

# DD2434 Machine Learning, Advanced Course - Assignment 1

Hannes Kindbom

## 1 The prior

### 1.1 Theory

#### Question 1

We often model the relationship between two single points  $\mathbf{x}_i \in R^q$  and  $\mathbf{t}_i \in R^D$  as  $\mathbf{t}_i = f(\mathbf{x}_i) + \epsilon$ , where the noise  $\epsilon$  follows a Gaussian distribution.

One argument for assuming a Gaussian distribution on the noise term is that the Gaussian distribution has certain properties which makes it easy to work with. The Gaussian choice must however of course be appropriate and not off.

Since we often do not have information about the noise distribution, we can assume that it consists of a large sum of independent and identically distributed random variables (measurement errors) with finite variance. According to the central limit theorem, this large sum will converge in distribution to Gaussian. We do also assume that all  $\epsilon$ 's and  $\epsilon_i$ 's are independent. The likelihood is a linear combination ( $f(\mathbf{x}_i) + \epsilon$ ) of the noise random variables conditioned on  $f$  and  $\mathbf{x}_i$ , which is thus also (multivariate) Gaussian. This is the rational argument behind the Gaussian assumption on the likelihood.

The components in  $\mathbf{t}_i$  of this multivariate Gaussian likelihood are independent (since the noise terms are independent), which implies that the covariances between them are 0. Thus, the covariance matrix is proportional to the identity matrix. To be specific, this spherical covariance matrix comes from the independence assumption on the noise terms.

#### Question 2

If we assume that each output point is conditionally dependent given the input and the mapping, the likelihood of the data takes the following form:

$$p(\mathbf{T}|f, \mathbf{X}) = p(\mathbf{t}_1|\mathbf{t}_{1+1}, \dots, \mathbf{t}_N, f, \mathbf{X}) \prod_{i=2}^N p(\mathbf{t}_i|f, \mathbf{X}) \quad (1)$$

where we have applied the product rule iteratively.

#### 1.1.1 Linear Regression

#### Question 3

Suppose we assume  $\mathbf{t}_i = \mathbf{W}\mathbf{x}_i + \epsilon$ ,  $i = 1, \dots, N$ , where  $\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ . Then the likelihood of the

data becomes:

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

where  $\mathbf{t}_i|\mathbf{x}_i, \mathbf{W} \sim N(\mathbf{W}\mathbf{x}_i, \sigma^2\mathbf{I})$ . Note that we get a product on the right hand side of equation 2, since the random variables  $\mathbf{t}_i$  in the joint distribution  $\mathbf{T}|\mathbf{X}, \mathbf{W}$  are conditionally independent given  $\mathbf{x}_i, \mathbf{W}$ .

#### Question 4

The conjugate prior for the Gaussian likelihood is a Gaussian prior. A Gaussian prior over each row of  $\mathbf{W}$  on the form  $p(\mathbf{w}) = N(\mathbf{w}_0, \tau^2\mathbf{I})$  will thus lead to a Gaussian posterior with the term  $-\frac{1}{2\tau^2}(\mathbf{w} - \mathbf{w}_0)^T(\mathbf{w} - \mathbf{w}_0)$  being part of the exponent, since *posterior*  $\propto$  *likelihood* \* *prior*. When searching for the MAP (maximum a posteriori estimator), the posterior can be logarithmised. Thereafter maximizing with respect to  $\mathbf{w}$  is equivalent to minimization of a sum of squares error function with a quadratic ( $L_2$  distance) penalization term:  $\frac{1}{2\tau^2}(\mathbf{w} - \mathbf{w}_0)^T(\mathbf{w} - \mathbf{w}_0)$  (this is the negative log-prior), where  $\frac{1}{2\tau^2}$  can be seen as the regularization coefficient. Regularization is used to reduce variance and overfitting. If  $\mathbf{w}_0 = 0$ , then the estimated weights  $\mathbf{w}_{MAP}$  will be encouraged to tend towards 0, if not supported by the data.

If we instead of a Gaussian prior, choose a Laplace prior, the preferred model would be encoded with  $L_1$  norm for model parameters. The one dimensional version of a Laplace distribution may be written as  $p(w) = \frac{1}{2b} \exp(-\frac{|w-w_0|}{b})$ , where  $w_0 \in R$  is the mean and  $b > 0$  is a scaling coefficient. The penalization term (part of the log posterior) in this case thus becomes  $\frac{|w-w_0|}{b}$ .

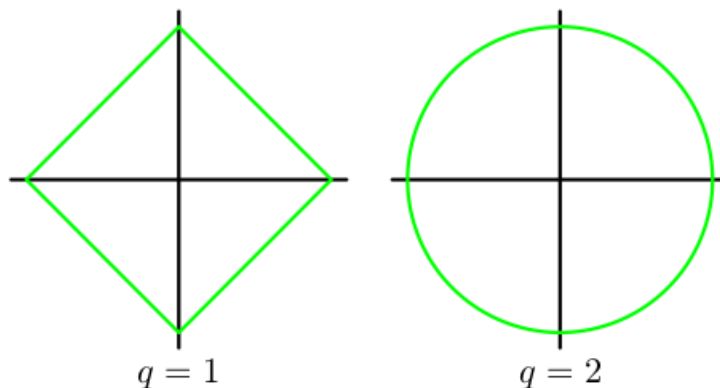


Figure 1: Contours of two different regularization terms when  $\mathbf{w}$  is two dimensional. The left plot is for  $L_1$  and the right is for  $L_2$  regularization.

