

This portfolio aims to cover the main concepts and results related to decision making and statistical learning. The notation to be used in the rest of this report will remain consistent and is described in Appendix A.

A statistical learning problem can typically be boiled down to a quantity y we want to predict based on a vector of explanatory variables \mathbf{x} , using a function $f(\mathbf{x})$. Often we are given a dataset $D := \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ which helps us determine a sensible prediction function.

1 Linear Regression

One way we can define $f(\mathbf{x})$ is as a linear function:

$$f(\mathbf{x}; \mathbf{w}) := \langle \mathbf{w}_1, \mathbf{x} \rangle + w_0 = [\mathbf{x}, 1] \mathbf{w},$$

where $\mathbf{w} := [\mathbf{w}_1, w_0]^T$

Then our problem is reduced to finding an estimate for \mathbf{w} .

1.1 Least Squares Estimation

A common approach to finding an appropriate \mathbf{w} is called least squares, where the general idea is to choose the \mathbf{w} that minimises the squared errors across the dataset. Mathematically, we define this as:

$$\mathbf{w}_{LS} := \arg \min_{\mathbf{w}} \sum_{i \in D} [y_i - f(\mathbf{x}_i; \mathbf{w})]^2$$

The main benefit of this is that we can obtain the closed form solution for this as

$$\mathbf{w}_{LS} = (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X}\mathbf{y}^T, \text{ where } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ 1 & \cdots & 1 \end{bmatrix}, \text{ and } \mathbf{y} = (y_1, \dots, y_n)$$

The derivation of this expression is provided in Appendix B.1

1.2 Maximum Likelihood

Another approach to estimating \mathbf{w} would be to choose the one that was most likely to produce our specific dataset, i.e. has the highest likelihood. In the linear regression context, we assume y_i are independent with distribution $\mathcal{N}_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2)$.

Then in regression the maximum likelihood function is given by

$$\mathbf{w}_{ML} := \arg \max_{\mathbf{w}} \log \prod_{i=1}^n \mathcal{N}_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2)$$

Rearranging this we actually see that $\mathbf{w}_{ML} = \mathbf{w}_{LS}$, this calculation is shown in Appendix B.2.

1.3 Transformed Linear Regression

Sometimes the response is not linear with respect to the predictors. We can combat this by transforming the data before fitting $f(\mathbf{x}; \mathbf{w})$ while still having a closed form solution using least squares.

$$\mathbf{w}_{LS} := \arg \min_{\mathbf{w}} \sum_{i \in D} [y_i - f'(\mathbf{x}_i, \mathbf{w})]^2$$

$$f'(\mathbf{x}_i, \mathbf{w}) := \langle \mathbf{w}_1, \phi(\mathbf{x}) \rangle + w_0,$$

where $\mathbf{w} := [\mathbf{w}_1, w_0]^T$ as before

We call the map $\phi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^b$ a feature transform. The closed form solution for \mathbf{w}_{LS} then becomes

$$\mathbf{w}_{LS} = \left(\phi(\mathbf{X})\phi(\mathbf{X})^T \right)^{-1} \phi(\mathbf{X})\mathbf{y}^T$$

Remark 1.1 If ϕ is symmetric we can simplify the $\phi(\mathbf{X})\phi(\mathbf{X})^T$ term to $\phi(\mathbf{X})^2$. Further, if ϕ is invertible, our expression for \mathbf{w}_{LS} becomes:

$$\mathbf{w}_{LS} = \left(\phi(\mathbf{X})^2 \right)^{-1} \phi(\mathbf{X})\mathbf{y}^T = \phi^{-1}(\mathbf{X}) \left(\phi^{-1}(\mathbf{X})\phi(\mathbf{X}) \right) \mathbf{y}^T = \phi^{-1}(\mathbf{X})\mathbf{y}^T$$

or the rest of this chapter, we will only consider polynomial transforms, i.e., functions of the form $\phi(x) := [x, x^2, x^3, \dots, x^b]^T$. To select b we refer to the following section.

