# Lecture Stat 461-561 Maximum Likelihood Estimation

A.D.

January 2008



- Maximum Likelihood Estimation
- Invariance
- Consistency
- Efficiency
- Nuisance Parameters

- Let f (x|θ) denote the joint pdf or pmf of the sample
   X = (X<sub>1</sub>, ..., X<sub>n</sub>) parametrized by θ ∈ Θ. Then given that X = x is observed, the function L (θ|x) = f (x|θ) is the likelihood function.
- The most common estimate is the Maximum Likelihood Estimate (MLE) given by

$$\widehat{ heta} = egin{arggammatrix} rgmax & L\left( \left. heta 
ight| {f x} 
ight) . \ heta \in \Theta \ \end{split}$$

• Example: Gaussian distribution

$$f(x_i \mid \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right).$$

Then we have with  $heta=\left(\mu,\sigma^2
ight)$ 

$$\log L(\theta | \mathbf{x}) = \sum_{i=1}^{n} \log f(x_i | \theta)$$
$$= -\frac{n}{2} \log (2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu)^2$$

• By taking the derivatives and setting them to zero

$$\frac{\partial \log L\left(\theta \mid \mathbf{x}\right)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0,$$
  
$$\frac{\partial \log L\left(\theta \mid \mathbf{x}\right)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2 = 0.$$

• By solving these equations, we obtain

$$\widehat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i,$$

$$\widehat{\sigma^2} = \frac{1}{n} \sum_{i=1}^n \left( x_i - \widehat{\mu} \right)^2.$$

• Note that  $\widehat{\mu}$  is an unbiased estimate but  $\widehat{\sigma^2}$  is biased.

• Example: Laplace Distribution (Double Exponential)

$$f(x_i|\theta) = \frac{1}{2} \exp\left(-|x_i - \theta|\right).$$

Then we have

$$\log L(\theta | \mathbf{x}) = -n \log 2 - \sum_{i=1}^{n} |x_i - \theta|.$$

• By taking the derivative, we obtain

$$\frac{d \log L\left(\left.\theta\right| \mathbf{x}\right)}{d \theta} = \sum_{i=1}^{n} \operatorname{sgn}\left(x_{i} - \theta\right)$$

hence

$$\widehat{\theta} = \mathsf{med} \{x_1, ..., x_n\}$$

for n = 2p + 1.

• *Example* (Uniform Distribution): Consider  $X_i \sim \mathcal{U}(0, \theta)$ , i.e.

$$f(x_i | \theta) = \begin{cases} 1/\theta & \text{if } 0 \le x < \theta, \\ 0 & \text{otherwise.} \end{cases}$$

We have

$$L(\theta | \mathbf{x}) = \prod_{i=1}^{n} f(x_i | \theta) = \begin{cases} (1/\theta)^n & \text{if } \theta \ge x_{(n)} \\ 0 & \text{if } \theta < x_{(n)} \end{cases}$$

• It follows that  $\widehat{ heta} = x_{(n)}$  where  $x_{(1)} < x_{(2)} < \cdots < x_{(n)}.$ 

.

• Example (Linear Regression): Let  $\{\mathbf{x}_i, y_i\}_{i=1}^n$  be a set of n data where  $\mathbf{x}_i = (x_1^i, x_2^i, ..., x_p^i)^T$  is a set of explanatory variables and  $y_i \in \mathbb{R}$  is the response. We assume

$$y_i = \mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta} + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}\left(\mathbf{0}, \sigma^2\right)$$

thus

$$f(y_i | \mathbf{x}_i, \beta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta})^2}{2\sigma^2}\right)$$

• We have for  $heta=ig(oldsymbol{eta},\!\sigma^2ig)$ 

$$\begin{split} \log L\left(\theta\right) &= \sum_{i=1}^{n} \log f\left(y_{i} | \mathbf{x}_{i}, \boldsymbol{\beta}\right) \\ &= -\frac{n}{2} \log\left(2\pi\sigma^{2}\right) - \frac{1}{2\sigma^{2}} \sum_{i=1}^{n} \left(y_{i} - \mathbf{x}_{i}^{\mathsf{T}} \boldsymbol{\beta}\right)^{2} \\ &= -\frac{n}{2} \log\left(2\pi\sigma^{2}\right) - \frac{1}{2\sigma^{2}} \left(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\right)^{\mathsf{T}} \left(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\right) \\ \end{split}$$
 where  $\mathbf{y} = (y_{1}, ..., y_{n})^{\mathsf{T}}$  and  $\mathbf{X} = (\mathbf{x}_{1}, ..., \mathbf{x}_{n})^{\mathsf{T}}$ .

• By taking the derivatives and setting them to zero

$$\frac{\partial \log L(\theta | \mathbf{x})}{\partial \boldsymbol{\beta}} = -\frac{1}{2\sigma^2} \left( -2\mathbf{X}^{\mathsf{T}} \boldsymbol{\beta} + 2\mathbf{X}^{\mathsf{T}} \mathbf{X} \boldsymbol{\beta} \right) = 0,$$

$$\frac{\partial \log L(\theta | \mathbf{x})}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \left( \mathbf{y} - \mathbf{X} \boldsymbol{\beta} \right)^{\mathsf{T}} \left( \mathbf{y} - \mathbf{X} \boldsymbol{\beta} \right) = 0.$$

• Thus we obtain

$$\widehat{\boldsymbol{\beta}} = \left(\mathbf{X}^{\mathsf{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y},$$
  
$$\widehat{\sigma^{2}} = \frac{1}{n}\left(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\right)^{\mathsf{T}}\left(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\right).$$

• Example (Time Series): Consider the following autoregression

$$X_0 = x_0, \ X_n = lpha X_{n-1} + \sigma V_n$$
 where  $V_n \stackrel{\mathrm{i.i.d.}}{\sim} \mathcal{N}\left(0,1
ight)$ 

where  $\theta = (\alpha, \sigma^2)$ .

• We have

$$L(\theta | \mathbf{x}) = f(\mathbf{x} | \theta) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \alpha x_{i-1})^2}{2\sigma^2}\right)$$

• Thus we have

$$\log L(\theta | \mathbf{x}) = cst - \frac{n}{2} \log \sigma^2 - \sum_{i=1}^{n} \frac{(x_i - \alpha x_{i-1})^2}{2\sigma^2}$$

#### • It follows that

$$\frac{\partial 2 \log L(\theta | \mathbf{x})}{\partial \sigma^2} = -\frac{n}{\sigma^2} + \sum_{i=1}^n \frac{(x_i - \alpha x_{i-1})^2}{\sigma^4},$$
$$\frac{\partial 2 \log L(\theta | \mathbf{x})}{\partial \alpha} = \frac{2}{\sigma^2} \sum_{i=1}^n x_{i-1} (x_i - \alpha x_{i-1})^2.$$

• Thus we have

$$\widehat{\alpha} = \frac{\sum_{i=1}^{n} x_{i-1} x_{i}}{\sum_{i=1}^{n} x_{i-1}^{2}}, \ \widehat{\sigma^{2}} = \frac{\sum_{i=1}^{n} (x_{i} - \widehat{\alpha} x_{i-1})^{2}}{n}.$$

#### Invariance

Consider  $\eta = g(\theta)$ . We introduce the induced likelihood function  $L^*$ 

$$L^{*}\left( \left. \boldsymbol{\eta} \right| \mathbf{x} \right) = \sup_{\{\boldsymbol{\theta}: \boldsymbol{g}(\boldsymbol{\theta}) = \boldsymbol{\eta}\}} L\left( \left. \boldsymbol{\theta} \right| \mathbf{x} \right).$$

• Invariance property: If  $\hat{\theta}$  is the MLE of  $\theta$  then for any function  $\eta = g(\theta)$  then  $g(\hat{\theta})$  is the MLE of  $\eta$ .