The difference between having  $L_1$  and  $L_2$  as penalizing terms is mainly that  $L_1$  works more as a feature selector, forcing some components in  $\mathbf{w}$  to be zero if the regularization coefficient is sufficiently large. The  $L_2$  norm will generally not make any feature exactly zero, only small. This can

be seen from the contours of the regularization term in figure 1<sup>1</sup> or from the fact that the density in the Laplace distribution is highly concentrated and peaked at its mean, while the Gaussian density is more wide spread. If we set the mean of the prior to be zero, then the Laplace prior will force weights to zero, unless the likelihood indicates something else.

### Question 5

The posterior when using a Gaussian prior and likelihood can be derived by using that  $p(\mathbf{W}|\mathbf{X}, \mathbf{T}) \propto p(\mathbf{T}|\mathbf{X}, \mathbf{W})p(\mathbf{W})$ . Gaussian is the conjugate prior to Gaussian, which means that the posterior will also be Gaussian. Let us firstly derive an expression for the likelihood, when assuming conditional independence of the target variables, as follows:

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

Equation 3 can be written as:

$$p(\mathbf{T}|\mathbf{X}, \mathbf{W}) = N(\mathbf{W}\mathbf{X}, \sigma^2 \mathbf{I}) \quad (4)$$

because a product of multiple Gaussian distributions is in turn Gaussian and each factor's covariance matrix is diagonal (spherical).

Let us now assume a Gaussian prior on the form  $p(\mathbf{W}) = N(\mathbf{W}_0, \tau^2 \mathbf{I})$ . Note thereafter that the probability density function of a multivariate Gaussian distribution with mean  $\mathbf{m}$  and covariance matrix  $\Sigma$  is proportional to  $\exp(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \Sigma^{-1}(\mathbf{x} - \mathbf{m}))$ . Hence, we can obtain an expression for the posterior:

$$p(\mathbf{W}|\mathbf{X}, \mathbf{T}) \propto \exp(-\frac{1}{2\sigma^2}(\mathbf{T} - \mathbf{W}\mathbf{X})^T(\mathbf{T} - \mathbf{W}\mathbf{X})) \exp(-\frac{1}{2\tau^2}(\mathbf{W} - \mathbf{W}_0)^T(\mathbf{W} - \mathbf{W}_0)) \quad (5)$$

As mentioned above, this will result in a Gaussian posterior. Let us take a closer look on the exponent to figure out the mean and covariance of the posterior. Adding the exponents in equation 5 yields:

$$-\frac{1}{2\sigma^2}(\mathbf{T} - \mathbf{W}\mathbf{X})^T(\mathbf{T} - \mathbf{W}\mathbf{X}) - \frac{1}{2\tau^2}(\mathbf{W} - \mathbf{W}_0)^T(\mathbf{W} - \mathbf{W}_0) \quad (6)$$

We often assume  $\mathbf{W}_0 = \mathbf{0}$ . Then, expression 6 simplifies to:

$$-\frac{1}{2}\left[\frac{1}{\sigma^2}(\mathbf{T}^T \mathbf{T} - 2\mathbf{X}^T \mathbf{W}^T \mathbf{T}) + \mathbf{W}^T \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{\tau^2} \mathbf{I}\right) \mathbf{W}\right] \quad (7)$$

We can now identify the covariance matrix  $\Sigma$  of the posterior immediately as:

$$\Sigma = \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{\tau^2} \mathbf{I}\right)^{-1} \quad (8)$$

The term  $-2\mathbf{X}^T \mathbf{W}^T \mathbf{T} \frac{1}{\sigma^2}$ , can be identified as the term  $-2\mathbf{W}^T \Sigma^{-1} \mathbf{m}$  obtained from expanding the general form of a multivariate Gaussian exponent  $(\mathbf{W} - \mathbf{m})^T \Sigma^{-1}(\mathbf{W} - \mathbf{m})$ . Rearranging a bit yields:

$$\mathbf{m} = \Sigma \mathbf{X}^T \mathbf{T} \frac{1}{\sigma^2} = \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{1}{\tau^2} \mathbf{I}\right)^{-1} \mathbf{X}^T \mathbf{T} \frac{1}{\sigma^2} \quad (9)$$

---

<sup>1</sup>C.M. Bishop. Pattern recognition and machine learning. 2nd ed. 2006. p.145

which is the desired mean of the posterior. We have thus found that the posterior is as follows:

$$p(\mathbf{W}|\mathbf{X}, \mathbf{T}) \sim N(\mathbf{m}, \Sigma) \quad (10)$$

Note that the least squares estimator  $\mathbf{W}_{LS}$  of  $\mathbf{W}$  is given by  $\mathbf{W}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}$  when  $\mathbf{X}$  has full rank. When using  $L_2$  regularization in ridge regression we get  $\mathbf{W}_{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{T}$ , where  $\lambda$  is a real constant. We can observe that the least squares estimate of  $\mathbf{W}$  is somewhat similar to the mean  $\mathbf{m}$  of the posterior but the ridge estimate is essentially on the same form. Since the Gaussian posterior is maximized at the mean value,  $\mathbf{m}$  is the MAP of  $\mathbf{W}$ . Hence, we can say that Bayesian linear regression has incorporated regularization and deals with overfitting, which can otherwise be a problem in the frequentist's approach.

Note that throughout the derivation of the posterior, we have omitted the normalizing constant  $Z$ . I.e. the posterior is actually on the form  $p(\mathbf{W}|\mathbf{X}, \mathbf{T}) = \frac{1}{Z} p(\mathbf{T}|\mathbf{X}, \mathbf{W}) p(\mathbf{W})$ .  $Z$  fills the purpose of making the posterior a probability distribution, which integral over the real line converges to 1.  $Z$  is furthermore often called the *evidence*, used for model selection.

