

Statistical Methods: Portfolio 1

By Henry Bourne

1 Introduction to Decision-Making

Computers are often used to answer complex questions, however, these methods are often a "black-box". We would like to have methods which let us conduct **rational decision-making**:

1. Predictions should be precise (no gibberish)
2. They should be data driven
3. They should take cost (of making the wrong decision) into consideration
4. They should take the random nature of data into consideration

In a **regression problem** we want to predict an outcome given some known inputs. We can use the following as an objective function:

Definition 1.1 *Least Squares (LS)*:

$$\min_f \sum_{i \in D_0} (y_i - f(x_i))^2 \quad (1)$$

where f is the function that gives our prediction for x , where x_i is the i -th input, y_i is the i -th (observed) output and $D_0 \subseteq D$ is the training dataset.

By minimizing this objective function wrt. f , we obtain a function f that minimizes the squared difference between its predictions and our observed values of the target variable.

Let \mathbf{w} be a vector parameterising f ¹, then the LS solution is $\mathbf{w}_{LS} := \operatorname{argmin}_{\mathbf{w}} \sum_{i \in D_0} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2$. We can prove that:

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

The proof can be found in section A.1.

Alternatively we can find \mathbf{w} in another data-driven way but also whilst taking the randomness of the data into account by using a probabilistic approach. We define a new objective function:

Definition 1.2 *Maximum likelihood estimation*:

$$\max_{\mathbf{w}} \log \mathbb{P}(D|\mathbf{w}) \quad (3)$$

we denote the parameters that maximize this as \mathbf{w}_{ML} , this is called the **Maximum Likelihood Estimator (MLE)**, we can write,

$$\mathbf{w}_{ML} := \operatorname{argmax}_{\mathbf{w}} \log \mathbb{P}(D|\mathbf{w}) \quad (4)$$

where D is the dataset and $\mathbb{P}(D|\mathbf{w})$ is called the **likelihood**.

Note that we can show $\mathbf{w}_{ML} = \mathbf{w}_{LS}$ ². We can also show that $\sigma_{ML}^2 = \frac{1}{n} \|\mathbf{y} - f(\mathbf{x}; \mathbf{w}_{ML})\|^2$.

1.1 LS with Feature Transform

It is possible to fit non-linear curves to our data using linear LS. All we do is augment our predictor variable, \mathbf{x} , using a function $\phi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times b}$, where:

$$\phi(\mathbf{x}) := (\mathbf{x}, \mathbf{x}^2, \dots, \mathbf{x}^b)^T \quad (5)$$

With a feature transform we have $\mathbf{w}_{LS} = (\phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^T \mathbf{y}$. A feature transform can let us fit more complex models for our data, however, it can lead to problems...

¹In the case of linear LS we have: $f(\mathbf{x}; \mathbf{w}) := \langle \mathbf{w}_1, \mathbf{x} \rangle + w_0$.

²This is true in cases where the underlying data-generating process is normal, ie. error terms Gaussian and iid.

2 Overfitting and the curse of dimensionality

In previous section we introduced the feature transform. By increasing the value of b we can fit increasingly complex models (models with terms with higher powers of x). We note that when we increase b our model can "bend" to better fit our data-points, however there reaches a point where if it can bend too much then our resultant model is going to be too specific to our training dataset and not generalize well when tested on more data.

2.1 Overfitting

Let's split our data into disjoint sets D_0 and D_1 ³. We'll use LS as our error function and denote it, E_{LS} , so $E_{LS}(D, \mathbf{w}) := \sum_{i \in D} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2$. The error function quantifies how well our model fits a dataset. We will use our partitioned datasets such that the **training error** is $E(D_0, \mathbf{w}_{LS})$ and the **testing error** is $E(D_1, \mathbf{w}_{LS})$.

The testing error tells us how well the model fits data it hasn't seen; how well it **generalizes**. Again if we consider fitting models with different values of b what we will see is that as we initially increase b the testing and training error decrease. However, at a certain point the testing error will begin to increase as the training error continues decreasing. What is happening here is called:

Definition 2.1 Overfitting:

Phenomenon where $f(\mathbf{x}; \mathbf{w}_{LS})$ fits too well on the training set, D_0 , while under-performing on unseen (testing) datasets, D_1 .

