

# Statistical Methods: Assessed Coursework 1

By Henry Bourne

## Question 0

### 0.1

Our prior knowledge says that the temperature changes periodically with time, we would like to add this behavior to our model. The trigonometric transform is good for approximating a generating function which has a temporal input, the resulting transform also consists of sinusoidal functions which would allow us to model the periodic behavior of temperature with time. Hence for  $\phi^{(1)}$  we choose the trigonometric transform.

From our prior knowledge we know that latitude changes linearly with temperature, to be able to model this behavior and incorporate it into our predictive function we should therefore choose a linear transform for  $\phi^{(3)}$ .

Finally, we have no prior knowledge on how CO2 emissions should affect temperature, this means we would like to select a basis function which isn't restrictive as it could lead to a highly inaccurate model being fit. We would therefore like to choose a very flexible basis function, so a good choice of basis function for  $\phi^{(4)}$  is the RBF.

### 0.2

B

## Question 1

### 1.1

We have,

$$p(\mathbf{f}|\mathbf{K}, \sigma) = N_{\mathbf{f}}(\mathbf{0}, \mathbf{K}), \quad (1)$$

$$p(\mathbf{y}|\mathbf{f}, \mathbf{K}, \sigma) = N_{\mathbf{y}}(\mathbf{f}, \sigma^2 \mathbf{I}) \quad (2)$$

By 2.115 in PRML we have,

$$p(\mathbf{y}|\sigma, \mathbf{K}) = N_{\mathbf{y}}(\mathbf{0}, \sigma^2 \mathbf{I} + \mathbf{K}) \quad (3)$$

### 1.2

Let us partition  $\mathbf{y}$  as follows,

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \mathbf{y}' \end{bmatrix} \quad (4)$$

Similarly we will partition,

$$\mathbf{0} = \begin{bmatrix} 0 \\ \mathbf{0} \end{bmatrix} \quad \text{and} \quad \sigma^2 \mathbf{I} + \mathbf{K} = \begin{bmatrix} \sigma^2 \mathbf{K}_{11} & \mathbf{K}_{12} \\ \mathbf{K}_{21} & \sigma^2 \mathbf{I} + \mathbf{K}_{22} \end{bmatrix} \quad (5)$$

Where we have partitioned  $\mathbf{K}$  such that  $\mathbf{K}_{11} = \mathbf{K}^{(1,1)}$ ,  $\mathbf{K}_{21} = \mathbf{K}^{(2:,1)}$ ,  $\mathbf{K}_{12} = \mathbf{K}^{(1,2)}$  and  $\mathbf{K}_{2,2} = \mathbf{K}^{(2:,2)}$ . By 2.81/2.82 in PRML we have that  $p(y_1|\mathbf{y}', \mathbf{K}, \sigma)$  has the following probability density function:

$$p(y_1|\mathbf{y}', \mathbf{K}, \sigma) = N_{y_1}(0 + \mathbf{K}_{12}(\sigma^2 \mathbf{I} + \mathbf{K}_{22})^{-1}(\mathbf{y}' - \mathbf{0}), \sigma^2 \mathbf{K}_{11} - \mathbf{K}_{12}(\sigma^2 \mathbf{I} + \mathbf{K}_{22})^{-1} \mathbf{K}_{22}) \quad (6)$$

$$= N_{y_1}(\mathbf{K}_{12}(\sigma^2 \mathbf{I} + \mathbf{K}_{22})^{-1} \mathbf{y}', \sigma^2 \mathbf{K}_{11} - \mathbf{K}_{12}(\sigma^2 \mathbf{I} + \mathbf{K}_{22})^{-1} \mathbf{K}_{22}) \quad (7)$$