"Proof": The MLE of  $\eta$  is defined by

$$\widehat{\eta} = \arg \sup_{\substack{\eta \quad \{\theta: g(\theta) = \eta\}}} L\left(\left.\theta\right| \mathbf{x}\right)$$

Define

$$g^{-1}(\eta) = \{\theta : g(\theta) = \eta\}.$$

Then clearly  $\widehat{\theta} \in g^{-1}(\widehat{\eta})$  and cannot be in any other preimage so  $\widehat{\eta} = g\left(\widehat{\theta}\right)$ .

A.D. ()

# Consistency

**Definition**. A sequence of estimators  $\widehat{\theta}_n = \widehat{\theta}_n(X_1, ..., X_n)$  is consistent for the parameter  $\theta$  if, for every  $\varepsilon > 0$  and every  $\theta \in \Theta$ 

$$\lim_{n\to\infty} P_{\theta}\left(\left|\widehat{\theta}_n - \theta\right| < \varepsilon\right) = 1 \text{ (equivalently } \lim_{n\to\infty} P_{\theta}\left(\left|\widehat{\theta}_n - \theta\right| \ge \varepsilon\right) = 0\text{)}.$$

• *Example*: Consider  $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, 1)$  and  $\widehat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$  then  $\widehat{\theta}_n \sim \mathcal{N}(\theta, 1/n)$  and

$$P_{\theta}\left(\left|\widehat{\theta}_{n}-\theta\right|<\varepsilon\right)=\int_{-\varepsilon\sqrt{n}}^{\varepsilon\sqrt{n}}\frac{1}{\sqrt{2\pi}}\exp\left(-\frac{u^{2}}{2}\right)du\rightarrow1.$$

It is possible to avoid this calculations and use instead Chebychev's inequality

$$P_{\theta}\left(\left|\widehat{\theta}_{n}-\theta\right| \geq \varepsilon\right) = P_{\theta}\left(\left|\widehat{\theta}_{n}-\theta\right|^{2} \geq \varepsilon^{2}\right) \leq \frac{\mathbb{E}_{\theta}\left(\left(\widehat{\theta}_{n}-\theta\right)^{2}\right)}{\varepsilon^{2}}$$
where  $\mathbb{E}_{\theta}\left(\left(\widehat{\theta}_{n}-\theta\right)^{2}\right) = var_{\theta}\left(\widehat{\theta}_{n}\right) + \left(\mathbb{E}_{\theta}\left(\widehat{\theta}_{n}-\theta\right)\right)^{2}$ .
A.D. ()
$$January 2008 = 13$$

63

• Example of *inconsistent* MLE (Fisher)

$$(X_i, Y_i) \sim \mathcal{N}\left(\left(\begin{array}{cc} \mu_i \\ \mu_i \end{array}\right), \left(\begin{array}{cc} \sigma^2 & 0 \\ 0 & \sigma^2 \end{array}\right)\right).$$

• The likelihood function is given by

$$L(\theta) = \frac{1}{(2\pi\sigma^2)^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n \left[ (x_i - \mu_i)^2 + (y_i - \mu_i)^2 \right] \right)$$

• We obtain

$$I(\theta) = cste - n \log \sigma^{2} \\ -\frac{1}{2\sigma^{2}} \left[ 2 \sum_{i=1}^{n} \left( \frac{x_{i} + y_{i}}{2} - \mu_{i} \right)^{2} + \frac{1}{2} \sum_{i=1}^{n} (x_{i} - y_{i})^{2} \right].$$

• We have

$$\widehat{\mu}_i = rac{x_i + y_i}{2}, \ \widehat{\sigma^2} = rac{\sum_{i=1}^n (x_i - y_i)^2}{4n} \to rac{\sigma^2}{2}.$$

# Consistency of the MLE

• Kullback-Leibler Distance: For any density f, g

$$D(f,g) = \int f(x) \log\left(\frac{f(x)}{g(x)}\right) dx$$

We have

$$D(f,g) \geq 0$$
 and  $D(f,f) = 0$ .

Indeed

$$-D(f,g) = \int f(x) \log\left(\frac{g(x)}{f(x)}\right) dx \le \int f(x) \left(\frac{g(x)}{f(x)} - 1\right) dx = 0$$

D (f, g) is a very useful 'distance' and appears in many different contexts.

# Alternative measures of similarity

• Hellinger distance

$$D(f,g) = \int \left(\sqrt{f(x)} - \sqrt{g(x)}\right)^2 dx$$

Generalized information

$$D(f,g) = \frac{1}{\lambda} \int \left( \left( \frac{f(x)}{g(x)} \right)^{\lambda} - 1 \right) f(x) \, dx$$

• L1-norm / Total variation

$$D(f,g) = \int |f(x) - g(x)| \, dx$$

• L2-norm / Total variation

$$D(f,g) = \int \left(f(x) - g(x)\right)^2 dx$$

### Example

• Suppose we have  $f(x) = \mathcal{N}(x; \xi, \tau^2)$  and  $g(x) = \mathcal{N}(x; \mu, \sigma^2)$ . • We have

$$\mathbb{E}_{f}\left[ (X-\mu)^{2} \right] = \mathbb{E}_{f}\left[ (X-\xi)^{2} + 2(X-\xi)(\xi-\mu) + (\xi-\mu)^{2} \right] \\ = \tau^{2} + (\xi-\mu)^{2}$$

• So it follows that

$$\mathbb{E}_{f} \left[ \log g \left( X \right) \right] = \mathbb{E}_{f} \left[ -\frac{1}{2} \log \left( 2\pi\sigma^{2} \right) - \frac{\left( X - \mu \right)^{2}}{2\sigma^{2}} \right]$$
$$= -\frac{1}{2} \log \left( 2\pi\sigma^{2} \right) - \frac{\tau^{2} + \left( \xi - \mu \right)^{2}}{2\sigma^{2}}$$

and

$$\mathbb{E}_{f}\left[\log f\left(X\right)\right] = -\frac{1}{2}\log\left(2\pi\tau^{2}\right) - \frac{1}{2}.$$

• It follows that

$$D(f,g) = \int f(x) \log\left(\frac{f(x)}{g(x)}\right) dx$$
  
=  $\mathbb{E}_f [\log f(X)] - \mathbb{E}_f [\log g(X)]$   
=  $\frac{1}{2} \left\{ \log\left(\frac{\sigma^2}{\tau^2}\right) + \frac{\tau^2 + (\xi - \mu)^2}{\sigma^2} - 1 \right\}$ 

• It can be easily checked that D(f, f) = 0 (less easy to show  $D(f, g) \ge 0$ ).