2.2 Cross-Validation

Consider the problem of how best to test the performance of our model given a dataset, D . We would like to both train our dataset on as much data as possible, but also evaluate its performance on unseen (testing) data such that the score quantifies well the model's ability to generalize. With our scenario earlier where we partitioned the dataset into D_0 (training) and D_1 (testing) datasets we can identify some shortcomings:

1. We have wasted D_1 for validation (testing), what if D_1 contains useful information for fitting a good model?
2. The selection of D_0 and D_1 is random, perhaps D_1 is better for training, or perhaps a different selection altogether would be better.

We now introduce a method which aims to solve these shortcomings:

Definition 2.2 K-fold Cross-Validation (CV)

Split D into disjoint D_0, \dots, D_k ,

for $i = 0, \dots, k$:

fit $f^{(i)}$ on all subsets but D_i

for all b compute: $E(D_i, f^{(i)})$

select b which minimizes: $\frac{\sum_i E(D_i, f^{(i)})}{k+1}$

*Note: the k picked must be $\leq n-1$ and if $k = n-1$ then we call this **leave-one-out-validation**.*

Although CV solves the problems stated earlier it does have its problems of its own:

1. Computational cost is high, as $f^{(i)}$ must be fitted and validated for all splits
2. The effectiveness of CV depends on the assumption that the data is iid. and often this assumption doesn't hold (eg. time series data)

2.3 Curse of dimensionality

Let's consider carrying out polynomial transform on higher dimensional datasets. When our input $\mathbf{x} \in \mathbb{R}^d$ then $\phi(\mathbf{x}) \in \mathbb{R}^{d \times b}$ and $\mathbf{w} \in \mathbb{R}^{b+1}$. And this is without considering pairwise cross-dimensional polynomials (eg. $x^{(1)}x^{(2)}$ or $x^{(1)}x^{(2)}x^{(3)}$). We can implement this by redesigning ϕ . Let $\phi(\mathbf{x}) := (h(x^{(1)}), \dots, h(x^{(d)}), \forall_{u < v} x^{(u)}x^{(v)})$, where $h(t) := (t^1, \dots, t^b)$, then $\phi(\mathbf{x}) \in \mathbb{R}^{d \times (b + \binom{d}{2})}$.

³We assume D contains iid. data-points

We can do the same for cross terms all the way up to d-plets and we know $\binom{d}{1} + \binom{d}{2} + \dots + \binom{d}{d} = 2^d$. The output dimension of $\phi(\mathbf{x})$ can grow exponentially with dimensionality⁴ and the number of observations much at least match the dimensionality, otherwise we cannot obtain \mathbf{w}_{LS} .

Definition 2.3 The Curse of Dimensionality:

A phenomenon where the number of observations needed to solve a problem grows exponentially with d , it 'forbids' us from solving high-dimensional problems.

3 Regularization and a Probabilistic View of Regression

3.1 Regularization

A trick for avoiding overfitting is regularization. Let's consider letting our regularization term be $\lambda \mathbf{w}^T \mathbf{w}$, where $\lambda > 0$. Note that $\mathbf{w}^T \mathbf{w}$ is the magnitude of \mathbf{w} , so by including this term in our objective function we are dissuading \mathbf{w} taking large values and hence reduces the chance of overfitting. Let \mathbf{w}_{LS-R} denote our regularized linear LS solution for the parameters, then:

$$\mathbf{w}_{LS-R} := \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i \in D} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 + \lambda \mathbf{w}^T \mathbf{w} \quad (6)$$

We can prove that if we are using regularization term $\lambda \mathbf{w}^T \mathbf{w}$ then:

$$\mathbf{w}_{LS-R} := (\phi(\mathbf{X})^T \phi(\mathbf{X}) + \lambda \mathbf{I})^{-1} \phi(\mathbf{X})^T \mathbf{y} \quad (7)$$

The proof is in section A.2.

