Lecture Stat 461-561 Information Criterion

Arnaud Doucet

February 2008



# Review of Maximum Likelihood Approach

• We have data 
$$X_i \stackrel{\mathrm{i.i.d.}}{\sim} g(x)$$
.

- We model the distribution of these data by a parametric model  $\{f(x|\theta); \theta \in \Theta \subseteq \mathbb{R}^p\}.$
- We estimate  $\theta$  by ML; that is

$$\widehat{ heta} = rgmax_{ heta \in \Theta} I\left( heta
ight) ext{ where } I\left( heta
ight) := \sum_{k=1}^n \log f\left(X_k \middle| heta
ight)$$

• We know that as  $n \to \infty$  then  $\widehat{\theta}$  converges towards

$$\theta^{*} = \operatorname*{arg\,min}_{\theta \in \Theta} \mathit{KL}\left(g\left(x\right); f\left(\left.x\right|\theta\right)\right)$$

where

$$KL(g(x); f(x|\theta)) := \int g(x) \log \frac{g(x)}{f(x|\theta)} dx$$

# A Central Limit Theorem

• Assuming the MLE is asymptotically consistent, then we have

$$\mathbf{0} = \frac{\partial I(\theta)}{\partial \theta} \Big|_{\widehat{\theta}} = \frac{\partial I(\theta)}{\partial \theta} \Big|_{\theta^*} + \frac{\partial^2 I(\theta)}{\partial \theta \partial \theta^{\mathsf{T}}} \Big|_{\theta^*} \left(\widehat{\theta} - \theta^*\right) + \dots$$

• The law of large numbers yields

$$-\frac{1}{n} \left. \frac{\partial I\left(\theta\right)}{\partial \theta \partial \theta^{\mathsf{T}}} \right|_{\theta^{*}} \to \left( J\left(\theta^{*}\right) := -\int g\left(x\right) \left. \frac{\partial^{2} \log f\left(x \mid \theta\right)}{\partial \theta \partial \theta^{\mathsf{T}}} \right|_{\theta^{*}} dx \right)$$

• The CLT yields

$$\sqrt{n}\frac{1}{n} \left. \frac{\partial I(\theta)}{\partial \theta} \right|_{\theta^*} \xrightarrow{\mathsf{D}} \mathcal{N}\left(\mathsf{0}, I(\theta^*)\right)$$

where

$$I\left(\theta^{*}\right) = \int g\left(x\right) \left.\frac{\partial \log f\left(x|\theta\right)}{\partial \theta}\right|_{\theta^{*}} \left.\frac{\partial \log f\left(x|\theta\right)}{\partial \theta^{\mathsf{T}}}\right|_{\theta^{*}} dx$$

• Using Slutzky's theorem we have

$$\sqrt{n}\left(\widehat{\theta}-\theta^{*}\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(\mathsf{0},J^{-1}\left(\theta^{*}\right)I\left(\theta^{*}\right)J^{-1}\left(\theta^{*}\right)\right).$$

- This is sometimes known as the sandwitch variance.
- When  $g(x) = f(x | \theta^*)$  then we have

$$J\left(\theta^{*}\right)=I\left(\theta^{*}\right)$$

and the 'standard' CLT follows

$$\sqrt{n}\left(\widehat{\theta}-\theta^{*}\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(\mathsf{0},I^{-1}\left(\theta^{*}\right)\right)$$

# **Evaluating The Statistical Model**

- We want to compare our model  $f\left(x|\hat{\theta}\right)$  to  $g\left(x\right)$ .
- One way consists of evaluating

$$\mathcal{K}L\left(g\left(x
ight);f\left(x|\,\widehat{ heta}
ight)
ight)=\int g\left(x
ight)\log g\left(x
ight)dx-\int g\left(x
ight)\log f\left(x|\,\widehat{ heta}
ight)dx.$$

and the larger  $\int g(x) \log f(x|\hat{\theta}) dx$ , the closer the model is to the true one.

• The crucial issue is to obtain an estimate of

$$\mathbb{E}_{g}\left[\log f\left(X|\widehat{\theta}\right)\right] = \int g(x)\log f\left(x|\widehat{\theta}\right) dx.$$

• Clearly an estimate of  $\mathbb{E}_{g}\left[\log f\left(X|\widehat{\theta}\right)\right]$  is  $n^{-1}l\left(\widehat{\theta}\right)$  and an estimate of  $n\mathbb{E}_{g}\left[\log f\left(X|\widehat{\theta}\right)\right]$  is  $l\left(\widehat{\theta}\right)$ .

