

# Coursework 2: Mathematics for Machine Learning

Romain REINERT

November, 17th 2020

## 1 Linear Regression

### 1.1 Question a: polynomial basis

As explained in the section "Basis Functions", we will use the notation  $\phi(x) = \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \phi_2(x) \\ \dots \\ \phi_M(x) \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \\ \dots \\ x^M \end{bmatrix}$

Using the theory of the book, we compute the maximum likelihood for  $\mathbf{w}$  using the formula :

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

where, for n points expressed in a polynomial basis of order M:

$$\Phi = [\phi(x_1) \quad \phi(x_2) \quad \dots \quad \phi(x_n)]^T = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^M \\ 1 & x_2 & x_2^2 & \dots & x_2^M \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^M \end{bmatrix} \in \mathbb{R}^{N \times (M+1)}$$

This formula is obtained by considering the log-likelihood  $\log p(\mathbf{y}|\mathbf{X}, \mathbf{w})$  and computing its gradient w.r.t  $\mathbf{w}$ . It is maximised by finding  $\mathbf{w}_{\text{ML}}$  such that this gradient is  $\mathbf{0}$  (necessary condition). For  $\sigma^2$ , the maximum likelihood estimate is computed using the gradient of the log-likelihood w.r.t  $\sigma^2$  and equaling it to 0. It is given by the formula:

$$\sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \phi^T(x_n) \mathbf{w})^2$$

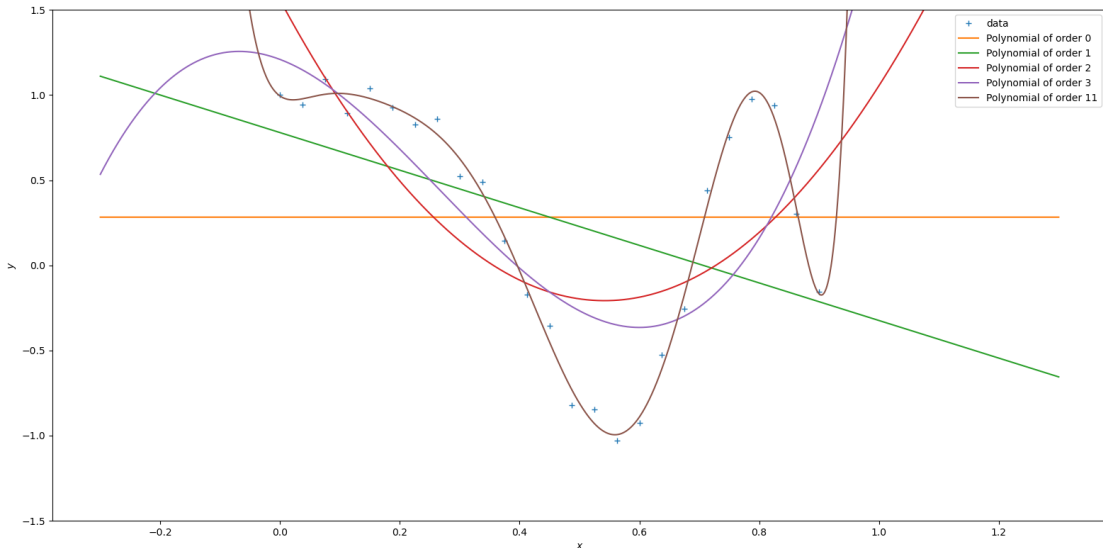


Figure 1: Predicted mean at test points in  $[-0.3, 1.3]$  for polynomial basis functions of different orders

In the python document, the matrix  $\Phi$  is created using *phi\_polynomial()*. The estimator of  $\mathbf{w}$  is computed in the *max\_lik\_estimate()* function. Then the function *plot\_polynomial()* plots the different predictions in  $[-0.3, 1.3]$  for basis of orders 0, 1, 2, 3 and 11.

As expected, higher orders fit better to the given data, which however may result in an over-fitting model (see question d).

## 1.2 Question b: trigonometric basis

In this question the new matrix  $\Phi$  is:

$$\Phi = [\phi(x_1) \quad \phi(x_2) \quad \dots \quad \phi(x_n)]^T = \begin{bmatrix} 1 & \sin 2\pi x_1 & \cos 2\pi x_1 & \dots & \sin 2\pi M x_1 & \cos 2\pi M x_1 \\ 1 & \sin 2\pi x_2 & \cos 2\pi x_2 & \dots & \sin 2\pi M x_2 & \cos 2\pi M x_2 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & \sin 2\pi x_n & \cos 2\pi x_n & \dots & \sin 2\pi M x_n & \cos 2\pi M x_n \end{bmatrix} \in \mathbb{R}^{N \times (2M+1)}$$

I used the functions *phi\_trigonometric()*, *max\_lik\_estimate()* and *plot\_trigonometric()* to plot the predictions for basis of orders 1 and 11.