When we increase λ the regularization term gets larger which means the coefficients of \mathbf{w}_{LS-R} must get smaller in order to decrease the value of the objective function, this in turn leads to a reduction in complexity of $f(\mathbf{x}; \mathbf{w}_{LS-R})$ and vice-versa. When we increase λ we expect to see at first a decrease in testing error (if the model was overfitted), but after a certain point the testing error to rise again. This is because increasing λ increases the generalization of the model which will reduce overfitting, however, there will reach a point where it begins to over generalize. We can also use different regularization terms, usually we use norms:

Definition 3.1 Norm:

To be a norm a positive function t must satisfy:

1. If $t(\mathbf{x}) = 0$ then $\mathbf{x} = 0$
2. $t(\mathbf{x}) + t(\mathbf{y}) \geq t(\mathbf{x} + \mathbf{y})$ (Triangle inequality)
3. $t(a \cdot \mathbf{x}) = |a| \cdot t(\mathbf{x})$

A special class of norm is the **L^p norm** which for a real $p \geq 1$ is $\|\mathbf{x}\|_p := (|x_1|^p + \dots + |x_d|^p)^{\frac{1}{p}}$.

3.2 A Probabilistic View of Regression

We are often faced with inverse problems. An **inverse problem** is where we have a dataset, D , of noisy observations, and we want to identify some latent unobserved data generating mechanism. Let's denote our latent function g , then the key to solving an inverse problem is inferring the posterior distribution $\mathbb{P}(g|D)$. If we let \mathbf{w} parameterize our latent function then all we must do is infer $\mathbb{P}(\mathbf{w}|D)$ or the:

Definition 3.2 Maximum A Posteriori (MAP)

The \mathbf{w} that maximizes $\mathbb{P}(\mathbf{w}|D)$, we denote the MAP estimator as:

$$\mathbf{w}_{MAP} := \underset{\mathbf{w}}{\operatorname{argmax}} \mathbb{P}(\mathbf{w}|D) \quad (8)$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{i \in D} N_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2) \cdot N_{\mathbf{w}}(0, \sigma_{\mathbf{w}}^2 \mathbf{I}) \quad (9)$$

⁴We could also include even more complex terms such as $(x^{(u)})^2 x^{(v)}$ for example

In section A.3 we prove that $\mathbf{w}_{MAP} = \mathbf{w}_{LS-R}$ using $\lambda = \frac{\sigma^2}{\sigma_w^2}$.

This approach of estimating \mathbf{w} is probabilistic, but we can go further, we can go fully probabilistic. What we can do is consider finding the distribution of $\mathbb{P}(\hat{y}|\mathbf{x}, D)$, of our target variable, the **predictive distribution**. We know that $\mathbb{P}(\hat{y}|\mathbf{x}, D) = \int \mathbb{P}(\hat{y}|\mathbf{x}, \mathbf{w}) \cdot \mathbb{P}(\mathbf{w}|D) d\mathbf{w}$. We can assume $\mathbb{P}(\hat{y}|\mathbf{x}, \mathbf{w})$ is distributed $N(f(\mathbf{x}; \mathbf{w}), \sigma^2)$, and we have that $\mathbb{P}(\mathbf{w}|D) \propto \prod_{i \in D} N_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2) \cdot N_{\mathbf{w}}(0, \sigma_w^2 \mathbf{I})$. In section A.4 we prove:

Suppose $f(\mathbf{x}; \mathbf{w}) = \langle \mathbf{w}, \phi(\mathbf{X}) \rangle$, prove the following,

$$\int \mathbb{P}(\hat{y}|\mathbf{x}, \mathbf{w}) \cdot \mathbb{P}(\mathbf{w}|D) d\mathbf{w} = N_{\hat{y}}(f(\mathbf{x}; \mathbf{w}_{LS-R}), \sigma^2 + \phi^T(\mathbf{x}) \sigma^2 (\phi \phi^T + \frac{\sigma^2}{\sigma_w^2} \mathbf{I}) \phi(\mathbf{x})) \quad (10)$$

where ϕ is short for $\phi(\mathbf{X})$ and \mathbf{w}_{LS-R} is the LS-R solution with $\lambda = \frac{\sigma^2}{\sigma_w^2}$

