

Choosing models with AIC & BIC (Zwick, 7 March 7, 2013)

Models are selected from the one of the measures that OCCAM outputs for different models applied to the training set data, namely the Bayesian Information Criterion (BIC) also known as the Schwartz Criterion (Schwartz 1978). BIC is a way of linearly integrating the error of a model and its complexity (DF) which differs from the Akaike Information Criterion (AIC) (Akaike 1994) by its inclusion of a factor which depends on the sample size, N:

$$\begin{aligned} \text{AIC} &= -2 N \sum p \ln q + 2 \text{DF}. \\ \text{BIC} &= -2 N \sum p \ln q + \ln(N) \text{DF} \end{aligned}$$

These measure are unaffected by adding the constant $N \sum p \ln p$, which gives

$$\begin{aligned} \text{AIC}' &= 2 N \sum p \ln (p/q) + 2 \text{DF}. \\ \text{BIC}' &= 2 N \sum p \ln (p/q) + \ln(N) \text{DF} \end{aligned}$$

The first term of AIC and BIC is now the familiar likelihood-ratio (LR, sometimes written equivalently as L^2) Chi-square measure of a model. **Good** models have **low** values of these measures, since LR, the model **error**, is ideally small and so is DF, the model **complexity**. In OCCAM, however, AIC and BIC are given *relative to a reference model*, usually taken to be the bottom (independence) model:

$$\begin{aligned} \Delta \text{AIC} &= \text{AIC}'(\text{ref}) - \text{AIC}'(\text{model}) = \text{AIC}(\text{ref}) - \text{AIC}(\text{model}) = \Delta \text{LR} + 2 * \Delta \text{DF} \\ \Delta \text{BIC} &= \text{BIC}'(\text{ref}) - \text{BIC}'(\text{model}) = \text{BIC}(\text{ref}) - \text{BIC}(\text{model}) = \Delta \text{LR} + \ln(N) * \Delta \text{DF} \end{aligned}$$

In this case (for reference=bottom), ΔAIC and ΔBIC have *high* values for good models, since ΔLR is the information *captured* in the model, and since $\Delta \text{DF} = \text{df}(\text{ref}) - \text{df}(\text{model})$, being negative, diminishes the measure the more complex the model is. Including the $\ln(N)$ factor in ΔBIC penalizes more complex models (for large enough N). BIC is more conservative than AIC in recommending departures from the reference independence model. In my experience, models picked by ΔBIC do better on generalization (test data) than the more complex models picked by ΔAIC .

If reference = top, it is *also* the case that good models have *high* ΔAIC and ΔBIC , since $\text{LR}(\text{reference}) = 0$, so $\Delta \text{LR} = -\text{LR}(\text{model})$, and thus negative. We want $\text{LR}(\text{model})$, the model error, to be as small as possible. We also want $\Delta \text{DF} = \text{df}(\text{reference}) - \text{df}(\text{model})$ to be as big as possible, which means a simple model. Having the smallest $\text{LR}(\text{model})$ and the biggest ΔDF will maximize ΔAIC or ΔBIC .

Akaike, H. (1994). "Implications of Informational Point of View on the Development of Statistical Science." In *Proceedings of the First US/Japan Conference on the Frontiers of Statistical Modeling: An Informational Approach*, H. Bozdogan, ed., pp. 27-38, Kluwer Academic Publishers, the Netherlands.

Schwartz, G. (1978). *Ann. Stat.* 6, pp. 461-464.