

Data Analysis

Christophe Ambroise

Model selection: How many clusters ?

The number of clusters K controls the model complexity.

Choosing K is an example of model selection.

The optimal Bayesian approach is to pick the model with the largest marginal likelihood,

$$K^* = \arg \max_k p(\mathcal{D}|K).$$

In practice,

- 1 Simple approximations, such as BIC, ICL can be used.
- 2 We can use the cross-validated likelihood as a performance measure
- 3 An alternative approach is to perform stochastic sampling in the space of models (MCMC)

The Laplace approximation

Gaussian approximation to a probability density defined over a set of continuous variables.

Considering the density

$$p(\mathbf{z}) = \frac{1}{Z} f(\mathbf{z})$$

The normalizing constant is

$$\begin{aligned} Z &= \int f(\mathbf{z}) d\mathbf{z} \\ &= f(\mathbf{z}_0) \int \exp -\frac{1}{2}(\mathbf{z}-\mathbf{z}_0)^T A(\mathbf{z}-\mathbf{z}_0) d\mathbf{z} \\ &\approx f(\mathbf{z}_0) \frac{(2\pi)^{p/2}}{|A|^{1/2}} \end{aligned}$$

where \mathbf{z}_0 is a mode of the distribution and A is the Hessian matrix of second derivatives of log-density $f(\mathbf{z})$ at $\mathbf{z} = \mathbf{z}_0$.

From the Bayes theorem the model evidence is

$$p(\mathcal{D}) = \int p(\boldsymbol{\theta})p(D|\boldsymbol{\theta})d\boldsymbol{\theta}$$

Using Laplace approximation in for $f(\boldsymbol{\theta}) = p(\boldsymbol{\theta})p(D|\boldsymbol{\theta})$ in $\boldsymbol{\theta} = \boldsymbol{\theta}_{MAP}$:

$$\ln p(\mathcal{D}) \approx \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) + \underbrace{\ln p(\boldsymbol{\theta}_{MAP}) + \frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln |A|}_{\text{Occamfactor}}$$

- $\ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP})$ represents the log-likelihood
- the Occam factor penalizes the model complexity

Assuming a simple Gaussian prior distribution over parameters, with full rank Hessian we can further approximate by

$$\ln p(\mathcal{D}) \approx \ln p(\mathcal{D}|\theta_{MAP}) - \frac{1}{2}p \ln n$$

which is known as the BIC (Bayesian Information Criterion) or the Schwartz criterion (1978).

BIC for choosing the number of clusters

$$K_{BIC} = \arg \max_k \ln p(\mathcal{D}|\theta_{MAP}^k) - \frac{1}{2}p_k \ln n$$

where p_k is the number of parameters of the model with k clusters and θ_{MAP}^k the MAP estimate of the model

Integrated Complete Likelihood (ICL)

$$\begin{aligned} BIC(k) &= p(\mathcal{D}|\boldsymbol{\theta}_{MAP}^k) - \frac{1}{2}p_k \ln n \\ &= \mathbb{E}_{Z|X;\boldsymbol{\theta}_{MAP}^k} [\ln p(X, Z; \boldsymbol{\theta}_{MAP}^k)] - \mathbb{E}_{Z|X;\boldsymbol{\theta}_{MAP}^k} [\ln p(Z|X; \boldsymbol{\theta}_{MAP}^k)] - \frac{1}{2}p_k \end{aligned}$$

Biernacki et al. (2000) proposed to favour clustering with high-confidence (low entropy) by removing entropy term to BIC.

ICL for choosing the number of clusters

$$K_{ICL} = \arg \max_k \ln \mathbb{E}_{Z|X;\boldsymbol{\theta}_{MAP}^k} [\ln p(X, Z; \boldsymbol{\theta}_{MAP}^k)] - \frac{1}{2}p_k \ln n$$

Using information theory Akaike (1974) derived an alternative criterion:

AIC

$$K_{AIC} = \arg \max_k \ln p(\mathcal{D}|\boldsymbol{\theta}_k) - p_k$$

Generally AIC chooses more complex models than BIC which chooses more complex models than ICL

$$K_{AIC} \geq K_{BIC} \geq K_{ICL}$$

Multivariate Gaussian Mixture models

Assumes K classes in proportion π_1, \dots, π_K with component densities

$$\mathbf{x}_i | z_i = k \sim \mathcal{N}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $\boldsymbol{\mu}_k \in \mathbb{R}^p$
- $\boldsymbol{\Sigma}_k \in \mathbb{R}^{p \times p}$

Number of parameters

$$p_k = p \cdot K + p(p-1)/2 + K - 1$$

Covariance matrix parametrization

Table 1: Parameterizations of the covariance matrix Σ_k currently available in `mclust` for hierarchical clustering (HC) and/or EM for multidimensional data. (‘•’ indicates availability).

identifier	Model	HC	EM	Distribution	Volume	Shape	Orientation
E		•	•	(univariate)	equal		
V		•	•	(univariate)	variable		
EII	λI	•	•	Spherical	equal	equal	NA
VII	$\lambda_k I$	•	•	Spherical	variable	equal	NA
EEI	λA		•	Diagonal	equal	equal	coordinate axes
VEI	$\lambda_k A$		•	Diagonal	variable	equal	coordinate axes
EVI	λA_k		•	Diagonal	equal	variable	coordinate axes
VVI	$\lambda_k A_k$		•	Diagonal	variable	variable	coordinate axes
EEE	$\lambda D A D^T$	•	•	Ellipsoidal	equal	equal	equal
EEV	$\lambda D_k A D_k^T$		•	Ellipsoidal	equal	equal	variable
VEV	$\lambda_k D_k A D_k^T$		•	Ellipsoidal	variable	equal	variable
VVV	$\lambda_k D_k A_k D_k^T$	•	•	Ellipsoidal	variable	variable	variable

Relation to kmeans algorithm I

Complete (Classification) log-likelihood

$$\begin{aligned} CL(\boldsymbol{\theta}; X, Z) &= \ln \prod_i p(x_i, z_i = k; \boldsymbol{\theta}_k) \\ &= \ln \prod_i \prod_k p(x_i, z_i = k; \boldsymbol{\theta}_k)^{\mathbb{1}(z_i=k)} \\ &= \sum_i \sum_k \mathbb{1}(z_i = k) \ln p(x_i, z_i = k; \boldsymbol{\theta}_k) \end{aligned}$$

CEM algorithm

$CL(\boldsymbol{\theta}; X, Z)$ can be maximized using CEM algorithm
if $\forall k$ we have $\pi_k = \frac{1}{K}$, $\boldsymbol{\Sigma}_k = \sigma I_p$ then $CEM \triangleq kmeans$