4 Risk and Bayes Optimal Prediction

Sometimes we need to make decisions, this is where classification comes into play. In **binary classification** we have an input, $\mathbf{x} \in \mathbb{R}^d$, an output, $y \in \{+1, -1\}$ ⁵, and a task where given \mathbf{x} we make a prediction y . In classification our aim is to find a **decision boundary**, which carves the input space into areas belonging to a certain prediction (or class). In binary classification we can define the kind of error we make: a **False Positive (FP)** is where a \mathbf{x} should be labelled “-1” but is labelled “+1”, a **False Negative (FN)** on the other hand is where \mathbf{x} is wrongly labelled “-1”, similarly we can define **True Positive (TP)** and **True Negative (TN)**. Let f be our **level-set function**, that is the value of $f(\mathbf{x})$ determines whether our prediction of \mathbf{x} is +1 or -1. The optimal level set function is:

Definition 4.1 Bayes Optimal Classifier

$$f(\mathbf{x}) = \mathbb{P}(\mathbf{x}, y = +1) - \mathbb{P}(\mathbf{x}, y = -1)$$

As $\mathbb{P}(\text{FP or FN}|f)$ is minimized when f defined as above (proof in section A.5). However, this only serves as an ideal optimal classifier, in reality we don’t have access to $\mathbb{P}(\mathbf{x}, y)$, we only have some data-points.

Sometimes wrong decisions may have different losses associated with them, we can define a **loss matrix**, \mathbf{L} , where entry i, j of the loss matrix is the loss associated with predicting j instead of the target i . To make an optimal decision whilst taking the risk into account we would like to minimize the expected loss of making a wrong decision: $\text{argmin}_{y_0} \mathbb{E}_{\mathbb{P}(y|\mathbf{x})}(L(y, y_0)|\mathbf{x}) = \sum_{y \in \{+1, -1\}} \mathbb{P}(y|\mathbf{x}) L(y, y_0)$, where y_0 is the decision we make (our prediction). Note, we don’t know $\mathbb{P}(y|\mathbf{x})$.

4.1 Inference of $\mathbb{P}(y|\mathbf{x})$

Instead of $\mathbb{P}(y|\mathbf{x})$ we can infer $\mathbb{P}(y|\mathbf{x}, D)$. There are two main approaches to doing this:

Definition 4.2 Discriminative Approach:

Model $\mathbb{P}(y|\mathbf{x})$ with $\mathbb{P}(y|\mathbf{x}; \mathbf{w})$. Note that with this method we cannot simulate a new \mathbf{x} given a class y , only tells us the difference between +1 and -1.

Definition 4.3 Generative Approach:

Note that $\mathbb{P}(y|\mathbf{x}, D) \propto \mathbb{P}(\mathbf{x}|y, D) \mathbb{P}(y)$. So we can model $\mathbb{P}(\mathbf{x}|y)$ with $\mathbb{P}(\mathbf{x}|y; \mathbf{w})$. With this approach we can generate a new input \mathbf{x} given an output y .

4.2 Risk in Regression

The output value in regression is a continuous variable, so we cannot have a loss matrix anymore. Instead, we have a loss function, such as squared loss. When we use squared-loss the optimal prediction is $\hat{y} = \mathbb{E}_{\mathbb{P}(y|\mathbf{x})}(y)$ (proved in section A.6). If we use the absolute loss then the optimal prediction is the median of $\mathbb{P}(y|\mathbf{x})$ (proved in section A.7).

⁵ Where $\{+1, -1\}$ are the class labels.