# Selecting a Model

- In practice, it is difficult to precisely capture the true structure of given phenomena from a limited number of data.
- Often several candidate statistical models are selected and we want to select the model that closely approximates the true distribution of the data.
- Given *m* candidate models {*f<sub>i</sub>* (*x*|*θ<sub>i</sub>*); *i* = 1, ..., *m*} and the associated ML estimates *θ̂<sub>i</sub>*, a simple solution would consist of selecting the model *j* where

$$j = rg\max_{i \in \{1,...,m\}} l_i\left(\widehat{ heta}_i
ight);$$

i.e. picking the model for which  $KL\left(g\left(x\right); f_{i}\left(x|\widehat{\theta}_{i}\right)\right)$  is the smallest.

- This approach does not provide a fair comparison of models.
- The quantity  $I_i\left(\widehat{\theta}_i\right)$  contains a bias as an estimator of  $n\mathbb{E}_g\left[\log f_i\left(X|\widehat{\theta}_i\right)\right]$ .
- The reason is that the data  $X_1, ..., X_n$  are used twice: to estimate  $\hat{\theta}_i$  and to approximate the expectation with respect to g.
- We will show that the size of the bias depends on the size p of the parameter  $\theta$ .

# Relationship between log-likelihood and expected log-likelihood

- Let us denote  $\theta_i^* = \arg \min KL(g(x); f_i(x|\theta_i))$  the 'true' parameter.
- We necessarily have

$$\mathbb{E}_{g}\left[\log f_{i}\left(X|\widehat{\theta}_{i}\right)\right] \leq \mathbb{E}_{g}\left[\log f_{i}\left(X|\theta_{i}^{*}\right)\right]$$

whereas

$$I_i\left(\widehat{\theta}_i\right) \geq I_i\left(\theta_i^*\right).$$

To correct for this bias, we want to estimate

$$b_{i}(g) = \mathbb{E}_{g}\left[l_{i}\left(\widehat{\theta}_{i}\right) - n\mathbb{E}_{g}\left[\log f_{i}\left(X|\widehat{\theta}_{i}\right)\right]\right]$$

where remember that

$$I_i\left(\widehat{\theta}_i\right) = \sum_{k=1}^n \log f_i\left(X_k | \, \widehat{\theta}_i\right)$$

and  $\widehat{\theta}_i = \widehat{\theta}_i (X_{1:n})$ .

• The information criterion for model *i* is defined as

$$\mathsf{IC}\left(i\right) = -2\sum_{k=1}^{n} \log f_{i}\left(X_{k} | \,\widehat{\theta}_{i}\right) + 2\left\{\mathsf{estimator \ for \ } b_{i}\left(g\right)\right\}$$

- IC(*i*) is a biased-corrected estimate of minus the expected log-likelihood.
- So given a collection of models, we will select the model

$$j = \underset{i \in \{1, \dots, m\}}{\arg\min} \operatorname{IC}(i)$$

• We want to estimate the bias b(g) which is given by

$$\mathbb{E}_{g}\left[I\left(\widehat{\theta}\left(X_{1:n}\right)\right) - n\mathbb{E}_{g}\left[\log f\left(X|\widehat{\theta}\left(X_{1:n}\right)\right)\right]\right]$$

$$= \underbrace{\mathbb{E}_{g}\left[I\left(\widehat{\theta}\left(X_{1:n}\right)\right) - I\left(\theta^{*}\right)\right]}_{D_{1}}$$

$$+ \underbrace{\mathbb{E}_{g}\left[I\left(\theta^{*}\right) - n\mathbb{E}_{g}\left[\log f\left(X|\theta^{*}\right)\right]\right]}_{D_{2}}$$

$$+ \underbrace{\mathbb{E}_{g}\left[n\mathbb{E}_{g}\left[\log f\left(X|\theta^{*}\right)\right] - n\mathbb{E}_{g}\left[\log f\left(X|\widehat{\theta}\left(X_{1:n}\right)\right)\right]\right]}_{D_{3}}$$

#### • We have

$$D_{2} = \mathbb{E}_{g} \left[ I(\theta^{*}) - n\mathbb{E}_{g} \left[ \log f(X|\theta^{*}) \right] \right]$$
  
=  $\mathbb{E}_{g} \left[ \sum_{k=1}^{n} \log f(X_{k}|\theta^{*}) - n\mathbb{E}_{g} \left[ \log f(X|\theta^{*}) \right] \right]$   
= 0

• No bias for this term...

