

# Statistical Methods: Portfolio 2

By Henry Bourne

## 1 Bias-Variance Decomposition

What is the mathematical explanation of overfitting? Why is cross-validation a good measurement of the generalization of a prediction? we introduce a frequentist analysis to answer these questions.

First we will assume an outcome,  $y_i$ , is generated by  $y_i = g(\mathbf{x}_i) + \epsilon_i$ . Where  $g_{\mathbf{x}} : \mathbb{R}^d \rightarrow \mathbb{R}$  is a deterministic function,  $\forall i, \epsilon_i$  independent of  $\mathbf{x}_i$  and  $\mathbb{E}(\epsilon_i) = 0$ . For now for the sake of simplicity we also assume  $\mathbf{x}_i$  are fixed. We also partition the dataset,  $D$ , into training,  $D_0$ , and testing,  $D_1$  datasets and will use least squares as our error/cost function, denoted  $E$ .

Whats better than the testing error on a specific dataset? One answer is the expectation of the testing error over the whole dataset. ie.  $\mathbb{E}_D(E(D_1, \mathbf{w}_{LS})) = \mathbb{E}_D(\sum_i [y_i - f(\mathbf{x}_i; \mathbf{w}_{LS})]^2) = \sum_i \mathbb{E}_D([y_i - f(\mathbf{x}_i; \mathbf{w}_{LS})]^2 | \mathbf{x}_i)$ . We can decompose this as follows (proof in section A.1):

**Definition 1.1 Bias-Variance Decomposition:**

$$\mathbb{E}_D([y_i - f(\mathbf{x}_i; \mathbf{w}_{LS})]^2 | \mathbf{x}_i) = \text{Var}(\epsilon_i) + [g(\mathbf{x}_i) - \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS})) | \mathbf{x}_i]^2 + \text{Var}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) \quad (1)$$

The 1<sup>st</sup> term measures the randomness of our data generating process (this is beyond our control), the 2<sup>nd</sup> term shows the accuracy of our expected prediction and the 3<sup>rd</sup> term shows how easily our fitted prediction function is affected by the randomness of the dataset. When we are using a polynomial feature transform and increase  $b$ , the prediction function becomes more complex and hence the bias decreases, but conversely the variance increases. Note that we use the  $\mathbf{w}_{LS}$  estimator here but the bias-variance decomposition holds true for any estimator, furthermore when we do use  $\mathbf{w}_{LS}$  the bias equals zero. Also note that we can derive bias-variance Decompositions for different loss functions. A balance between trying to minimize the bias and variance gives us the **minimum expected error**.

### 1.1 Justifying K-fold cross-validation

First we define the:

**Definition 1.2 In-Sample Error:**

*This is the collective error of the training set, we can calculate it by averaging the expected error obtained over all the  $\mathbf{x}_i$ 's in the training set,*

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

In practice the in-sample error is not useful as we cannot calculate it, because we do not know  $g(\mathbf{x})$  or the distribution of the  $\epsilon_i$ , instead we use the:

**Definition 1.3 Out-sample Error:**

*The error over the whole distribution of  $\mathbf{x}$ ,*

$$\mathbb{E}_{\mathbf{x}} \mathbb{E}[(\mathbf{y} - f(\mathbf{x}; \mathbf{w}_{LS}))^2 | \mathbf{x}] \quad (3)$$

Now that  $\mathbf{x}$  is being treated as a random quantity, we have,

$$\mathbb{E}_{\mathbf{x}} \mathbb{E}[(\mathbf{y} - f(\mathbf{x}; \mathbf{w}_{LS}))^2 | \mathbf{x}] = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{D_1} \mathbb{E}_{D_0}[(y - f(\mathbf{x}; \mathbf{w}_{LS}))^2 | \mathbf{x}] \quad (4)$$

$$= \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{p(y|\mathbf{x})} \mathbb{E}_{D_0}[(y - f(\mathbf{x}; \mathbf{w}_{LS}))^2] \quad (5)$$

$$= \mathbb{E}_{D_0} \mathbb{E}_{p(y, \mathbf{x})}[(y - f(\mathbf{x}; \mathbf{w}_{LS}))^2] \quad (6)$$