Appendices

A Proofs

A.1 Proof of Equation (2)

Let \mathbf{X} be,

$$\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad (11)$$

we can write our linear model in matrix form,

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon} \quad (12)$$

note that,

$$\sum_{i \in D_0} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \quad (13)$$

we can find the minimum by differentiating wrt. \mathbf{w} and finding the solution when the gradient equals zero. However, first we expand our expression,

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \quad (14)$$

$$= \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \quad (15)$$

we now can find the derivative wrt. \mathbf{w} ,

$$\frac{\partial}{\partial \mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} \quad (16)$$

now setting this to zero and solving,

$$\Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (17)$$

$$\Rightarrow \mathbf{w} = \mathbf{X}^{-1} \mathbf{X}^{-T} \mathbf{X}^T \mathbf{y} \quad (18)$$

$$\Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (19)$$

hence proven.

A.2 Proof of Equation (7)

We have that,

$$\mathbf{w}_{LS-R} := (\phi(\mathbf{X})^T \phi(\mathbf{X}) + \lambda \mathbf{I})^{-1} \phi(\mathbf{X})^T \mathbf{y} \quad (20)$$

we can rewrite our objective function as:

$$(\mathbf{y} - \phi(\mathbf{X})\mathbf{w})^T (\mathbf{y} - \phi(\mathbf{X})\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w} \quad (21)$$

$$= \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \phi(\mathbf{X})^T \mathbf{y} - \mathbf{y}^T \phi(\mathbf{X}) \mathbf{w} + \mathbf{w}^T \phi(\mathbf{X})^T \phi(\mathbf{X}) \mathbf{w} + \lambda \mathbf{w}^T \mathbf{w} \quad (22)$$

now taking the partial derivative $\frac{\partial}{\partial \mathbf{w}}$ and setting the derivative to zero we have,

$$0 = -\phi(\mathbf{X})^T \mathbf{y} - \mathbf{y}^T \phi(\mathbf{X}) + 2\phi(\mathbf{X})^T \phi(\mathbf{X}) \mathbf{w} + 2\lambda \mathbf{w} \quad (23)$$

$$\rightarrow 2\lambda \mathbf{w} + 2\phi(\mathbf{X})^T \phi(\mathbf{X}) \mathbf{w} = \phi(\mathbf{X})^T \mathbf{y} + \mathbf{y}^T \phi(\mathbf{X}) \quad (24)$$

$$\rightarrow (\phi(\mathbf{X})^T \phi(\mathbf{X}) + \lambda \mathbf{I}) \mathbf{w} = \phi(\mathbf{X})^T \mathbf{y} \quad (25)$$

$$\rightarrow \mathbf{w} = (\phi(\mathbf{X})^T \phi(\mathbf{X}) + \lambda \mathbf{I})^{-1} \phi(\mathbf{X})^T \mathbf{y} \quad (26)$$

hence proven.

A.3 Proof of $\mathbf{w}_{MAP} = \mathbf{w}_{LS-R}$ when $\lambda = \frac{\sigma^2}{\sigma_w^2}$

Writing down the objective function of regularized LS, we have,

$$\prod_{i \in D} N_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2) \cdot N_{\mathbf{w}}(0, \sigma_w^2 \mathbf{I}) \quad (27)$$

$$\propto \prod_{i \in D} N_{y_i}(f(\mathbf{x}_i; \mathbf{w}), \sigma^2) \cdot \prod_{i \in \mathbf{w}} N_{\mathbf{w}_i}(0, \sigma_w^2) \quad (28)$$

$$= \frac{1}{2\sigma^2} \sum_{i \in D} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 + \frac{1}{2\sigma_w^2} \sum_{i \in \mathbf{w}} w_i^2 + c \propto \sum_{i \in D} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 + \lambda \sum_{i \in \mathbf{w}} w_i^2 + c \quad (29)$$

Note: for the second step we took the negative log probability and for the third step we multiplied by $2\sigma^2$. We note that what we ended up with is the objective function for regularized LS.