1.4 Testing

The basis for most testing methods involves the following step:

1. We randomly split our data D into a training set D_0 and testing set D_1
2. We fit \mathbf{w}_{LS} using D_0
3. Evaluate the testing error $E(D_1, \mathbf{w}_{LS}) = \sum_{i \in D_1} [y_i - f(\mathbf{x}_i; \mathbf{w})]^2$

The aim of the last step is to test the performance of our fit on unseen data. A low testing error would imply that the model is generalisable. An obvious way of selecting b then is to choose the value that produces the lowest testing error, then fit f_{LS} on the full dataset for that value of b . The issue with this is that the D_1 might contain pertinent information to the fit of the data that isn't incorporated with this fitting method.

1.4.1 Cross-Validation

A natural extension of this is to use k -fold cross validation, where k go up to $n - 1$. The broad steps for this are:

1. Form a disjoint union of D, D_0, \dots, D_k
2. For $i = 0, \dots, k$ and for each b , fit $f_{LS}^{(i)}(b)$ on $D \setminus D_i$ and compute $E(D_i, f_{LS}^{(i)}(b))$
3. Choose b such that it minimises $\frac{\sum_i E^{(i)}(b)}{k+1}$

While cross-validation is easy to implement, the computational cost is high. It also fails when the data is not iid, which occurs often in real world data.

1.5 High Dimensional Data

Often we have a situation where we have multiple variables contributing to an outcome y . For this we can introduce some new notation, we write $\mathbf{h}(t) := [t^1, \dots, t^b]$ as a polynomial transform and $\phi(\mathbf{x}) := [h(x^{(1)}), \dots, h(x^{(d)})]^T$ for a vector of transformed \mathbf{x}

Remark 1.2 In this case $\phi(\mathbf{x}) \in \mathbb{R}^{db}$ is still a vector since we are just increasing the number of covariates we are considering, and still only relates to one observation.

Other terms we might consider are interaction terms, as the response might depend on the joint value of multiple inputs. In the case where we want pairwise cross-dimensional polynomials, our $\phi(\mathbf{x})$ becomes:

$$\phi(\mathbf{x}) := \left[\mathbf{h}\left(x^{(1)}\right), \dots, \mathbf{h}\left(x^{(d)}\right), x^{(u)} x_{u < v}^{(v)} \right] \in \mathbb{R}^{db + \binom{d}{2}}$$

The dimensionality of ϕ actually increases exponentially with d , and since we need $n > d$ to perform least squares methods, we also need n to increase exponentially with d , which is often not reasonable so we look to a probabilistic approach to find other ways to increase the flexibility of $f(\mathbf{x})$.

1.6 Regularisation

If we want the more flexible fit to our predictions provided by high order polynomials, we can consider using a regularisation term. This is a penalty of the magnitude of the vector $\mathbf{w}^T \mathbf{w}$ (other examples of terms to use are in Appendix B.3):

$$\mathbf{w}_{LS-R} := \arg \min_{\mathbf{w}} \sum_{i \in D} [y_i - f(\mathbf{x}_i, \mathbf{w})]^2 + \lambda \mathbf{w}^T \mathbf{w}$$

which has the solution

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

for which the proof is given in Appendix B.3. To see how this regularisation term helps reduce over-fitting, when $\lambda \rightarrow \infty$ we have $f(x; \mathbf{w}_{LS-R}) \rightarrow 0$, i.e. our prediction function reduces in complexity and gets flatter as λ increases. Note, it still remains to choose a suitable λ . If we have a large iid data set, we can choose λ using cross validation. If however we have limited data or non iid data, we need to consider probabilistic methods.

