

COURSEWORK 2

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

CO496

Author:

Chengyi Zhang (CID: 01604256)

Date: November 5, 2018

1

1.1 Question 1a

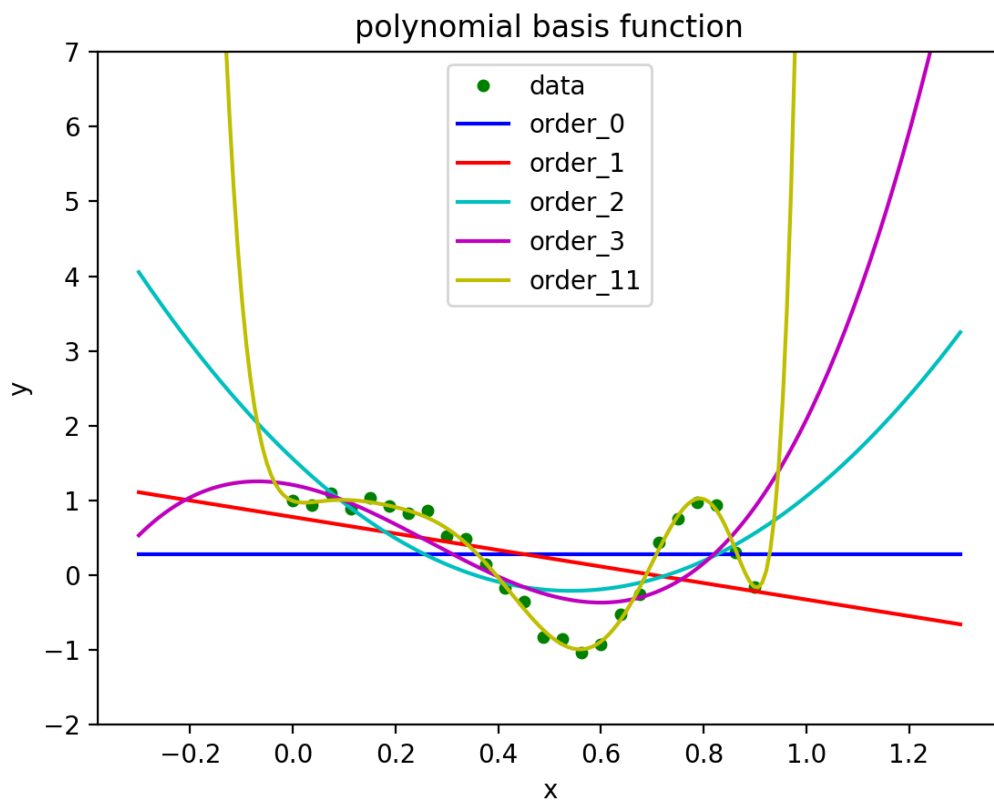


Figure 1: Fit the data using polynomial basis functions with various orders

1.2 Question 1b

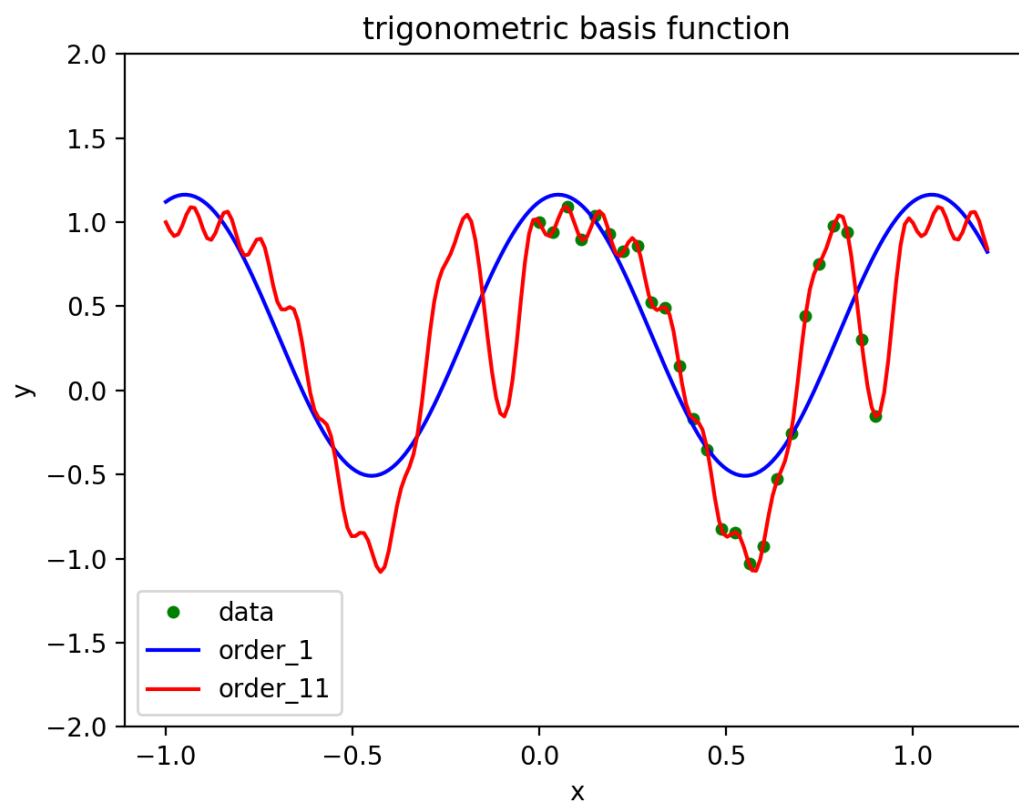


Figure 2: Fit the data using trigonometric basis functions with various orders

1.3 Question 1c

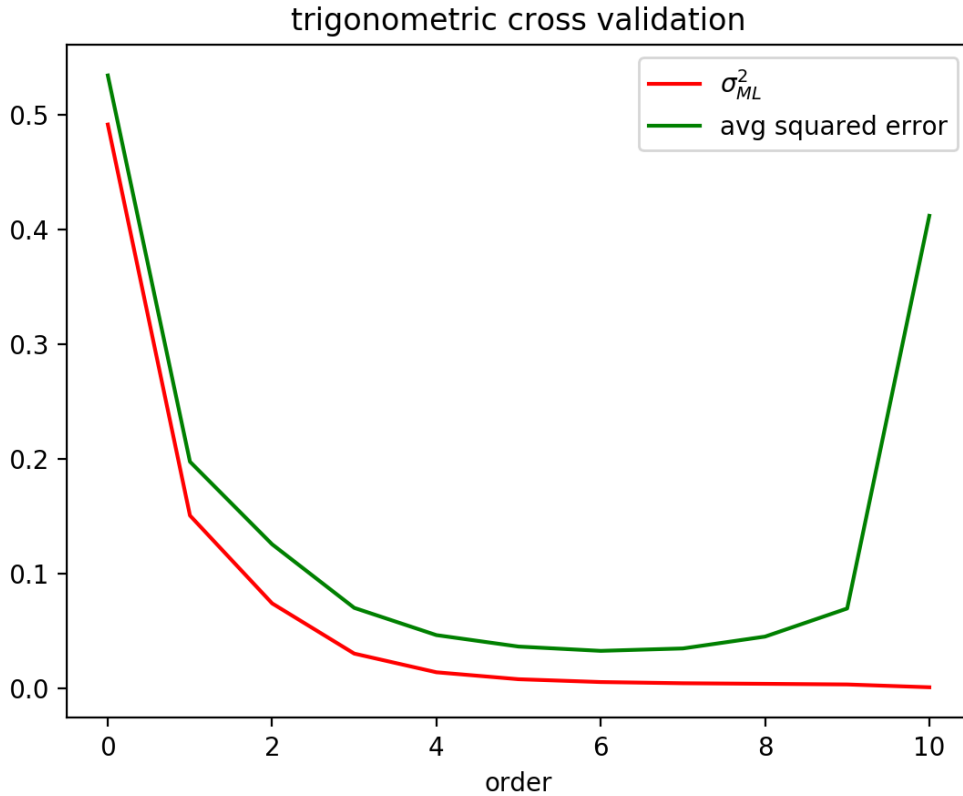


Figure 3: σ_{ML}^2 and average squared error of order 0 to 10

1.4 Question 1d

When the model increases the order of parameters to fit the training data set better, σ_{ML}^2 decreases. Meanwhile, however, the influence of noise is also getting stronger, causing the model to perform worse on testing data. When the model has more parameters than the data can express, the model starts to over-fit the training data set. Therefore, the average squared error on test data decreases at first, and then increases with the order.