# Calculation of third term

• Let us denote

$$\eta\left(\widehat{\theta}\right) = \mathbb{E}_{g}\left[\log f\left(X|\widehat{\theta}\right)\right].$$

 $\bullet$  By performing a Taylor expansion around  $\theta^*,$  we obtain

$$\eta\left(\widehat{\theta}\right) = \eta\left(\theta^{*}\right) + \frac{\partial\eta\left(\theta\right)}{\partial\theta}\Big|_{\theta^{*}}^{\mathsf{T}}\left(\widehat{\theta} - \theta^{*}\right) + \frac{1}{2}\left(\widehat{\theta} - \theta^{*}\right)^{\mathsf{T}}\left.\frac{\partial^{2}\eta\left(\theta\right)}{\partial\theta\partial\theta^{\mathsf{T}}}\Big|_{\theta^{*}}\left(\widehat{\theta} - \theta^{*}\right)^{2}\right)$$

• As 
$$\frac{\partial \eta(\theta)}{\partial \theta}\Big|_{\theta^*} = \mathbf{0}$$
 then  
 $\eta\left(\widehat{\theta}\right) \approx \eta\left(\theta^*\right) - \frac{1}{2}\left(\widehat{\theta} - \theta^*\right)^{\mathsf{T}} J\left(\theta^*\right)\left(\widehat{\theta} - \theta^*\right)$ 

where

$$J\left(\theta^{*}\right) = -\left.\frac{\partial^{2}\eta\left(\theta\right)}{\partial\theta\partial\theta^{\mathsf{T}}}\right|_{\theta^{*}} = -\int g\left(x\right)\left.\frac{\partial^{2}\log f\left(x\right|\theta)}{\partial\theta\partial\theta^{\mathsf{T}}}\right|_{\theta^{*}}dx$$

• It follows that

$$D_{3} = \mathbb{E}_{g} \left[ n \left( \eta \left( \theta^{*} \right) - \eta \left( \widehat{\theta} \right) \right) \right]$$
  

$$\approx \frac{n}{2} \mathbb{E}_{g} \left[ \left( \widehat{\theta} - \theta^{*} \right)^{\mathsf{T}} J \left( \theta^{*} \right) \left( \widehat{\theta} - \theta^{*} \right) \right]$$
  

$$= \frac{n}{2} \mathbb{E}_{g} \left[ \operatorname{tr} \left\{ J \left( \theta^{*} \right) \left( \widehat{\theta} - \theta^{*} \right) \left( \widehat{\theta} - \theta^{*} \right)^{\mathsf{T}} \right\} \right]$$
  

$$= \frac{n}{2} \operatorname{tr} \left\{ J \left( \theta^{*} \right) \mathbb{E}_{g} \left[ \left( \widehat{\theta} - \theta^{*} \right) \left( \widehat{\theta} - \theta^{*} \right)^{\mathsf{T}} \right] \right\}$$

• We have from the CLT established previously we have

$$\mathbb{E}_{g}\left[\left(\widehat{\theta}-\theta^{*}\right)\left(\widehat{\theta}-\theta^{*}\right)^{\mathsf{T}}\right]\approx\frac{1}{n}J\left(\theta^{*}\right)^{-1}I\left(\theta^{*}\right)J\left(\theta^{*}\right)^{-1}$$

so

$$D_{3} \approx \frac{1}{2} \operatorname{tr} \left\{ I\left(\theta^{*}\right) J\left(\theta^{*}\right)^{-1} \right\}$$

