

# 1 Bayesian Linear Regression

We use the model

$$p(\mathbf{y}, \mathbf{w} | \mathbf{x}, \alpha, \beta) = \left( \prod_{i=1}^N \mathcal{N}(y_i | \mathbf{w}^T \phi_i, \beta) \right) \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha \mathbf{I}) \quad (1)$$

and again the data set as defined in the instructions.

The marginal likelihood is  $\log \int p(\mathbf{y}, \mathbf{w} | \mathbf{x}, \alpha, \beta) d\mathbf{w}$

- a) Write a python function that `lml(alpha, beta, Phi, Y)` that returns the log marginal likelihood, and also a function `grad_lml(alpha, beta, Phi, Y)` that returns the gradient of the log marginal likelihood with respect to the vector  $(\alpha, \beta)$ .<sup>2</sup> The function should return a numpy vector with the gradient with respect to  $\alpha$  in the first component and gradient with respect to  $\beta$  in the second.

4 marks for the marginal likelihood correct and 8 for the gradients (4 per component)

- b) For the given data set and the linear basis functions (i.e. polynomial of order 1), maximize the log marginal likelihood with respect to  $\alpha$  and  $\beta$  using gradient descent. Show your steps on a contour plot as you did in previous questions. It is up to you where you start, but be careful that the log marginal likelihood varies over several orders of magnitude so you may have to start fairly close. You may have to clip your contours to show anything interesting on the plot. Don't use a log scale for  $\alpha$  and  $\beta$  (though this would be sensible). Report your results for the maximum.

- Correct value for  $\alpha$  and  $\beta$
- Contour with sensible scales showing the maximum clearly
- Gradient descent steps shown, with sensible starting position and step size indicated

- c) In the case of trigonometric basis function, compute the maximum of the log marginal likelihood for orders 0 to 11 inclusive using gradient descent (make sure you choose good starting values and a small step size with plenty of iterations). Plot these values on a graph against the order of the basis functions. Compare your answer to your cross validation graph from the first coursework (1c) and describe briefly the merits of the two approaches

---

<sup>2</sup>It is more straightforward if you do this in  $N \times N$  form. That is, write the likelihood as  $\mathcal{N}(\mathbf{y} | \Phi \mathbf{w}, \beta \mathbf{I})$  and use the standard results for Gaussians. Leaving the matrices in this form is very inefficient for large  $N$ , but if you like you can use the Woodbury identity to rewrite it in a way that only requires the determinant and inverse of an  $M \times M$  matrix, where  $M$  is the dimension of the basis functions. Alternatively you can complete the square and get the result out directly, following e.g. Bishop PRML p167

3 marks for the correct graph

- A point in favour of the Bayesian approach
- A point in favour of the cross validation approach

d) For  $\alpha = 1.$  and  $\beta = 0.1$  take 5 samples from the posterior distribution over the weights in the case of 10 Gaussian basis functions equally spaced between -0.5 and 1 (inclusive) with scale 0.1. Use these samples to plot the noise-free predicted function values at the test points (i.e. with y-values  $\Phi^* \mathbf{w}$ , where  $\Phi^*$  is the matrix of stacked basis functions evaluated at the test inputs  $x^*$ ). Plot also the predictive mean and 2 standard deviation error bars as a shaded region. Don't include the noise in your shaded region, but do add also dotted curves indicated two standard deviations including the noise (i.e. dotted for  $\mathbf{y}^*$  and shaded for  $\Phi^* \mathbf{w}$ ). Use test points in the interval -1 to 1.5 to show the behaviour away from the data and away from the basis function centers. Plot the samples in a different colour and use a low alpha (in the sense of opacity!) for the shaded region. Plot also the data.

- Predictive mean
- Shaded region for noise-free prediction
- Error bars including noise
- 5 samples, clearly shown in different colour
- All correct

e) (Extension: not to be graded). Use a large number of basis functions in a wider interval and experiment with different values of  $\alpha$  and  $\beta$ . Use gradient descent to find the best  $\alpha$  and  $\beta$  (you will probably have to use a log scale to get this to work effectively), or alternatively use a more sophisticated algorithm like conjugate gradients (you will certainly need to use a log scale for this work). Plot your results.