A.4 Proof of Equation (10)

First we note,

$$\mathbb{P}(\hat{y}|\mathbf{x}, \mathbf{w}) = N_{\hat{y}}(f(\mathbf{x}; \mathbf{w}), \sigma^2) \quad (30)$$

also note that,

$$\mathbb{P}(\mathbf{w}|D) \propto \mathbb{P}(D|\mathbf{w}) \cdot \mathbb{P}(\mathbf{w}) \quad (31)$$

$$= N_{\hat{y}}(\mathbf{f}(\mathbf{x}; \mathbf{w}), \sigma^2 \mathbf{I}) \cdot N_{\mathbf{w}}(0, \sigma_w^2 \mathbf{I}) \quad (32)$$

$$\propto (\mathbf{y} - \phi \mathbf{w})^T (\sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \phi \mathbf{w}) + \mathbf{w}^T (\sigma_w^2 \mathbf{I})^{-1} \mathbf{w} \quad (33)$$

$$= \frac{1}{\sigma^2} \mathbf{y}^T \mathbf{I} \mathbf{y} - \frac{2}{\sigma^2} \mathbf{y}^T \phi \mathbf{w} + \mathbf{w}^T \phi^T \frac{1}{\sigma^2} \phi \mathbf{w} + \frac{1}{\sigma_w^2} \mathbf{w}^T \mathbf{I} \mathbf{w} \quad (34)$$

$$\propto -\frac{2}{\sigma^2} \mathbf{y}^T \phi \mathbf{w} + \frac{1}{\sigma^2} \mathbf{w}^T \phi^T \phi \mathbf{w} + \frac{1}{\sigma_{bmw}^2} \mathbf{w}^T \mathbf{w} \quad (35)$$

$$= \mathbf{w}^T \left(\frac{1}{\sigma_w^2} \mathbf{I} + \frac{1}{\sigma^2} \phi^T \phi \right) \mathbf{w} - \frac{2}{\sigma^2} \mathbf{w}^T \phi^T \mathbf{y} \quad (36)$$

$$= \frac{1}{\sigma^2} (\mathbf{w}^T (\frac{\sigma^2}{\sigma_w^2} \mathbf{I} + \phi^T \phi) \mathbf{w} - 2 \mathbf{w}^T \phi^T \mathbf{y}) \quad (37)$$

$$= \frac{1}{\sigma^2} ((\mathbf{w}^T - (\lambda \mathbf{I} + \phi^T \phi)^{-1} \phi^T \mathbf{y})^T (\frac{\sigma^2}{\sigma_w^2} \mathbf{I} + \phi^T \phi) (\mathbf{w} - (\lambda \mathbf{I} + \phi^T \phi)^{-1} \phi^T \mathbf{y}) \quad (38)$$

$$- (\phi^T \mathbf{y})^T (\frac{\sigma^2}{\sigma_w^2} \mathbf{I} + \phi^T \phi)^{-1} \phi^T \mathbf{y}) \quad (39)$$

$$= N_{\mathbf{w}}(\phi \mathbf{w}_{LS-R}, \frac{\sigma^2}{\sigma_w^2} \mathbf{I} + \phi^T \phi) \quad (40)$$

therefore, with the help of 2.115 in PRML, we have,

$$\int \mathbb{P}(\hat{y}|\mathbf{x}, \mathbf{w}) \cdot \mathbb{P}(\mathbf{w}|D) d\mathbf{w} \quad (41)$$

$$= \mathbb{P}(\hat{y}|\mathbf{x}, D) \quad (42)$$

$$= N_{\hat{y}}(f(\mathbf{x}; \mathbf{w}_{LS-R}), \sigma^2 + \phi^T(\mathbf{x}) \sigma^2 (\phi \phi^T + \frac{\sigma^2}{\sigma_w^2} \mathbf{I}) \phi(\mathbf{x})) \quad (43)$$

A.5 $\mathbb{P}(\text{FP or FN}|f)$ minimized when $f(\mathbf{x}) = \mathbb{P}(\mathbf{x}, y = +1) - \mathbb{P}(\mathbf{x}, y = -1)$

We can write our objective as,

$$\int_{\{\mathbf{x}: f(\mathbf{x}) \geq 0\}} \mathbb{P}(\mathbf{x}, y = -1) d\mathbf{x} + \int_{\{\mathbf{x}: f(\mathbf{x}) \leq 0\}} \mathbb{P}(\mathbf{x}, y = +1) d\mathbf{x} \quad (44)$$

