

Advanced Machine Learning, Assignment 1

Youssef Taoudi, yousseft@kth.se

December 24, 2019

1 The Prior

1.1 Theory

1.1.1 Question 1

Insert statement about Gaussian form of likelihood

The assumption we make about the data by choosing a spherical covariance matrix is that the covariance becomes a diagonal matrix, meaning all data points are assumed to be independent. The assumption also means that all features will share the same variance.

1.1.2 Question 2

The likelihood will be a multivariate normal distribution with a mean vector as a collection of the function of all data points of the distribution. If the data points are not assumed to be independent, then the covariance will be used instead of the variance in a joint distribution.

$$P(\mathbf{T}|f, \mathbf{X}) \sim \mathcal{N} \left(\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \dots & \Sigma_{1n} \\ \vdots & & \\ \Sigma_{n1} & \dots & \Sigma_{nn} \end{bmatrix} \right) \quad (1)$$

1.1.3 Question 3

$$p(\mathbf{T}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{T}|f(\mathbf{x}_i, \mathbf{W}), \beta^{-1}) \quad (2)$$

Equation 2 shows the probability $p(\mathbf{T}|\mathbf{X}, \mathbf{W})$ as a normal distribution with the mean $\mathbf{T}|f(\mathbf{x}_i, \mathbf{W})$ and variance β . Where t is the target value that is conditioned on the deterministic function $f(\mathbf{x}_i, \mathbf{W}) = \mathbf{W}\mathbf{x}_i$ and β is the precision (inverse of the variance) for the error function ϵ [1, p. 142]. In our case, the error function has precision $\beta = (\sigma^2 \mathbf{I})^{-1}$. The inverse precision β^{-1} can thus be written as $(\sigma^2 \mathbf{I})$ which gives us the Gaussian in equation 3. Note that the product

of the probabilities hold true only because the data samples are assumed to be conditionally independent.

$$p(\mathbf{T}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{T}|\mathbf{W}\mathbf{x}_i, \sigma^2\mathbf{I}) \quad (3)$$

1.1.4 Question 4

Regularizing with L1 norm is known as the lasso. One property of the lasso is that as the shrinkage term λ increases, an increasing number of coefficients w_i are driven to zero which leads to a sparse model where some features will not be used for predictions [1, p. 145]. Regularization terms are used to reduce complexity on models that are trained on limited data [1, p. 145] which will lead to a smoother distribution.

The distribution of the prior is shown in equation 4 as a normal distribution [3] where the denominator is a constant value for all vectors \mathbf{w} .

$$p(\mathbf{w}) = \frac{\exp -\frac{1}{2\tau^2}(\mathbf{w} - \mathbf{w}_0)(\mathbf{w} - \mathbf{w}_0)^T}{C} \quad (4)$$

The negative logarithmic prior of one row takes the form of equation 5 derived from the logarithm of the normal distribution of the prior in equation 4.

$$\ln p(\mathbf{w}) = \frac{1}{2\tau^2}(\mathbf{w} - \mathbf{w}_0)(\mathbf{w} - \mathbf{w}_0)^T + \ln C \quad (5)$$

Equation 6 and 7 derived by Bishop [1, p. 153] represent the penalization terms for L_1 and L_2 regularization respectively. In our case, $\lambda = \frac{1}{\tau^2}$ and \mathbf{w} is one row in our \mathbf{W} matrix. Naturally, as L_2 takes a very similar form as our log prior, an L_2 regularization would be more suitable for the problem.

$$\frac{\lambda}{2} \sum_{j=1}^D |w_j| \quad (6)$$

$$\frac{\lambda}{2} \sum_{j=1}^D |w_j|^2 = \frac{\lambda}{2} \mathbf{w}\mathbf{w}^T \quad (7)$$

Figure 1 illustrates the general contour forms in two dimensions when using L_1 and L_2 regularization respectively.

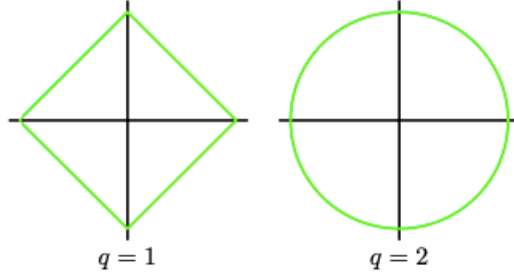


Figure 1: The figure illustrates the contours of the regularization terms where $q=1$ (L_1) on the left and $q=2$ (L_2) on the right. Image taken from Bishop[1, p. 145].

1.1.5 Question 5

Our posterior for each row in matrix W_i will be a Gaussian distribution in the form 8. Because output points are assumed independent, we can formulate a normal distribution over each output feature by dividing the W matrix in rows and the T matrix in columns where \mathbf{W}_i is one row in \mathbf{W} and \mathbf{T}_i is one column in \mathbf{T} .

$$p(\mathbf{W}_i | \mathbf{X}, \mathbf{T}_i) = \mathcal{N}(\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) \quad (8)$$

The posterior 8 is made up of three factors: the prior, the likelihood and the evidence shown in expression 10. The first exponent is quadratic in terms of \mathbf{W}_i , the second exponent is linear and the third exponent is constant.

$$p(\mathbf{W}_i | \mathbf{X}, \mathbf{T}_i) \propto P(\mathbf{T}_i | \mathbf{W}_i, \mathbf{X}) \times P(\mathbf{W}_i) = \exp\left(-\frac{1}{2} \mathbf{W}_i \boldsymbol{\Sigma}_w^{-1} \mathbf{W}_i^T\right) \times \exp(\mathbf{W}_i \boldsymbol{\Sigma}_w^{-1} \boldsymbol{\mu}_w) \quad (9)$$