Now, can we approximate the out-sample error? Suppose we have datasets  $D^{(1)}, \dots, D^{(K)}$  containing pairs  $(y, \mathbf{x})$  from  $p(y, \mathbf{x})$ . Let  $D^{(k)} := D_0^{(k)} \cup D_1^{(k)}$ , the following holds under mild conditions (conditions that aren't too hard to obtain):

$$\mathbb{E}_{D_0} \mathbb{E}_{p(y, \mathbf{x})} [(y - f(\mathbf{x}; \mathbf{w}_{LS}))^2] \approx \frac{1}{K} \sum_{k=1, \dots, K} \frac{1}{n_k} \sum_{(y, \mathbf{x}) \in D_1^{(k)}} (y - f(\mathbf{x}; \mathbf{w}_{LS})^{(k)})^2 \quad (7)$$

Where  $f(\mathbf{x}; \mathbf{w}_{LS})^{(k)}$  denotes the prediction function trained on  $D_0^{(k)}$  and where  $n_k$  is the size of  $D_1^{(k)}$ . If we suppose that  $D_0^{(k)}$  is the  $k$ -th split of an iid. dataset and  $D_1^{(k)}$  is the rest of the dataset, then the above result provides some justification for the use of K-fold cross-validation.

## 2 Feature Transforms and Kernel Methods

We will start by discussing different types of feature transforms before diving into kernel methods.

### 2.1 Feature transforms

We've already seen polynomial feature transforms, but why do they work? Let's get some intuition from the 1-dimensional case. The Taylor series of the data generating function,  $g$ , at zero is:  $g(x) = g(0)(x - 0)^0 + g'(0)(x - 0)^1 + \frac{g''(0)}{2!}(x - 0)^2 + \dots$ . The Taylor series tells us that we can approximate a smooth function using polynomial terms (at some cost), this gives some intuition for how a polynomial transform could allow us to create a good model of the data generating function. Another way we can decompose  $g$  is using the **Fourier series**:  $g(x) = a_0 + \sum_{i=1}^{\infty} [a_i \sin(ix) + b_i \cos(ix)]$ . We can use the Fourier series as intuition for another type of transform:

**Definition 2.1 Trigonometric transform:**

Used to approximate  $g(x)$  over the time domain, defined as follows,

$$\phi(\mathbf{x}) := [\sin(\mathbf{x}), \cos(\mathbf{x}), \sin(2\mathbf{x}), \cos(2\mathbf{x}), \dots, \sin(b\mathbf{x}), \cos(b\mathbf{x})] \quad (8)$$

polynomial and trigonometric transforms are based on the idea a function can be approximated by a:

**Definition 2.2 Linear basis expansion:**

A linear basis expansion of  $g(\mathbf{x})$  is,

$$g(\mathbf{x}) \approx f(\mathbf{x}; \mathbf{w}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_{i=1}^b w^{(i)} \phi^{(i)}(\mathbf{x}) \quad (9)$$

where  $\phi^{(i)}$  are called **basis functions**

A widely used basis function for regression tasks is the,

**Definition 2.3 Radial Basis Function (RBF):**

$$\phi^{(i)}(\mathbf{x}) := \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2\sigma^2}\right) \quad (10)$$

where  $\sigma > 0$  is called the **bandwidth**. Note that  $\sigma$  is determined before fitting, a common practice is setting it equal to the median of all pairwise distances of  $\mathbf{x}$  in your dataset.

$\mathbf{x}_i$  are called **RBF centroids** (often chosen randomly from dataset). We will use  $\phi(\mathbf{x})$  to denote  $[\phi^{(1)}(\mathbf{x}), \dots, \phi^{(b)}(\mathbf{x})]$

If we imagine the prediction function starting as a flat line, if  $w^{(i)} > 0$  then the RBF at  $\mathbf{x}_i$  lifts up the value of the prediction function around  $\mathbf{x}_i$ , we have the converse is true if  $w^{(i)} < 0$ . Each RBF defines a ball on which the prediction function is supported (with radius  $\sigma$ ), if  $g(\mathbf{x})$  has a wide support then the prediction function must be supported almost everywhere and need many centroids. The packing number (of the number of balls)  $b = O(c^d)$ , where  $d$  is dimension of  $\mathbf{x}$ .

## 2.2 Kernel methods

The larger the value of  $b$ , the higher the dimensionality of the feature space and the more flexible our prediction function is. Can we have an  $\infty$ -dimensional feature space?