### 1.3

Our prediction for  $\mathbf{x}_1$  is,

$$f(\mathbf{x}_1; \mathbf{w}_{LS}) := \mathbf{K}_{21}(\mathbf{K}_{22} + \lambda \mathbf{I})\mathbf{y}'^T \quad (8)$$

$$= \mathbf{K}_{12}(\mathbf{K}_{22} + \lambda \mathbf{I})\mathbf{y}' \quad (9)$$

As  $\mathbf{K}$  is a covariance matrix. Notice that this is exactly the mean of the distribution for  $y_1|\mathbf{y}', \sigma, \mathbf{K}$ , if we set the regularization parameter  $\lambda = \sigma^2$ , hence performing the kernel regression with  $\lambda$  equal to the variance we have that our prediction for a given  $\mathbf{x}_i$  will be equal to the expected value of the target variable  $y_i$  given the rest of our target variables,  $\mathbf{K}$  and  $\sigma$ .

### 1.4

We don't assume the data points are independent. If we don't need to assume that the datapoints are independent then it opens us up to be able to work on other types of data such as time series data.

## Question 2

### 2.1

Let  $\mathbf{h} := \phi(\xi)^T(\phi(\mathbf{X})\phi(\mathbf{X})^T)^{-1}\phi(\mathbf{X})$ , then we have,

$$\text{Var}(f(\mathbf{x}_i; \mathbf{w}_{LS})|\mathbf{x}_i) = \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS})^2|\mathbf{x}_i) - \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS})|\mathbf{x}_i)^2 \quad (10)$$

$$= \mathbb{E}(\langle \mathbf{h}, \mathbf{y} \rangle^2|\mathbf{x}_i) - \mathbb{E}(\langle \mathbf{h}, \mathbf{y} \rangle|\mathbf{x}_i)^2 \quad (11)$$

$$= \mathbb{E}(\langle \mathbf{h}, g(\mathbf{x}) + \epsilon \rangle^2|\mathbf{x}_i) - \mathbb{E}(\langle \mathbf{h}, g(\mathbf{x}) + \epsilon \rangle|\mathbf{x}_i)^2 \quad (12)$$

$$= \langle \mathbf{h}, g(\mathbf{x}) \rangle^2 + \mathbb{E}(\langle \mathbf{h}, \epsilon \rangle^2) - \mathbb{E}(\langle \mathbf{h}, g(\mathbf{x}) \rangle)^2 \quad (13)$$

$$= \mathbb{E}(\langle \mathbf{h}, \epsilon \rangle^2) \quad (14)$$

$$= \mathbb{E}(\mathbf{h}^T \epsilon \epsilon^T \mathbf{h}) \quad (15)$$

$$= \mathbf{h} \mathbb{E}(\epsilon \epsilon^T) \mathbf{h} \quad (16)$$

$$= \mathbf{h}^T \mathbf{h} \sigma^2 \quad (17)$$

$$= \langle \mathbf{h}, \mathbf{h} \rangle \cdot \sigma^2 \quad (18)$$

Therefore,

$$\frac{1}{n} \sum_{i=1}^n \text{Var}(f(\mathbf{x}_i; \mathbf{w}_{LS})|\mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n \langle \mathbf{h}, \mathbf{h} \rangle \cdot \sigma^2 \quad (19)$$

Note that,

$$\langle \mathbf{h}, \mathbf{h} \rangle = \text{tr}(\mathbf{h} \mathbf{h}^T) \quad (20)$$

$$= \text{tr}(\phi(\mathbf{x}_i)^T (\phi(\mathbf{X})\phi(\mathbf{X})^T)^{-1} \phi(\mathbf{x}_i)) \quad (21)$$

$$= \text{tr}(\phi(\mathbf{x}_i)^T (\phi(\mathbf{X})\phi(\mathbf{X})^T)^{-1} \phi(\mathbf{x}_i)) \quad (22)$$