# Calculation of first term

• We have

$$I(\theta) = I\left(\widehat{\theta}\right) + \frac{\partial I(\theta)}{\partial \theta}\Big|_{\widehat{\theta}}^{\mathsf{T}}\left(\theta - \widehat{\theta}\right) + \frac{1}{2}\left(\theta - \widehat{\theta}\right)^{\mathsf{T}}\left.\frac{\partial^{2}I(\theta)}{\partial \theta \partial \theta^{\mathsf{T}}}\Big|_{\widehat{\theta}}\left(\theta - \widehat{\theta}\right) + \dots$$

so

٢

$$I(\theta^*) - I(\widehat{\theta}) \approx -\frac{n}{2} \left(\theta^* - \widehat{\theta}\right)^{\mathsf{T}} J(\theta^*) \left(\theta^* - \widehat{\theta}\right)$$

$$D_{1} = \mathbb{E}_{g} \left[ I \left( \widehat{\theta} \left( X_{1:n} \right) \right) - I \left( \theta^{*} \right) \right]$$

$$\approx \frac{n}{2} \mathbb{E}_{g} \left[ \left( \theta^{*} - \widehat{\theta} \right)^{\mathsf{T}} J \left( \theta^{*} \right) \left( \theta^{*} - \widehat{\theta} \right) \right]$$

$$= \frac{n}{2} \mathbb{E}_{g} \left[ J \left( \theta^{*} \right) \operatorname{tr} \left\{ \left( \theta^{*} - \widehat{\theta} \right)^{\mathsf{T}} \left( \theta^{*} - \widehat{\theta} \right) \right\} \right]$$

$$= \frac{1}{2} \operatorname{tr} \left( I \left( \theta^{*} \right) J \left( \theta^{*} \right)^{-1} \right)$$

## Estimate of the Bias

We have

$$\begin{array}{lll} b\left(g\right) &=& D_{1}+D_{2}+D_{3}\\ &\approx& \frac{1}{2}\mathrm{tr}\left\{I\left(\theta^{*}\right)J\left(\theta^{*}\right)^{-1}\right\}+0+\frac{1}{2}\mathrm{tr}\left\{I\left(\theta^{*}\right)J\left(\theta^{*}\right)^{-1}\right\}\\ &=& \mathrm{tr}\left\{I\left(\theta^{*}\right)J\left(\theta^{*}\right)^{-1}\right\} \end{array}$$

• Let  $\widehat{I}$  and  $\widehat{J}$  be consistent estimate of  $I\left(\theta^{*}\right)$  and  $J\left(\theta^{*}\right)$ , say

$$\widehat{I} = \frac{1}{n} \sum_{k=1}^{n} \frac{\partial \log f(X_{k} | \theta)}{\partial \theta} \Big|_{\widehat{\theta}} \frac{\partial \log f(X_{k} | \theta)}{\partial \theta^{\mathsf{T}}} \Big|_{\widehat{\theta}},$$

$$\widehat{J} = -\frac{1}{n} \sum_{k=1}^{n} \frac{\partial^{2} \log f(X_{k} | \theta)}{\partial \theta \partial \theta^{\mathsf{T}}} \Big|_{\widehat{\theta}}$$

then we can estimate the bias through

$$\widehat{b}\left( g
ight) = ext{tr}\left( \widehat{I}\widehat{J}^{-1}
ight)$$

• We have for the information criterion for model *i* 

$$IC\left(i
ight)=-2\sum_{k=1}^{n}\log f_{i}\left(\left.X_{k}
ight|\widehat{ heta}_{i}
ight)+2 ext{tr}\left(\widehat{I}_{i}\widehat{J}_{i}^{-1}
ight)$$

• Assuming  $g(x) = f_i(x|\theta_i^*)$  then  $I_i(\theta_i^*) = J_i(\theta_i^*)^{-1}$  so  $b_i(g) = p_i$ (where  $p_i$  is the dimension of  $\theta_i$ ) so in this case

$$AIC(i) = -2\sum_{k=1}^{n} \log f_i\left(X_k | \widehat{\theta}_i\right) + 2p_i.$$

# Checking the Equality of Two Discrete Distributions

• Assume we have two sets of data each having k categories and

Category	1	2	• • •	k
Data set 1	$n_1$	<i>n</i> 2	• • •	n <sub>k</sub>
Data set 2	$m_1$	$m_2$	• • •	$m_k$

where we have  $n_1 + \cdots + n_k = n$  observations for the first dataset and  $m_1 + \cdots + m_k = m$  for the second.