1.7 Maximum A Posteriori

Note this section makes use of the Bayesian approach to linear regression, details for which are given in Appendix B.4. One way of choosing \mathbf{w} is to select the value most likely given our dataset, i.e. the point \mathbf{w}_{MAP} at which $p(\mathbf{w}|D)$ is highest

$$\mathbf{w}_{MAP} := \arg \max_{\mathbf{w}} \prod_{i \in D} \mathcal{N}_{y_i}(f(\mathbf{x}_i, \mathbf{w}), \sigma^2) \cdot \mathcal{N}_{\mathbf{w}}(0, \sigma_{\mathbf{w}}^2 \mathbf{I})$$

As shown in B.4, this is actually equivalent to finding \mathbf{w}_{LS-R} for $\lambda = \frac{\sigma^2}{\sigma_{\mathbf{w}}^2}$.

1.8 Full Bayesian Approach

One flaw of MAP is that the uncertainty $p(\mathbf{w}|D)$ is not captured. We can remedy this by providing a predictive distribution $p(\hat{y}|\mathbf{x}, D)$ rather than a point estimate of y for any given \mathbf{x} . Using the definition of conditional distributions, we can write

$$\begin{aligned} p(\hat{y}|\mathbf{x}, D) &= \int p(\hat{y}|\mathbf{x}, \mathbf{w}) p(\mathbf{w}|D) d\mathbf{w} \\ &= \mathcal{N}_{\hat{y}} \left[f(\mathbf{x}, \mathbf{w}_{LS-R}), \sigma^2 + \phi^T(\mathbf{x}) \sigma^2 \left(\phi \phi^T + \frac{\sigma^2}{\sigma_{\mathbf{w}}^2} \phi(\mathbf{x}) \right) \right] \end{aligned}$$

With the proof for the latter equality provided in Appendix B.4

2 Binary Classification

With this type of problem, we are given an observed vector $\mathbf{x} \in \mathbb{R}^d$ and want to predict a binary output $y \in \{-1, +1\}$. We do this by choosing an f such that $f(x) = 0$ is a boundary separating the space of \mathbf{x} into two regions R_+ and R_- , relating to a prediction of $+1$ and -1 respectively.

We can choose the decision boundary by minimising the number of misclassifications. For any observed \mathbf{x} we can calculate the probability of such of errors as

$$\mathbb{P}(\mathbf{x} \text{ is misclassified} | f) = \int_{R_+} p(\mathbf{x}, y = -1) d\mathbf{x} + \int_{R_-} p(\mathbf{x}, y = +1) d\mathbf{x}$$

The f that minimises this is $f = p(\mathbf{x}, y = +1) - p(\mathbf{x}, y = -1)$, and is known as the Bayes optimal classifier. The proof for which is given in Appendix C.

2.1 Loss in Classification

The consequences of a false positive (FP) and a false negative (FN) might be different, which we can model using a loss matrix. Then we might instead want to minimise the expected loss of a decision, i.e. the optimal decision is given by

$$\arg \min_{y_0} \mathbb{E}_{p(y|\mathbf{x}, D)} [L(y, y_0) | \mathbf{x}] = \arg \min_{y_0} \sum_y p(y | \mathbf{x}, D) L(y, y_0); \text{ where } L := \begin{bmatrix} 0 & \text{Cost of FN} \\ \text{Cost of FP} & 0 \end{bmatrix}$$

The issue with this is that we do not have direct access to $p(y | \mathbf{x}, D)$, we can take two approaches with this:

1. **Discriminative:** this takes the direct approach where we model $p(y | \mathbf{x})$ with $p(y | \mathbf{x}, \mathbf{w})$
2. **Generative:** this is an indirect approach where we generate an input \mathbf{x} given output y from $p(\mathbf{x} | y)$ and use this to obtain $p(y | \mathbf{x})$

Remark 2.1 *It is easy to see we will always have $p(y = +1 | \mathbf{x}) + p(y = -1 | \mathbf{x}) = 1$ but in some cases we may find both these individual probabilities are quite low. We may then want to define a region of \mathbf{x} for which we reject making a decision. For a threshold α , we can define a corresponding region*

$$R = \{\mathbf{x} : \max\{p(y = +1 | \mathbf{x}), p(y = -1 | \mathbf{x})\} < \alpha\}$$