### 1.1.2 Non-parametric Regression

#### Question 6

Assume we have the relationship  $\mathbf{t}_i = \mathbf{f}(\mathbf{x}_i) + \epsilon$ , where  $\mathbf{f}_i$  is considered as a random variable. When we take a non-parametric approach and assign  $\mathbf{f}$  to be a Gaussian Process on the space of real valued functions on the input space, we put a prior on the space of functions instead of on some parameters. The problem of "manually" choosing an appropriate parametric form of the mapping is therefore avoided. Samples from a Gaussian process prior are illustrated in figure 8, where it can be observed that the curves do not follow a specified parametric mapping. We can formally define the Gaussian Process prior as:

$$p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) = N(\mathbf{0}, k(\mathbf{X}, \mathbf{X})) \quad (11)$$

where  $k$  is the covariance function and  $\boldsymbol{\theta}$  the hyper-parameters. The joint distribution of any  $\mathbf{f}_1(x_1), \dots, \mathbf{f}_n(x_n)$  is normally distributed per definition in a Gaussian Process. The marginal distribution  $p(\mathbf{f})$  is also normally distributed, which motivates the choice of this prior; The mapping between  $\mathbf{X}$  and  $\mathbf{T}$  can result in any real value in the output space, which is possible with the normal distribution on  $\mathbf{f}$ .

Regarding the mean and covariance function, we often set the mean to zero in the prior due to lack of knowledge about  $\mathbf{f}$ . The covariance function can be chosen on the form, such that similar  $x_i$ 's produce corresponding  $\mathbf{f}_i$ 's that are highly correlated and vice versa. This may be a desired behaviour. One such covariance function is the squared exponential covariance function. The hyper-parameters  $\boldsymbol{\theta}$  furthermore allow for customization of the covariance function to the specific dataset.

#### Question 7

The joint likelihood of the full model above can without making any assumptions (only applying the product rule) be formulated as follows:

$$p(\mathbf{T}, \mathbf{X}, \mathbf{f}, \boldsymbol{\theta}) = p(\mathbf{T}|\mathbf{f}, \mathbf{X}, \boldsymbol{\theta}) p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) p(\mathbf{X}, \boldsymbol{\theta}) \quad (12)$$

Using the assumptions that  $\mathbf{X}$  and  $\boldsymbol{\theta}$  are independent and that  $\mathbf{T}$  is conditionally independent of  $\mathbf{X}$  and  $\boldsymbol{\theta}$  given  $\mathbf{f}$ , equation 12 may be written as:

$$p(\mathbf{T}, \mathbf{X}, \mathbf{f}, \boldsymbol{\theta}) = p(\mathbf{T}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})p(\mathbf{X})p(\boldsymbol{\theta}) \quad (13)$$

Figure 2 represents the right-hand side of (13) in terms of a simple graphical model (a so called Bayesian Network).

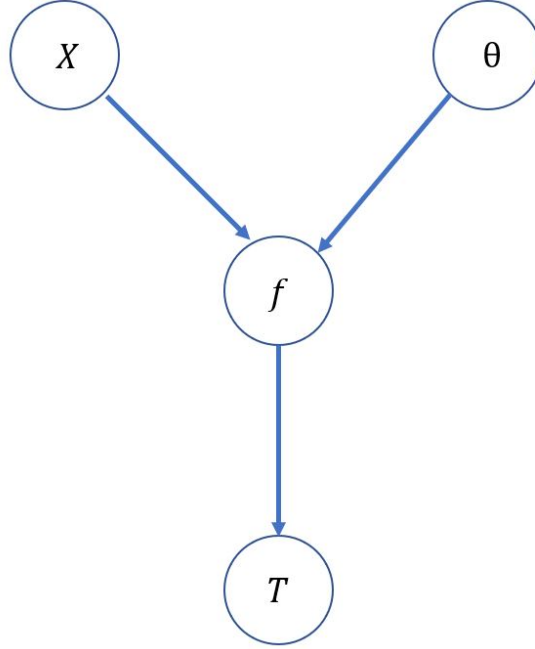


Figure 2: Bayesian Network for (13), where directed arrows are drawn from the nodes corresponding to the variables on which the distribution is conditioned.

### Question 8

An expression for  $p(\mathbf{T}|\mathbf{X}, \boldsymbol{\theta})$  is derived below, using the result from (13).

$$p(\mathbf{T}, \mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) = p(\mathbf{T}, \mathbf{X}, \mathbf{f}, \boldsymbol{\theta}) / (p(\mathbf{X}, \boldsymbol{\theta})) = p(\mathbf{T}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) \implies$$

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

Here  $\mathbf{f}$  is marginalized out since we are not interested in it directly.  $p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})$  in the integral above can be identified as the prior over instantiation of functions in (11).  $p(\mathbf{T}|\mathbf{f})$  is the likelihood. The data in  $\mathbf{T}$  and  $\mathbf{X}$  is thus connected and the uncertainty filter through, since the uncertainty in  $\mathbf{T}$  remains.  $\boldsymbol{\theta}$  is held constant during the integration, wherefore we still condition on it after the marginalization.

## 1.2 Practical

### 1.2.1 Linear Regression

#### Question 9

A synthetic dataset according to (14) was first generated.

$$t_i = w_0 x_i + w_1 + \epsilon = 0.5x_i - 1.5 + \epsilon, \text{ where } \epsilon \sim N(0, \sigma^2) \text{ and } \mathbf{x} = [-1, -0.99, \dots, 0.99, 1] \quad (14)$$

The generated dataset for standard deviation on the noise  $\sigma = 0.2$  is shown in figure 3. A Gaussian prior on  $W = [w_0, w_1]$  was thereafter set to  $W \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \tau^2 I\right)$ , where  $\tau = 0.5$ . The prior is illustrated in figure 4.