### Example

• Assume we have  $f(x) = \frac{1}{2} \exp(-|x|)$  and  $g(x) = \mathcal{N}(x; \mu, \sigma^2)$ . • We obtain

$$\mathbb{E}_{f} \left[ \log f(X) \right] = -\log 2 - \frac{1}{2} \int_{-\infty}^{\infty} |x| \exp(-|x|) dx$$
  
=  $-\log 2 - \int_{0}^{\infty} x \exp(-|x|) dx$   
=  $-\log 2 - 1$ ,

$$\mathbb{E}_{f}\left[\log g\left(X\right)\right] = -\frac{1}{2}\log\left(2\pi\sigma^{2}\right) - \frac{1}{4\sigma^{2}}\left(4 + 2\mu^{2}\right)$$

It follows that

$$D\left(f,g
ight)=rac{1}{2}\log\left(2\pi\sigma^{2}
ight)+rac{1}{2\sigma^{2}}\left(2+\mu^{2}
ight)-\log2-1.$$

- Assume the pdfs  $f(x|\theta)$  have common support for all  $\theta$  and  $f(x|\theta) \neq f(x|\theta')$  for  $\theta \neq \theta'$ ; i.e.  $S_{\theta} = \{x : f(x|\theta) > 0\}$  is independent of  $\theta$ .
- Denote

$$M_{n}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \log \frac{f(X_{i}|\theta)}{f(X_{i}|\theta_{*})}$$

- As the MLE  $\hat{\theta}_{n}$  maximises  $L(\theta | \mathbf{x})$ , it also maximizes  $M_{n}(\theta)$ .
- Assume  $X_i \stackrel{\text{i.i.d.}}{\sim} f(x | \theta_*)$ . Note that by the law of large numbers  $M_n(\theta)$  converges to

$$\mathbb{E}_{\theta_*}\left(\log\frac{f\left(X|\theta\right)}{f\left(X|\theta_*\right)}\right) = \int f\left(x|\theta_*\right)\log\frac{f\left(x|\theta\right)}{f\left(x|\theta_*\right)}dx$$
$$= -D\left(f\left(\cdot|\theta_*\right), f\left(\cdot|\theta\right)\right) := M\left(\theta\right).$$

• Hence,  $M_n(\theta) \approx -D(f(\cdot|\theta_*), f(\cdot|\theta))$  which is maximized for  $\theta^*$  so we expect that its maximizer will converge towards  $\theta_*$ .

# Example

- Assume  $f(x) = g(x) = \mathcal{N}(x; 0, 1)$ .
- We approximate

$$\mathbb{E}_{f}\left[\log g\left(X\right)\right] = -\frac{1}{2}\log\left(2\pi\right) - \frac{1}{2} = -1.4189$$

through

$$\mathbb{E}_{\widehat{f}}\left[\log g\left(X\right)\right] = -\frac{1}{2}\log\left(2\pi\right) - \frac{1}{2n}\sum_{i=1}^{n}X_{i}^{2}$$

Numerical examples

n	10	100	1,000	10,000	$\mathbb{E}_{f}\left[\log g\left(X ight) ight]$
Mean	-1.4188	-1.4185	-1.4191	-1.4189	-1.4189
Variance	0.05079	0.00497	0.00050	0.00005	-
Standard deviation	0.22537	0.07056	0.02232	0.00696	-

Mean, variance and standard deviation by running 1,000 Monte Carlo trials

Theorem. Suppose

$$\sup_{\theta\in\Theta}\left|M_{n}\left(\theta\right)-M\left(\theta\right)\right|\xrightarrow{\mathsf{P}}0$$

and that, for every  $\varepsilon > 0$ ,

then

$$\widehat{\theta}_n \xrightarrow{\mathsf{P}} \theta_*$$

*Proof.* Since  $\widehat{\theta}_n$  maximizes  $M_n(\theta)$ , we have  $M_n(\widehat{\theta}_n) \ge M_n(\theta_*)$ . Thus,

Thus it implies that for any  $\delta > 0$ , we have

$$\Pr\left(M\left(\widehat{\theta}_{n}\right) < M\left(\theta_{*}\right) - \delta\right) \to 0.$$

Now for any  $\varepsilon > 0$ , there exists  $\delta > 0$  such that  $|\theta - \theta^*| \ge \varepsilon$  implies  $M(\theta) < M(\theta_*) - \delta$ . Hence,

$$\Pr\left(\left|\widehat{\theta}_{n}-\theta_{*}\right|>\varepsilon\right)\leq\Pr\left(M\left(\widehat{\theta}_{n}\right)< M\left(\theta_{*}\right)-\delta\right)\rightarrow0.$$

#### Asymptotic Normality

- Assuming we have  $\widehat{\theta}_n \xrightarrow{\mathsf{P}} \theta_*$ , what can we say about  $\sqrt{n} \left( \widehat{\theta}_n \theta_* \right)$ ?
- Lemma. Let  $s(x|\theta) := \frac{\partial \log f(x|\theta)}{\partial \theta}$  be the score function, then we have for any  $\theta$

 $\mathbb{E}_{\theta}\left[s\left(X|\theta\right)\right]=0.$ 

Proof. We have

$$\int \frac{\partial \log f(x|\theta)}{\partial \theta} f(x|\theta) dx$$

$$= \int \frac{\frac{\partial f(x|\theta)}{\partial \theta}}{f(x|\theta)} f(x|\theta) dx = \int \frac{\partial f(x|\theta)}{\partial \theta} dx$$

$$= \frac{\partial}{\partial \theta} \underbrace{\int f(x|\theta) dx}_{=1} = 0.$$

• Lemma. We also have

$$\mathit{var}_{\theta}\left[\mathit{s}\left(\left. X \right| \theta \right)\right] = \mathbb{E}_{\theta}\left[\mathit{s}\left(\left. X \right| \theta \right)^{2}\right] = -\mathbb{E}_{\theta}\left[\frac{\partial^{2} \log \mathit{f}\left(\left. X \right| \theta \right)}{\partial \theta^{2}}\right] := \mathit{I}\left(\theta\right)$$

• Proof. This follows from

$$\int \frac{\partial \log f(x|\theta)}{\partial \theta} f(x|\theta) \, dx = 0$$

thus by taking the derivative once more with respect to  $\boldsymbol{\theta}$ 

$$0 = \frac{\partial}{\partial \theta} \int \frac{\partial \log f(x|\theta)}{\partial \theta} f(x|\theta) dx$$
  
= 
$$\int \frac{\partial^2 \log f(x|\theta)}{\partial \theta^2} f(x|\theta) + \int \frac{\partial \log f(x|\theta)}{\partial \theta} \frac{\partial f(x|\theta)}{\partial \theta} dx$$

• Heuristic Derivation. We have for  $I(\theta) := \log L(\theta | \mathbf{x})$ 

$$0 = l'\left(\widehat{\theta}_{n}\right) \approx l'\left(\theta_{*}\right) + \left(\widehat{\theta}_{n} - \theta_{*}\right)l''\left(\theta_{*}\right)$$
$$\Rightarrow \left(\widehat{\theta}_{n} - \theta_{*}\right) = -\frac{l'\left(\theta_{*}\right)}{l''\left(\theta_{*}\right)}$$

That is

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta_{*}\right)=\frac{\frac{1}{\sqrt{n}}l^{\prime}\left(\theta_{*}\right)}{-\frac{1}{n}l^{\prime\prime}\left(\theta_{*}\right)}$$

• Now remember that  $I'(\theta_*) = \sum_{i=1}^n s(X_i | \theta_*)$  where  $\mathbb{E}_{\theta_*}[s(X_i | \theta_*)] = 0$  and  $var_{\theta_*}[s(X_i | \theta_*)] = I(\theta_*)$  so the CLT tells us that

$$\frac{1}{\sqrt{n}}I'(\theta_*) \xrightarrow{\mathsf{D}} \mathcal{N}(\mathsf{0}, I(\theta_*))$$

Now the law of large number yields

$$-\frac{1}{n}I''\left(\theta_{*}\right)\xrightarrow{\mathsf{P}}I\left(\theta_{*}\right)$$

so by Slutsky's theorem

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

- Note that you have already seen this expression when establishing the Cramer-Rao bound.
- It is important to remember that depending on θ<sub>\*</sub> the parameter can be more or less easy to estimate.