2.2 Link to Regression

In regression, we don't have a loss matrix, but we can use a loss function to be minimised. One such function, the square loss function, is given below (with the proof of the result given in Appendix C), :

$$\hat{y} := \arg \min_{y_0} \mathbb{E}_{p(y|\mathbf{x})} [(y - y_0)^2 | \mathbf{x}] = \mathbb{E}_{p(y|\mathbf{x})} [y]$$

A Notation

The notation used throughout this document is defined below

- x, y, z : scalars,
- $\mathbf{x}, \mathbf{y}, \mathbf{z}$: vectors,
- $\mathbf{x} \in \mathbb{R}^d$: a vector in d dimensional real space,
- $x^{(i)}$: the i^{th} entry in \mathbf{x} ,
- $\mathbf{x}_i \in X$: the i^{th} element of the set X ,
- $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$: takes vector \mathbf{x} and maps it to m dimensional space,
- $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \in \mathbb{R}^{b \times d}$: matrices with b rows and d columns,
- $\langle \mathbf{x}, \mathbf{y} \rangle$: the inner product of vectors \mathbf{x} and \mathbf{y}

B Linear Regression

B.1 Least Squares

Proof of closed form LS solution:

Supposing $\mathbf{x}_1, \dots, \mathbf{x}_n$ are column vectors of length d we can construct the matrix

$$\mathbf{X} := \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \\ 1 & \dots & 1 \end{bmatrix} \in \mathbb{R}^{(d+1) \times n}$$

We can easily see that applying f to $[\mathbf{x}_1 \dots \mathbf{x}_n]$ column-wise gives us

$$\mathbf{f}([\mathbf{x}_1 \dots \mathbf{x}_n]; \mathbf{w}) = \mathbf{X}^T \mathbf{w}$$

We can also rewrite the expression for \mathbf{w}_{LS} by the definition of a vector norm:

$$\begin{aligned} \mathbf{w}_{LS} &= \arg \min_{\mathbf{w}} \|\mathbf{y}^T - \mathbf{X}^T \mathbf{w}\| = \arg \min_{\mathbf{w}} [(\mathbf{y}^T - \mathbf{X}^T \mathbf{w})^T (\mathbf{y}^T - \mathbf{X}^T \mathbf{w})] \\ &= \arg \min_{\mathbf{w}} [\mathbf{y} \mathbf{y}^T - \mathbf{y} \mathbf{X}^T \mathbf{w} - \mathbf{w}^T \mathbf{X} \mathbf{y}^T + \mathbf{w}^T \mathbf{X} \mathbf{X}^T \mathbf{w}] \\ &= \arg \min_{\mathbf{w}} [\mathbf{y} \mathbf{y}^T - 2 \mathbf{w}^T \mathbf{X} \mathbf{y}^T + \mathbf{w}^T \mathbf{X} \mathbf{X}^T \mathbf{w}] \end{aligned}$$

Since we are trying to minimise a function of \mathbf{w} , we can simply differentiate the function w.r.t \mathbf{w} and set the result to 0. This yields:

$$-2 \mathbf{X} \mathbf{y}^T + \mathbf{X} \mathbf{X}^T \mathbf{w}_{LS} = 0 \implies \mathbf{w}_{LS} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{y}^T$$

Remark B.1 If $n < d$ then the matrix \mathbf{X} does not achieve full column rank, and thus $\mathbf{X} \mathbf{X}^T$ cannot be inverted, rendering our expression for \mathbf{w}_{LS} invalid. One way we might rectify this is by removing covariates that are strongly correlated with each other. The intuition behind this is that the information these two rows contribute to the model is similar enough that only one is enough. We want to remove enough rows such that we have full column rank.