We can determine the content above in the figures. In Fig. 1, the curves with higher order are obviously fitting the training data better. However, the curve with order 3 increases rapidly when $x > 0$, and the curve with order 11, which almost perfectly fits the training data, changes extremely fast when x is not in the domain of training data. Similarly, in Fig. 2, the curve with order 11 is strongly influenced by the noise and thus it fits the data perfectly.

2

2.1 Question 2a

Set $\mathbf{Y} = \begin{bmatrix} y_1 \\ \dots \\ y_N \end{bmatrix}_{N \times 1}$, $\mathbf{\Phi} = \begin{bmatrix} \phi(x_1) \\ \dots \\ \phi(x_N) \end{bmatrix}_{N \times m}$

Prior:

$$p(\mathbf{w}) \sim N(0, b^2 \mathbf{I}_{m \times m}), \quad b^2 = \frac{\sigma^2}{\lambda} \quad (1)$$

(Negative) log likelihood:

$$\mathbb{L}(y_i | x_i, \mathbf{w}) = -\log\left(\prod_{i=1}^N p(y_i | x_i, \mathbf{w})\right) \quad (2)$$

$$= -\log\left(\prod_{i=1}^N p(y_i | \phi(x_i), \mathbf{w})\right) \quad (3)$$

$$= -\log\left(\prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \mathbf{w}^\top \phi(x_i))^2}{2\sigma^2}\right\}\right) \quad (4)$$

$$= \frac{N}{2} \log 2\pi + \frac{N}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \phi(x_i))^2 \quad (5)$$