We want to pick $f(\mathbf{x})$ that minimizes the above, the $f(\mathbf{x})$ that minimizes the above objective is the one such that,

$$f(\mathbf{x}) \geq 0 \text{ if } \mathbb{P}(\mathbf{x}, y = -1) \leq \mathbb{P}(\mathbf{x}, y = +1) \quad (45)$$

and

$$f(\mathbf{x}) \leq 0 \text{ if } \mathbb{P}(\mathbf{x}, y = +1) \leq \mathbb{P}(\mathbf{x}, y = -1) \quad (46)$$

an $f(\mathbf{x})$ that satisfies these statements is,

$$f(\mathbf{x}) = \mathbb{P}(\mathbf{x}, y = +1) - \mathbb{P}(\mathbf{x}, y = -1) \quad (47)$$

if $f(\mathbf{x}) \geq 0$, then $\mathbb{P}(\mathbf{x}, y = +1) \geq \mathbb{P}(\mathbf{x}, y = -1)$,

if $f(\mathbf{x}) \leq 0$, then $\mathbb{P}(\mathbf{x}, y = +1) \leq \mathbb{P}(\mathbf{x}, y = -1)$. Hence, proven.

A.6 The optimal prediction using squared loss is $\hat{y} = \mathbb{E}_{\mathbb{P}(y|\mathbf{x})}(y)$

We can rewrite the objective function as,

$$\int \mathbb{P}(y|\mathbf{x}) L(y, y_0) dy = \int \mathbb{P}(y|\mathbf{x}) (y - y_0)^2 dy \quad (48)$$

$$= \int \mathbb{P}(y|\mathbf{x}) (y^2 + y_0^2 - 2yy_0) dy \quad (49)$$

taking the derivative of $\frac{\partial}{\partial y_0}$ and setting to 0 we get,

$$0 = \int \mathbb{P}(y|\mathbf{x}) (2y_0 - 2y) dy \quad (50)$$

$$\Rightarrow 0 = 2y_0 \int \mathbb{P}(y|\mathbf{x}) dy - 2 \int \mathbb{P}(y|\mathbf{x}) y dy \quad (51)$$

$$\Rightarrow 0 = 2y_0 - 2\mathbb{E}_{\mathbb{P}(y|\mathbf{x})}(y) \quad (52)$$

$$\Rightarrow y_0 = \mathbb{E}_{\mathbb{P}(y|\mathbf{x})}(y) \quad (53)$$

so $\hat{y} = \mathbb{E}_{\mathbb{P}(y|\mathbf{x})}(y)$.

A.7 The optimal prediction using absolute loss is the median of $\mathbb{P}(y|\mathbf{x})$

The objective function is,

$$\int \mathbb{P}(y|\mathbf{x}) |y - y_0| dy \quad (54)$$

taking the derivative of $\frac{\partial}{\partial y_0}$ and setting to 0 we get,

$$\int -\frac{\mathbb{P}(y|\mathbf{x})(y - y_0)}{|y - y_0|} dy = 0 \quad (55)$$

note that,

$$\frac{y - y_0}{|y - y_0|} = \begin{cases} 1, & y > y_0 \\ -1, & y < y_0 \end{cases} \quad (56)$$

so we have,

$$\Rightarrow \int_{y > y_0} \mathbb{P}(y|\mathbf{x}) dy - \int_{y < y_0} \mathbb{P}(y|\mathbf{x}) dy = 0 \quad (57)$$

$$\Rightarrow \int_{y_0}^{\infty} \mathbb{P}(y|\mathbf{x}) dy = \int_{-\infty}^{y_0} \mathbb{P}(y|\mathbf{x}) dy \quad (58)$$

so $\hat{y} = m$ where m is the medium of $\mathbb{P}(y|\mathbf{x})$

B Homeworks

B.1 For Section 1

Question B.1 *Prove $w_{LS} = (X^T X)^{-1} X^T y$.*

Answer is section A.1

Question B.2 *Why is the solution of w_{LS} useless if $n < d$?*

If $n < d$ then the columns of the model matrix don't have full rank, we know that for a matrix $A \in \mathbb{R}^{n \times d}$, $\text{rank}(A) \leq \min(n, d)$. Therefore, for model matrix X with n, d we have $\text{rank}(X) \leq n < d$ and so X is not full rank. Since X is not full rank we have that $X^T X$ is not invertible.