B.2 Maximum Likelihood

Proof that least squares and maximum likelihood are equivalent:

Under the independence assumption, the ML estimate of \mathbf{w} becomes:

$$\begin{aligned}\mathbf{w}_{ML} &:= \arg \max_{\mathbf{w}} \log \prod_{i=1}^n \mathcal{N}_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2) \\ &= \arg \max_{\mathbf{w}} \left[\sum_{i=1}^n -\frac{(y_i - f(\mathbf{x}_i; \mathbf{w}))^2}{2\sigma^2} \right] - n \log \sigma \sqrt{2\pi} \\ &= \arg \min_{\mathbf{w}} \sum_{i \in D_0} [y_i - f(\mathbf{x}_i; \mathbf{w})]^2 = \mathbf{w}_{LS}\end{aligned}$$

Remark B.2 *The normality assumption on the outcome might not always be valid. For example:*

- *The response might be bounded rather than able to take any values in the real numbers (e.g. exponentially distributed data)*
- *The response might take discrete integer values rather than being continuous (e.g. Poisson distributed data)*

B.3 Regularisation

Proof of expression of regularised least squares estimate:

We replace \mathbf{X} in the standard least squares expression with $\phi(\mathbf{X})$ since we are considering the transformed \mathbf{X}

$$\sum_{i \in D} [y_i - f(\mathbf{x}_i; \mathbf{w})]^2 = \mathbf{y}\mathbf{y}^T - 2\mathbf{w}^T \phi(\mathbf{X})\mathbf{y}^T + \mathbf{w}^T \phi(\mathbf{X})\phi(\mathbf{X})^T \mathbf{w},$$

so we can rewrite the objective function as

$$\begin{aligned}\mathbf{w}_{LS-R} &= \arg \min_{\mathbf{w}} [\mathbf{y}\mathbf{y}^T - 2\mathbf{w}^T \phi(\mathbf{X})\mathbf{y}^T + \mathbf{w}^T (\phi(\mathbf{X})\phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{w}] \\ &= \arg \min_{\mathbf{w}} [\mathbf{w}^T (\phi(\mathbf{X})\phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{w} - 2\mathbf{w}^T \phi(\mathbf{X})\mathbf{y}^T]\end{aligned}$$

Differentiating with respect to \mathbf{w} and rearranging we reach the closed form solution for regularised least squares as:

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

Other Regularisation Terms

In fact, we can use any norm on \mathbf{w} . A common class of these is L^p norms, the term we used is in fact one of these; the L^2 norm. The general form for these is:

$$\|\mathbf{x}\|_p = (|x_1|^p + \dots + |x_d|^p)^{\frac{1}{p}}$$

These all penalise the size of \mathbf{w} slightly differently, in fact we have for any p $\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_{p+1}$.

B.4 Bayesian Approach

In a generalised scenario, we have a random data set D , and we want to quantify the underlying process. In regression, we observe the output y_i , which we assume is generated by $y = g(x_i) + \epsilon$ where ϵ models the noisiness of the data. The unobserved process g is what we are trying to model. We do this by trying to infer the probability distribution of g given the dataset, $p(g|D)$, which we call the posterior.

To obtain the posterior we use Bayes' Rule to rewrite this conditional probability as

$$p(g|D) = \frac{p(D|g)p(g)}{p(D)}$$

with the necessary functions defined below

- $p(D|g)$: the likelihood of observing a specified data set D given
- $p(g)$: the prior distribution - this contains our beliefs about the data
generating process before observing the data
- $p(D)$: the marginal probability of observing the dataset D

Proof for expression of maximum a posteriori estimate

Since the log is an increasing function we can apply this to get the \mathbf{w}_{MAP} . First considering each of the terms we see that

$$\begin{aligned} \log \left[\prod_{i \in D} \mathcal{N}_{y_i}(f(\mathbf{x}_i, \mathbf{w}), \sigma^2) \right] &= -\frac{1}{2\sigma^2} \sum_{i \in D} (y_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 \\ \log [\mathcal{N}_{\mathbf{w}}(0, \sigma_{\mathbf{w}}^2 \mathbf{I})] &= -\frac{1}{2} \mathbf{w}^T \left(\frac{1}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right) \mathbf{w} = -\frac{1}{2\sigma_{\mathbf{w}}^2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