$$= \frac{N}{2} \log 2\pi + \frac{N}{2} \log \sigma^2 + \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{\Phi} \mathbf{w})^\top (\mathbf{Y} - \mathbf{\Phi} \mathbf{w}) \quad (6)$$

(Negative) log prior:

$$\mathbb{L}(\mathbf{w}) = -\log p(\mathbf{w}) \quad (7)$$

$$= -\log((2\pi b^2)^{-\frac{m}{2}}) \exp\left\{-\frac{\mathbf{w}^\top \mathbf{w}}{2b^2}\right\} \quad (8)$$

$$= \frac{m}{2} \log 2\pi + \frac{m}{2} \log b^2 + \frac{\mathbf{w}^\top \mathbf{w}}{2b^2} \quad (9)$$

Differentiate the product of both log and solve:

$$\frac{\partial \mathbb{L}(y_i | x_i, \mathbf{w}) \mathbb{L}(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{2\sigma^2} (-2\mathbf{\Phi}^\top \mathbf{Y} + 2\mathbf{\Phi}^\top \mathbf{\Phi} \mathbf{w}) + \frac{\mathbf{w}}{b^2} = 0 \quad (10)$$

$$(\mathbf{\Phi}^\top \mathbf{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I}) \mathbf{w} = \mathbf{\Phi}^\top \mathbf{Y} \quad (11)$$

$$\mathbf{w} = (\mathbf{\Phi}^\top \mathbf{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I})^{-1} \mathbf{\Phi}^\top \mathbf{Y} \quad (12)$$