Equation 9 can be written in terms of one quadratic, one linear and one constant exponent in terms of W as shown in equation 10.

$$p(\mathbf{W}_i | \mathbf{X}, \mathbf{T}_i) \propto \exp\left(-\frac{1}{2} \mathbf{W}_i \boldsymbol{\Sigma}_w^{-1} \mathbf{W}_i^T\right) \times \exp(\mathbf{W}_i \boldsymbol{\Sigma}_w^{-1} \boldsymbol{\mu}_w) \times \exp\left(-\frac{1}{2} \boldsymbol{\mu}_w \boldsymbol{\Sigma}_w^{-1} \boldsymbol{\mu}_w\right) \quad (10)$$

Our task is thus to match posterior 10 with our prior 4 and likelihood 1. The evidence can be ignored for now as it is constant and does not need to be used to find the parameters $\boldsymbol{\mu}_w$ and $\boldsymbol{\Sigma}_w$. If we only focus on the exponents, we can write a tidy expression as shown in 11

$$-\frac{1}{2\sigma^2} (\mathbf{T}_i - \mathbf{X} \mathbf{W}_i^T)^T (\mathbf{T}_i - \mathbf{X} \mathbf{W}_i^T) - \frac{1}{2} \mathbf{W}_i \boldsymbol{\Sigma}^{-1} \mathbf{W}_i^T \quad (11)$$

Expression 11 can be written into constant, linear and quadratic terms in the perspective of \mathbf{W}_i as shown in 12 where the first term is constant, the second term is linear and the last two terms are quadratic.

$$-\frac{1}{2\sigma^2}\mathbf{T}_i^T\mathbf{T}_i + \frac{1}{\sigma^2}\mathbf{T}_i^T\mathbf{X}\mathbf{W}_i^T - \frac{1}{2\sigma^2}(\mathbf{X}\mathbf{W}_i^T)^T(\mathbf{X}\mathbf{W}_i^T) - \frac{1}{2}\mathbf{W}_i\Sigma^{-1}\mathbf{W}_i^T \quad (12)$$

The quadratic terms in 12 can be used to complete the square by setting them equal to the quadratic exponent in 10.

$$-\frac{1}{2\sigma^2}(\mathbf{X}\mathbf{W}_i^T)^T(\mathbf{X}\mathbf{W}_i^T) - \frac{1}{2}\mathbf{W}_i\Sigma^{-1}\mathbf{W}_i^T = -\frac{1}{2}\mathbf{W}_i\Sigma_w^{-1}\mathbf{W}_i^T \quad (13)$$

The expression on the left can be re-written to the same form as the right side expression.

$$-\frac{1}{2}\mathbf{W}_i\left(\frac{1}{\sigma^2}\mathbf{X}^T\mathbf{X} + \Sigma^{-1}\right)\mathbf{W}_i^T = -\frac{1}{2}\mathbf{W}_i\Sigma_w^{-1}\mathbf{W}_i^T \quad (14)$$

The covariance Σ_w can be identified easily through 14.

$$\Sigma_w = \left(\frac{1}{\sigma^2}\mathbf{X}^T\mathbf{X} + \Sigma^{-1}\right)^{-1} \quad (15)$$

To find the mean μ_w , we have to look at the linear terms of equation 12 and the exponents in equation 10.

$$\frac{1}{\sigma^2}\mathbf{T}_i^T\mathbf{X}\mathbf{W}_i^T = (\mathbf{W}_i\Sigma_w^{-1}\mu_w) \quad (16)$$

Once again, the left side expression can be tweaked to make the calculation easier.

$$\frac{1}{\sigma^2}\mathbf{W}_i^T\mathbf{X}^T\mathbf{T}_i = (\mathbf{W}_i\Sigma_w^{-1}\mu_w) \quad (17)$$

Since we know Σ_w , it can be substituted into equation 17.

$$\frac{1}{\sigma^2}\mathbf{W}_i^T\mathbf{X}^T\mathbf{T}_i = (\mathbf{W}_i\left(\frac{1}{\sigma^2}\mathbf{X}^T\mathbf{X} + \Sigma^{-1}\right)\mu_w) \quad (18)$$

The mean $\boldsymbol{\mu}_w$ can thus be rewritten as shown in expression 19.

$$\boldsymbol{\mu}_w = \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}^{-1} \right)^{-1} \mathbf{X}^T \mathbf{T}_i \quad (19)$$

Since we now know the parameters for the distribution, we can now express the posterior for each output feature i as a multivariate gaussian as shown in 20. Note that the covariance $\boldsymbol{\Sigma} = \tau^2 \mathbf{I}$

$$p(\mathbf{W}_i | \mathbf{X}, \mathbf{T}_i) = \mathcal{N} \left(\frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{\tau^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{T}_i, \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{\tau^2} \mathbf{I} \right)^{-1} \right) \quad (20)$$

Note that the constant term corresponds to the evidence (Z) and is not interesting for finding the gaussian distributions for the posteriors.

Looking at 11, it is clear that we can identify the log-posterior as a least square likelihood with L2-regularization as prior.

1.1.6 Question 6

The prior must be a gaussian distribution in order for it to be conjugate to the posterior gaussian process. The gaussian prior $p(f | \mathbf{X}, \boldsymbol{\theta})$ with a covariance function is chosen to express the property that the values for similar points will be more strongly correlated than for dissimilar points [1, p. 306].

The reason for using a covariance function (kernel) is because it can be used to adjust smoothing in the distribution. Figure 2 illustrates how smoothness can be adjusted using a squared exponential kernel $k(x, x') = \exp \left(-\frac{1}{2\tau^2} \|x - x'\|^2 \right)$ for different values of τ [2].