• Similarly, we can prove that

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

• We can also prove that

$$\sqrt{n}\left|g'\left(\widehat{\theta}_{n}\right)\right|\sqrt{I\left(\widehat{\theta}_{n}\right)}\left(g\left(\widehat{\theta}_{n}\right)-g\left(\theta_{*}\right)\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(0,1\right).$$

• This allows us to derive some confidence intervals.

# Making the proof more rigourous

We have

$$I'\left(\widehat{\theta}_{n}\right) = I'\left(\theta_{*}\right) + \left(\widehat{\theta}_{n} - \theta_{*}\right)I''\left(\theta_{*}\right) + \frac{1}{2}\left(\widehat{\theta}_{n} - \theta_{*}\right)^{2}I'''\left(\theta_{*,n}\right)$$

where  $\theta_{*,n}$  lies between  $\widehat{\theta}_n$  and  $\theta_*$  so that

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta_{*}\right)=\frac{\frac{1}{\sqrt{n}}I^{\prime\prime}\left(\theta_{*}\right)}{-\frac{1}{n}I^{\prime\prime}\left(\theta_{*}\right)-\frac{1}{2n}\left(\widehat{\theta}_{n}-\theta_{*}\right)^{2}I^{\prime\prime\prime}\left(\theta_{*,n}\right)}$$

• To proof the result, we need to check that  $\frac{1}{2n} \left(\widehat{\theta}_n - \theta_*\right)^2 I'''(\theta_{*,n}) \xrightarrow{P} 0. \text{ As } \widehat{\theta}_n \xrightarrow{P} \theta_*, \text{ we just need to prove that } \frac{1}{n} I'''(\theta_{*,n}) \text{ is bounded (in probability). So we need an additional condition of the form say for any <math>\theta$ 

$$\left|\frac{\partial^{3}\log f\left(\left.x\right|\theta\right)}{\partial\theta^{3}}\right| \leq C\left(x\right)$$

with  $\mathbb{E}_{\theta_*} \left[ C(X) \right] < \infty$ .

# Multiparameter Case

• The extension to the multiparameter case  $\theta = (\theta_1, ..., \theta_d)$  is straightforward

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta_{*}\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(0,J\left(\theta_{*}\right)\right)$$

where  $J\left( heta_{*}
ight)=I\left( heta_{*}
ight)^{-1}$  where

$$\left[I\left(\theta_{*}\right)\right]_{k,l} = -\mathbb{E}_{\theta_{*}}\left[\frac{\partial^{2}\log f\left(x|\theta\right)}{\partial\theta_{k}\partial\theta_{l}}\right]$$

• We define  $\nabla g := \left(\frac{\partial g}{\partial \theta_1}, \dots, \frac{\partial g}{\partial \theta_d}\right)^{\mathsf{T}}$  then if  $\nabla g\left(\theta_*\right) \neq 0$  $\sqrt{n}\left(g\left(\widehat{\theta}_n\right) - g\left(\theta_*\right)\right) \xrightarrow{\mathsf{D}} \mathcal{N}\left(0, \nabla g\left(\theta_*\right) J\left(\theta_*\right) \nabla g\left(\theta_*\right)^{\mathsf{T}}\right)$  • Example: If  $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$  with  $\theta = (\mu, \sigma)$  then  $\log f(x|\theta) = cst - \log \sigma - \frac{(x-\mu)^2}{2\sigma^2}$   $s(x|\theta) = \begin{pmatrix} \frac{(x-\mu)}{\sigma^2} \\ -\frac{1}{\sigma} + \frac{(x-\mu)^2}{\sigma^3} \end{pmatrix}$ ,  $I(\theta) = \begin{pmatrix} \frac{1}{\sigma^2} & \mathbb{E}_{\theta} \left[ \frac{2(x-\mu)}{\sigma^3} \right] \\ \mathbb{E}_{\theta} \left[ \frac{2(x-\mu)}{\sigma^3} \right] & \mathbb{E}_{\theta} \left( -\frac{1}{\sigma^2} + \frac{3(x-\mu)^2}{\sigma^4} \right) \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{pmatrix}$ 

• The MLE of  $\mu$  is given by

$$\widehat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i \Rightarrow var\left[\widehat{\mu}\right] = \frac{\sigma^2}{n}$$

and the MLE is indeed efficient (it reaches Cramer-Rao lower bound).

# • Assume we observe a vector $X=(X_1,...,X_k)$ where $X_j\in\{0,1\}$ , $\sum_{j=1}^k X_j=1$ with

$$f(x|p_1,...,p_{k-1}) = \left(\prod_{j=1}^{k-1} p_j^{x_j}\right) \left(1 - \sum_{j=1}^{k-1} p_j\right)^{x_k}$$

where  $p_j > 0$  and  $p_k := 1 - \sum_{j=1}^{k-1} p_j < 1$ . We have  $\theta = (p_1, ..., p_{k-1})$ .

We have

$$\frac{\partial \log f(x|\theta)}{\partial p_j} = \frac{x_j}{p_j} - \frac{x_k}{p_k},$$
  
$$\frac{\partial^2 \log f(x|\theta)}{\partial p_j^2} = -\frac{x_j}{p_j^2} - \frac{x_k}{p_k^2},$$
  
$$\frac{\partial^2 \log f(x|\theta)}{\partial p_j \partial p_l} = -\frac{x_k}{p_k^2}, j \neq l < k.$$

• Recall that  $X_j$  has a Bernoulli distribution with mean  $p_j$  so

$$I(\theta) = \begin{bmatrix} p_1^{-1} + p_k^{-1} & p_k^{-1} & p_k^{-1} \\ p_k^{-1} & p_2^{-1} + p_k^{-1} & \vdots \\ \vdots & p_k^{-1} & \vdots \\ \vdots & \vdots & p_{k-1}^{-1} + p_k^{-1} \end{bmatrix}$$

• One can check that

$$I(\theta)^{-1} = \begin{bmatrix} p_1(1-p_1) & -p_1p_2 & \cdots & -p_1p_{k-1} \\ -p_1p_2 & p_2(1-p_2) & & -p_2p_{k-1} \\ \vdots & \vdots & & \vdots \\ -p_1p_{k-1} & -p_2p_{k-1} & \cdots & p_{k-1}(1-p_{k-1}) \end{bmatrix}$$

• Now assume we observe  $X^1, X^2, ..., X^n$  then

$$\log L\left(\left.\theta\right|\mathbf{x}\right) = \sum_{j=1}^{k} t_j \log p_j = \sum_{j=1}^{k-1} t_j \log p_j + t_k \log \left(1 - \sum_{j=1}^{k-1} p_j\right)$$

where 
$$t_j = \sum_{i=1}^n x_j^i$$
.