Hence,

$$\frac{1}{n} \sum_{i=1}^n \langle \mathbf{h}, \mathbf{h} \rangle \cdot \sigma^2 = \text{tr}(\sum_{i=1}^n \phi(\mathbf{x}_i)^T (\phi(\mathbf{X})\phi(\mathbf{X})^T)^{-1} \phi(\mathbf{x}_i)) \cdot \frac{\sigma^2}{n} \quad (23)$$

$$= \text{tr}((\phi(\mathbf{X})\phi(\mathbf{X})^T)(\phi(\mathbf{X})\phi(\mathbf{X})^T)^{-1}) \cdot \frac{\sigma^2}{n} \quad (24)$$

$$= \frac{b\sigma^2}{n} \quad (25)$$

Hence, as  $b$  increases  $\frac{1}{n} \sum_{i=1}^n \text{Var}(f(\mathbf{x}_i; \mathbf{w}_{LS})|\mathbf{x}_i)$  grows.

## 2.2

The in-sample error is,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}((y_i - f_{LS}(\mathbf{x}_i))^2 | \mathbf{x}_i) \quad (26)$$

Using bias-variance decomposition we know this equals,

$$\frac{1}{n} \sum_{i=1}^n \text{Var}(\epsilon_i) + [g(\mathbf{x}_i) - \mathbb{E}(f_{LS}(\mathbf{x}_i) | \mathbf{x}_i)]^2 + \text{Var}(f_{LS}(\mathbf{x}_i) | \mathbf{x}_i) \quad (27)$$

Looking at the bias we have that its equal to,

$$g(\mathbf{x}_i)^2 - 2g(\mathbf{x}_i) \cdot \mathbb{E}(f_{LS}(\mathbf{x}_i) | \mathbf{x}_i) + \mathbb{E}(f_{LS}(\mathbf{x}_i) | \mathbf{x}_i)^2 \quad (28)$$

$$= g(\mathbf{x}_i)^2 - 2g(\mathbf{x}_i) \cdot \mathbb{E}(\phi(\mathbf{x}_i)^T (\phi(\mathbf{X})\phi(\mathbf{X})^T)^{-1} \phi(\mathbf{X})(g(\mathbf{x}) + \epsilon_i) | \mathbf{x}_i) \quad (29)$$

$$+ \mathbb{E}(\phi(\mathbf{x}_i)^T (\phi(\mathbf{X})\phi(\mathbf{X})^T)^{-1} \phi(\mathbf{X})(g(\mathbf{x}) + \epsilon_i) | \mathbf{x}_i)^2 \quad (30)$$

$$= g(\mathbf{x}_i)^2 - 2g(\mathbf{x}_i) \cdot \mathbb{E}(\phi(\mathbf{x}_i)^T \mathbf{w}^*) + \mathbb{E}(\phi(\mathbf{x}_i)^T \mathbf{w}^*)^2 \quad (31)$$

$$= g(\mathbf{x}_i)^2 - 2g(\mathbf{x}_i) \cdot \mathbb{E}(g(\mathbf{x}_i)) + \mathbb{E}(g(\mathbf{x}_i))^2 \quad (32)$$

$$= 0 \quad (33)$$

We also have that  $\text{Var}(\epsilon_i)$  is constant as n increases, finally we have from the previous question that,

$$\frac{1}{n} \sum_{i=1}^n \text{Var}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) = \frac{b\sigma^2}{n} \quad (34)$$

Hence as n increases, this above decreases and the in-sample error decreases

## 2.3

Because the least squares estimator is unbiased if we reduce the variance associated with our prediction we obtain an estimate closer to the true data generating function. Therefore, if we increase n (by question 2.1/2.2) we have that the average variance of the prediction function given the dataset decreases, this means that we are less likely to overfit as our prediction is likely to stray from the expected value of the prediction function which is the data generating function. Conversely if we decrease n then we are more likely to overfit the data as the variance of our prediction given our dataset increases.

If we increase b, then the variance increases and using the same argument as before we are more likely to overfit, conversely if we decrease b the variance is more likely to decrease and we are less likely to overfit.