Suppose  $\phi(\mathbf{x})$  maps  $\mathbf{x}$  to an  $\infty$ -dim. feature space, we will have that  $\mathbf{w}$  will be  $\infty$ ly long. However, we can show that  $\mathbf{w}_{LS-R} := (\phi\phi^T + \lambda\mathbf{I})^{-1}\phi\mathbf{y}^T = \phi(\phi^T\phi + \lambda\mathbf{I})^{-1}\mathbf{y}^T$ , where  $\phi$  is short for  $\phi(\mathbf{X})$  (proof in section A.2). We now no longer need to compute  $\phi\phi^T$  (intractable), we instead can compute  $\phi^T\phi \in \mathbb{R}^{n \times n}$ .

Is there a way we can compute the inner-product without computing  $\phi(\cdot)$ ? Define  $k(\mathbf{a}, \mathbf{b}) := \langle \phi(\mathbf{a}), \phi(\mathbf{b}) \rangle$ , denote  $\mathbf{K} := \phi^T\phi$ . We have that  $K^{(i,j)} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = k(\mathbf{x}_i, \mathbf{x}_j)$ . Going back to our prediction function we have that,  $f(\mathbf{x}; \mathbf{w}_{LS-R}) = \langle \mathbf{w}_{LS-R}, \phi(\mathbf{x}) \rangle = \langle \phi(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}^T, \phi(\mathbf{x}) \rangle = \langle (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}^T, \phi(\mathbf{x})^T\phi \rangle$ . Let's denote  $\mathbf{k} := \phi(\mathbf{x})^T\phi$ , where  $k^{(i)} = \langle \phi(\mathbf{x}), \phi(\mathbf{x}_i) \rangle = k(\mathbf{x}, \mathbf{x}_i)$ . We can then rewrite our prediction functions as  $f(\mathbf{x}; \mathbf{w}_{LS-R}) := \mathbf{k}(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}^T$ , note that  $\phi(\mathbf{x})$  only appears in inner products! <sup>1</sup>

All we need now is a function,  $k$ , that mimics the behavior of the inner product (without actually having to compute the inner product) and if our  $k(\mathbf{a}, \mathbf{b})$  is positive definite  $\exists \phi$  such that  $k(\mathbf{a}, \mathbf{b}) = \langle \phi(\mathbf{a}), \phi(\mathbf{b}) \rangle$ . Such a function is called a:

### Definition 2.4 Kernel function:

Function  $k(\mathbf{a}, \mathbf{b})$  that mimics the inner product, if explicit  $\phi(\mathbf{x})$  can be derived from  $k$  we say  $k$  induces feature transform  $\phi(\mathbf{x})$ . Many known choices of  $k$  exist.

## 2.3 Choosing a kernel function, $k$

First we give some examples of kernel functions,

### Definition 2.5 Linear Kernel Function:

$k(\mathbf{a}, \mathbf{b}) := \langle \mathbf{a}, \mathbf{b} \rangle$ , induces feature transform  $\phi(\mathbf{x}) = \mathbf{x}$

### Definition 2.6 Polynomial Kernel Function with Degree $b$ :

$k(\mathbf{a}, \mathbf{b}) := (\langle \mathbf{a}, \mathbf{b} \rangle + 1)^b$ , this induces feature transform  $\phi(\mathbf{x}) = [\mathbf{x}, \dots, \mathbf{x}^b]$ .

### Definition 2.7 RBF (or Gaussian) Kernel <sup>2</sup>:

$k(\mathbf{a}, \mathbf{b}) := \exp(-\frac{\|\mathbf{a}-\mathbf{b}\|^2}{2\sigma^2})$ , induces feature transform  $\phi(\mathbf{x})$  that is infinite-dimensional,  $\sigma$  is chosen before fitting (often as the median of the pairwise distances of all inputs).

How do we choose an appropriate kernel function? it depends on the learning task and dataset. The RBF kernel is a good all-round choice for  $\mathbf{x} \in \mathbb{R}^d$ .

## 3 Probabilistic Model Selection in Regression

In frequentist model selection we wanted to minimize the expected squared error, however, this couldn't be done in practice, so we used an approximation called the out-sample error (definition 1.3). This approach, however, requires us to hold out part of our sample during training meaning:

1. We lose information (the data held out for testing).
2. Cross-validation helps combat the above point, however, leads to heavy calculation.
3. Our dataset may not be iid.