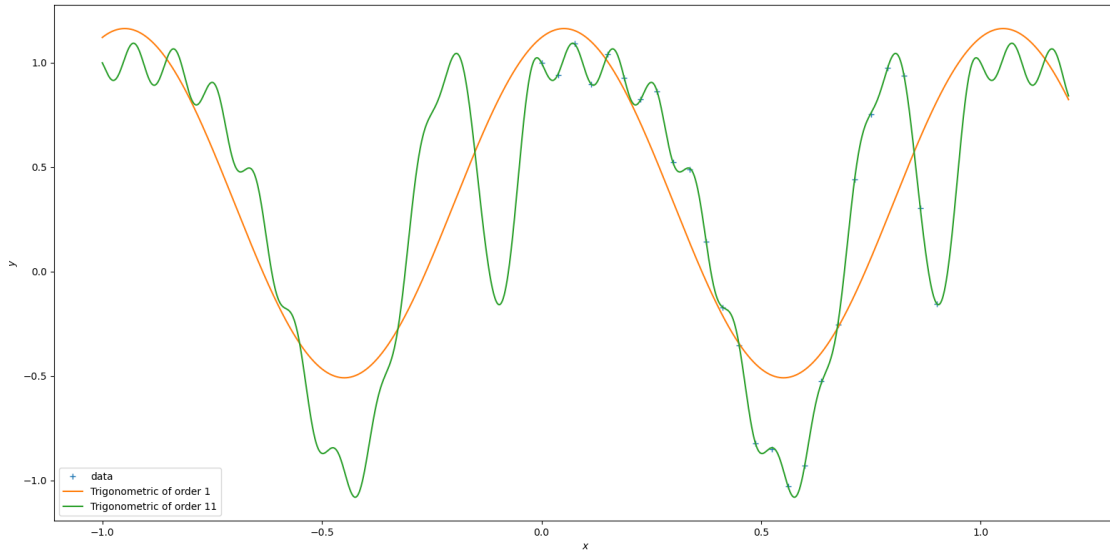


Figure 2: Predicted mean at test points in  $[-1, 1.2]$  for trigonometric basis functions of orders 1 and 11

### 1.3 Question c: leave-one-out cross-validation

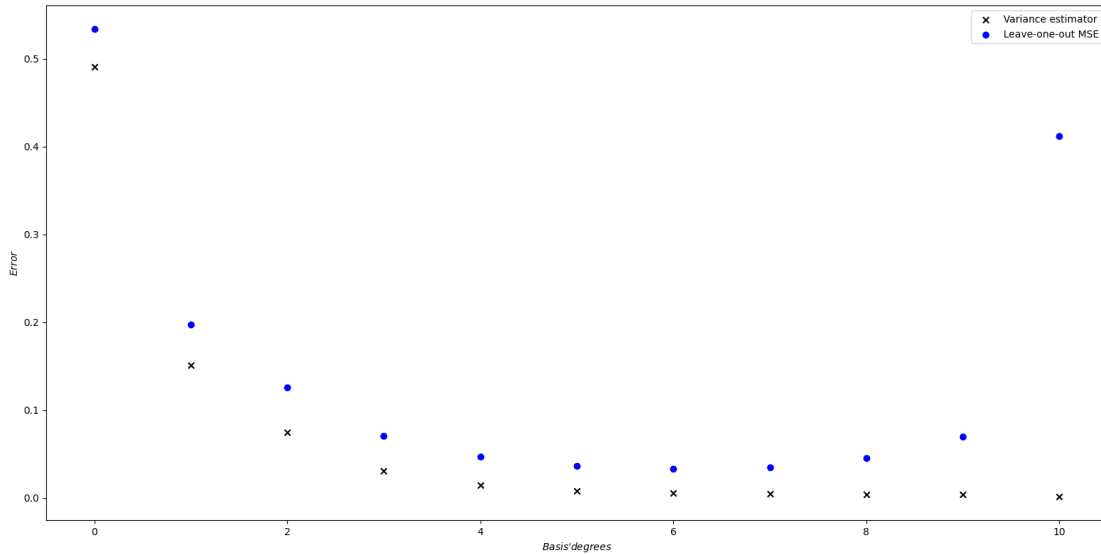


Figure 3: Evolution of the averaged error and the variance estimation using maximum likelihood, against the degree of the trigonometric basis

### 1.4 Question d: over-fitting

Over-fitting appears when a model fits to a given training data, but doesn't generalize well to unknown data. This can clearly be seen on the previous figure: using leave-one-out cross-validation, we see that while the mean-squared error keeps decreasing on training data as orders of the basis increase from 0 to 10, errors on test data starts to increase for basis of orders superior to 6.

It can be illustrated by the plots of the first two questions, where high degrees polynomials and trigonometric functions show high oscillations to fit exactly to each data-point. The following figure shows that between polynomial basis of orders 4 and 11, weights are increasing widely and successively in opposite signs, in order to allow these strong oscillations fitting perfectly to the data. It suggests that despite a very good fitting to our current points, the model will probably not generalize well to unseen data.

```
w_ml for a polynomial basis of order 4:  
[[ 0.54204609]  
 [ 16.83338977]  
 [-103.88164731]  
 [ 177.58671124]  
 [-92.61235966]]  
  
w_ml for a polynomial basis of order 11:  
[[ 9.94382437e-01]  
 [-2.42250063e+00]  
 [ 8.26455789e+01]  
 [-9.73074488e+02]  
 [ 5.40927681e+03]  
 [-1.55296378e+04]  
 [ 1.78595422e+04]  
 [ 1.60023399e+04]  
 [-7.74854209e+04]  
 [ 1.03589628e+05]  
 [-6.59935728e+04]  
 [ 1.70551217e+04]]
```

Figure 4: Increase in the squared values of the weights

Some solutions can be put in place to try to delay over-fitting. One way to do it is to modify the loss function to penalize the increase in the squared values of weights: this is called regularisation.

An other approach, the Maximum A Posteriori estimation (MAP), uses a prior Gaussian probabilistic distribution on  $\mathbf{w}$  to represent the knowledge we have on our weights before having observed any data, thus regularizing their value. Finally, Bayesian Linear Regression (BLR) can be used to give an idea of our confidence in the predictions we are making.