

Data analysis and data Mining

Report 1

Mauro Angelini
Alessio Gilardi

October 26, 2018

Contents

First exercise	2
Second exercise	4

First exercise

In this first exercise we implemented the regression of a function in one-dimensional case, for example: the quadratic function

$$y = x^2, \forall x \in [0, 1]$$

Vengono definite:

- Function to be rebuilt
- The number of samples
- The samples disturbed by the Gaussian additive noise
- The variance of the noise
- The degree of the polynomial regressor

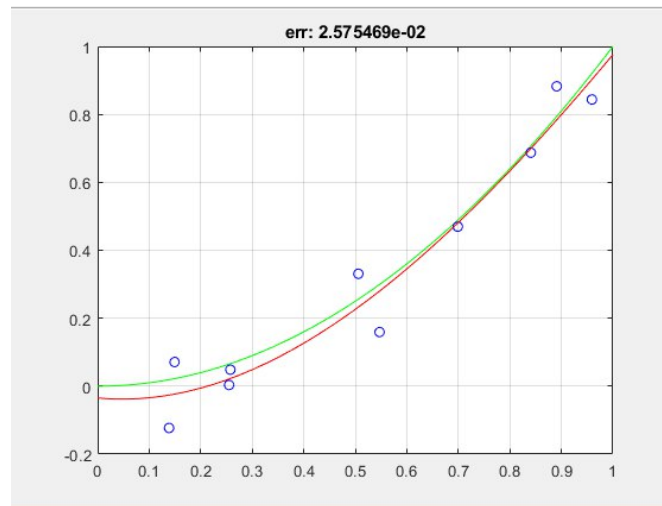
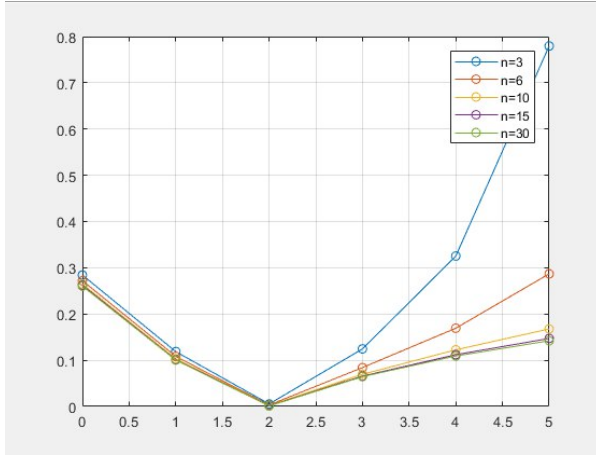


Figure 1: Regression Function

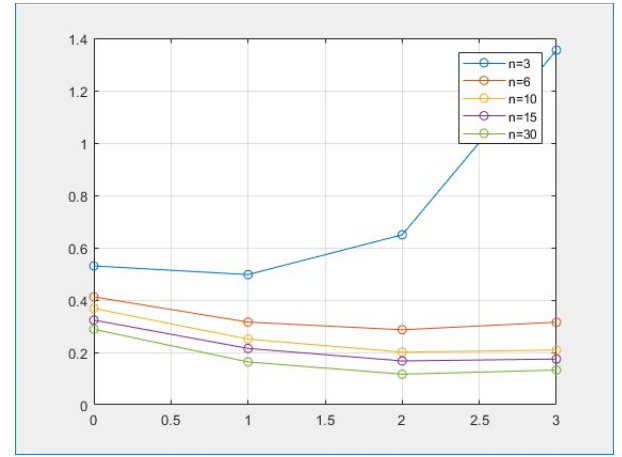
Figure 1 Regression function of a quadratic function

Starting from this first test, we notice that increasing the complexity of the polynomial too much produce a great loss of precision in the reconstruction. It's therefore necessary to evaluate the variation of reconstruction, by changing n (Number of samples), p (Degree of the polynomial regressor) and σ (Noise variance). So the error is calculated between YT (True function) and YP (Prediction function) and plotted in the graph.

Now we notice that in order to make the measurement of error more precise, it's better to repeat the experiment a sufficiently high number of times, for example, we make 30 repetitions of the prediction.