The posterior was thereafter computed, using the theory described under question 5 above. The results are displayed in figure 5. As expected, the posterior gets more and more concentrated around the true value of  $W$ , as more data is added to the likelihood. This is reasonable, since it gets increasingly influenced by the dataset (generated with the true  $W$ ).

Figure 6 shows the resulting functions, after drawing random samples from the seven different posteriors. The graphs in 6 are ordered in the same way as in figure 5. Again as expected, the function lines get more and more concentrated around the true blue fat line as more data is added to the likelihood. It is especially difficult to get a reasonable function line with only one data point, as can be observed in the first graph.

The same experiment was thereafter performed with  $\sigma = 0.4$  as well as  $\sigma = 0.8$ . The results are displayed in figure 7. It may be observed that the increased variance  $\sigma$  is reflected in the variance of the posterior as well as in the sampled function lines. This is consistent with the theory and a higher variance in the data is intuitively reflected in our parameter estimates.

### 1.2.2 Non-parametric Regression

#### Question 10

10 samples per length scale  $l$  were drawn from a Gaussian Process prior with four different length scales in the squared exponential kernel. The results are shown in figure 8. The Gaussian Process prior is defined as in (11). The squared exponential kernel is defined as in (15).

$$k(x_i, x_j) = \sigma_f^2 \exp(-(x_i - x_j)^T(x_i - x_j)/l^2), \sigma_f^2 = 4 \quad (15)$$

As can be observed in figure 8,  $f(x_i)$  becomes increasingly smooth/stable as  $l$  increases. The explanation for this is that, as  $l$  increases,  $k(x_i, x_j)$  increases for fixed signal variance  $\sigma_f^2$ . A larger  $k(x_i, x_j)$ , means that  $f(x_i)$  and  $f(x_j)$  will have a stronger correlation, hence the smoother graphs. The opposite holds for small  $l$ .

#### Question 11

Suppose  $\mathbf{x}_*$  is a collection of new inputs  $\mathbf{x}_1^*, \dots, \mathbf{x}_m^*$ . Also, let  $\mathbf{T}$  be a collection of training target variables  $\mathbf{t}_1, \dots, \mathbf{t}_n$  and  $\mathbf{X}$  a collection of training inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . If the same relationship as stated under question 6 is assumed, then the joint distribution between  $\mathbf{T}$  and  $\mathbf{f}_*$  is Gaussian as

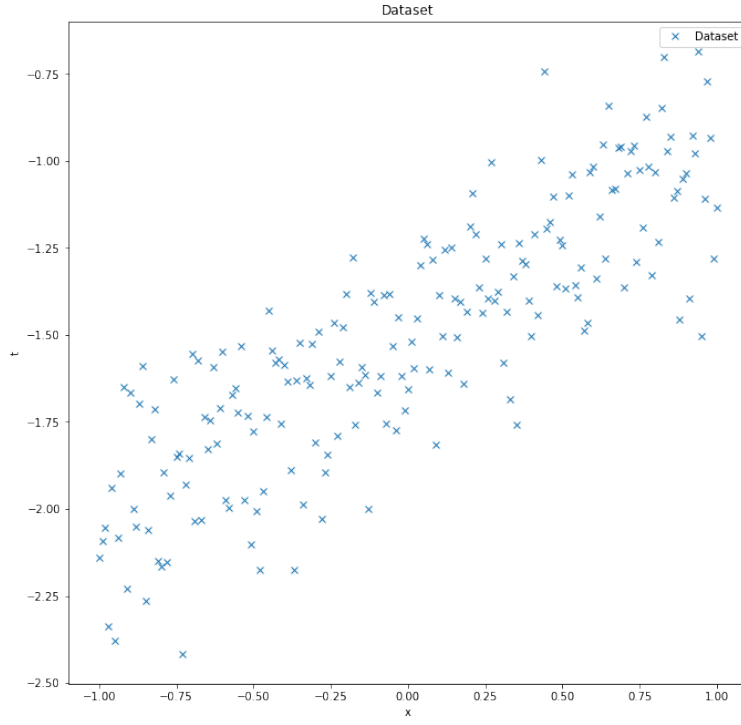


Figure 3: Generated dataset from equation 14 with  $\sigma = 0.2$

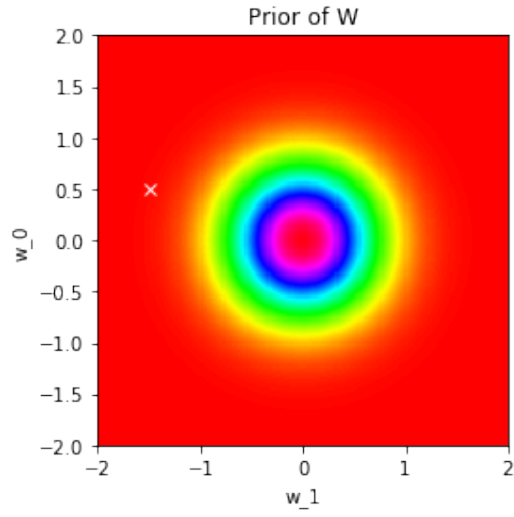


Figure 4: The prior on  $W$ , where the white X marks the true values of  $W$ .