So we have

$$\frac{\partial \log L\left(\left.\theta\right|\mathbf{x}\right)}{\partial p_{j}} = \frac{t_{j}}{p_{j}} - \frac{t_{k}}{p_{k}} \text{ for } j = 1, ..., k - 1 \Rightarrow \widehat{p}_{j} = \frac{t_{j}}{n}.$$

• Clearly  $t_j$  is Binomial $(n, p_j)$  with variance  $np_j (1 - p_j)$  so  $\hat{p}_j$  is efficient.

- Assume θ = (θ<sub>1</sub>, ..., θ<sub>d</sub>) is the parameter vector but only the scalar θ<sub>1</sub> is of interest whereas (θ<sub>2</sub>, ..., θ<sub>d</sub>) are nuisance parameters.
- We want to assess how the asymptotic precision with which we estimate θ<sub>1</sub> is influenced by the presence of nuisance parameters; i.e. if θ̂ is an efficient estimate for θ, then how does θ̂<sub>1</sub> as an estimator of θ<sub>1</sub> compare to an efficient estimation of θ<sub>1</sub>, say θ̂<sub>1</sub>, which would assume that all the nuisance parameters are known.
- Intuitively, we should have  $var\left[\widetilde{\theta}_{1}\right] \leq var\left[\widehat{\theta}_{1}\right]$ ; i.e. ignorance cannot bring you any advantage.

- Asymptotic variance of  $\sqrt{n} \left( \widehat{\theta}_n \theta_* \right)$  is  $I^{-1} \left( \theta_* \right)$  whose (i, j) parameter is denoted  $\alpha_{i,j}$ .
- Asymptotic variance of  $\sqrt{n} \left( \widetilde{\theta}_1 \theta_{*,1} \right)$  is  $I^{-1} \left( \theta_{*,1} \right) = 1/\gamma_{1,1}$ .
- **Theorem**. We have  $\alpha_{1,1} \ge 1/\gamma_{1,1}$ , with equality if and only if  $\alpha_{1,2} = \cdots = \alpha_{1,d} = 0$ .

• Partition  $I(\theta)$  as follows

$$I(\theta) = \left( egin{array}{cc} \gamma_{1,1} & \rho^{\mathsf{T}} \\ \rho & \Sigma \end{array} 
ight).$$

Now we use the fact that

$$I^{-1}(\theta) = \frac{1}{\tau} \left( \begin{array}{cc} 1 & -\rho^{\mathsf{T}} \Sigma^{-1} \\ -\Sigma^{-1} \rho & \Sigma^{-1} \rho \rho^{\mathsf{T}} \Sigma^{-1} + \tau \Sigma^{-1} \end{array} \right)$$

where  $\tau = \gamma_{1,1} - \rho^{\mathsf{T}} \Sigma^{-1} \rho$ .

• As I  $(\theta)$  is definite positive then  $\Sigma^{-1}$  is definite positive and

$$\alpha_{1,1} = \frac{1}{\tau} \geq 1/\gamma_{1,1}$$

with equality iff  $\rho = 0$ .

• To show that  $\tau > 0$  we use the fact that  $I(\theta)$  is p.d. and that

$$au = v^{\mathsf{T}} I\left( heta
ight) v$$
 where  $v = \left(egin{array}{c} 1 \ -
ho^{\mathsf{T}} \Sigma^{-1} \end{array}
ight)$ 

# Beyond Maximum Likelihood: Method of Moments

- MLE estimates can be difficult to compute, the method of moments is a simple alternative. The obtained estimators are typically not optimal but can be used as starting values for more sophisticated methods.
- For  $1 \leq j \leq d$ , define the  $j^{\text{th}}$  moment of  $f(x|\theta)$  where  $\theta = (\theta_1, ..., \theta_d)$

$$\alpha_{j}(\theta) = \mathbb{E}_{\theta}\left[X^{j}\right] = \int x^{j}f(x|\theta) dx$$

and, given  $old X=(X_1,...,X_n)$  , the  $j^{ ext{th}}$  sample as

$$\widehat{\alpha}_j = \frac{1}{n} \sum_{i=1}^n X_i^j.$$

 The idea of the method of moments method is to match the theoretical moments to the sample moments; that is we defined θ<sub>n</sub> as the value of θ such that

$$lpha_j\left(\widehat{ heta}_n
ight)=\widehat{lpha}_j ext{ for } j=1,...,d.$$

• Example: Let 
$$X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$$
 with  $\theta = (\mu, \sigma^2)$  then  
 $\alpha_1(\theta) = \mu, \ \alpha_2(\theta) = \sigma^2 + \mu^2,$   
 $\widehat{\alpha}_1 = \frac{1}{n} \sum_{i=1}^n X_i^1, \ \widehat{\alpha}_2 = \frac{1}{n} \sum_{i=1}^n X_i^2$ 

• Thus we obtain

$$\widehat{\mu} = \widehat{lpha}_1$$
 and  $\widehat{\sigma^2} = \widehat{lpha}_2 - \left(\widehat{lpha}_1
ight)^2$ .

• Note that  $\widehat{\sigma^2}$  is not unbiased.

• Assume  $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}(\theta_1, \theta_2)$  where  $-\infty < \theta_1 < \theta_2 < +\infty$  then

$$\alpha_1(\theta) = \frac{\theta_1 + \theta_2}{2}, \ \alpha_2(\theta) = \frac{\theta_1^2 + \theta_2^2 + \theta_1 \theta_2}{3}$$

Now we solve and obtain

$$\begin{aligned} \theta_1 &= 2\widehat{\alpha}_1 - \theta_2, \\ 3\widehat{\alpha}_2 &= (2\widehat{\alpha}_1 - \theta_2)^2 + \theta_2^2 + (2\widehat{\alpha}_1 - \theta_2) \theta_2 \Leftrightarrow (\theta_2 - \widehat{\alpha}_1)^2 = 3\left(\widehat{\alpha}_2 - \widehat{\alpha}_1^2\right) \end{aligned}$$

Since  $\theta_2 > \mathbb{E}(X)$  then

$$\widehat{ heta}_2 = \widehat{lpha}_1 + \sqrt{3\left(\widehat{lpha}_2 - \widehat{lpha}_1^2
ight)}, \ \widehat{ heta}_1 = \widehat{lpha}_1 - \sqrt{3\left(\widehat{lpha}_2 - \widehat{lpha}_1^2
ight)}.$$

• Note that  $(\hat{\theta}_1, \hat{\theta}_2)$  is NOT a function of the sufficient statistics  $X_{(1)}, X_{(n)}$ .

Assume X<sub>i</sub> <sup>i.i.d.</sup> ∼ Bi (p, k) with parameters k ∈ N and p ∈ (0, 1).
We have

$$\alpha_1(\theta) = kp, \ \alpha_2(\theta) = kp(1-p) + k^2p^2.$$

Thus we obtain

$$\widehat{\boldsymbol{\rho}} = \left( \widehat{\alpha}_1 + \widehat{\alpha}_1^2 - \widehat{\alpha}_2 \right) / \widehat{\alpha}_1,$$

$$\widehat{\boldsymbol{k}} = \widehat{\alpha}_1^2 / \left( \widehat{\alpha}_1 + \widehat{\alpha}_1^2 - \widehat{\alpha}_2 \right).$$

• The estimator  $\widehat{p} \in (0, 1)$  but  $\widehat{k}$  is generally not an integer.

