[\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})]]^2$$

and the expectations refer to different sets of training data.



Underfitting/Overfitting and learning parameters

- For very low p , the model is very simple, and so cannot capture the full complexity of the data (Underfitting! also called **bias error**)
- For very high p , the model is complex, and so tends to overfit to spurious properties of the data (Overfitting! also called **variance error**)

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 the value of p ?



Model Selection and Hold-out

We can keep some of our training data out

Hold-out procedure

- 1 A small subset of Tr , called the validation set (or hold-out set), denoted Va , is kept apart;
- 2 A classifier/regressor is trained using examples in $Tr - Va$;
- 3 Step 2 is performed with different values of the hyperparameter(s) (in our example, the value p), and tested against the hold-out sample.

It is possible to show (Hoeffding's Inequality) that the evaluation performed in step 2 gives an **unbiased estimate** of the error made by a classifier trained with the same values of the hyperparameter(s) and with training set of cardinality $|Tr| - |Va| < |Tr|$.

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!).



K-fold Cross Validation

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

K-fold CV procedure

- 1 K different classifiers/regressors h_1, h_2, \dots, h_k are built by partitioning the initial corpus Tr into k disjoint sets Va_1, \dots, 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, \dots, h_k , and then averaging the results

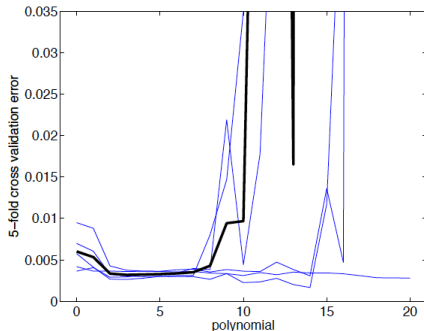
The above procedure is repeated for different values of the hyperparameter(s) and the setting (model) with the smallest final error is selected

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.



- What happens varying k ?
- For higher 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!



Almost all learning algorithms have (hyper)parameters!

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

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



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)

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

Alternative measures (precision and recall):

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



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:
 $P(\text{RELEVANT}|\text{RETRIEVED})$
- **Recall** is the “degree of completeness” of the system:
 $P(\text{RETRIEVED}|\text{RELEVANT})$



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

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



Multiclass classification consists of a classification task with more than two classes; e.g., classify a set of images of fruits which may be oranges, apples, or pears. Multiclass classification makes the assumption that each sample is assigned to one and only one label: a fruit can be either an apple or a pear but not both at the same time.

How can the multiclass problem be reduced to a set of binary problems?



The **one-vs-rest** strategy, is implemented in `OneVsRestClassifier` in `sklearn`. The strategy consists in fitting one classifier per class.

For each classifier, the class is fitted against all the other classes. In addition to its computational efficiency (only `n_classes` classifiers are needed), one advantage of this approach is its interpretability. Since each class is represented by one and only one classifier, it is possible to gain knowledge about the class by inspecting its corresponding classifier. This is the most commonly used strategy and is a fair default choice.



The `one-vs-one` strategy, is implemented in `OneVsOneClassifier` in `sklearn`. The strategy consists in fitting one classifier for each pair of classes.

At prediction time, the class which received the most votes is selected.

Since it requires to fit $n_classes * (n_classes - 1) / 2$ classifiers, this method is usually slower than `one-vs-the-rest`. However, this method may be advantageous for algorithms such as kernel algorithms which don't scale well with `n_samples`. This is because each individual learning problem only involves a small subset of the data whereas, with `one-vs-the-rest`, the complete dataset is used `n_classes` times.



Evaluation in Multiclass classification

A very intuitive way to show the results of a multiclass predictor is using the **confusion matrix**.

Example:

Prediction : A B A B C C C A B

True labels: A A A A C C C B B

		Predicted			
		A	B	C	
True labels	A	2	2	0	4
	B	1	2	0	3
	C	0	0	3	3
		3	4	3	Total

Precision and Recall for multiclass problems



Precision can be calculated separately for each class. For each row, we take the number on the diagonal, and divide it by the sum of all the elements in the column.

$$\text{prec_A: } 2/3 = 0.67$$

$$\text{prec_B: } 2/4 = 0.50$$

$$\text{prec_C: } 3/3 = 1.00$$

Recall can be calculated separately for each class. It is the value on the diagonal, divided by the sum of the values in the row.

$$\text{recall_A: } 2/4 = 0.50$$

$$\text{recall_B: } 2/3 = 0.67$$

$$\text{recall_C: } 3/3 = 1.00$$



Micro and Macro Averaging

To extract a single number from the class precision, recall or F1-score, it is common to adopt two different kinds of average.

The first, is to compute the score separately for each class and then taking the average value — this is called **macro averaging**. This kind of averaging is useful when the dataset is unbalanced and we want to put the same emphasis on all the classes.

$$\frac{r_A + r_B + r_C}{3} = 0.7233$$

The second is to calculate the measure from the grand total of the numerator and denominator, this is called **micro averaging** and is useful when you want to bias your classes towards the most populated class. For our 'recall' example:

$$\frac{2 + 2 + 3}{4 + 3 + 3} = 0.636$$



Notions

- Bias and Variance
- Under-fitting and Over-fitting
- Model Selection (Hold-out and Cross Validation)
- Unbalanced data
- Multi-class classification and evaluation

Exercises

- Bias and Variance evaluation in a controlled regression task
- Create a large regression dataset and demonstrate that by taking the mean of the models obtained with large p and relatively small training sets, we are able to break down the variance and thus significantly improve performance. Why can't this methodology actually be used? how to reduce the variance of the models even with small training sets?