Differentiate the loss function and solve:

$$L(\mathbf{x}) = (\mathbf{Y} - \mathbf{\Phi} \mathbf{w})^\top (\mathbf{Y} - \mathbf{\Phi} \mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w} \quad (13)$$

$$\frac{\partial L}{\partial w} = -2\Phi^T Y + 2\Phi^T \Phi w + 2\lambda w = 0 \quad (14)$$

$$(\Phi^T \Phi + \lambda I)w = \Phi^T Y \quad (15)$$

$$w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y \quad (16)$$

With $b^2 = \frac{\sigma^2}{\lambda}$, the solutions are the same. Note that the scaling coefficient, $\frac{1}{2\sigma^2}$, and other constants in the product of log likelihood and log prior do not affect the solution.

The loss function has two parts. First, we calculate the sum of the squared error between the actual value y_i and the expected value of the corresponding input data x_i and our parameters w . However, if we have parameters with very high order, it is likely that the model can fit training data and result perfectly while it may have bad performance on test data, and such situation is called over-fitting. Therefore, we add the second part, named regularization, as a penalty to the order of parameters. Since the square of any entry of parameter, w_i^2 , is equal to or larger than zero, then with $\lambda > 0$, the value of the second part is non-decreasing (increasing in most circumstances where $w_i \neq 0$) when the order is increasing.

2.2 Question 2b

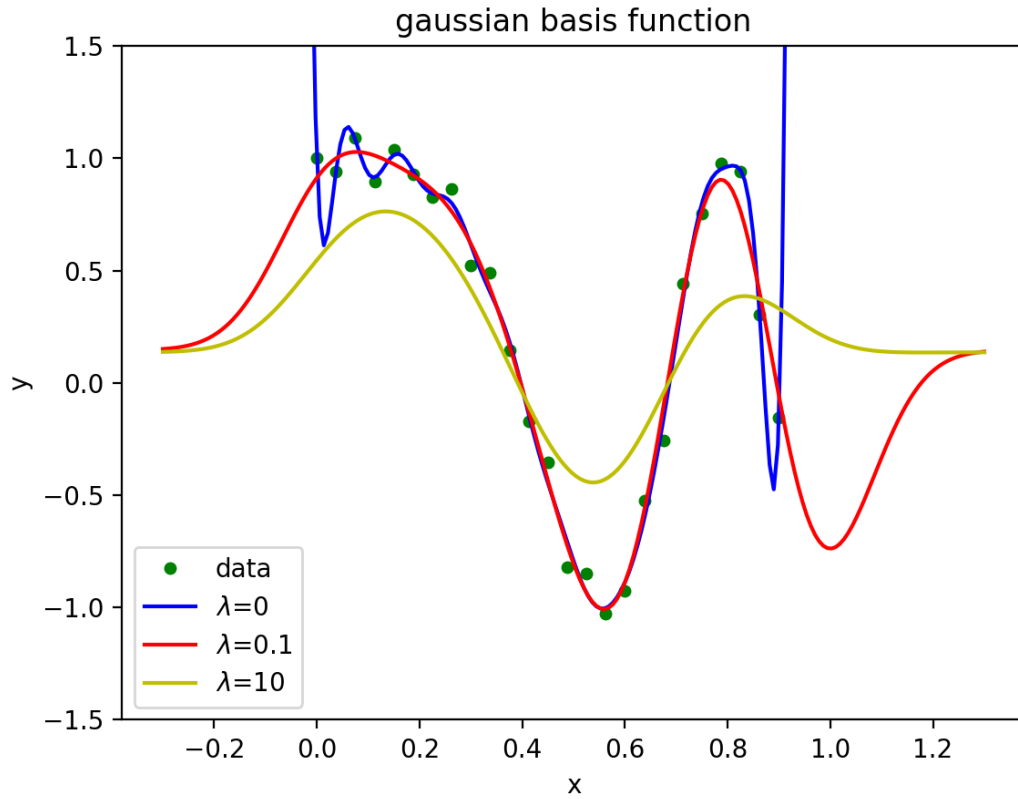


Figure 4: Ridge regression with different λ