

Exercise 1

Linear regression can be transformed into a Gaussian process by assuming a Gaussian prior over the weights $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma)$. Assuming that $f(\mathbf{x})$ is linear, i.e. $f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$, derive the posterior distribution over \mathbf{w} .

Solution Exercise 1 The linear regression model can be written as

$$y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \epsilon$$

where $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$ is a Gaussian additive noise. It follows that $y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \epsilon \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}^{(i)}, \sigma_n^2)$, or equivalently

$$p(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T \mathbf{x}^{(i)}, \sigma_n^2)$$

For simplicity, let us stack the input data into the matrix \mathbf{X} and the output data into the vector \mathbf{y} . Assume a zero mean Gaussian prior over the weights $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma)$. From the Bayes' rule we can write the posterior distribution over \mathbf{w} as

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}) p(\mathbf{w})}{p(\mathbf{y} | \mathbf{X})}$$

Considering that the marginal likelihood $p(\mathbf{y} | \mathbf{X})$ is constant we have

$$\begin{aligned} p(\mathbf{w} | \mathbf{y}, \mathbf{X}) &\propto p(\mathbf{y} | \mathbf{X}, \mathbf{w}) p(\mathbf{w}) \\ &= \exp \left(-\frac{1}{2\sigma_n^2} (\mathbf{y} - \mathbf{X}^T \mathbf{w})^T (\mathbf{y} - \mathbf{X}^T \mathbf{w}) \right) \exp \left(-\frac{1}{2} \mathbf{w}^T \Sigma^{-1} \mathbf{w} \right) \\ &= \exp \left(-\frac{1}{2\sigma_n^2} (\mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{X} \mathbf{y} - \mathbf{y}^T \mathbf{X}^T \mathbf{w} + \mathbf{w}^T \mathbf{X} \mathbf{X}^T \mathbf{w}) - \frac{1}{2} \mathbf{w}^T \Sigma^{-1} \mathbf{w} \right) \\ &= \exp \left(-\frac{1}{2\sigma_n^2} (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X} \mathbf{y}) - \frac{1}{2} \mathbf{w}^T \left(\frac{\mathbf{X} \mathbf{X}^T}{\sigma_n^2} + \Sigma^{-1} \right) \mathbf{w} \right) \end{aligned}$$

Choosing $\mathbf{A} = \left(\frac{\mathbf{X} \mathbf{X}^T}{\sigma_n^2} + \Sigma^{-1} \right)$ and considering that $\frac{1}{2\sigma_n^2} \mathbf{y}^T \mathbf{y} = \frac{1}{2\sigma_n^2} \|\mathbf{y}\|^2 = c$ is constant we have

$$\begin{aligned} p(\mathbf{w} | \mathbf{y}, \mathbf{X}) &\propto \exp \left(-c + \frac{1}{\sigma_n^2} \mathbf{w}^T \mathbf{X} \mathbf{y} - \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} \right) \\ &= \exp(-c) \exp \left(\frac{1}{\sigma_n^2} \mathbf{w}^T \mathbf{X} \mathbf{y} - \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} \right) \\ &= \exp(-c) \exp \left(\frac{1}{\sigma_n^2} \mathbf{w}^T \mathbf{A} \mathbf{A}^{-1} \mathbf{X} \mathbf{y} - \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} \right) \\ &= \exp(-c) \exp \left(-\frac{1}{2} \left(\mathbf{w} - \frac{1}{\sigma_n^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \right)^T \mathbf{A} \left(\mathbf{w} - \frac{1}{\sigma_n^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \right) - \frac{1}{2(\sigma_n^2)^2} (\mathbf{X} \mathbf{y})^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \right) \\ &= \exp(-c - \frac{1}{2(\sigma_n^2)^2} (\mathbf{X} \mathbf{y})^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y}) \exp \left(-\frac{1}{2} \left(\mathbf{w} - \frac{1}{\sigma_n^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \right)^T \mathbf{A} \left(\mathbf{w} - \frac{1}{\sigma_n^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \right) \right) \\ &\propto \exp \left(-\frac{1}{2} \left(\mathbf{w} - \frac{1}{\sigma_n^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \right)^T \mathbf{A} \left(\mathbf{w} - \frac{1}{\sigma_n^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \right) \right) = \mathcal{N} \left(\frac{1}{\sigma_n^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{y}, \mathbf{A}^{-1} \right) \end{aligned}$$

Exercise 2

An autonomous flying robot, while performing a transportation task, encounters a storm. The odometry module of the robot computes and transmits position values (y) every $0.2s$. Due to the electromagnetic interference of the storm, odometry data are influenced by a Gaussian noise.

Being unable to complete the task, the robot activates an emergency landing procedure. To work properly, the landing procedure requires the prediction of the robot position one step ($0.2s$) in the future. The limited amount of memory of the robots allows to store only the last 3 positions.

To make predictions, the landing strategy uses Gaussian processes, with the radial basis kernel function:

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2l^2}\right) + \sigma_n^2 \delta_{xx'}$$

where x are the position sampling time instants, l defines the lengthscale of the kernel, σ_n^2 the variance of the noise component, and $\delta_{xx'} = 1$ if $x = x'$, 0 otherwise. Given the last 3 positions $\mathbf{y} = [1.5, 1, 0.8]$ with associated times $\mathbf{x} = [1, 1.2, 1.4]$, compute the position at $x^{(3)} = 1.4s$ and $x_t = 1.6s$ assuming:

- Variance $\sigma_n^2 = 0.1$ and lengthscale $l = 0.1$.
- Variance $\sigma_n^2 = 0.1$ and lengthscale $l = 0.001$ (close to zero).
- Variance $\sigma_n^2 = 0.1$ and lengthscale $l = 100$ (very high).
- State your conclusions regarding the relationship between lengthscale and obtained results.

Solution Exercise 2

- a) Given a query point x_t , the predictive mean \hat{y}_t is:

$$\hat{y}_t = \mathbf{K}(x_t, \mathbf{x})[\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}^T$$

For $\sigma_n^2 = 0.1$ and $l = 0.1$, we have:

$$\hat{y}(x^{(3)} = 1.4) = [0.45 \quad 0.82 \quad 1] \begin{bmatrix} 1.1 & 0.82 & 0.45 \\ 0.82 & 1.1 & 0.82 \\ 0.45 & 0.82 & 1.1 \end{bmatrix}^{-1} \begin{bmatrix} 1.5 \\ 1 \\ 0.8 \end{bmatrix} \approx 0.73$$

$$\hat{y}(x_t = 1.6) \approx 0.44$$

- b) For $\sigma_n^2 = 0.1$ and $l = 0.001$, we have:

$$\hat{y}(x^{(3)} = 1.4) \approx 0.73$$

$$\hat{y}(x_t = 1.6) \approx 0.0$$

Good approximation of the training points, poor predictions with new targets \rightarrow over-fitting.

- c) For $\sigma_n^2 = 0.1$ and $l = 100$, we have:

$$\hat{y}(x^{(3)} = 1.4) \approx 1.06$$

$$\hat{y}(x_t = 1.6) \approx 1.06$$

The predicted value is almost constant.

- d) The variation of the lengthscale parameter is key to obtain optimal data fitting. For very low values of l the model has high complexity and is prone to over-fitting, while for high values the prediction is almost independent from the data.