

Statistical Methods: Portfolio 1

By Henry Bourne

Abstract

In this document we will summarise content from the first 4 lectures from the Statistical methods course (at Compass, University of Bristol). These lectures cover content on using statistical methods for decision making.

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 the appendix (CITE).

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**.

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

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

Appendices

A Proofs

A.1 Proof of Equation (2)

Let \mathbf{X} be,

$$\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad (5)$$

we can write our linear model in matrix form,

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon \quad (6)$$

note that,

$$\sum_{i \in D_0} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \quad (7)$$

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,

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \quad (8)$$

$$= \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \quad (9)$$

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

$$\frac{\partial}{\partial \mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} \quad (10)$$

now setting this to zero and solving,

$$\Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (11)$$

$$\Rightarrow \mathbf{w} = \mathbf{X}^{-1} \mathbf{X}^{-T} \mathbf{X}^T \mathbf{y} \quad (12)$$

$$\Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (13)$$

hence proven.

References

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