#### Statistical Properties of the Estimate

• Let  $\widehat{\alpha} = (\widehat{\alpha}_1, ..., \widehat{\alpha}_d)$ , we have

$$\widehat{\alpha} = h\left(\widehat{\theta}\right)$$

and if the inverse function  $g = h^{-1}$  exists, then

$$\widehat{ heta}=oldsymbol{g}\left(\widehat{lpha}
ight)$$
 .

- If g is continuous at  $\alpha = (\alpha_1, ..., \alpha_d)$  then  $\widehat{\theta}$  is a consistent estimate of  $\theta$  as  $\widehat{\alpha}_j \rightarrow \alpha_j$ .
- Moreover if g is differentiable at lpha and  $\mathbb{E}_{ heta}\left[X^{2d}
  ight]<\infty$  then

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta\right)\overset{\mathrm{D}}{\rightarrow}\mathcal{N}\left(0,\nabla g\left(\alpha\right)^{\mathsf{T}}V_{\alpha}\nabla g\left(\alpha\right)\right)$$

where

$$V_{\alpha}[i,j] = \alpha_{i+j} - \alpha_i \alpha_j.$$

• The result follows from

$$\sqrt{n}\left(\widehat{\alpha}_{n}-\alpha\right)\xrightarrow{\mathrm{D}}\mathcal{N}\left(0,\,V_{\alpha}\right)$$

as

$$\mathbb{E}_{\theta}\left[\widehat{\alpha}_{j,n}\right] = \alpha_{j}, \ cov\left[\widehat{\alpha}_{i,n}\widehat{\alpha}_{j,n}\right] = \frac{\alpha_{i+j} - \alpha_{i}\alpha_{j}}{n}$$

• We have

$$\widehat{\theta}_{n}-\theta=g\left(\widehat{\alpha}_{n}\right)-g\left(\alpha_{n}\right)$$

so using the delta method

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(0,\nabla g\left(\alpha\right)^{\mathsf{T}}V_{\alpha}\nabla g\left(\alpha\right)\right)$$

• We can also establish the  $n^{-1}$  order asymptotic bias of the estimate as

$$g(\widehat{\alpha}_{n}) = g(\alpha) + \nabla g(\alpha)^{\mathsf{T}} (\widehat{\alpha}_{n} - \alpha) + \frac{1}{2} (\widehat{\alpha}_{n} - \alpha)^{\mathsf{T}} \nabla^{2} g(\alpha) (\widehat{\alpha}_{n} - \alpha) + o\left(\frac{1}{n}\right)$$

where  $\sqrt{n} \left( \widehat{\alpha}_n - \alpha \right) \xrightarrow{\mathsf{D}} Z_{\Sigma}$  with  $Z_{\Sigma} \sim \mathcal{N} \left( 0, \Sigma \right)$  so

$$n\left(\widehat{\alpha}-\alpha\right)^{\mathsf{T}}\nabla^{2}g\left(\alpha\right)\left(\widehat{\alpha}-\alpha\right)\xrightarrow{\mathsf{D}}Z_{\Sigma}^{\mathsf{T}}\nabla^{2}g\left(\alpha\right)Z_{\Sigma}$$

as recall that  $X_n \xrightarrow{\mathsf{D}} X$  implies  $\varphi(X_n) \xrightarrow{\mathsf{D}} \varphi(X)$ .

Thus we have

$$\mathbb{E}\left[g\left(\widehat{\alpha}_{n}\right)\right] = g\left(\alpha\right) + \frac{1}{2n}\mathbb{E}\left[Z_{\Sigma}^{\mathsf{T}}\nabla^{2}g\left(\alpha\right)Z_{\Sigma}\right] + o\left(\frac{1}{n}\right)$$
$$= g\left(\alpha\right) + \frac{tr\left(\nabla^{2}g\left(\alpha\right)\Sigma\right)}{2n} + o\left(\frac{1}{n}\right).$$

# Beyond Maximum Likelihood: Pseudo-Likelihood

- Assume  $X = (X_1, ..., X_q) \sim f(x|\theta)$ . Given *n* observations  $X^i \sim f(x|\theta)$ , the MLE requires maximizing  $L(\theta|\mathbf{x})$ .
- However in some problems, it might be difficult to specify  $f(x|\theta)$  and we may be only able to specify say

$$f(x_k, x_l | \theta)$$
 for  $1 \le k < l \le q$ 

• Based on this information and *n* observations, we could define the pseudo-log-likelihood functions

$$l_{1}(\theta | \mathbf{x}) = \sum_{i=1}^{n} l_{1}(\theta | x^{i}) = \sum_{i=1}^{n} \sum_{s=1}^{q} \log f(x_{s} | \theta),$$
  
$$l_{2}(\theta | \mathbf{x}) = \sum_{i=1}^{n} l_{2}(\theta | x^{i}) = \sum_{i=1}^{n} \sum_{s=1}^{q} \sum_{t=s+1}^{q} \log f(x_{s}, x_{t} | \theta) + \alpha l_{1}(\theta | \mathbf{x}).$$

• These pseudo-likelihood functions are simpler that the full likelihood.

 Under regularity conditions very similar to the ones for the MLE, solving

$${\it l}_k^\prime \left( \left. heta 
ight| {f x} 
ight) = 0$$
 for  $k=1,2$ 

will provide unbiased estimates.

• To derive the asymptotic variance, we use

$$0 = l'_{k}\left(\widehat{\theta}_{n}\right) \approx l'_{k}\left(\theta_{*}\right) + \left(\widehat{\theta}_{n} - \theta_{*}\right)l''_{k}\left(\theta_{*}\right)$$
$$\Rightarrow \sqrt{n}\left(\widehat{\theta}_{n} - \theta_{*}\right) = -\frac{\frac{1}{\sqrt{n}}l'_{k}\left(\theta_{*}\right)}{\frac{1}{n}l''_{k}\left(\theta_{*}\right)}$$

where 
$$\frac{1}{n}I_k''(\theta_*) \xrightarrow{\mathsf{P}} \mathbb{E}_{\theta_*}[I_k'']$$
 and  $\frac{1}{\sqrt{n}}I'(\theta_*) \xrightarrow{\mathsf{D}} \mathcal{N}\left(0, \mathbb{E}_{\theta_*}\left[I_k'^2\right]\right)$ , thus  $\sqrt{n}\left(\widehat{\theta}_n - \theta_*\right) \xrightarrow{\mathsf{D}} \mathcal{N}\left(0, \mathbb{E}_{\theta_*}\left[-I_k''\right]^{-2} \mathbb{E}_{\theta_*}\left[I_k'^2\right]\right)$ .