What about a probabilistic approach? In probabilistic model selection we model the uncertainty of a model,  $m \in \{m_1, \dots, m_k\}$ , using a prior,  $p(m)$ , over the possible models. Then we can write the posterior of the model using Bayes rule:  $p(m|D) \propto p(D|m)p(m)$ . This expresses the preference over all models,  $\{m_1, \dots, m_k\}$ , given the data,  $D$ .

But now that we have a model of preference over all models, how do we select a model? As Bayesians always do we marginalize! Our prediction function is  $p(\hat{y}|D) = \sum_{m \in \{m_1, \dots, m_k\}} p(\hat{y}|D, m) \cdot p(m|D)$ , which is a weighted sum of the prediction functions for each model, where the weights

<sup>1</sup>  $\mathbf{K} = O(n^2)$  and  $(\mathbf{K} + \lambda\mathbf{I})^{-1}$  is usually  $O(n^3)$ , so computational cost gets very large for a large  $n$

<sup>2</sup> Note: RBF basis function and the RBF kernel are **not** the same thing.

are the model posteriors. Note that using Bayes rule we can write:  $p(m|D) \propto p(D|m)p(m)$ , where  $p(D|m)$  is called the **Model Evidence**.

Now, how do we calculate  $p(m|D)$ ? We have the model priors so we need to calculate  $p(D|m)$ . Suppose our model is parameterized by  $\mathbf{w}$ , then,  $p(D|m) = \int p(D|\mathbf{w}, m) \cdot p(\mathbf{w}|m) d\mathbf{w}$ . Note that we can write,

$$p(\mathbf{w}|D, m) = \frac{p(D|\mathbf{w}, m)p(\mathbf{w}|m)}{p(D|m)} \quad (11)$$

We notice that the model evidence ( $p(D|m)$ ) normalizes the posterior of the parameters ( $p(\mathbf{w}|D, m)$ ).

Is there a way of approximating the model evidence? Let's consider the case where we have only one parameter and suppose that  $p(w|D, m) \propto p(D|w, m) \cdot p(w|m)$  plateaus at  $w_{MAP}$ . Then,  $\int p(D|w, m) \cdot p(w|m) dw \approx p(D|w_{MAP}, m) \cdot p(w_{MAP}|m) \cdot \Delta_{posterior}$ <sup>3</sup>. If the prior of the parameters is flat as well:  $p(w|m) = \frac{1}{\Delta_{prior}}$ , then,

$$p(D|m) \approx p(D|w_{MAP}, m) \frac{\Delta_{posterior}}{\Delta_{prior}} \quad (12)$$

To further analyze the above let's consider its log,  $\log p(D|m) \approx \log p(D|w_{MAP}, m) + \log \frac{\Delta_{posterior}}{\Delta_{prior}}$ . Note that  $\Delta_{posterior}/\Delta_{prior} < 1$ <sup>4</sup> therefore the second term is always negative and it increases in magnitude as the ratio gets smaller. So if parameters are finely tuned to the data in the posterior distribution the second term (the **penalty term**) gets large. The first term corresponds to the log likelihood (because the prior is flat).

Furthermore, for a model with  $b$  parameters we have:  $\log p(D|m) \approx \log p(D|\mathbf{w}_{MAP}, m) + b \log \frac{\Delta_{posterior}}{\Delta_{prior}}$ , assuming  $\Delta_{posterior}/\Delta_{prior}$  the same for all  $w_i$  and that the  $w_i$  are independent (proof in section A.3). If we increase  $b$  the second term increases in magnitude (gets more negative), however, the first term will increase as a more complex model can better fit the data. So we have a trade-off where the model-evidence is optimal for some intermediate model-complexity.

### 3.1 Tuning hyper-parameters

In most scenarios we normally have that before trying to fit a model we must first parameterize the model such as choosing the regularization parameter or the degree ( $b$ ) of the polynomial transform, these parameters are called **hyper parameters**. Can we use probabilistic model selection to help us determine a hyper-parameter? With hyper parameters we write the predictive distribution as:

$$p(\hat{y}|D) = \int p(\hat{y}|D, \alpha) \cdot p(\alpha|D) d\alpha = \int \int p(\hat{y}|\mathbf{w}\alpha) \cdot p(\mathbf{w}|D, \alpha) \cdot p(\alpha|D) d\mathbf{w} d\alpha \quad (13)$$