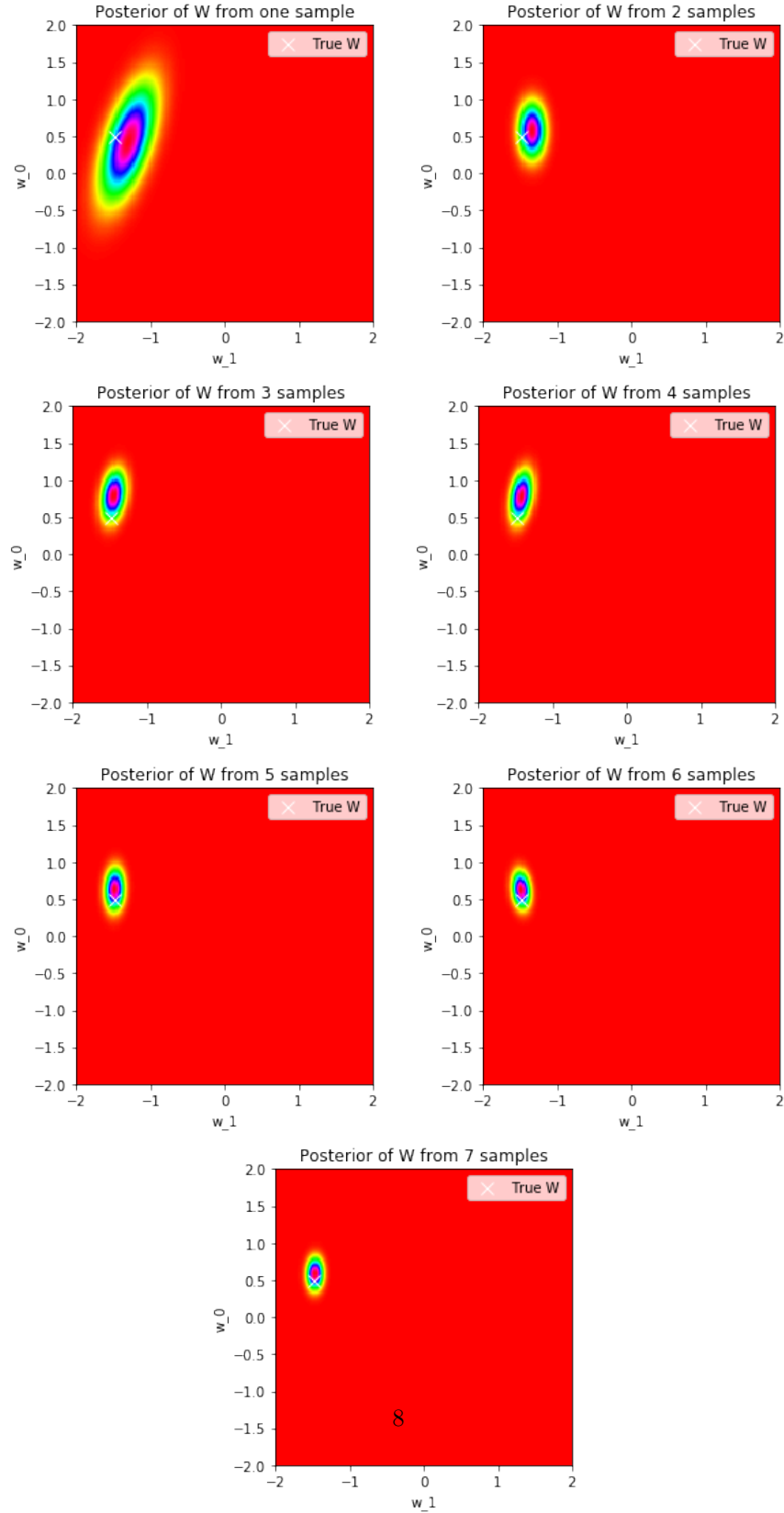


Figure 5: Posteriors for  $W$ , where random data points from the dataset have been added successively (starting at 1 ending at 7 data points). As before, the white cross marks the true values of  $W$ .