• We assume these datasets follow the multinomial distributions

$$p(n_{1}, \dots, n_{k} | p_{1}, \dots, p_{k}) = \frac{n!}{n_{1}! \cdots n_{k}!} p_{1}^{n_{1}} \cdots p_{k}^{n_{k}},$$
  
$$p(m_{1}, \dots, m_{k} | q_{1}, \dots, q_{k}) = \frac{m!}{m_{1}! \cdots m_{k}!} q_{1}^{m_{1}} \cdots q_{k}^{m_{k}}$$

• We want to check whether  $p_1, \cdots, p_k \neq q_1, \cdots, q_k$  or  $p_1, \cdots, p_k = q_1, \cdots, q_k$ 

• Assume the two distributions are different then

$$I(p_{1:k}, q_{1:k}) = \log n! - \sum_{j=1}^{k} (\log n_j! + n_j \log p_j) + \log m! - \sum_{j=1}^{k} (\log m_j! + m_j \log q_j)$$

• We have the MLE  $\widehat{p}_j = \frac{n_j}{n}, \ \widehat{q}_j = \frac{m_j}{n}$  • So

$$I\left(\widehat{p}_{1:k}, \widehat{q}_{1:k}\right) = C + \sum_{j=1}^{k} \left(n_j \log \frac{n_j}{n} + m_j \log \frac{m_j}{n}\right)$$

with  $C = \log n! + \log m! - \sum_{j=1}^{k} (\log n_j! + \log m_j!)$  and

$$AIC = -2\left(C + \sum_{j=1}^{k} \left(n_j \log \frac{n_j}{n} + m_j \log \frac{m_j}{n}\right)\right) + 4\left(k - 1\right)$$

• If we assume the two distributions are equal then

$$I(r) = C - \sum_{j=1}^{k} (n_j + m_j) \log r_j$$

 $\mathsf{and}$ 

$$\widehat{r}_j = rac{n_j + m_j}{n + m}$$

• We have

$$I(\widehat{r}_{1:k}) = C + \sum_{j=1}^{k} (n_j + m_j) \log \left(\frac{n_j + m_j}{n + m}\right)$$

and

$$AIC = -2\left(C + \sum_{j=1}^{k} (n_j + m_j) \log\left(\frac{n_j + m_j}{n + m}\right)\right) + 2(k - 1)$$



#### • Consider the following data

Category	1	2	3	4	5
Data set 1	304	800	400	57	323
Data set 2	174	509	362	80	214

• From this table we can deduce

Category	1	2	3	4	5
$\widehat{p}_j$	0.16	0.42	0.21	0.03	0.17
$\widehat{q}_j$	0.13	0.38	0.27	0.06	0.16
$\hat{r}_j$	0.15	0.41	0.24	0.04	0.17

 Ignoring the constant C, the maximum log-likelihood of the models are

Model 1 : 
$$\sum_{j=1}^{k} \left( n_j \log \frac{n_j}{n} + m_j \log \frac{m_j}{n} \right) = -4567.36,$$
  
Model 2 :  $\sum_{j=1}^{k} \left( n_j + m_j \right) \log \left( \frac{n_j + m_j}{n + m} \right) = -4585.61.$ 

 The number of free parameters of the models is 2 (k - 1) = 8 in model 1 and 4 in model 2. So it follows that AIC for model 1 and model 2 is respectively 9150,73 and 9179,22 so AIC selects Model 1.

# Determining the number of bin size of an histogram

- Histograms are use for representing the properties of a set of observations obtained from either a discrete or continuous distribution.
- Assume we have a histogram  $\{n_1, n_2, ..., n_k\}$  where k is refer to as the bin size.
- If k is too large or too small, then it is difficult to capture the characteristics of the true distribution.
- We can think of an histogram as a model specified by a multinomial distribution

$$p(n_1, \cdots, n_k | p_1, \cdots, p_k) = \frac{n!}{n_1! \cdots n_k!} p_1^{n_1} \cdots p_k^{n_k}$$

where  $n_1 + \cdots + n_k = n$ ,  $p_1 + \cdots + p_k = 1$ .

• In this case, we have seen that the MLE is  $\hat{p}_j = \frac{n_j}{n}$  and

$$I(\widehat{p}_{1:k}) = C + \sum_{j=1}^{k} n_j \log \frac{n_j}{n}$$

with 
$$C = \log n! - \sum_{j=1}^{k} \log n_j!$$

• We have k-1 free parameters so

$$AIC = -2\left\{C + \sum_{j=1}^{k} n_j \log \frac{n_j}{n}\right\} + 2(k-1).$$

• We want to compare this model to a model with lower resolutions.