• We have the estimates

$$\mathbb{E}_{\theta_*}\left[I_k''\right] \approx \frac{1}{n} \sum_{i=1}^n I_k''(\theta | x_i), \ \mathbb{E}_{\theta_*}\left[I_k'^2\right] \approx \frac{1}{n} \sum_{i=1}^n I_k'^2(\theta | x_i).$$

- Example. Assume that  $X = (X_1, ..., X_q) \sim \mathcal{N}(0, \Sigma)$  where  $[\Sigma](i, j) = 1$  if i = j and  $\rho$  otherwise. We are interested in estimating  $\theta = \rho$ .
- There is no information about  $\rho$  in  $I_1(\theta | \mathbf{x})$  so we use  $I_2(\theta | \mathbf{x})$  for  $\alpha = 0$ . For *n* observations  $(X^1, ..., X^n)$ , we have

$$l_{2}(\theta | \mathbf{x}) = -\frac{nq(q-1)}{4} \log(1-\rho^{2}) - \frac{q-1+\rho}{2(1-\rho^{2})} SS_{W} - \frac{(q-1)(1-\varrho)}{2(1-\rho^{2})} \frac{SS_{B}}{q}$$

where

$$SS_W = \sum_{i=1}^n \sum_{s=1}^q \left( X_s^i - \left( \sum_{t=1}^q X_t^i \right) \right)^2$$
,  $SS_B = \sum_{i=1}^n X_t^{i^2}$ 

• After simple but tiedous calculations, we obtain for the asymptotic variance

$$\frac{2}{nq(q-1)} \frac{(1-\rho^{2}) c(q,\rho)}{(1+\rho^{2})^{2}}$$

where

$$c(q,\rho) = (1-\rho^2)(1+3\rho^2) + q\rho(-3\rho^3+8\rho^2-3\rho+2) +q^2\rho^2(1-\rho^2)$$

whereas for MLE we have

$$\frac{2}{nq(q-1)} \frac{\left(1 + (q-1)\rho\right)^2 \left(1 - \rho\right)^2}{1 + (q-1)\rho^2}$$

• The ratio is 1 for q = 2 as expected and also 1 if  $\rho = 0$  or 1. For any other values, there is a loss of efficiency for  $l_2(\theta | \mathbf{x})$  which increases as  $q \to \infty$ .

Consider the following time series

$$X_0 \sim \mathcal{N}\left(0, \frac{\sigma^2}{1-\alpha^2}\right)$$
,  $X_n = \alpha X_{n-1} + \sigma V_n$  where  $V_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ 

where  $\theta = \sigma^2$ .

• We can show that we have for any i = 0, 1, ..., n

$$f(x_i | \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(1-\alpha^2) x_i^2}{2\sigma^2}\right)$$

and we consider

$$2I_{1}\left(\left.\theta\right|\mathbf{x}\right) = 2\sum_{i=1}^{n}\log f\left(\left.x_{i}\right|\theta\right) = cste - n\log\sigma^{2} - \frac{\left(1-\alpha^{2}\right)}{2\sigma^{2}}\sum_{i=1}^{n}x_{i}^{2}$$

This pseudo-likelihood can easily be maximized

$$\widehat{\sigma^2} = \frac{\left(1 - \alpha^2\right)\sum_{i=1}^n x_i^2}{n}.$$

• If one is interested in estimating  $\alpha$ , it would be necessary to introduce  $f(x_i, x_{i+1} | \theta)$ .

• Pseudo-likelihood is widely used for Markov random fields since its introduction by Besag (1975). In the Gaussian context, we have  $X = (X_1, ..., X_n)$  where d is extremely large and Gaussian and the model is specified by

$$\mathbb{E}_{ heta}\left[\left.X_{i}\right|x_{-i}
ight]=\lambda\sum_{i=1}^{n}H_{ij}x_{j}, \; \textit{var}_{ heta}\left[\left.X_{i}\right|x_{-i}
ight]=\kappa.$$

• Computing the likelihood for  $\theta = (\lambda, \kappa)$  can be too computatonally intensive so the pseudo-likelihood is defined through

$$\widetilde{I}(\theta | \mathbf{x}) = \sum_{i=1}^{n} \log f(x_i | \theta, x_{-i})$$

thus

$$\widehat{\lambda} = \frac{\mathbf{x}^{\mathsf{T}} \mathbf{H} \mathbf{x}}{\mathbf{x}^{\mathsf{T}} \mathbf{H}^{2} \mathbf{x}}, \ \kappa = d^{-1} \left( \mathbf{x}^{\mathsf{T}} \mathbf{x} - \frac{\left( \mathbf{x}^{\mathsf{T}} \mathbf{H} \mathbf{x} \right)^{2}}{\mathbf{x}^{\mathsf{T}} \mathbf{H}^{2} \mathbf{x}} \right)$$

 In this context, it can be show that the estimate is consistent and has a reasonable asymptotic variance.

- In many applications, the log-likelihood *l*(θ; y<sub>1:n</sub>) is very complex to compute.
- Instead we maximize a surrogate function  $I_S(\theta; y_{1:n})$ .
- If possible, we pick this function such that if  $\theta^*$  is the 'true' parameter then

$$\mathsf{E}_{\theta^*}\left(I_{\mathcal{S}}\left(\theta; Y_{1:n}\right)\right)$$

is maximized for  $heta= heta^*$  and solving

$$abla I_{\mathcal{S}}\left(\widehat{\theta}_{n};Y_{1:n}\right)=0$$

is 'easy'.

• Under regularity assumptions, we have

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta^{*}\right)\Rightarrow\mathcal{N}\left(0,\,G_{n}^{-1}\left(\theta^{*}\right)\right)$$

where

$$G_{n}^{-1}\left(\theta\right)=H_{n}^{-1}\left(\theta\right)J_{n}\left(\theta\right)H_{n}^{-\mathsf{T}}\left(\theta\right)$$

with

$$J_{n}(\theta) = \mathbb{V} \{ \nabla I_{S}(\theta; Y_{1:n}) \}, H_{n}(\theta) = \mathbb{E} \{ \nabla^{2} I_{S}(\theta; Y_{1:n}) \}.$$

- When  $I_S(\theta; Y_{1:n}) = I(\theta; Y_{1:n})$  and the model is correctly specified then  $G_n(\theta)$  is the Fisher information matrix.
- When  $I_{S}(\theta; Y_{1:n}) \neq I(\theta; Y_{1:n})$ , we typically lose in terms of efficiency.

# Application to General State-Space Models

 Consider the following general state-space model. Let {X<sub>k</sub>}<sub>k≥1</sub> be an Markov process defined by

$$X_1 \sim \mu_{ heta}$$
 and  $X_k | (X_{k-1} = x_{k-1}) \sim f_{ heta} (\cdot | x_{k-1})$ .

• Then we have that for any n > 0

$$p_{\theta}(x_{1:n}) = p_{\theta}(x_{1}) \prod_{k=2}^{n} p_{\theta}(x_{k} | x_{1:k-1})$$
$$= \mu_{\theta}(x_{1}) \prod_{k=2}^{n} f_{\theta}(x_{k} | x_{k-1}).$$

• We are interested in estimating  $\theta$  from the data but we do not have access to  $\{X_k\}_{k\geq 1}$ . We only have access to a process  $\{Y_k\}_{k\geq 1}$  such that, conditional upon  $\{X_k\}_{k\geq 1}$ , the observations are statistically independent and

$$Y_k | (X_k = x_k) \sim g_\theta (\cdot | x_k).$$

That is we have for any n > 0

$$p_{\theta}(y_{1:n}|x_{1:n}) = \prod_{k=1}^{n} p_{\theta}(y_{k}|x_{k}) = \prod_{k=1}^{n} g_{\theta}(y_{k}|x_{k}).$$

