

An introduction to linear regression and variations

Anna Aasen, Carl Martin Fevang, Håkon Kvernmoen

September 14, 2022

1 Introduction



Figure 1: Some figure caption

2 Theory

2.1 Statistical introduction to linear regression

A way to motivate the process of statistical learning is that we want to develop a process by which we use observed data to inform our belief in a future outcome. Such a process can be developed from the ideas of Bayes' theorem, which we can use to inform us what confidence we should have in the parameters of a given model, given data points we have observed.

Let \mathbf{y} denote a vector of a series of measured values y_i at points X , where X is a matrix where the rows x_i correspond to the input values for the measurement y_i . Further let $\boldsymbol{\theta}$ denote the parameters of our given model – then Bayes' theorem tells us that our *posteriori* confidence in the parameters $\boldsymbol{\theta}$ given the available data follows

$$P(\boldsymbol{\theta}|X, \mathbf{y}) = \frac{P(\mathbf{y}|X, \boldsymbol{\theta})P(\boldsymbol{\theta})}{P(X, \mathbf{y})}. \quad (1)$$

Here $P(X, \mathbf{y})$ plays the role of a normalisation, and can safely be ignored – either because we look for values of $\boldsymbol{\theta}$ that maximise the probability, or because we assume the probabilities $P(\mathbf{y}|X, \boldsymbol{\theta})$ and $P(\boldsymbol{\theta})$ which we then can properly normalise. Now we have a framework through which we can use assumptions about the distributions $P(\mathbf{y}|X, \boldsymbol{\theta})$ and $P(\boldsymbol{\theta})$ to give us a probability distribution for the parameters $\boldsymbol{\theta}$ that we can maximise.

2.1.1 Ordinary least squares

Now we can define the ordinary least squares method by assuming that the data y that we want to make a fit to has a linear noise term that is normally distributed with mean nought. That is $y(x) = f(x) + \epsilon$, where $f(x)$ is some analytic,

non-stochastic function of the data input x , and $\epsilon \sim N(0, \sigma^2)$. This means that if we make a model fitting $y(x)$ by a $\tilde{y}(x)$, we should expect the error at every data point (y_i, x_i) to be normally distributed with the same variance. In mathematical terms, we should expect

$$P(\mathbf{y}|X, \boldsymbol{\theta}) = \prod_{i=0}^{n-1} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \tilde{y}_i)^2}{2\sigma^2}}. \quad (2)$$

Furthermore, we assume the parameters $\boldsymbol{\theta}$ to be distributed evenly, and as such, $P(\boldsymbol{\theta})$ just contributed to an overall normalisation. This means that maximising $P(\mathbf{y}|X, \boldsymbol{\theta})$ amounts to the same as maximising $P(\boldsymbol{\theta}|X, \mathbf{y})$, and can be done analytically.

Knowing that the Gaussian distribution has a single extremum that is a maximum, we can find the values for $\boldsymbol{\theta}$ maximising the probability as

$$\begin{aligned} \frac{dP(\mathbf{y}|X, \boldsymbol{\theta})}{d\theta_k} &= 0 \\ \Rightarrow \prod_{i=0}^{n-1} \frac{1}{\sqrt{2\pi\sigma^2}} \frac{1}{\sigma^2} X_{ki}(y_i - X_{ij}\theta_j) e^{-\frac{(y_i - X_{ij}\theta_j)^2}{2\sigma^2}} &= 0, \end{aligned} \quad (3)$$

which as we can see, is maximisable point by point. This means that the optimal parameters $\hat{\boldsymbol{\theta}}_{\text{OLS}}$ that maximise $P(\mathbf{y}|X, \boldsymbol{\theta})$ can be found as

$$\hat{\boldsymbol{\theta}}_{\text{OLS}} = (X^T X)^{-1} X^T \mathbf{y} \quad (4)$$

2.2 Ridge regression and Lasso

So far, we have assumed that the parameters $\boldsymbol{\theta}$ are uniformly distributed, and as such have no *bias* towards any particular value. In practice, this means that the OLS model will contort to fit itself to all values in the data. As such, a motivation for adding a bias could be for the sake of stability in the fit. This is done in ridge regression, where one assumes that the parameters are normally distributed, such that parameter values far from the mean are thought less likely to occur. This in practice means that the model is less willing to contort to outliers in the dataset, trading it for a bias towards certain parameter values. Assuming the probability distribution for $\boldsymbol{\theta}$

$$P(\boldsymbol{\theta}) = \prod_{i=0}^{n-1} \frac{1}{\sqrt{2\pi\tau^2}} e^{-\frac{\theta_i^2}{2\tau^2}} \quad (5)$$

$$\hat{\boldsymbol{\theta}}_{\text{ridge}} = (X^T X + 2\lambda\sigma^2)^{-1} X^T \mathbf{y} \quad (6)$$

$$\mathbb{E}(y_i) = \mathbb{E}[X_{ij}\boldsymbol{\beta}_j + \epsilon_i] = X_{ij}\boldsymbol{\beta}_j + \mathbb{E}(\epsilon) = X_{ij}\boldsymbol{\beta}_j \quad (7)$$

$$\begin{aligned} \text{Var}(y_i) &= \mathbb{E}[(y_i - \mathbb{E}(y_i))^2] = \mathbb{E}[(X_{ij}\boldsymbol{\beta}_j + \epsilon_i - X_{ij}\boldsymbol{\beta}_j)^2] \\ &= \mathbb{E}(\epsilon_i^2) = \text{Var}(\epsilon_i) = \sigma^2 \end{aligned} \quad (8)$$

$$\begin{aligned} \mathbb{E}(\hat{\boldsymbol{\beta}}) &= \mathbb{E}((X^T X)^{-1} X^T \mathbf{y}) = \mathbb{E}((X^T X)^{-1} X^T [X\boldsymbol{\beta} + \boldsymbol{\epsilon}]) \\ &= \mathbb{E}(\boldsymbol{\beta}) + (X^T X)^{-1} X^T \mathbb{E}(\boldsymbol{\epsilon}) = \boldsymbol{\beta} \end{aligned} \quad (9)$$

$$\begin{aligned} \text{Var}(\hat{\boldsymbol{\beta}}) &= \mathbb{E}(\hat{\boldsymbol{\beta}}\hat{\boldsymbol{\beta}}^T) - \mathbb{E}(\hat{\boldsymbol{\beta}})\mathbb{E}(\hat{\boldsymbol{\beta}}^T) \\ &= \mathbb{E}((X^T X)^{-1} X^T \mathbf{y} \mathbf{y}^T X (X^T X)^{-1}) - \boldsymbol{\beta}\boldsymbol{\beta}^T \\ &= (X^T X)^{-1} X^T [X\boldsymbol{\beta}\boldsymbol{\beta}^T X^T + \sigma^2 \mathbb{I}] X (X^T X)^{-1} - \boldsymbol{\beta}\boldsymbol{\beta}^T \\ &= \sigma^2 (X^T X)^{-1}, \end{aligned} \quad (10)$$

where we have used that $\mathbb{E}(\mathbf{y} \mathbf{y}^T) = \mathbb{E}((X\boldsymbol{\beta} + \boldsymbol{\epsilon})(X\boldsymbol{\beta} + \boldsymbol{\epsilon})^T) = X\boldsymbol{\beta}\boldsymbol{\beta}^T X^T + \sigma^2 \mathbb{I}.$

A Appendix entry

some appendix things