• To compare the histogram model with a simpler one, we may assume  $p_{2j-1} = p_{2j}$  for j = 1, ..., m and for sake of simplicity, we consider here k = 2m. In this case, the new MLE is  $\hat{p}_{2j-1} = \hat{p}_{2j} = \frac{n_{2j-1}+n_{2j}}{2n}$ 

In this case, the AIC is given by

$$AIC = -2\left\{C + \sum_{j=1}^{m} (n_{2j-1} + n_{2j}) \log \frac{n_{2j-1} + n_{2j}}{2n}\right\} + 2(m-1).$$

• Similarly we can consider the histogram with m bins where k = 4m.

• We apply this to galaxy data for which AIC selects 14

bin size 28		14	7	
log-like	-189.19	-197.71	-209.52	
AIC	432.38	421.43	431.03	

# Application to Polynomial Regression

- We are given *n* observations  $\{x_i, y_i\}$ .
- We want to use the following model

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \ldots + \beta_k x^k + \varepsilon, \ \varepsilon \sim \mathcal{N}\left(0, \sigma^2\right)$$

where k is unknown.

$$ullet$$
 We have  $heta=ig(eta_{0:k},\sigma^2ig)$  and

$$I(\theta) = -\frac{n}{2}\log\left(2\pi\sigma^2\right) - \frac{1}{2\sigma^2}\sum_{j=1}^n\left(y_j - \left(\sum_{m=0}^k \beta_m x_j^m\right)\right)^2$$

• We can easily establish that

$$I\left(\widehat{\theta}\right) = -\frac{n}{2}\log\left(2\pi\widehat{\sigma}^2\right) - \frac{n}{2}$$

 In this case for the polynomial model of order k we have k + 2 unknown parameters so

$$AIC(k) = n(\log 2\pi + 1) + n\log \widehat{\sigma}^{2} + 2(k+2)$$

• This yields

k	0	1	2	3	4	5	6	7
$I\left(\widehat{\theta}_{k}\right)$	22.41	31.19	41.51	42.52	43.75	44.44	45.00	45.45
AIC	-40.81	-56.38	-75.03	-75.04	-75.50	-74.89	-74.00	-72.89

• AIC selects model 4.

- Daily minimum temperatures y<sub>i</sub> in January averaged from 1971 through 2000 for city i
- x<sub>i1</sub> latitudes, x<sub>i2</sub> longitudes and x<sub>i3</sub> altitudes.
- A standard model is

$$y_i = a_0 + a_1 x_{i1} + a_2 x_{i2} + a_3 x_{i3} + \varepsilon_i$$

where  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ .

• But we would also like to compare models where some of the explanatory variables are omitted.

### Results

We have

 $AIC (model) = n (\log 2\pi + 1) + n \log \hat{\sigma}^2 + 2 (nb. explanatority var. + 2)$ 

#### • This yields

model	<i>x</i> <sub>1</sub> , <i>x</i> <sub>3</sub>	<i>x</i> <sub>1</sub> , <i>x</i> <sub>2</sub> , <i>x</i> <sub>3</sub>	<i>x</i> <sub>1</sub> , <i>x</i> <sub>2</sub>	<i>x</i> <sub>1</sub>	<i>x</i> <sub>2</sub> , <i>x</i> <sub>3</sub>	<i>x</i> <sub>2</sub>	<i>x</i> <sub>3</sub>
$\widehat{\sigma}^2$	1.49	1.48	5.11	5.54	5.69	7.81	19.96
AIC	88.92	90.81	119.71	119.73	122.43	128.35	151.88

• We select the model

$$y_i = 40.49 - 1.11x_{i1} - 0.010x_{i3} + \varepsilon_i$$

with  $\varepsilon_{i} \sim \mathcal{N}\left(0, 1.49
ight)$  .

