

# OUTLINE

1 MAXIMUM LIKELIHOOD REGRESSION

2 BAYESIAN REGRESSION

## BAYESIAN REGRESSION: POSTERIOR DISTRIBUTION

- Let's assume the regression weights have a Gaussian prior  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$  and that the bias is zero
- The log posterior is

$$\log p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \alpha, \beta) = -\frac{\beta}{2} \sum_{j=1}^N [t_j - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_j)]^2 - \alpha \mathbf{w}^T \mathbf{w} + \text{const}$$

- Hence the posterior is Gaussian

$$p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

with mean and variance

$$\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

## BAYESIAN REGRESSION: PREDICTIVE DISTRIBUTION

- Given the posterior, one can find the MAP estimate. However, in a fully Bayesian treatment, one makes predictions by integrating out the parameters via their posterior distribution.
- The predictive distribution

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{t}, \mathbf{w}, \alpha, \beta)p(\mathbf{w}|\mathbf{t}, \alpha, \beta)d\mathbf{w}$$

is Gaussian

$$p(t|\mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

with mean  $\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x})$  and variance

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x})$$

## MARGINAL LIKELIHOOD

- We can find  $\alpha$  and  $\beta$  by maximising the marginal likelihood:  
 $p(\mathbf{t}|\alpha, \beta)$
- The log-marginal likelihood is:

$$\rightsquigarrow \left[ \log p(\mathbf{t}|\alpha, \beta) = \frac{M}{2} \log \alpha + \frac{N}{2} \log \beta - E(\mathbf{m}_N) - \frac{1}{2} \log |\mathbf{S}_N^{-1}| - \frac{N}{2} \log 2\pi \right]$$

with

$$\rightsquigarrow E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N$$

## OPTIMISING THE MARGINAL LIKELIHOOD

- We will present a fix-point algorithm: we will write the gradient equations equal to zero as fix-point equations and iterate until convergence.
- In taking the derivative w.r.t  $\alpha$  or  $\beta$ , the most challenging term is the log of the determinant of  $\mathbf{S}_N^{-1} = (\alpha \mathbf{I} + \beta \Phi^T \Phi)$
- To deal with it, let  $\lambda_j$  be the eigenvalues of  $\beta \Phi^T \Phi$ , so that  $|\mathbf{S}_N^{-1}| = \prod_{i=0}^{M-1} (\alpha + \lambda_i)$
- We then have that

$$\partial \log |\mathbf{S}_N^{-1}| / \partial \alpha = \sum_i \frac{1}{\alpha + \lambda_i}$$

$$\Phi^T \Phi$$

eigenvalues  
 $\lambda_i = \beta \cdot \mu_i$

- Moreover,  $\lambda_i$  are proportional to  $\beta$ , so that  $\partial \lambda_i / \partial \beta = \lambda_i / \beta$

## OPTIMISING THE MARGINAL LIKELIHOOD

- Now, define

$$\gamma = \sum_i \frac{\lambda_i}{\alpha + \lambda_i}$$

(which measures the number of well determined parameters)

- By deriving the log-marginal w.r.t.  $\alpha$  and setting derivative to zero, we obtain:

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N} = g_\alpha(\alpha, \beta)$$

- By deriving the log-marginal w.r.t.  $\beta$  and setting derivative to zero, we obtain:

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N [t_n - \mathbf{m}_n^T \phi(\mathbf{x}_n)]^2 = g_\beta(\alpha, \beta)$$

- We start from an initial value  $\alpha_0$  and  $\beta_0$  and iterate  $\alpha_{n+1} = g_\alpha(\alpha_n, \beta_n)$ ,  $\beta_{n+1} = g_\beta(\alpha_n, \beta_n)$  until convergence.

## TASK 3

- Implement Bayesian regression, with type II likelihood optimisation of  $\alpha$  and  $\beta$ .
- For the 1d non-linear dataset, use polynomial model of degree 12.
- Plot predictions and 95% confidence intervals, from the predictive distribution.
- For the 2d non-linear dataset, use the Gaussian functions models. How can we set the lengthscale  $\gamma$ ?