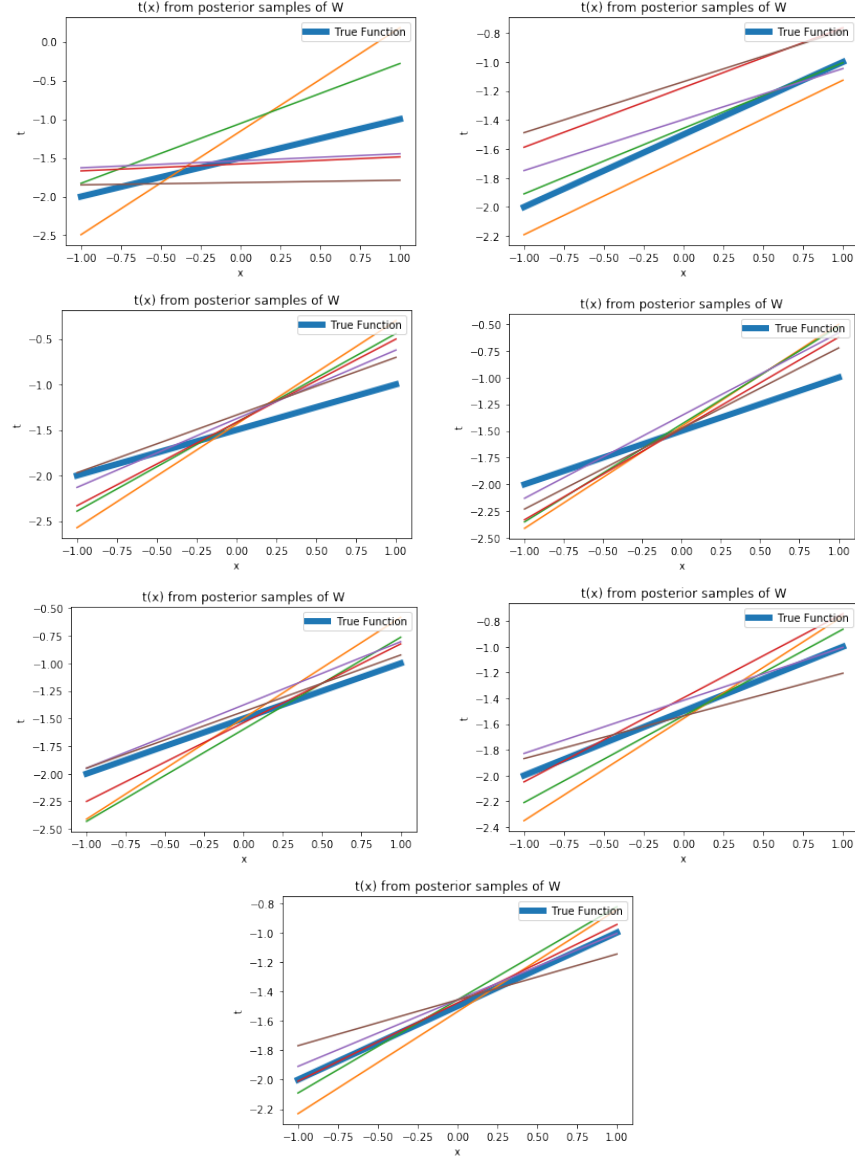


Figure 6: The function  $t = w_0x + w_1$  with five different  $W$ , sampled from the seven different posteriors in figure 5 (plots are ordered the same way). The "fat" blue line shows the true function.

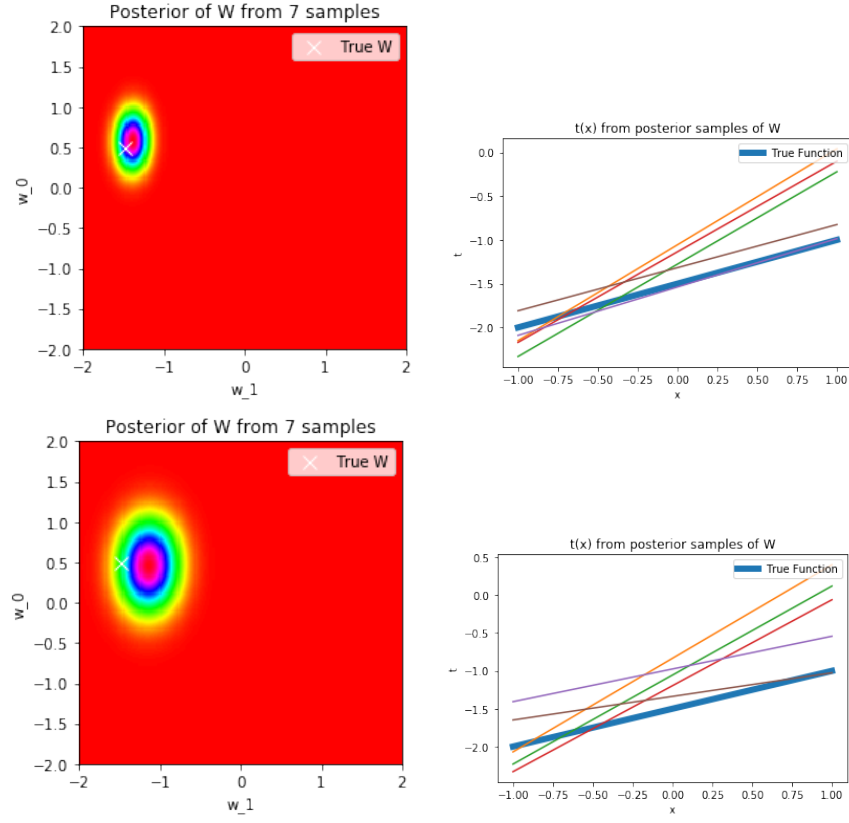


Figure 7: Posteriors for  $W$  and corresponding function lines, where 7 random data points from the data are included in the likelihood. As before, the white cross marks the true values of  $W$  and the blue line is the true one. The standard deviation on the noise term in the data is now  $\sigma = 0.4$  in the first row and  $\sigma = 0.8$  in the second row

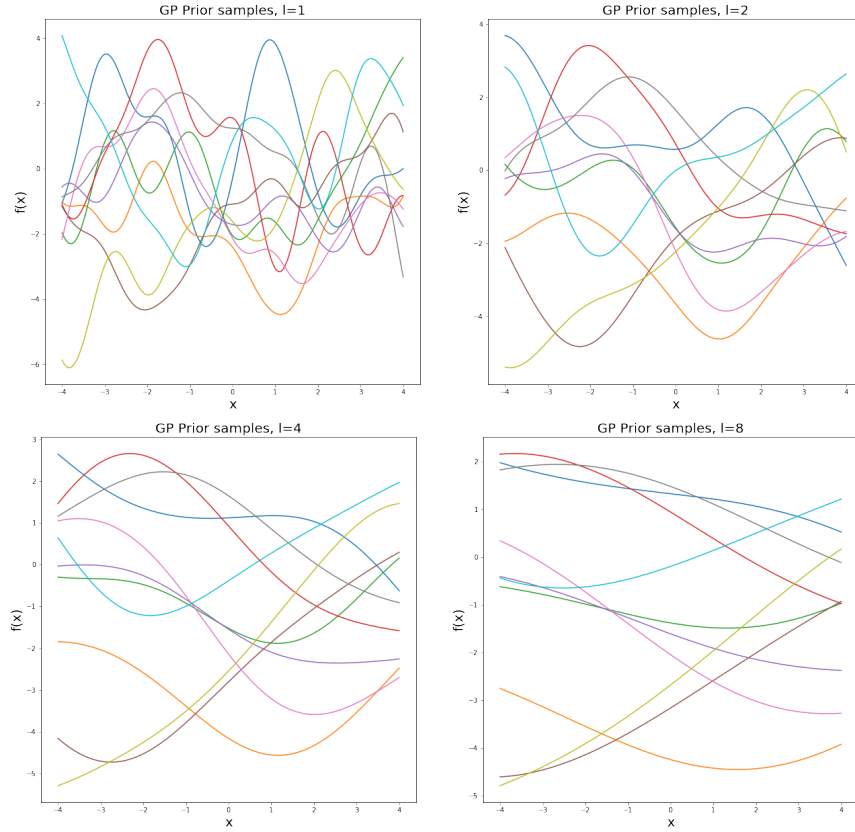


Figure 8: 10 samples from a Gaussian Process prior with varying length scale  $l$  in the squared exponential kernel.