Question B.3 *In what scenarios is the use of the Normal dist. to model $\mathbb{P}(y|x, w, \sigma)$ a bad idea?*

When the errors aren't normally distributed, which could arise when the predictor or target variables are non-normal or when outliers disrupt the model prediction. When this is the case it can cause lots of problems as we often use this assumption of normality when computing confidence intervals and to carry out hypothesis testing for example. So when our assumption of normality is incorrect it can lead to incorrect analyses.

Question B.4 *Prove $w_{LS} = (\phi(X))^{-1} y$.*

We have,

$$w_{LS} = (\phi(X)^T \phi(X))^{-1} \phi(X)^T y \quad (59)$$

$$= \phi(X)^{-1} \phi(X)^{-T} \phi(X)^T y \quad (60)$$

$$= \phi(X)^{-1} y \quad (61)$$

Question B.5 *If we increase b of $\phi(x)$ by 2-fold, by how many folds will the computation time of w_{LS} increase?*

Let's consider computing the solution of w_{LS} using the shortened form we found in question B.4. Increasing b by 2-fold will increase the number of rows in the model matrix by 2-fold, finding the inverse of a matrix is $O(n^3)$ (if using Gaussian elimination), so this would lead to a 2^3 -fold increase in computation of the inverse. We also will have a 2-fold increase in the number of computations during the matrix multiplication.

B.2 For Section 2

Question B.6 *Why do machine learning algorithms still work on high-dimensional datasets (such as images), despite the curse of dimensionality telling us that the number of observations needed for solving high dimensional problems should grow exponentially with dimensionality?*

This is down to machine learning algorithms being able to learn something about the underlying structure of the data. In learning about the underlying structure it reduces the dimensionality of the problem, as it can use features within the data. For example in convolutional neural networks it's been shown that the network uses sets of features built up hierarchically in order to help it classify images. Instead of having to directly learn what every combination of pixels contained within an image should be classified as it simply learns that certain low level-patterns of pixels are important for classification, it then learns that these low-level patterns put together in different ways create other larger features which are important and so on...

B.3 For Section 3

Question B.7 *The proof in section A.2*

Question B.8 *Is the statement "the solution of w_{LS} is useless if $n < d$ " still true for w_{LS-R} ?*

If $n < d$ then X not full rank and neither is $X^T X$, denote $A := X^T X$. A is not full rank and therefore its columns are not linearly independent, that means there exist μ_1, \dots, μ_d such that there exists $i \in \{1, \dots, d\}$ with $\mu_i \neq 0$ and $\mu_1 A^{(1)} + \dots + \mu_d A^{(d)} = 0$, where $A^{(i)}$ denotes the i -th column of A . Let μ_1, \dots, μ_d be any reals such that the above is true, then we have for any

$j \in \{1, \dots, n\}$, $\sum_{i \in \{1, \dots, d\}} \mu_i \mathbf{A}_{ji} = 0$. Where \mathbf{A}_{ji} denotes the entry found at the j -th row and i -th column of \mathbf{A} . Now consider the same linear combination for any row, j , of $\mathbf{A} + \lambda \mathbf{I}$, we have $(\sum_{i \in \{1, \dots, d\} \setminus \{j\}} \mu_i \mathbf{X}_{ji}) + \mu_j \mathbf{X}_{jj} + \lambda = \lambda$. As $\lambda > 0$ we have that no linear combination of columns equals zero, unless $\mu_1 = \dots = \mu_d = 0$. Hence, the columns of $(\mathbf{A} + \lambda \mathbf{I}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})$ are linearly independent and is therefore also invertible. So there does exist a solution to \mathbf{w}_{LS-R} even if $n < d$.

Question B.9 *The proof in section [A.3](#)*

Question B.10 *The proof in section [A.4](#)*