

Model Selection Criteria: AIC & BIC

This document discusses some of the mathematical details of Akaike's Information Criterion (AIC) and Schwarz's Bayesian Information Criterion (BIC). We assume the reader knowledge of the matrix form for multiple linear regression. Please see Matrix Form of Linear Regression for a review.

Maximum Likelihood Estimation of β and σ

To understand the formulas for AIC and BIC, we will first briefly explain the likelihood function and maximum likelihood estimates for regression.

Let \mathbf{Y} be $n \times 1$ matrix of responses, \mathbf{X} , the $n \times (p+1)$ matrix of predictors, and β , $(p+1) \times 1$ matrix of coefficients. If the multiple linear regression model is correct then,

$$\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2) \quad (1)$$

When we do linear regression, our goal is to estimate the unknown parameters β and σ^2 from (1). In Matrix Form of Linear Regression, we showed a way to estimate these parameters using matrix algebra. Another approach for estimating β and σ^2 is using *maximum likelihood estimation*.

A **likelihood function** is used to summarise the evidence from the data in support of each possible value of a model parameter. Using (1), we will write the likelihood function for linear regression as

$$L(\mathbf{X}, \mathbf{Y}|\beta, \sigma^2) = \prod_{i=1}^n (2\pi\sigma^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mathbf{X}_i\beta)^T (Y_i - \mathbf{X}_i\beta) \right\} \quad (2)$$

where Y_i is the i^{th} response and \mathbf{X}_i is the vector of predictors for the i^{th} observation. One approach estimating β and σ^2 is to find the values of those parameters that maximize the likelihood in (2), i.e. **maximum likelihood estimation**. To make the calculations more manageable, instead of maximizing the likelihood function, we will instead maximize its logarithm, i.e. the log-likelihood function. The values of the parameters that maximize the log-likelihood function are those that maximize the likelihood function. The log-likelihood function we will maximize is

$$\begin{aligned} \log L(\mathbf{X}, \mathbf{Y}|\beta, \sigma^2) &= \sum_{i=1}^n -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mathbf{X}_i\beta)^T (Y_i - \mathbf{X}_i\beta) \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) \end{aligned} \quad (3)$$

[–insert details MLES–]

The maximum likelihood estimate of β and σ^2 are

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad \hat{\sigma}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\beta})^T (\mathbf{Y} - \mathbf{X}\hat{\beta}) = \frac{1}{n} RSS \quad (4)$$

where RSS is the residual sum of squares. Note that the maximum likelihood estimate is not exactly equal to the estimate of σ^2 we typically use $\frac{RSS}{n-p-1}$. This is because the maximum likelihood estimate of σ^2 in (4) is a *biased* estimator of σ^2 . When n is much larger than the number of predictors p , then the differences in these two estimates are trivial.

AIC

Akaike's Information Criterion (AIC) is

$$AIC = -2 \log L + 2(p + 1) \quad (5)$$

where $\log L$ is the log-likelihood. This is the general form of AIC that can be applied to a variety of models, but for now, let's focus on AIC for multiple linear regression.

$$\begin{aligned} AIC &= -2 \log L + 2(p + 1) \\ &= -2 \left[-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right] + 2(p + 1) \\ &= n \log \left(2\pi \frac{RSS}{n} \right) + \frac{1}{RSS/n} RSS \\ &= n \log(2\pi) + n \log(RSS) - n \log(n) + 2(p + 1) \end{aligned} \quad (6)$$

BIC

[—]