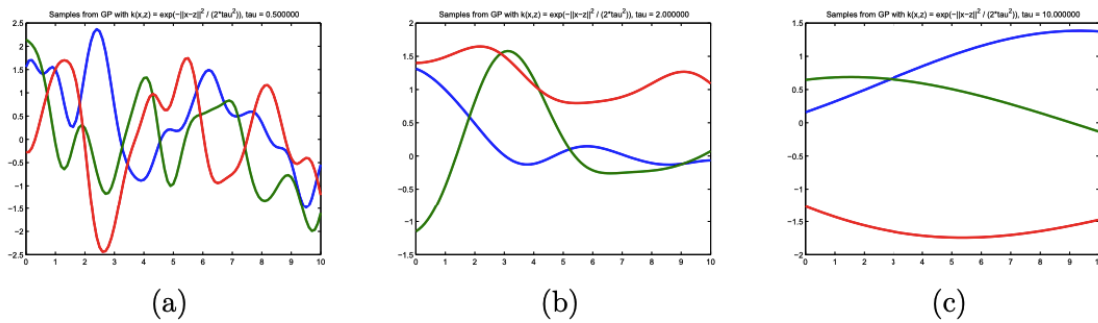


Figure 2: The figure illustrates samples from three different gaussian process priors with squared exponential covariance functions. a) has $\tau = 0.5$, b) has $\tau = 2$ and c) has $\tau = 10$. Image taken from Chuong B. Do[2].

1.1.7 Question 7

Figure 21 is the formula for the joint probability of the four variables. The graphical model 3 displays how the variables are related. f is conditioned on \mathbf{X} and θ and \mathbf{T} is in turn conditioned on f .

$$p(\mathbf{t}, \mathbf{f}, \mathbf{X}, \theta) = p(\mathbf{t}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \theta)p(\mathbf{X})p(\theta) \quad (21)$$

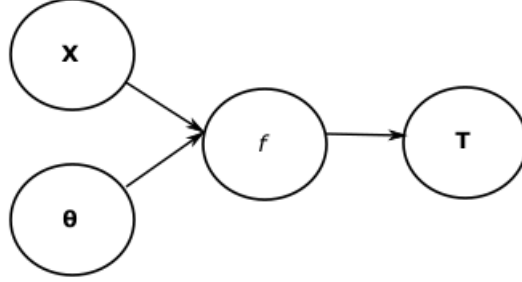


Figure 3: Graphical model displaying the relationships between the variables

1.1.8 Question 8

$$p(\mathbf{T}|\mathbf{X}, \theta) = \int p(\mathbf{T}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \theta)d\mathbf{f} \quad (22)$$

The prior gives all possible functions \mathbf{f} generated by data points from the input space \mathbf{X} and hyperparameters θ . The resulting functions from the prior (\mathbf{f}) are then conditioned on by the output points \mathbf{T} in the likelihood to find the posterior.

There are uncertainties in our data because of the noise (ϵ) and in our functions (\mathbf{f}). Because they are Gaussian sources of uncertainty, their covariances will simply be added and thus pass through the marginalisation [1, p. 307].

The θ is still needed because it is a constant in our covariance function which is needed to calculate the probabilities inside the integral.

1.2 Practical

1.2.1 Question 9

The prior over \mathbf{W} is displayed in figure 4.

It is clear from plots 5 and 6 that the more data points that are used for the model, the more certain it will be in predictions. For example, in the plots 5 the uncertainty of the posterior becomes smaller the more datapoints are used in training (smaller circles in the graph). The same can be said for the lineplots 6, the more datapoints that are used the closer (on average) we get to the true values of $w_0 = 0.5$ and $w_1 = -1.5$. Varying the noise

levels did not affect the uncertainty (variance) of the posterior but it did change the average values (center) of w_0 and w_1 especially for lower amounts of data points.

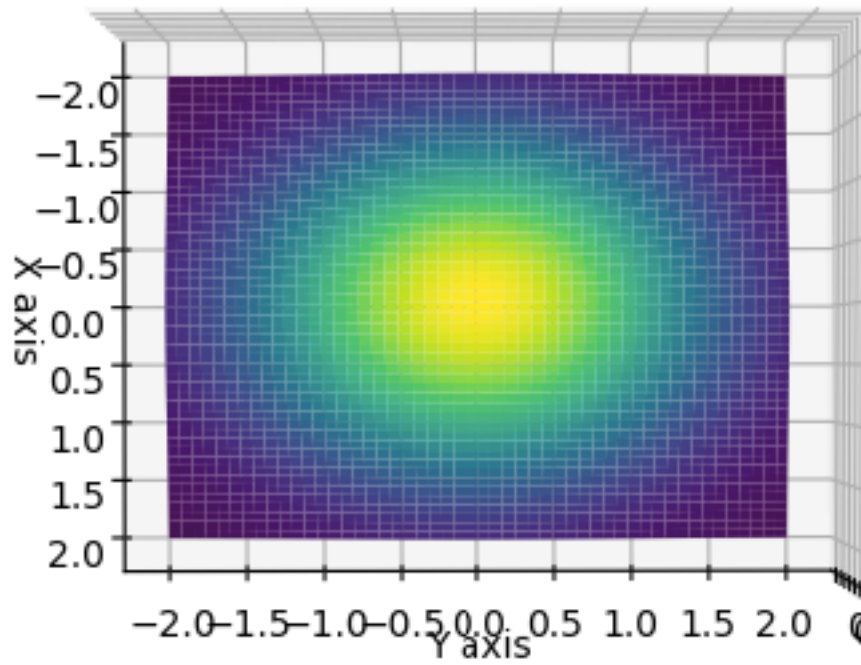


Figure 4: Vizualization of the prior over the \mathbf{W} . The mean is assumed to be 0 and the variance is assumed to be 1.