However, the integral wrt.  $\alpha$  is intractable. If  $p(\alpha|D)$  is very sharp and centred on  $\hat{\alpha}$  then  $\int \int p(\hat{y}|\mathbf{w}\alpha) \cdot p(\mathbf{w}|D, \alpha) \cdot p(\alpha|D) d\mathbf{w} d\alpha \approx \int p(\hat{y}|\mathbf{w}, \hat{\alpha}) \cdot p(\mathbf{w}|D, \hat{\alpha}) d\mathbf{w}$ . However, to find  $\hat{\alpha}$  we need to maximize  $p(\alpha|D)$ . Note that  $p(\alpha|D) \propto p(D|\alpha)p(\alpha) = p(\alpha) \int p(D|\mathbf{w}, \alpha)p(\mathbf{w}|\alpha) d\mathbf{w}$ . If  $p(\alpha)$  is relatively flat then we just set  $\hat{\alpha} := \operatorname{argmax}_{\alpha} \int p(D|\mathbf{w}, \alpha) \cdot p(\mathbf{w}|\alpha) d\mathbf{w}$ , this is called **Marginalized Likelihood Maximization** or **Evidence Approximation**.

<sup>3</sup> As  $\int f(x)dx \approx f(x_0) \cdot \Delta x$ , if  $f$  can be approximated by a step function with length  $\Delta x$  that peaks at  $x_0$ .

<sup>4</sup> As the posterior is almost always sharper than the prior.

# Appendices

## A Proofs

### A.1 Proof of Bias-Variance Decomposition

By our data-generating assumption:

$$\mathbb{E}_D([y_i - f(\mathbf{x}_i; \mathbf{w}_{LS})]^2 | \mathbf{x}_i) = \mathbb{E}_D([g_{\mathbf{x}_i} + \epsilon_i - f(\mathbf{x}_i; \mathbf{w}_{LS})]^2 | \mathbf{x}_i) \quad (14)$$

$$= \mathbb{E}_D([g_{\mathbf{x}_i} + \epsilon_i - f(\mathbf{x}_i; \mathbf{w}_{LS})]^2) + \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) \quad (15)$$

$$- \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i)) \quad (16)$$

$$= \mathbb{E}_D(\epsilon_i^2) + \mathbb{E}_D(\epsilon_i(g_{\mathbf{x}_i} - f(\mathbf{x}_i; \mathbf{w}_{LS}))) + \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) \quad (17)$$

$$- \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i)) + \mathbb{E}_D((g_{\mathbf{x}_i} - f(\mathbf{x}_i; \mathbf{w}_{LS})) \quad (18)$$

$$+ \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) - \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i))) \quad (19)$$

$$= \mathbb{E}_D(\epsilon_i^2) + \mathbb{E}_D((g_{\mathbf{x}_i} - f(\mathbf{x}_i; \mathbf{w}_{LS}) + \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) \quad (20)$$

$$- \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i))^2) \quad (21)$$

$$= \text{Var}(\epsilon_i) + \mathbb{E}_D((g_{\mathbf{x}_i} - \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i))^2) \quad (22)$$

$$+ \mathbb{E}_D((\mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) - f(\mathbf{x}_i; \mathbf{w}_{LS}))^2) \quad (23)$$

$$+ 2 \cdot \mathbb{E}_D((g_{\mathbf{x}_i} - \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i)) \cdot (\mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) \quad (24)$$

$$- f(\mathbf{x}_i; \mathbf{w}_{LS}))) \quad (25)$$

note that,

$$\mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) - f(\mathbf{x}_i; \mathbf{w}_{LS}) = f(\mathbf{x}_i; \mathbf{w}_{LS}) - f(\mathbf{x}_i; \mathbf{w}_{LS}) = 0 \quad (26)$$

so,

$$\mathbb{E}_D([y_i - f(\mathbf{x}_i; \mathbf{w}_{LS})]^2 | \mathbf{x}_i) = \text{Var}(\epsilon_i) + (g_{\mathbf{x}_i} - \mathbb{E}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i))^2 + \text{Var}(f(\mathbf{x}_i; \mathbf{w}_{LS}) | \mathbf{x}_i) \quad (27)$$

hence, proven.

