

Foundations Statistics

Wafaa Mohammed

Exercise 4.4. [BIC for Gaussians]

(Source: Jaakkola.) The Bayesian information criterion (BIC) is a penalized log-likelihood function that can be used for model selection. It is defined as

$$BIC = \log p(\mathcal{D}|\hat{\theta}_{ML}) - \frac{d}{2}\log(N)$$

where d is the number of free parameters in the model and N is the number of samples. In this question, we will see how to use this to choose between a full covariance Gaussian and a Gaussian with a diagonal covariance. Obviously a full covariance Gaussian has higher likelihood, but it may not be “worth” the extra parameters if the improvement over a diagonal covariance matrix is too small. So we use the BIC score to choose the model.

We can write

$$\log p(\mathcal{D}|\hat{\Sigma}, \hat{\mu}) = -\frac{N}{2}\text{tr}(\hat{\Sigma}^{-1}\hat{S}) - \frac{N}{2}\log(|\hat{\Sigma}|)$$
$$\hat{S} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

where \hat{S} is the scatter matrix (empirical covariance), the trace of a matrix is the sum of its diagonals, and we have used the trace trick.

- Derive the BIC score for a Gaussian in D dimensions with full covariance matrix. Simplify your answer as much as possible, exploiting the form of the MLE. Be sure to specify the number of free parameters d .
- Derive the BIC score for a Gaussian in D dimensions with a diagonal covariance matrix. Be sure to specify the number of free parameters d . Hint: for the diagonal case, the ML estimate of Σ is the same as $\hat{\Sigma}_{ML}$ except the off-diagonal terms are zero:

$$\hat{\Sigma}_{diag} = \text{diag}(\hat{\Sigma}_{ML}(1,1), \dots, \hat{\Sigma}_{ML}(D,D))$$

Solution. a. The MLE for the covariance matrix is the empirical covariance matrix: $\hat{\Sigma}_{MLE} = \hat{S}$. The log likelihood becomes:

$$\log p(\mathcal{D}|\hat{\Sigma}, \hat{\mu}) = -\frac{N}{2}\text{tr}(\hat{S}^{-1}\hat{S}) - \frac{N}{2}\log(|\hat{S}|)$$

$$= -\frac{N}{2}tr(I) - \frac{N}{2}log(|\hat{S}|) = -\frac{N}{2}D - \frac{N}{2}log(|\hat{S}|)$$

The number of free parameters:

- for the symmetric covariance matrix = $\frac{D(D+1)}{2}$
- for the mean = D

Thus, the total number of free parameters:

$$d = \frac{D(D+1)}{2} + D$$

The BIC for a Gaussian with full covariance matrix:

$$\begin{aligned} BIC &= log p(\mathcal{D}|\hat{\theta}_{ML}) - \frac{d}{2}log(N) \\ &= -\frac{N}{2}D - \frac{N}{2}log(|\hat{S}|) - \frac{1}{2}\left(\frac{D(D+1)}{2} + D\right)log(N) \end{aligned}$$

b. for a Gaussian with a diagonal covariance matrix

$$\hat{\Sigma}_{MLE} = \hat{\Sigma}_{diag} = \text{diag}(\hat{S}(1,1), \dots, \hat{S}(D,D)).$$

The log likelihood becomes:

$$\begin{aligned} log p(\mathcal{D}|\hat{\Sigma}, \hat{\mu}) &= -\frac{N}{2}tr(\hat{\Sigma}_{diag}^{-1}\hat{S}) - \frac{N}{2}log(|\hat{\Sigma}_{diag}|) \\ &= -\frac{N}{2} \sum_{i=1}^D \frac{\hat{S}(i,i)}{\hat{\Sigma}_{diag}(i,i)} - \frac{N}{2} \sum_{i=1}^D log \hat{\Sigma}_{diag}(i,i) = -\frac{N}{2}D - \frac{N}{2} \sum_{i=1}^D log \hat{S}(i,i) \end{aligned}$$

The number of free parameters:

- for the diagonal covariance matrix = D
- for the mean = D

Thus, the total number of free parameters:

$$d = 2D$$

The BIC for a Gaussian with full covariance matrix:

$$\begin{aligned} BIC &= log p(\mathcal{D}|\hat{\theta}_{ML}) - \frac{d}{2}log(N) \\ &= -\frac{N}{2}D - \frac{N}{2} \sum_{i=1}^D log \hat{S}(i,i) - D log(N) \end{aligned}$$

□