

Lecture 19:

31/03/2023

Lecturer: Abir De

Scribe: Group 21

1 Recap

Consider the following process :

$$y = w^T \phi(x) + \epsilon \quad (1)$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$, $w \sim \mathcal{N}(0, \Sigma_p)$ is a $d \times 1$ weight vector, and $\phi(x)$ is the $d \times 1$ feature vector, then the probability of observing y^* given x^* and D , where $D = \{(x_i, y_i) | i = 1, 2, \dots, n\}$ is given by:

$$P(y^* | x^*, D) \sim \mathcal{N}(\mu_{y^*}, \sigma_{y^*}^2)$$

$$\mu_{y^*} = \phi(x^*)^T \left[\frac{\Phi \Phi^T}{\sigma^2} + \Sigma_p^{-1} \right]^{-1} \frac{\Phi y}{\sigma^2} \quad (2)$$

$$\mu_{y^*} = \phi(x^*)^T \Sigma_p \Phi (\Phi \Sigma_p \Phi + \sigma^2 I)^{-1} y \quad (3)$$

Note that these two expressions of μ_{y^*} will be equal only when both the inverses exist

$$\sigma_{y^*}^2 = \phi(x^*)^T [\Sigma_p - \Sigma_p \Phi (\Phi^T \Sigma_p \Phi + \sigma^2 I)^{-1} \Phi^T \Sigma_p] \phi(x^*) + \sigma^2 \quad (4)$$

Recall ,

$$\mu_{A|B} = \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (y_B - \mu_B)$$

Hence here $\mu_A = \mu_B = 0$, $\Sigma_{AB} = \phi(x^*)^T \Sigma_p \Phi$, $\Sigma_{BB} = \Phi^T \Sigma_p \Phi$

For $\sigma = 0$ we have shown that,

$$\mu_{y^*} = y_i \quad \text{if } (x^*, y^*) = (x_i, y_i) \quad (5)$$

2 Analyzing the Variance

Similar to calculation of μ_{y^*} , for $\sigma = 0$ and $(x^*, y^*) = (x_i, y_i)$,

$$\sigma_{y^*}^2 = \phi(x^*)^T \Sigma_p \phi(x^*) - \phi(x^*)^T \Sigma_p \Phi (\Phi^T \Sigma_p \Phi + \sigma^2 I)^{-1} \Phi^T \Sigma_p \phi(x^*)$$

As shown in previous lectures, $\phi(x^*)^T \Sigma_p \Phi (\Phi^T \Sigma_p \Phi + \sigma^2 I)^{-1} = [1 \ 0 \ 0 \dots 0]$, when $(x^*, y^*) \in D$

$$\implies \sigma_{y^*}^2 = \phi(x^*)^T \Sigma_p \phi(x^*) - [1 \ 0 \ 0 \dots 0] \Phi^T \Sigma_p \phi(x^*)$$

$$\begin{aligned}\implies \sigma_{y^*}^2 &= \phi(x^*)^T \Sigma_p \phi(x^*) - \phi(x^*)^T \Sigma_p \phi(x^*) \\ &\implies \sigma_{y^*}^2 = 0\end{aligned}$$

Similarly we can also show that if (x^*, y^*) are such that,

$$\phi(x^*) = \Sigma \alpha_i \phi(x_i) \quad , \quad y^* = \Sigma \alpha_i y_i \quad , \quad \text{then } \sigma_{y^*} = 0 \quad (6)$$

3 Finding Variance of Gaussian process when the data belongs to the training data itself

We have to find $\sigma_{y^*}^2$ of $P(y^*|x^*, D)$ for $x^* = x_i$

We know that,

$$\sigma_{y^*}^2 = \phi(x^*)^T [\Sigma_p - \Sigma_p \Phi (\Phi^T \Sigma_p \Phi + \sigma^2 I) \Phi^T \Sigma_p] \phi(x^*) + \sigma^2$$

$$\begin{aligned}\Sigma_p - \Sigma_p \Phi (\Phi^T \Sigma_p \Phi + \sigma^2 I) &= \Sigma_p - \Sigma_p \Phi (I + \sigma^2 (\Phi^T \Sigma_p \Phi)^{-1})^{-1} (\Phi^T \Sigma_p \Phi)^{-1} \Phi^T \Sigma_p \\ &= \Sigma_p - \Sigma_p \Phi (I - \sigma^2 (\Phi^T \Sigma_p \Phi)^{-1}) (\Phi^T \Sigma_p \Phi)^{-1} \Phi^T \Sigma_p \\ &= \Sigma_p - \Sigma_p \Phi (\Phi^T \Sigma_p \Phi)^{-1} \Phi^T \Sigma_p + \sigma^2 \Sigma_p \Phi (\Phi^T \Sigma_p \Phi)^{-2} \Phi^T \Sigma_p\end{aligned}$$

Hence,

$$\begin{aligned}\sigma_{y^*}^2 &= \Phi(x^*)^T \Sigma_p \Phi(x^*) - \Phi(x^*)^T \Sigma_p \Phi (\Phi^T \Sigma_p \Phi)^{-1} \Phi^T \Sigma_p \Phi(x^*) + \sigma^2 \Phi(x^*)^T \Sigma_p \Phi (\Phi^T \Sigma_p \Phi)^{-2} \Phi^T \Sigma_p \Phi(x^*) + \sigma^2 \\ &\quad \Phi(x^*)^T \Sigma_p \Phi(x^*) - \Phi(x^*)^T \Sigma_p \Phi (\Phi^T \Sigma_p \Phi)^{-1} \Phi^T \Sigma_p \Phi(x^*) = 0 \quad \text{for } x^* = x_i\end{aligned}$$

Hence,

$$\begin{aligned}\sigma_{y^*}^2 &= \sigma^2 \Phi(x^*)^T \Sigma_p \Phi (\Phi^T \Sigma_p \Phi)^{-2} \Phi^T \Sigma_p \Phi(x^*) + \sigma^2 \\ &= 2\sigma^2\end{aligned}$$

4 Writing the above variance in terms of kernel function

Let

$$\phi(x^*)^T \Sigma_p \Phi = K(x^*, x)$$

$$\Phi^T \Sigma_p \Phi = K(X, X) = K$$

Hence,

$$\sigma_{y^*}^2 = K(x^*, x + \sigma^2 - K(x^*, x)) [K + \sigma^2 I]^{-1} K(x, x^*) y$$

5 Some points to think about

If $\sigma \neq 0$ then for what instance the variance σ_{y^*} will be least?

Lets say we have $x_1, x_2, \dots, x_{1million}$ unlabelled points. $x'_1, x'_2, \dots, x'_{1000}$ are labelled as $y'_1, y'_2, \dots, y'_{1000}$. Using these we want to pick some $x_{i's}$ for labelling, but which ones to pick, because picking all of them is not practically possible?

We will want those $x_{i's}$ which are dissimilar with the given $x_{j's}$. Hence they will have more variance. Therefore we find low variance $x_{i's}$ and just label them using nearest neighbour and discard them from getting picked up to label the hard way.