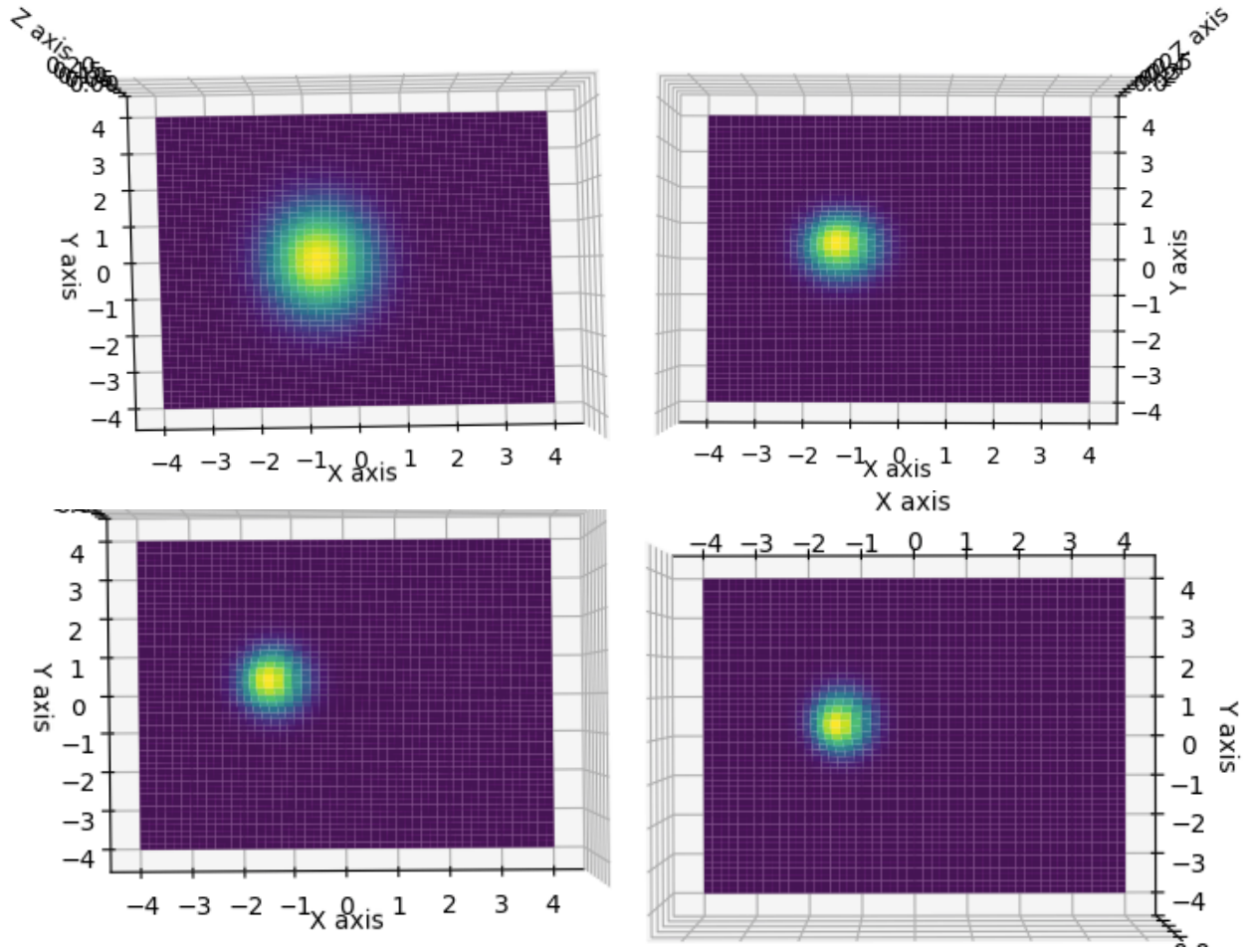


Figure 5: Posterior distributions over W for up to 7 different datapoints. The top-left is for 1 datapoint, top-right for 3 datapoints, bottom-left for 5 datapoints and bottom-right for 7 datapoints. Note that the Y-axis corresponds to w_0 and the X-axis to w_1