follows:

$$\begin{bmatrix} \mathbf{T} \\ \mathbf{f}_* \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma^2 I & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right) \quad (16)$$

where  $\sigma^2 I$  is the covariance matrix for the error term  $\epsilon$ . The posterior  $\mathbf{f}_* | \mathbf{T}, \mathbf{X}, \mathbf{x}_*$  is then also a Gaussian process with mean function  $m_{post} = k(\mathbf{x}_*, \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \sigma^2 I)^{-1} \mathbf{T}$  and covariance function  $cov_{post} = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})[k(\mathbf{X}, \mathbf{X}) + \sigma^2 I]^{-1} k(\mathbf{X}, \mathbf{x}_*)$ . However if no data is observed and  $\mathbf{T}$  and  $\mathbf{X}$  do not exist, we have  $\mathbf{f}_* | \mathbf{x}_* \sim N(\mathbf{0}, k(\mathbf{x}_*, \mathbf{x}_*))$ , which is merely the Gaussian Process prior in (11) for  $\mathbf{x}_*$ .

A synthetic dataset was then generated according to the following:

$$t_i = (2 + (0.5x_i - 1)^2) * \sin(3x_i) + \epsilon_i \text{ where } \epsilon_i \sim N(0, 3) \text{ and } \mathbf{x} = [-4, -3, -2, -1, 0, 2, 3, 5]^T \quad (17)$$

Samples from the posterior using this dataset, where noise is included, is illustrated in figure 9. Figure 10 shows the result when ignoring the Gaussian noise in the model. A length scale of 2 and  $\sigma_f^2 = 4$  was used throughout all experiments. It can be observed in both cases, that the posterior samples vary more the further away from the observed data. This is expected since  $\mathbf{f}_*$  is conditioned on the observed data and the squared exponential kernel assigns larger covariance to nearby points. The higher uncertainty reflects our lack of knowledge in those regions.

Note how samples in figure 10 pass exactly through the observed data points, which is undesirable in this case, as we have noise in the observations and that uncertainty should be included in the model to avoid overfitting. Hence, the diagonal covariance matrix  $\sigma^2 I$  should be added to the squared exponential as in (16). Such posterior is, as mentioned, shown in figure 9.

The predictive posterior mean with and without considering the noise in the data are furthermore shown in figure 11 and 12 respectively. The same behaviour as in figure 9 and 10 is illustrated here. Comparing the posterior to the prior samples, one may observe that the posterior mean converges towards the prior mean of 0 and the posterior standard deviation towards that of the prior (i.e  $\sigma = 0.2$ ), as we move away from the observed data. This is desirable, since we do not want the observed data to affect our prior belief at a far distance from the observed data.

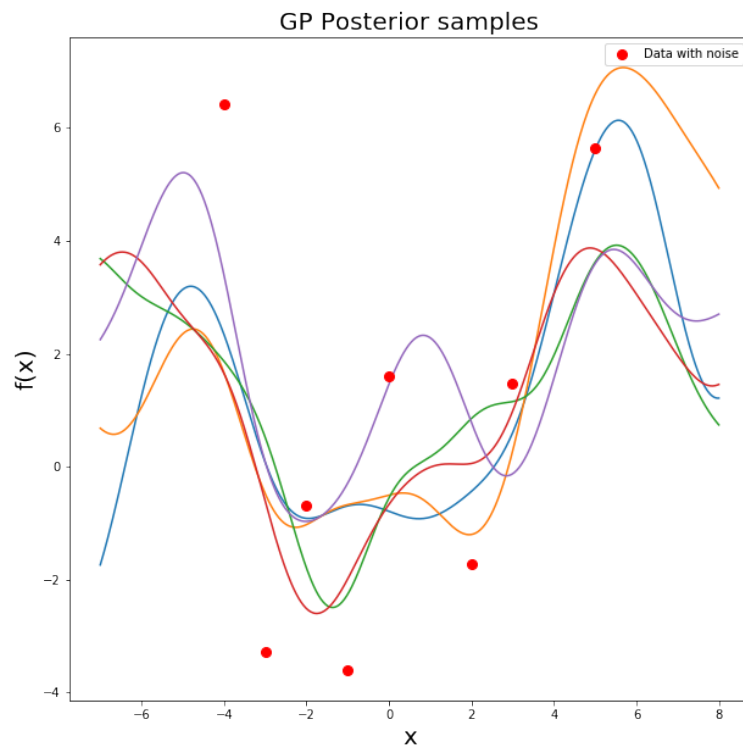


Figure 9: Samples from the Gaussian process posterior with noise term. The red dots mark the data used.