(a) Error by varying n and p (from 0 to 5)



(b) Error by varying n and p (from 0 to 3)

Figure 2: Comparing error with different values of p

	p = 0	p = 1	p = 2	p = 3	p = 4	p = 5
n = 3	2.842e-01	1.186e-01	5.349e-03	1.245e-01	3.253e-01	7.796e-01
n = 6	2.719e-01	1.084e-01	3.093e-03	8.410e-02	1.693e-01	2.869e-01
n = 10	2.647e-01	1.030e-01	2.056e-03	6.957e-02	1.225e-01	1.677e-01
n = 15	2.622e-01	1.012e-01	1.514e-03	6.572e-02	1.120e-01	1.471e-01
n = 30	2.603e-01	1.004e-01	1.095e-03	6.449e-02	1.092e-01	1.421e-01

(a) Error by varying n and p (from 0 to 5)

	p = 0	p = 1	p = 2	p = 3
n = 3	5.303e-01	4.972e-01	6.493e-01	1.354e+00
n = 6	4.121e-01	3.157e-01	2.865e-01	3.154e-01
n = 10	3.678e-01	2.507e-01	2.010e-01	2.089e-01
n = 15	3.235e-01	2.152e-01	1.674e-01	1.746e-01
n = 30	2.882e-01	1.636e-01	1.168e-01	1.325e-01

(b) Error by varying n and p (from 0 to 3)

Figure 3: Comparing error with different values of p

We performed tests with some values of n and p , from which it becomes clear how the error always maintains the same trend. In particular, observing the error plots, it's evident that the maximum degree of precision is obtained for p values around the degree of the original function. Obviously, using a second degree polynomial the error is minimal since the function to be reconstructed is a parabola. Utilizing a large number of samples n , in general, also increases accuracy, while an increase of p results in a loss of precision. It's evident that the variation of the error between one execution of the exercise and the other is minimal for the values greater than n . While using a small number of samples n , it's better to use low-grade polynomials, since an excessive number of degrees of freedom results in an attempt of the regression function to pass for each sample, making the reconstruction less precise, in particular in presence of a lot of noise, as it will try to pass for the samples even if noisy.

Second exercise

Ripeto la regressione polinomiale dell'esercizio precedente con l'aggiunta di regolarizzatore (bias).

Con l'uso del *bias* λ cerco di regolarizzare il mio polinomio: in sostanza aumentando λ semplifico il mio risultato mentre diminuendolo utilizzo una funzione più complessa per la regressione, quindi modificare λ è concettualmente simile a cambiare il grado del polinomio, offrendo, però, una maggior granularità nella scelta del valore ottimale. λ rappresenta il livello di fiducia nella qualità dei dati in ingresso: più l'ingresso è rumoroso più conviene utilizzare λ grandi e quindi adottare soluzioni semplici in modo to not fit the noise, meno è rumoroso e più possiamo diminuire λ ed utilizzare soluzioni più complesse in order to fit data.

In this exercise it's repeated the polynomial regression of the previous one with the addition of a regularizer *bias*. With the use of *bias* λ we try to regularize our polynomial: so by stretching λ we simplify my result while using it to use a more complex function for regression, so modifying λ is conceptually similar to changing the degree of polynomial, offering, however, a greater granularity in the choice of the optimal value. λ represents the level of confidence in the quality of input data: the more noisy the input, the better it's to use λ large and then adopt simple solutions to not fit the noise, the less it is noisy and the more we can decrease λ and use more complex solutions in order to fit data.