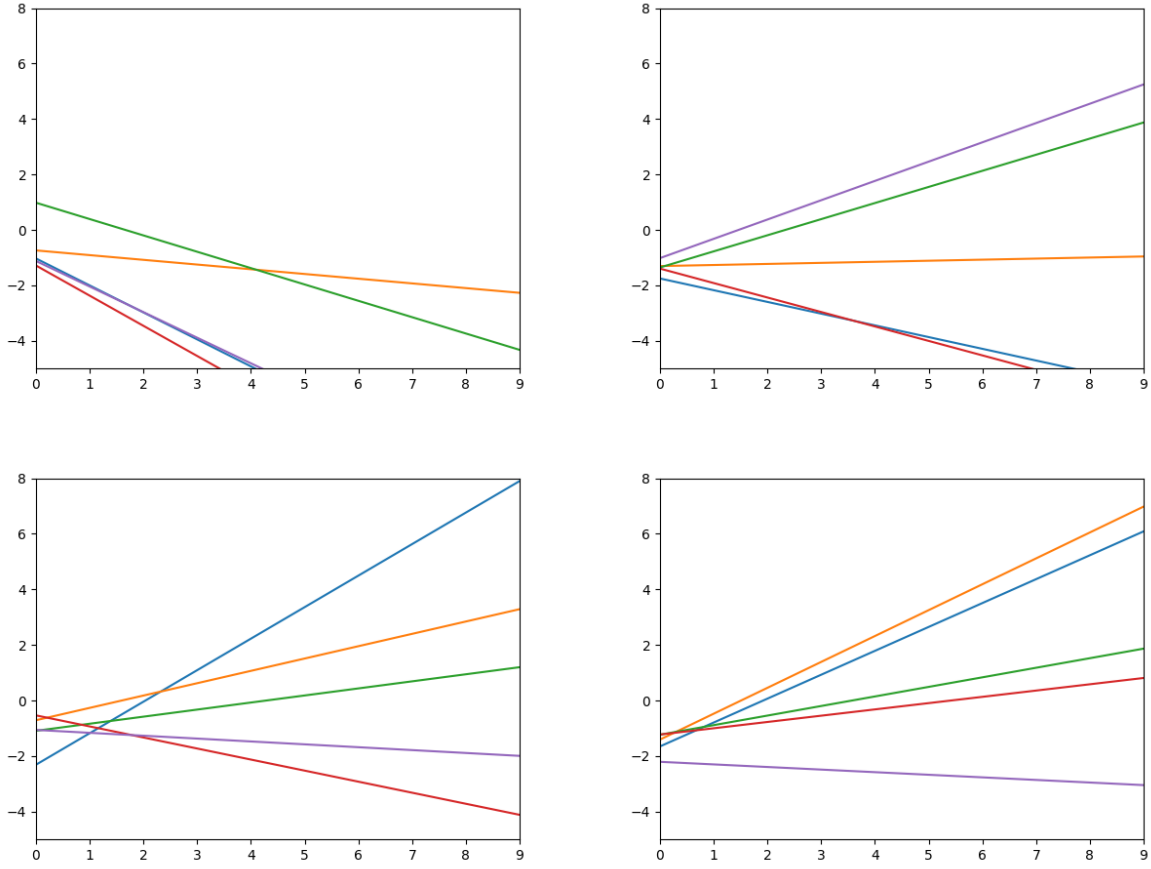


Figure 6: 5 Sample lines drawn from different amount datapoints. Top left is sampled from a distribution with 1 datapoint, top-right from 3 datapoints, bottom-left from 5 datapoints and bottom-right from 7 datapoints. Note that the intersection in the vertical axis when the horizontal axis equals 0 is the predicted value of w_1 while the slope of the line is the predicted value of w_0

1.2.2 Question 10

It is clear from the plots in figure 7 that the length scale adjusts the smoothness of the function curve. For a high value (for example top-left image in 7), the lines will be fairly flat. For lower values on the length scale on the other hand, the plots will be more complex.

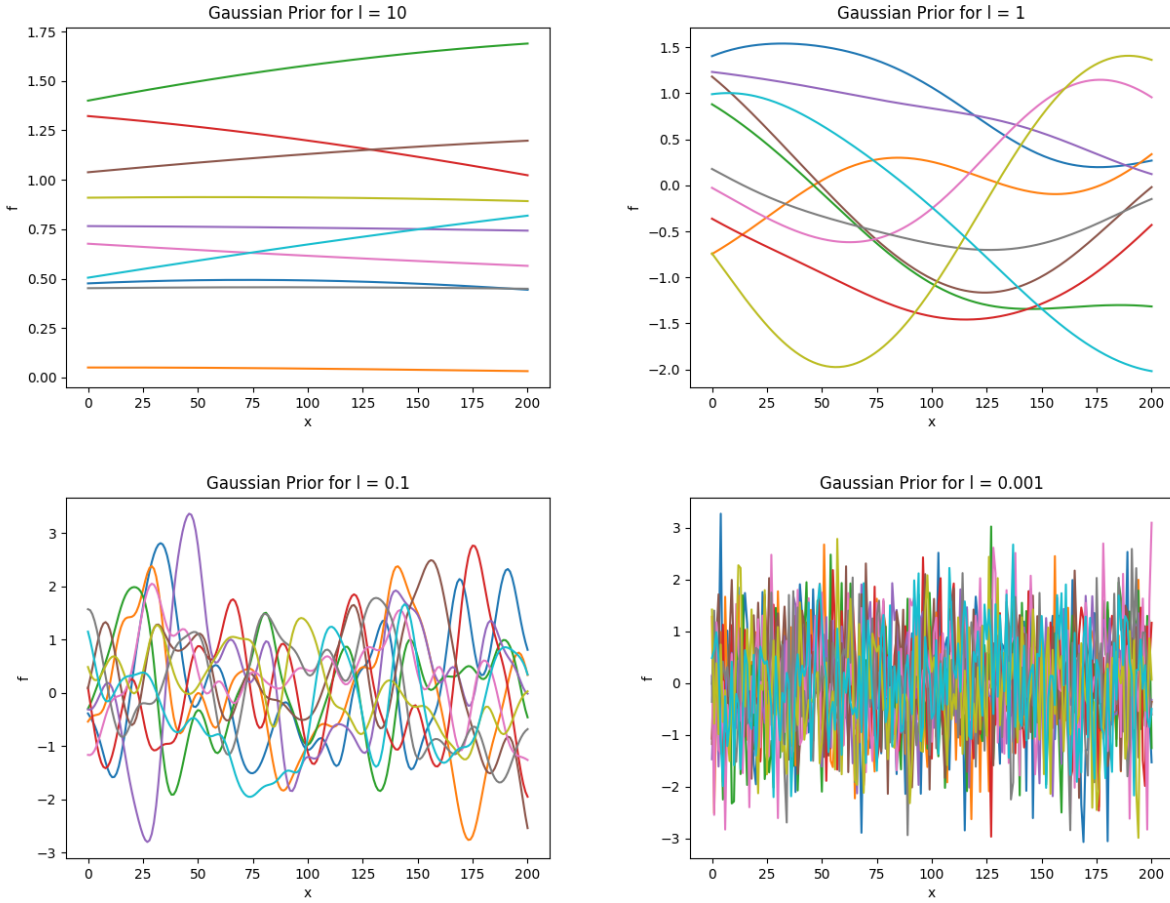


Figure 7: The prior from 10 samples for 4 different values of the length scale l .