• We have the following model

$$y_{k} = \sum_{i=1}^{m} a_{i} y_{k-i} + \varepsilon_{k}, \quad \varepsilon_{k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, \sigma^{2}\right).$$

 Assuming we have data y<sub>1</sub>, y<sub>2</sub>, ..., y<sub>n</sub> and that y<sub>0</sub>, y<sub>-1</sub>, ..., y<sub>1-m</sub> are deterministic then we can easily obtain the MLE estimate and check that

$$I\left(\widehat{\theta}\right) = -\frac{n}{2}\log\left(2\pi\widehat{\sigma}^2\right) - \frac{n}{2}$$

• Since the model with order m has m+1 unknown parameters then

$$AIC(m) = n(\log 2\pi + 1) + n\log \hat{\sigma}^2 + 2(m+1)$$

- This corresponds to the logarithms of the annual numbers of lynx trapped from 1821 to 1934; *n* = 114 observations.
- We try 20 candidate models and m = 11 is selected using AIC.

- We some times encounter the situation in which the stochastic structure of the data changes at a certain time or location.
- Let us start with a simple model where

$$Y_n \sim \mathcal{N}\left(\mu_n, \sigma^2\right)$$

and

$$\mu_n = \begin{cases} \theta_1 & \text{if } n < k \\ \theta_2 & \text{if } n \ge k \end{cases}$$

where k is the unknown so-called changepoint.

• We have N data and, for the model k, the likelihood is

$$L\left(\theta_{1},\theta_{2},\sigma^{2}\right)=\prod_{n=1}^{k-1}\mathcal{N}\left(y_{n};\theta_{1},\sigma^{2}\right)\prod_{n=k}^{N}\mathcal{N}\left(y_{n};\theta_{2},\sigma^{2}\right)$$

• It can be easily established that

$$\widehat{\theta}_{1} = \frac{1}{k-1} \sum_{n=1}^{k-1} y_{n}, \ \widehat{\theta}_{2} = \frac{k}{N-k+1} \sum_{n=k}^{N} y_{n},$$
$$\widehat{\sigma}^{2} = \frac{1}{N} \left\{ \sum_{n=1}^{k-1} \left( y_{n} - \widehat{\theta}_{1} \right)^{2} + \sum_{n=k}^{N} \left( y_{n} - \widehat{\theta}_{2} \right)^{2} \right\}$$

So we have the maximum log-likelihood given by

$$I\left(\widehat{\theta}_{1},\widehat{\theta}_{2},\widehat{\sigma}^{2}\right) = -\frac{N}{2}\log\left(2\pi\widehat{\sigma}^{2}\right) - \frac{N}{2}$$

where we emphasize that  $\hat{\sigma}^2$  is a function of k,

• The AIC criterion is given by

$$AIC = N \log \left(2\pi \widehat{\sigma}^2\right) + N + 2 \times 3$$

and the changepoint k can be automatically determined by finding the value of k that gives the smallest AIC.

Arnaud Doucet ()

# Application to Factor Analysis

- Suppose we have  $x = (x_1, ..., x_p)^T$  a vector of mean  $\mu$  and variance-covariance  $\Sigma$ .
- The factor analysis model is

$$x = \mu + Lf + \varepsilon$$

where *L* is a  $p \times m$  matrix of factor loadings whereas  $f = (f_1, ..., f_m)^{\mathsf{T}}$ and  $\varepsilon = (\varepsilon_1, ..., \varepsilon_p)^{\mathsf{T}}$  are unobserved random vectors assumed to satisfy

$$\begin{split} \mathbb{E}\left[f\right] &= 0, \ \textit{Cov}\left[f\right] = \textit{I}_{m}, \\ \mathbb{E}\left[\varepsilon\right] &= 0, \ \textit{Cov}\left[\varepsilon\right] = \Psi = \mathsf{diag}\left(\Psi_{1}, ..., \Psi_{p}\right), \\ \textit{Cov}\left[f, \varepsilon\right] &= 0. \end{split}$$

It then follows that

$$\Sigma = LL^{\mathsf{T}} + \Psi.$$

- Let S denote the sample covariance.
- Under a normal assumption for f and ε, it can be shown that the ML estimates of L and Ψ can be obtained by minimizing

$$\log |\Sigma| - \log |S| + \operatorname{tr} (\Sigma^{-1}S) - p$$

subject to  $L^{\mathsf{T}}\Psi^{-1}L$  being a diagonal matrix.

In this case, we have

$$AIC(m) = n\left\{ p \log(2\pi) + \log\left|\widehat{\Sigma}\right| + tr\left(\Sigma^{-1}S\right) \right\} \\ + 2\left\{ p(m+1) - \frac{1}{2}m(m-1) \right\}$$