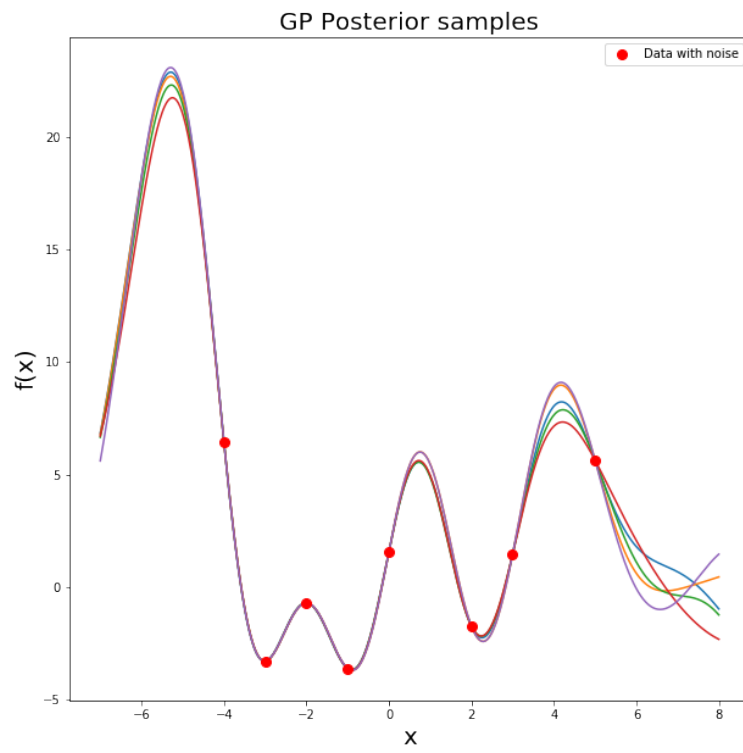


Figure 10: Samples from the Gaussian process posterior without noise term. The red dots mark the data used.

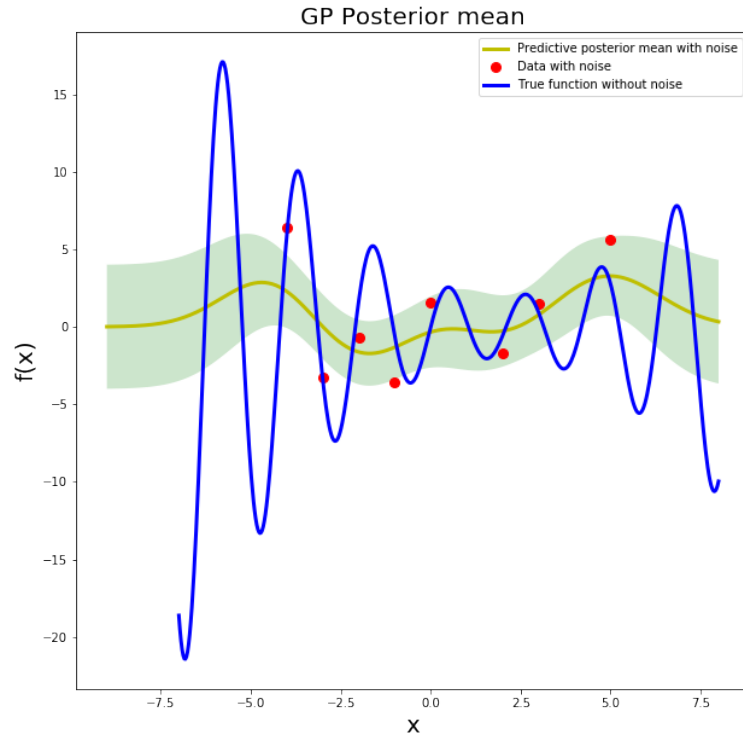


Figure 11: The Gaussian process posterior mean with noise term. The red dots mark the data used and the blue curve is the true function without noise term. The green shaded area illustrates mean plus minus two standard deviations.

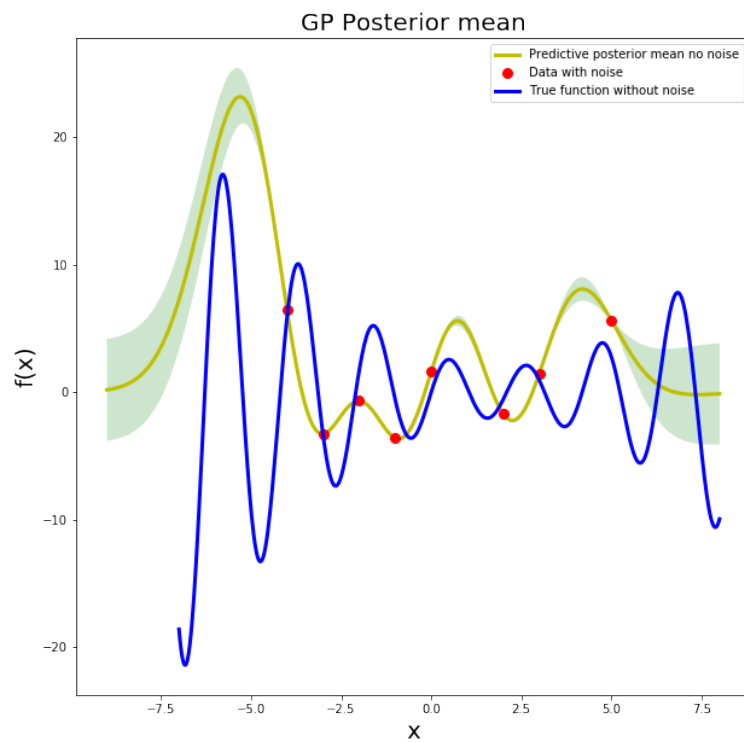


Figure 12: The Gaussian process posterior mean without noise term. The red dots mark the data used and the blue curve is the true function without noise term. The green shaded area illustrates mean plus minus two standard deviations.