1.2.3 Question 11

$$p(\mathbf{f}^* | \mathbf{X}^*, \mathbf{X}, \mathbf{f}) \sim \mathcal{N} \left(\mathbf{k}(\mathbf{X}^*, \mathbf{X})^\top K(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f}, \mathbf{k}(\mathbf{X}^*, \mathbf{X}^*) - \mathbf{k}(\mathbf{X}^*, \mathbf{X})^\top K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{X}^*) \right) \quad (23)$$

The predictive posterior is a gaussian distribution as shown in equation 23 [5, p. 16]. \mathbf{X} is the training data, \mathbf{X}^* is the data that is to be predicted and \mathbf{f} are the functions of the training data. If we were to have no training data available, the only term that would remain would be $\mathbf{k}(\mathbf{X}^*, \mathbf{X}^*)$. Thus the equation can be re-written as 24 which is the same as the gaussian prior (as in the assignment).

$$p(\mathbf{f}^* | \mathbf{X}^*, \mathbf{X}, \mathbf{f}) \sim \mathcal{N}(0, \mathbf{k}(\mathbf{X}^*, \mathbf{X}^*)) \quad (24)$$

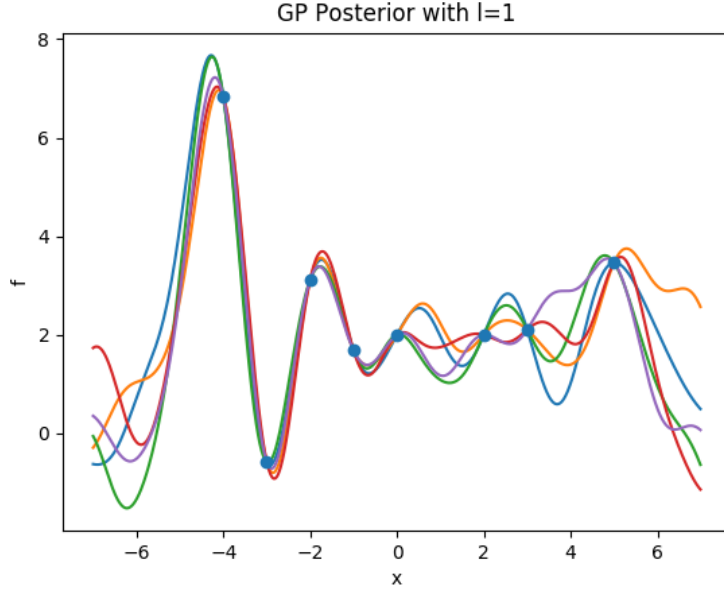


Figure 8: GP Posterior over 1400 observations of x^* where $\{x^* \in \mathbf{Q} \mid -7 \leq x \leq 7\}$. The blue dots are the true values of \mathbf{f} for the training data \mathbf{x} .

In figure 8 you can see the predictive distribution of the posterior where each line is a sample of f . Note that all lines at the true values of the observed data \mathbf{x} .

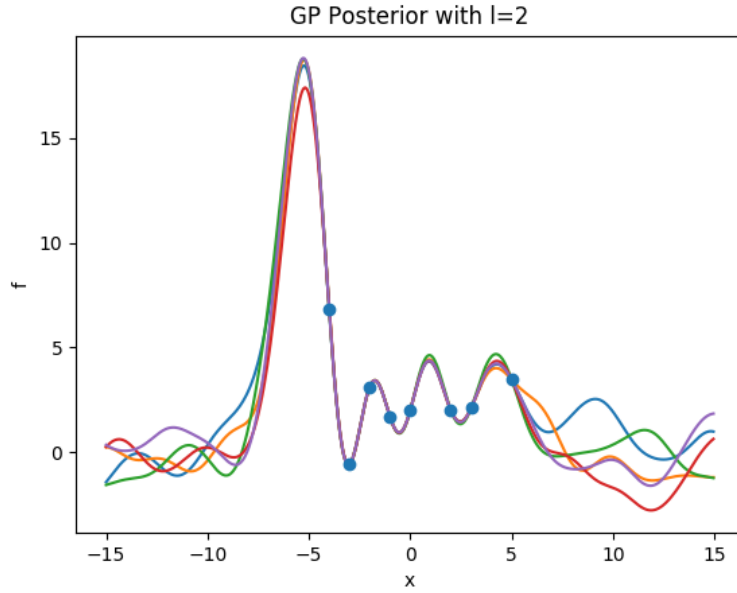


Figure 9: GP Posterior over the x-axis where $\{x^* \in \mathbf{Q} \mid -15 \leq x \leq 15\}$. The blue dots are the true values of \mathbf{f} for the training data \mathbf{x} .

In figure 9 it is clear that the uncertainty grows the further away from the observed data while points close to the observed data generally has lower variance.