### A.2 Proof of $\mathbf{w}_{LS-R} = \phi(\phi^T \phi + \lambda \mathbf{I})^{-1} \mathbf{y}^T$

Let  $\phi(\mathbf{X})$  be denoted by  $\phi$ , recall that,

$$\mathbf{w}_{LS-R} := (\phi \phi^T + \lambda \mathbf{I})^{-1} \phi \mathbf{y}^T \quad (28)$$

the Woodbury identity says,

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T = \mathbf{P} \mathbf{B}^T (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{I})^{-1} \quad (29)$$

let  $\mathbf{P} = \frac{1}{\lambda} \mathbf{I}$  and  $\mathbf{B} = \phi^T$ , then,

$$(\phi \phi^T + \lambda \mathbf{I})^{-1} \phi \mathbf{y}^T = \frac{1}{\lambda} \mathbf{I} \phi (\phi^T \frac{1}{\lambda} \mathbf{I} \phi + \mathbf{I})^{-1} \quad (30)$$

$$= \phi (\phi^T \phi + \lambda \mathbf{I})^{-1} \quad (31)$$

hence, proven.

### A.3 Proof that $\log p(D|m) \approx \log p(D|\mathbf{w}_{MAP}, m) + b \log \frac{\Delta_{\text{posterior}}}{\Delta_{\text{prior}}}$ , assuming that $\frac{\Delta_{\text{posterior}}}{\Delta_{\text{prior}}}$ same $\forall w_i$ and $w_i$ independent

We have that,

$$p(D|m) = \int p(D|\mathbf{w}, m) \cdot p(\mathbf{w}|m) d\mathbf{w} \quad (32)$$

Note that for all i,

$$p(w^{(i)}|m) \approx \frac{\Delta_{\text{posterior}}}{\Delta_{\text{prior}}} \quad (33)$$

so we have,

$$p(\mathbf{w}|m) = \prod_i p(w_i|m) \approx \left(\frac{\Delta_{posterior}}{\Delta_{prior}}\right)^b \quad (34)$$

so,

$$p(D|m) \approx p(D|\mathbf{w}_{MAP}, m) \cdot \left(\frac{\Delta_{posterior}}{\Delta_{prior}}\right)^b \quad (35)$$

Taking the log,

$$\log p(D|m) \approx \log p(D|\mathbf{w}_{MAP}, m) + b \log \frac{\Delta_{posterior}}{\Delta_{prior}} \quad (36)$$

hence, proven.

## B Homeworks

### B.1 For Section 1

**Question B.1.1** *Prove the bias-variance decomposition (definition 1.1):*

*Proof in section A.1*

### B.2 For Section 2

**Question B.2.1** *Prove that  $w_{LS-R} = \phi(\phi^T \phi + \lambda I)^{-1} y^T$ :*

*Proof in section A.2*

### B.3 For Section 3

**Question B.3.1** *Prove that  $\log p(D|m) \approx \log p(D|w_{MAP}, m) + b \log \frac{\Delta_{posterior}}{\Delta_{prior}}$ , assuming  $\Delta_{posterior}/\Delta_{prior}$  the same for all  $w_i$  and that the  $w_i$  are independent.*

*Proof in section A.3.*

**Question B.3.2** *Example of using marginalized likelihood maximization/ evidence approximation in linear regression:*

*Suppose we have a likelihood model:*

$$p(y_1 \cdots y_n | x_1 \cdots x_n; w, b) := \prod_{i \in D} N_{y_i}(\langle w, \phi(x_i) \rangle, \sigma^2) \quad (37)$$

$$p(w; \sigma_w, b) := N_w(0, \sigma_w^2 I) \quad (38)$$

*we would like to find an expression for:*

$$p(y_1 \cdots y_n | x_1 \cdots x_n; b, \sigma, \sigma_w) \quad (39)$$

*We have that,*

$$p(y_1 \cdots y_n | x_1 \cdots x_n; b, \sigma, \sigma_w) = \int p(y_1 \cdots y_n | x_1 \cdots x_n; w, b, \sigma, \sigma_w) \cdot p(w) dw \quad (40)$$

*We can also rewrite,*

$$\prod_{i \in D} N_{y_i}(\langle w, \phi(x_i) \rangle, \sigma^2) = N_y(\phi(X)^T w, \sigma^2 I) \quad (41)$$

*By 2.115 in PRML we have,*

$$p(y_1 \cdots y_n | x_1 \cdots x_n; b, \sigma, \sigma_w) = N_y(\phi(X)^T 0 + 0, \sigma^2 I + \phi(X)^T \sigma_w^2 I \phi(X)) \quad (42)$$

$$= N_y(0, \sigma^2 I + \sigma_w^2 \phi(X)^T \phi(X)) \quad (43)$$