# Model Selection — Akaike Information Criterion

PkGu

10/06/2021

## Viewpoint from Homo-Frequentistus

In Akaike's work (1979), he generated a reasonable "metric" from the likelihood principle, which is initially denoted by

$$I(f_1, f_0; \Phi) := \int_X \Phi(\frac{f_1}{f_0}) f_0(x) \ dx = \mathbb{E}_{f_0}[\Phi(\frac{f_1}{f_0})]$$

where $f_0$ means the "real" PDF and $f_1$ is a candidate. For $\mathcal{P}$ parametrized by $\Theta$, and assuming that $\theta_0$ stands for the "real" one, we naturally need the $I$ and $\nabla I$ take 0 at $\theta_0$ and the Hessian matrix at $\theta_0$ is positive-definite for sensitivity of $I(\theta, \theta_0; \Phi)$ near $\theta_0$, *i.e.* equations:

$$I(\theta_0, \theta_0; \Phi) = 0;$$

$$\partial_i I(\theta, \theta_0; \Phi)|_{\theta=\theta_0} = 0;$$

$$\partial_{i,j}^2 I(\theta, \theta_0; \Phi)|_{\theta=\theta_0} = \Phi^{''}(1) \cdot \int_X \left[ \left( \frac{\partial_i f_\theta}{f_\theta} \right) \left( \frac{\partial_j f_\theta}{f_\theta} \right) \right]\Bigg|_{\theta=\theta_0} dx;$$

And consider the multi-variate case: Suppose there are $N$ *i.i.d.* variables for which we do the same calculation while expecting that

$$I_N(\theta_0, \theta_0; \Phi) = 0;$$

$$\partial_i I_N(\theta, \theta_0; \Phi)|_{\theta=\theta_0} = 0;$$

$$\partial_{i,j}^2 I_N(\theta, \theta_0; \Phi)|_{\theta=\theta_0} = N \cdot \partial_{i,j}^2 I(\theta, \theta_0; \Phi)|_{\theta=\theta_0};$$

It suggests the adoption of $\Phi(r) := log \ r$ for the amount of "information-distance". And it coincides with the definition of *KL divergence* or *relative entropy*.

For an imperfect sub-model, which means that the "real" parameter $\theta \notin \Theta_{sub}$, we can still look for an MLE solution $\hat{\theta_{sub}}$ in $\Theta_{sub}$. For "real" parameter $\theta_0$, we define the loss and risk functions:

$$W(\theta_0; \hat{\theta}) := 2 \cdot KL(\theta_0; \hat{\theta});$$

$$R(\theta_0; \hat{\cdot}) := \mathbb{E}_X[W(\theta_0; \hat{\theta}(x))];$$

Here we encounter the main problem: suppose $\Theta$ has local dimension L, from merely $i.i.d.$ dataset $X = (x_1, ..., x_N)$ and the model itself to select the "best" submodel $\Theta_k$ in which the setting of coordinates is $\theta_{k;k+1} = ... = \theta_{k;L} = 0$. We are able to compute:

$$\omega_{k;L} := -2 \cdot \left( \frac{1}{N} \sum_1^N log \left( \frac{f_{\hat{\theta}_k}(x_i)}{f_{\hat{\theta}_L}(x_i)} \right) \right);$$

as an estimate of $W(\theta_L; \theta_k)$, where $\theta_k$ minimizes $W(\theta_L; \_)$ in $\Theta_k$ and especially $\theta_L = \theta_0$ is the "real" parameter. And $\eta_{k;L} := N \cdot \omega_{k;L}$, by large sample theory we've already known that

$$\eta_{k;L} \xrightarrow{N \to \infty} \chi^2_{L-k};$$

For simulation of $W(\theta_L; \hat{\theta}_k)$, assuming some reasonable smoothness of $W$, we rather compute the order-2 simulator around $\theta_L$:

$$W_2(\theta_L; \hat{\theta}_k) := (\hat{\theta}_k^{\ i} - \theta_L^{\ i})(\hat{\theta}_k^{\ j} - \theta_L^{\ j})I_{ij};$$

where $I_{ij}$ is exactly the $i, j$-th element of Fisher information matrix at $\theta_L$. We regard $I$ as an inner product, by redefining $\theta_k$ to be the I-projection of $\theta_L$ to $\Theta_k$, we have

$$W_2(\theta_L; \hat{\theta}_k) = \|\hat{\theta}_k - \theta_L\|_I^2 = \|\hat{\theta}_k - \theta_k\|_I^2 + \|\theta_k - \theta_L\|_I^2;$$

By doing order-2 Taylor expansion of log-likelihood function around respectively $\hat{\theta}_L$ and $\hat{\theta}_k$ and substitute the variable $\theta$ by $\theta_k$, we get that up to an order-3 term,

$$\begin{aligned} \eta_{k;L} &= N(\|\hat{\theta}_L - \theta_k\|_I^2 - \|\theta_k - \hat{\theta}_k\|_I^2) \\ &= N[\|\hat{\theta}_L - \theta_L\|_I^2 + \|\theta_L - \theta_k\|_I^2 - (\hat{\theta}_L - \theta_L, \theta_L - \theta_k)_I] - N\|\theta_k - \hat{\theta}_k\|_I^2 \\ &= A + B - C - D; \end{aligned}$$

And under the assumption of smoothness of $I$, by the same Taylor expansions, noticing that $\sqrt{N}(\hat{\theta}_L - \theta_L)$ have tendency $N(0, I^{-1})$ as $N \to \infty$, we have that the estimate and also the order-1 residue

$$\frac{1}{N} \sum_1^N \partial^2_{ij} log \ f_{\lambda\hat{\theta}_k + (1-\lambda)\theta_k}(x_n) - I_{ij} = O(\frac{1}{\sqrt{N}})$$

so we have the equation as $N \to \infty$

$$\sqrt{N}(\hat{\theta}_k^{\ j} - \theta_k^{\ j})I_{ij} = \sqrt{N}(\hat{\theta}_L^{\ j} - \theta_L^{\ j})I_{ij}$$

which means that $\hat{\theta}_k - \theta_k$ is almost the I-projection of $\hat{\theta}_L - \theta_L$ onto $\Theta_k$, hence independently and asymptotically $A - D \sim \chi^2_{L-k}$ and $D \sim \chi^2_k$. And $B$ has no relation with $N$ or $X$ while $C$ has expectation 0 and its variation is equal to $B$. So when $N$ significantly larger than $L$, $B$ and $C$ are both insignificant compared to $\eta_{k;L}$. So asymptotically

$$N \cdot W_2(\theta_L; \hat{\theta}_k) = \eta_{k;L} - (A - D) + D$$

After taking expectation, we find a good estimator for $W(\theta_L; \hat{\theta}_k)$:

$$\hat{W}(\theta_L; \hat{\theta}_k) = N^{-1}(\eta_{k;L} + 2k - L)$$

and we just need to find a best $\hat{\theta}_k$ whose k coordinates minimize $\eta_{k;L} + 2k$.