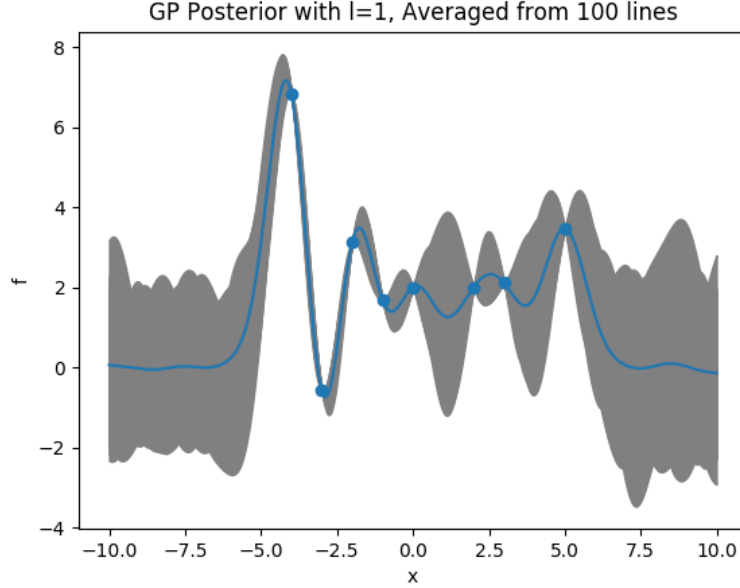


Figure 10: The blue line is the averaged posterior over 100 different samples from many observations $x^* \in \mathbf{Q}$ where $\{x^* \in \mathbf{Q} \mid -10 \leq x \leq 10\}$. The shaded areas is the predictive variance of the posterior over the 100 lines.

Comparing all posterior plots to the prior plot 7 the main difference is that the posterior is fit after our labeled data \mathbf{t} whereas the prior has no knowledge of the data. This is a desirable effect as it leads to better predictions, especially for points close to the training data.

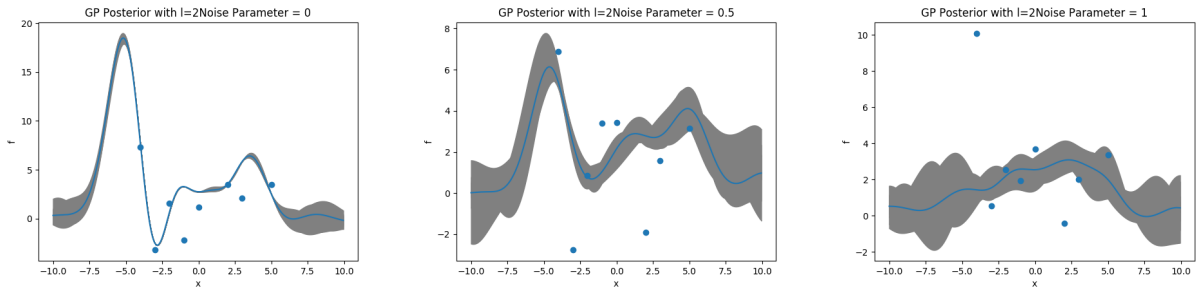


Figure 11: The average predictive posterior distribution with increasing values for the noise variance (σ_n^2) going from left to right.

Figure 11 shows how noisy data is predicted using different values for the noise variance[5] on the diagonal covariance matrix. Higher values for the noise variance lead to higher

uncertainties (variance) for f and it also leads to a smoother function of the predictive mean.

2 The Posterior

2.1 Theory

2.1.1 Question 12

The preferred prior is defined as a Gaussian because of conjunction, with independence over all latent variables \mathbf{X} which are properties that simplifies our derivations and calculations.

2.1.2 Question 13

The prior as given in the assignment 25.

$$p(\mathbf{X}) = \mathcal{N}(0, \mathbf{I}) \quad (25)$$

We want to calculate the marginal likelihood 26.

$$p(\mathbf{Y}|\mathbf{W}) = \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{X})d\mathbf{X} \quad (26)$$

Due to the independence in \mathbf{X} and \mathbf{Y} , the marginal posterior for the entire \mathbf{Y} can be written as a product through Naive Bayes:

$$p(\mathbf{Y}|\mathbf{W}) = \prod p(\mathbf{Y}_i|\mathbf{W}) \quad (27)$$

Thus, we can simply calculate the posterior over each row of \mathbf{Y} by using the likelihood 28 as given from Bishop[1, p. 571] where \mathbf{Y}_i is a single output vector and \mathbf{X}_i is corresponding \mathbf{X} .

$$p(\mathbf{Y}_i|\mathbf{X}_i) = \mathcal{N}(\mathbf{Y}_i|\mathbf{W}\mathbf{X}_i + \boldsymbol{\mu}, \sigma^2\mathbf{I}) \quad (28)$$

Since we know that the marginal is Gaussian, we can derive the mean and covariance from every instance of $p(\mathbf{Y}_i|\mathbf{X}_i)$ defined in equation 28[1, p. 573].

$$\begin{aligned} E[\mathbf{Y}_i] &= E[\mathbf{W}\mathbf{X}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon}] = \boldsymbol{\mu} \\ \text{cov}[\mathbf{Y}_i] &= E[(\mathbf{W}\mathbf{X}_i + \boldsymbol{\epsilon})(\mathbf{W}\mathbf{X}_i + \boldsymbol{\epsilon})^T] \\ &= E[\mathbf{W}\mathbf{X}_i^T\mathbf{X}_i\mathbf{W}^T] + E[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I} \end{aligned} \quad (29)$$

Following 29, we may express our Gaussian as:

$$p(\mathbf{Y}_i|\mathbf{W}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}) \quad (30)$$

2.1.3 Question 14

We know that maximizing the likelihood is the same as minimizing the sum of square error in the log probability [1, p. 23]:

$$\ln p(\mathbf{Y}|\mathbf{W}, \mathbf{X}) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (\mathbf{Y}_n - \mathbf{W}\mathbf{X}_n)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \quad (31)$$

MAP on the other hand is a means of finding the parameters that maximise the posterior. In this case, the log-prior is added onto the maximum likelihood as a constraint[1, p. 30]:

$$\ln p(\mathbf{Y}|\mathbf{W}, \mathbf{X}) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (\mathbf{Y}_n - \mathbf{W}\mathbf{X}_n)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) - \frac{1}{2\sigma_w} \mathbf{W}\mathbf{W}^T \quad (32)$$

