

Base form

$$X \sim f$$

$$f(x) = \sum_{k=1}^K \omega_k g_k(x|\theta_k)$$

with $\sum_{k=1}^K \omega_k = 1$

Hierarchical representation

$$X|C \sim g_c$$

$$C \sim \text{Pr}(C = k) = \omega_k$$

$$p(x) = \sum_{k=1}^K p(x|C = k)p(C = k)$$

Observed data likelihood

$$L(X|\omega, \theta) = \prod_{i=1}^n \sum_{k=1}^K \omega_k g_k(x_i|\theta_k)$$

Complete data likelihood

$$L(X|\omega, \theta, c) = \prod_{i=1}^n \prod_{k=1}^K [\omega_k g_k(x_i|\theta_k)]^{\mathbb{1}_{c_i=k}}$$

Probability of a data point x_i being generated by a component g_k

$$\text{Pr}(c_i = k|\omega, \theta) = \frac{\omega_k g_k(x_i|\theta_k)}{\sum_l \omega_l g_l(x_i|\theta_l)} =: v_{i,k}(\omega, \theta)$$

Maximum likelihood framework

Bayesian framework

MODEL FITTING

Expectation-maximization algorithm

MCMC algorithm/Gibbs sampling

output: maximum likelihood estimates of ω and θ
 procedure:

output: a sample from the posterior distribution
 $p(\omega, \theta, c|X)$

$\hat{\omega}, \hat{\theta} \leftarrow$ initial values

procedure:

repeat

set priors $p(\omega)$ and $p(\theta)$ on ω and θ

$\hat{\omega}, \hat{\theta} \leftarrow \text{argmax}_{\omega, \theta} Q(\omega, \theta|\hat{\omega}, \hat{\theta})$

convenient: $(\omega_1, \dots, \omega_K) \sim \text{Dirichlet}(a_1, \dots, a_K)$

until convergence

repeat

where

$\omega \leftarrow \text{Dirichlet}(a_1^*, \dots, a_K^*),$

$Q(\omega, \theta|\hat{\omega}, \hat{\theta}) = E_{c|X, \hat{\omega}, \hat{\theta}} [\log L(X|\omega, \theta, c)]$

with $a_k^* = a_k - 1 + \sum_{i=1}^n \mathbb{1}_{c_i=k}$

$$= \sum_{i=1}^n \sum_{k=1}^K v_{i,k}(\hat{\omega}, \hat{\theta}) [\log \omega_k + \log g_k(x_i|\theta_k)]$$

$c_i \leftarrow p(c_i = k|\dots) = v_{i,k}(\omega, \theta)$

$\theta_k \leftarrow p(\theta_k|\dots) \propto \left[\prod_{\{i: c_i=k\}} g_k(x_i|\theta_k) \right] p_k(\theta_k)$

DETERMINING THE NUMBER OF COMPONENTS

For each candidate model \mathcal{M} compute its Bayesian information criterion value

With

$$\text{BIC}(\mathcal{M}) = -2 \log L(\hat{\eta}_{\mathcal{M}}) + r_{\mathcal{M}} \log n,$$

K – upper bound on # components,

K^* – best guess for actual # components \bar{K} ,

solve

$$K^* = \alpha \log \frac{n+\alpha-1}{\alpha}$$

for α .

where

$$L(\hat{\eta}_{\mathcal{M}}) = L(X|\hat{\omega}_{\mathcal{M}}, \hat{\theta}_{\mathcal{M}}),$$

$r_{\mathcal{M}}$ – # independent parameters in the model.

Set the prior $p(\omega) \sim \text{Dirichlet}(\frac{\alpha}{K}, \dots, \frac{\alpha}{K})$.

Fit a model with K components.

Select the model with the lowest BIC value.

Posterior distribution for \bar{K} obtained by counting the number of unique values assumed by indicators c_i per iteration of Gibbs sampling.

CLASSIFICATION WITH MIXTURE MODELS

Unsupervised: fit a model, classify as $c_i = \text{argmax}_k v_{i,k}(\hat{\omega}, \hat{\theta})$

Semi-supervised: fit a model to the whole dataset, keeping $v_{i,k}$ constant on the training data;

for the test data, classify as $c_i = \text{argmax}_k v_{i,k}(\hat{\omega}, \hat{\theta})$

Supervised: estimate $\hat{\theta}_k$ by fitting g_k to training data from class k ; set $\hat{\omega}_k = (\text{\#obs. in class } k)/n$;

for the test data, classify as $c_i = \text{argmax}_k v_{i,k}(\hat{\omega}, \hat{\theta})$