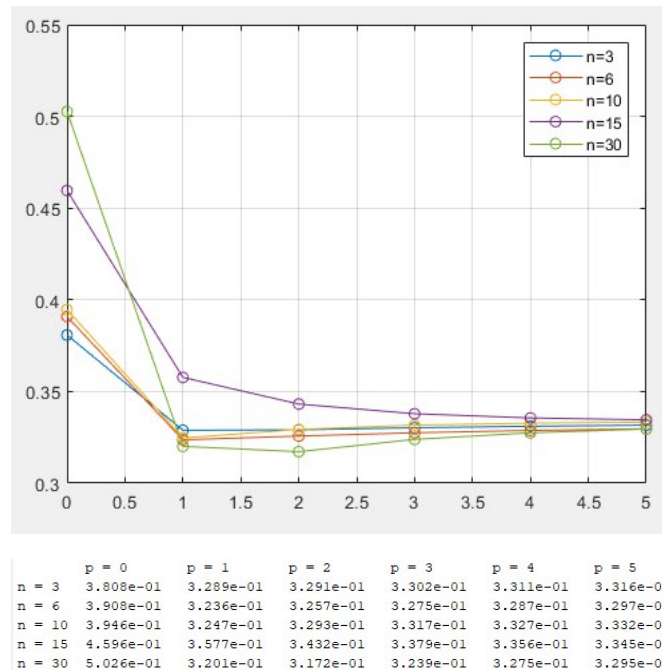


Figure 4: $\lambda = 100$ $\sigma = 10$

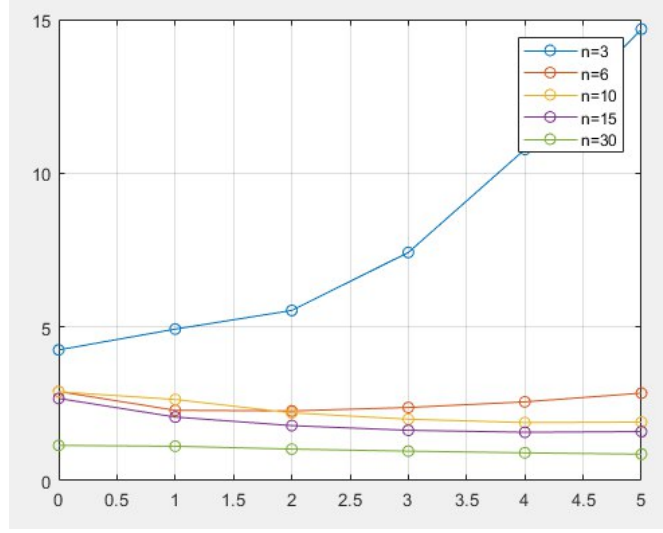


Figure 5: $\lambda = 0,001$ $\sigma = 10$

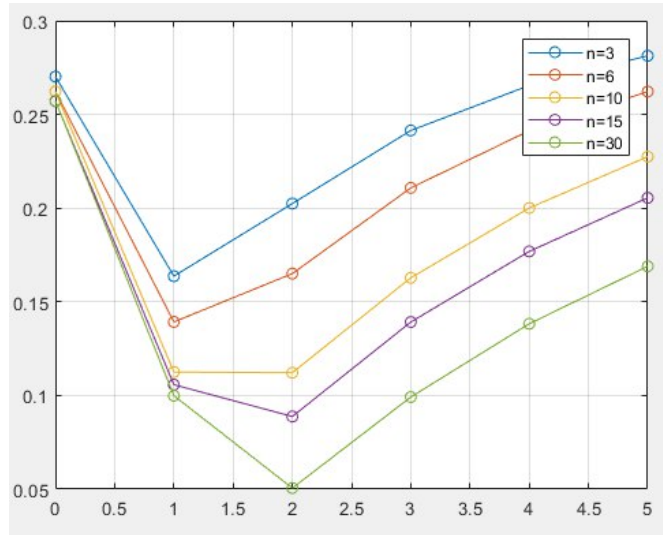


Figure 6: $\lambda = 1$ $\sigma = 0,01$

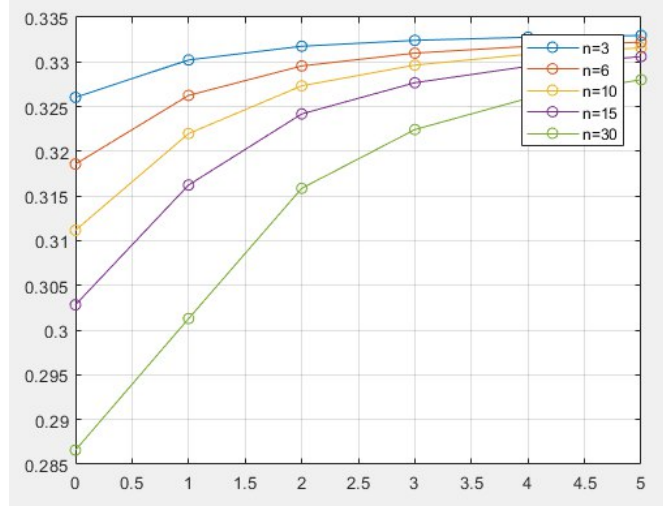


Figure 7: $\lambda = 100$ $\sigma = 0,01$

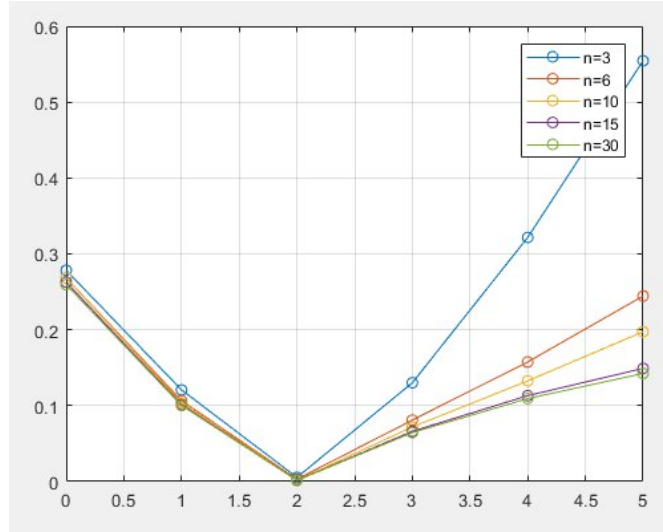


Figure 8: $\lambda = 0,00001$ $\sigma = 0,01$

List of Figures

1	Regression Function	2
2	Comparing error with different values of p	3
3	Comparing error with different values of p	3
4	$\lambda = 100$ $\sigma = 10$	4
5	$\lambda = 0,001$ $\sigma = 10$	5
6	$\lambda = 1$ $\sigma = 0,01$	5
7	$\lambda = 100$ $\sigma = 0,01$	6
8	$\lambda = 0,00001$ $\sigma = 0,01$	6