Note that for ML, there is no prior knowledge about the observations, so far smaller number of observations, there can occur overfitting. This in means that ML is more useful when a high amount of data has been observed and less useful for low amount of observed data. MAP on the other hand has a regularization term that counteracts the problem of overfitting making it useful even for low amounts of observed data [1, p. 29]. Also worth noting is that MAP will share the same behaviour as ML when sufficient data has been observed (i.e when the observed data is very large). For type 2 maximum likelihood, the observed variables are marginalised out, meaning there is no need for observed data as long as there is some belief of the prior distribution of the data. Thus there is no need for actual data when using Type 2 ML.

$$\hat{\mathbf{W}} = \operatorname{argmax}_{\mathbf{W}} \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W})}{\int p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W})d\mathbf{W}} = \operatorname{argmax}_{\mathbf{W}} p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W}) \quad (33)$$

The integral can be re-written in terms of only X and Y .

$$\hat{\mathbf{W}} = \operatorname{argmax}_{\mathbf{W}} \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W})}{p(\mathbf{Y}|\mathbf{X})} = \operatorname{argmax}_{\mathbf{W}} p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W}) \quad (34)$$

Type 2 ML is the best suitable choice in our case because we have no observed data to use, but we have an assumption of how the data is distributed. Thus we can simply marginalize out the observed data.

2.1.4 Question 15

For convinience sake, we will call the covariance \mathbf{C}

$$p(\mathbf{Y}) = \mathcal{N}(\mu, \mathbf{C}) \quad (35)$$

This is of course the same covariance that we calculated in 30.

$$\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I} \quad (36)$$

The inverse of \mathbf{C} is needed for the next step:

$$\mathbf{C}^{-1} = \sigma^{-1}\mathbf{I} - \sigma^{-2}\mathbf{W}(\mathbf{W}^T\mathbf{W} + \sigma^2\mathbf{T})^{-1}\mathbf{W}^T \quad (37)$$

Now on to our log-posterior which of course is made up of the sum of all individual probabilities for \mathbf{Y}_i . Note that N is the amount of observations and D is the amount of features.

$$\ln p(\mathbf{Y}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \ln p(\mathbf{Y}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{Y}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{Y}_n - \boldsymbol{\mu}) \quad (38)$$

$$\mathcal{L}(\mathbf{W}) = \frac{N}{2} \left\{ D \ln(2\pi) + \ln |\mathbf{C}| + \text{Tr} \left(\mathbf{C}^{-1} \frac{1}{N} \sum_{n=1}^N (\mathbf{Y}_n - \boldsymbol{\mu}) (\mathbf{Y}_n - \boldsymbol{\mu})^T \right) \right\} \quad (39)$$

For the gradient we follow rules derived from The Matrix Cookbook[4].

$$\begin{aligned} \partial(\det(X)) &= \text{Tr}(X^{-1}\partial X) \\ \partial \text{Tr}(X) &= \text{Tr}(\partial X) \\ \partial X^{-1} &= -X^{-1}(\partial X)X^{-1} \end{aligned} \quad (40)$$

We can write the derivative of \mathbf{C} in terms of \mathbf{W} :

$$\frac{\partial \mathbf{C}}{\partial W_{ij}} = \frac{\partial \mathbf{W}\mathbf{W}^T}{\partial W_{ij}} \quad (41)$$

We choose to solve the determinant and the trace separately. Starting with the determinant.

$$\frac{\partial(\det(\mathbf{C}))}{\partial W_{ij}} = \text{Tr}(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial W_{ij}}) \quad (42)$$

Now for the trace, firstly the partial can be moved into the trace because of the rules in 40 meaning we want to solve the derivative inside the trace, note that constants with respect to \mathbf{W} can be removed from the partial:

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{Y}_n - \boldsymbol{\mu}) (\mathbf{Y}_n - \boldsymbol{\mu})^T \frac{\partial(\mathbf{C}^{-1})}{\partial W_{ij}} \quad (43)$$

$$\frac{\partial(\mathbf{C}^{-1})}{\partial W_{ij}} = -\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial W_{ij}} \mathbf{C}^{-1} \quad (44)$$

Thus our final derivative with respects to \mathbf{W} is:

$$\frac{\mathcal{L}(\mathbf{W})}{\partial W_{ij}} = \frac{N}{2} \text{Tr}(\mathbf{C}^{-1} \frac{\partial \mathbf{W}\mathbf{W}^T}{\partial W_{ij}}) - \text{Tr}((- \mathbf{C}^{-1} \frac{\partial \mathbf{W}\mathbf{W}^T}{\partial W_{ij}} \mathbf{C}^{-1}) \frac{1}{N} \sum_{n=1}^N (\mathbf{Y}_n - \boldsymbol{\mu}) (\mathbf{Y}_n - \boldsymbol{\mu})^T) \quad (45)$$

2.1.5 Question 16

References

- [1] Christopher M. Bishop. *Pattern recognition and machine learning*. en. Information science and statistics. New York: Springer, 2006. ISBN: 978-0-387-31073-2.
- [2] Chuong B. Do. *Gaussian Processes*. en. 2008. URL: http://cs229.stanford.edu/section/cs229-gaussian_processes (visited on 12/04/2019).
- [3] *Multivariate normal distribution*. en. Page Version ID: 922344737. Oct. 2019. URL: https://en.wikipedia.org/w/index.php?title=Multivariate_normal_distribution&oldid=922344737 (visited on 11/18/2019).
- [4] Kaare Brandt Petersen and Michael Syskind Pedersen. “[<http://matrixcookbook.com>]”. en. In: (), p. 72.
- [5] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. en. Adaptive computation and machine learning. OCLC: ocm61285753. Cambridge, Mass: MIT Press, 2006. ISBN: 978-0-262-18253-9.