Then returning this into our expression for \mathbf{w}_{MAP} we get

$$\begin{aligned} \mathbf{w}_{MAP} &= \arg \max_{\mathbf{w}} \left[-\frac{1}{2\sigma^2} \sum_{i \in D} (y_i^2 - 2y_i \mathbf{w}^T \phi(\mathbf{x}_i) + \mathbf{w}^T \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T \mathbf{w}) - \frac{1}{2\sigma_{\mathbf{w}}^2} \mathbf{w}^T \mathbf{w} \right] \\ &= \arg \max_{\mathbf{w}} \left[-\frac{1}{2} \mathbf{w}^T \left(\sum_{i \in D} \frac{1}{\sigma^2} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T + \frac{1}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right) \mathbf{w} - \frac{1}{2} \left(\sum_{i \in D} -\frac{2}{\sigma^2} y_i \phi(\mathbf{x}_i)^T \mathbf{w} \right) \right] \\ &= \arg \max_{\mathbf{w}} \left[\mathbf{w}^T \left(\phi(\mathbf{X}) \phi(\mathbf{X})^T + \frac{\sigma^2}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right) \mathbf{w} - 2\mathbf{w}^T \phi(\mathbf{X}) \mathbf{y}^T \right] \end{aligned}$$

Which if we take $\lambda = \frac{\sigma^2}{\sigma_{\mathbf{w}}^2}$, we can recognise this as the expression for \mathbf{w}_{LS-R} , hence we conclude $\mathbf{w}_{LS-R} = \mathbf{w}_{MAP}$ as desired.

Proof for distribution of a prediction

If we assume $g(\mathbf{x}) = f(\mathbf{x}; \mathbf{w})$, we can replace g with the \mathbf{w} vector we are trying to estimate. Further assuming $\epsilon \sim \mathcal{N}(0, \sigma^2)$ and D is iid, we can write:

$$p(D|\mathbf{w}) = \prod_{i \in D} p(y_i|\mathbf{x}_i, \mathbf{w}, \sigma^2) = \prod_{i \in D} \mathcal{N}_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2)$$

We also need the prior $p(\mathbf{w})$, which we take to be a normal distribution: $\mathbf{w} \sim \mathcal{N}(0, \sigma_{\mathbf{w}}^2 \mathbf{I})$. Our expression for the posterior then becomes

$$\begin{aligned} p(\mathbf{w}|D) &= \frac{\prod_{i \in D} \mathcal{N}_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2) \cdot \mathcal{N}_{\mathbf{w}}(0, \sigma_{\mathbf{w}}^2 \mathbf{I})}{p(D)} \\ &\propto \prod_{i \in D} \mathcal{N}_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2) \cdot \mathcal{N}_{\mathbf{w}}(0, \sigma_{\mathbf{w}}^2 \mathbf{I}) \end{aligned}$$

where we can disregard $p(D)$ since it does not depend on \mathbf{w} . We can ignore any constants not dependent on \mathbf{w} . We can also take the log for simplicity. Since we have already performed this simplification above for \mathbf{w}_{MAP} , we know

$$\log(p(\mathbf{w}|D)) = \mathbf{w}^T \left(\phi(\mathbf{X})\phi(\mathbf{X})^T + \frac{\sigma^2}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right) \mathbf{w} - 2\mathbf{w}^T \phi(\mathbf{X})\mathbf{y}^T$$

Then by matching powers we can see (and denoting $\phi = \phi(\mathbf{X})$)

$$p(\mathbf{w}|D) = \mathcal{N} \left(\mathbf{y}^T \phi \sigma^2 \left[\phi \phi^T + \frac{\sigma^2}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right]^{-1}, \sigma^2 \left[\phi \phi^T + \frac{\sigma^2}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right]^{-1} \right)$$

Using results regarding conditional and marginal distributions given in Section 2.3 of Bishop (2006)¹, which states that given distributions:

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{\Lambda}), \quad p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L})$$

the marginal distribution of \mathbf{y} is

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L} + \mathbf{A}\mathbf{\Lambda}\mathbf{A}^T)$$

Since we have access to $p(\mathbf{w}|D)$ as above and by assumption $p(\hat{y}|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\phi(\mathbf{x})^T \mathbf{w}, \sigma^2)$, we can see

$$\begin{aligned} \boldsymbol{\mu}_{\hat{y}} &= \phi(\mathbf{x})^T \mathbf{y}^T \phi \sigma^2 \left[\phi \phi^T + \frac{\sigma^2}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right]^{-1} = \phi(\mathbf{x})^T \mathbf{w}_{LS-R} = f(\mathbf{x}, \mathbf{w}_{LS-R}) \\ \boldsymbol{\Sigma}_{\hat{y}} &= \sigma^2 + \phi(\mathbf{x})^T \mathbf{y}^T \phi \sigma^2 \left[\phi \phi^T + \frac{\sigma^2}{\sigma_{\mathbf{w}}^2} \mathbf{I} \right]^{-1} \phi(\mathbf{x}) \end{aligned}$$

and so we have obtained the marginal distribution of \hat{y} .

C Binary Classification

Proof for Bayes Optimal Classifier

We want to minimise the following expression with respect to f

$$\mathbb{P}(\mathbf{x} \text{ is misclassified} | f) = \int_{R_+} p(\mathbf{x}, y = -1) d\mathbf{x} + \int_{R_-} p(\mathbf{x}, y = +1) d\mathbf{x}$$

¹Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.

We can re-write this as

$$\begin{aligned}
 \mathbb{P}(\mathbf{x} \text{ is misclassified} | f) &= \int_{R_+} p(\mathbf{x}, y = -1) d\mathbf{x} + \int_{R_-} p(\mathbf{x}, y = +1) d\mathbf{x} \\
 &= \int_{\mathbf{x}} \mathbb{I}[f(x) \geq 0] p(\mathbf{x}, y = -1) d\mathbf{x} + \int_{\mathbf{x}} \mathbb{I}[f(x) < 0] p(\mathbf{x}, y = +1) d\mathbf{x} \\
 &= \int_{\mathbf{x}} \mathbb{I}[f(x) \geq 0] (p(\mathbf{x}, y = -1) d\mathbf{x} - p(\mathbf{x}, y = +1)) + \int_{\mathbf{x}} p(\mathbf{x}, y = +1) d\mathbf{x} \\
 &= \int_{\mathbf{x}} \mathbb{I}[f(x) \geq 0] (p(\mathbf{x}, y = -1) - p(\mathbf{x}, y = +1)) d\mathbf{x} + p(y = +1)
 \end{aligned}$$

Since we are interested in finding the function f that minimises the expression, we can disregard the second term and focus only on the integral term. To minimise the integral term, we want to choose an f such that whenever $(p(\mathbf{x}, y = -1) - p(\mathbf{x}, y = +1)) > 0$ the indicator function is 0 and 1 when the expression is negative. This occurs when $f = p(\mathbf{x}, y = +1) - p(\mathbf{x}, y = -1)$.

Proof of regression predictions with respect to loss function:

In the following I will use \mathbb{E} in place of $\mathbb{E}_{p(y|\mathbf{x})}$ for simplicity since this is the only expectation we consider. In the case of the square loss, our objective function is

$$\mathbb{E}[(y - y_0)^2 | \mathbf{x}] = \mathbb{E}[y^2] - 2y_0\mathbb{E}[y] + y_0^2$$

Which we can differentiate with respect to y_0 to obtain the optimal \hat{y}

$$-2\mathbb{E}[y] + 2\hat{y} = 0 \implies \hat{y} = \mathbb{E}[y]$$