• Linear Gaussian model. Consider say for |lpha| < 1

$$\begin{array}{lll} X_{1} & \sim & \mathcal{N}\left(0, \frac{\sigma^{2}}{1-\alpha^{2}}\right), \ X_{k} = \alpha X_{k-1} + \sigma V_{k} \ \text{where} \ V_{k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right), \\ Y_{k} & = & \beta + X_{k} + \tau W_{k} \ \text{where} \ W_{k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right). \end{array}$$

ullet In this case we have say  $heta=\left(eta,\sigma^2,lpha, au^2
ight)$  and

$$\begin{aligned} f_{\theta}\left(\left.x_{k}\right|x_{k-1}\right) &= \mathcal{N}\left(x_{k};\alpha x_{k-1},\sigma^{2}\right), \\ g_{\theta}\left(\left.y_{k}\right|x_{k}\right) &= \mathcal{N}\left(y_{k};\beta+x_{k},\tau^{2}\right). \end{aligned}$$

• Stochastic Volatility model. Consider say for  $|\alpha| < 1$ 

$$\begin{array}{lll} X_{1} & \sim & \mathcal{N}\left(0, \frac{\sigma^{2}}{1-\alpha^{2}}\right), \ X_{k} = \alpha X_{k-1} + \sigma V_{k} \ \text{where} \ V_{k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right), \\ Y_{k} & = & \beta \exp\left(X_{k}/2\right) W_{k} \ \text{where} \ W_{k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right). \end{array}$$

ullet In this case we have say  $heta=\left(eta,\sigma^2,lpha
ight)$  and

$$\begin{array}{rcl} f_{\theta}\left(\left.x_{k}\right|x_{k-1}\right) & = & \mathcal{N}\left(x_{k};\alpha x_{k-1},\sigma^{2}\right), \\ g_{\theta}\left(\left.y_{k}\right|x_{k}\right) & = & \mathcal{N}\left(y_{k};\mathbf{0},\beta^{2}\exp\left(x_{k}\right)\right). \end{array}$$

• In this case, the likelihood of the observations  $y_{1:n}$  is given by

$$p_{\theta}(y_{1:n}) = \int p_{\theta}(x_{1:n}, y_{1:n}) dx_{1:n}$$

$$= \int p_{\theta}(y_{1:n} | x_{1:n}) p_{\theta}(x_{1:n}) dx_{1:n}$$

$$= \int \left( \prod_{k=1}^{n} g_{\theta}(y_{k} | x_{k}) \right) \left( \mu_{\theta}(x_{1}) \prod_{k=2}^{n} f_{\theta}(x_{k} | x_{k-1}) \right) dx_{1:n}.$$

 If the model is linear Gaussian or finite state-space, we can compute the likelihood in closed-form but the maximization is not trivial. Otherwise, we cannot.

### Pairwise likelihood for state-space models

• We consider the following pseudo-likelihood for  $m \geq 1$ 

$$L_{S}\left(\theta; y_{1:n}\right) = \prod_{i=1}^{n-1} \prod_{j=i+1}^{\min\{i+m,n\}} p_{\theta}\left(y_{i}, y_{j}\right)$$

where

$$p_{\theta}(y_i, y_j) = \int g_{\theta}(y_i | x_i) g_{\theta}(y_j | x_j) p_{\theta}(x_i, x_j) dx_i dx_j.$$

• As an alternative, if n = pm, we could maximize

$$L_{\mathcal{S}}(\theta; y_{1:n}) = \prod_{i=1}^{p} p_{\theta}\left(y_{(i-1)m+1:im}\right)$$

• For the two models discussed earlier, it is possible to compute exactly  $p_{\theta}\left(x_{i},x_{j}
ight)$  as

$$p_{\theta}(x_i, x_j) = p_{\theta}(x_i) p_{\theta}(x_j | x_i)$$

where

$$p_{\theta}(x_i) = \mathcal{N}\left(x_i; 0, \frac{\sigma^2}{1-\alpha^2}\right)$$

$$p_{\theta}(x_j | x_i) = \mathcal{N}\left(x_j; \alpha^{j-i} x_i, \sigma^2 \sum_{k=0}^{j-i-1} \alpha^{2k}\right).$$

• In a general case, we could approximate  $p_{\theta}(y_i, y_j)$  through Monte Carlo

$$\widehat{p}_{\theta}\left(y_{i}, y_{j}\right) = \frac{1}{N} \sum_{l=1}^{N} g_{\theta}\left(y_{i} | X_{i}^{l}\right) g_{\theta}\left(y_{j} | X_{j}^{l}\right)$$

where 
$$\left(X_{i}^{I}, X_{j}^{I}\right) \sim p_{\theta}\left(x_{i}, x_{j}\right)$$
.

• Under regularity assumptions, we have

$$\int I_{S}\left(\theta; y_{1:n}\right) p_{\theta^{*}}\left(y_{1:n}\right) dy_{1:n}$$

which is maximum in  $\theta^*$  so maximizing this pseudo-likelihood method makes sense.

• To prove it, note that

$$I_{S}\left(\theta; y_{1:n}\right) = \sum_{i=1}^{n-1} \sum_{j=i+1}^{\min\{i+m,n\}} \log p_{\theta}\left(y_{i}, y_{j}\right)$$

and

$$\int \log p_{\theta} (y_i, y_j) . p_{\theta^*} (y_{1:n}) dy_{1:n}$$
$$= \int \log p_{\theta} (y_i, y_j) . p_{\theta^*} (y_i, y_j) dy_i dy_j$$

which is maximum in  $heta= heta^{*.}$ 

# Application

Consider

$$X_{1} \sim \mathcal{N}\left(0, \frac{\sigma^{2}}{1-\alpha^{2}}\right), X_{k} = \alpha X_{k-1} + \sigma V_{k} \text{ where } V_{k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right),$$
  

$$Y_{k} = \beta + X_{k} + \tau W_{k} \text{ where } W_{k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right).$$
  
where  $\theta = (\beta, \sigma^{2}, \alpha, \tau^{2}).$   
In this case, we can directly establish not only  $p_{\theta}\left(x_{i}, x_{j}\right)$  but  $p_{\theta}\left(y_{i}, y_{j}\right)$ 

$$\begin{pmatrix} Y_i \\ Y_j \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} \beta \\ \beta \end{pmatrix}, \begin{pmatrix} \tau^2 + \frac{\sigma^2}{1-\alpha^2} & \alpha^{j-i}\frac{\sigma^2}{1-\alpha^2} \\ \alpha^{j-i}\frac{\sigma^2}{1-\alpha^2} & \tau^2 + \frac{\sigma^2}{1-\alpha^2} \end{pmatrix}\right)$$

- For m = 2, ..., 20 we compare the performance of  $\hat{\theta}_{MLE}$  and  $\hat{\theta}_{MPL}$  where the likelihood and pseudo-likelihood are maximized using a simple gradient algorithm (EM could be used).
- 1,000 time series of length n = 500 with  $\beta^* = 0.1$ ,  $\tau^* = 1.0$ ,  $\alpha^* = 0.95$  and  $\sigma^* = 0.55$  are simulated.

true	$\widehat{\theta}_{MPL}^{(2)}$	$\widehat{\theta}_{MPL}^{(6)}$	$\widehat{\theta}_{MPL}^{(12)}$	$\widehat{ heta}_{MPL}^{(20)}$	$\widehat{\theta}_{ML}$
$\beta$ 0.1	0.108	0.108	0.109	0.109	0.102
	(0.488)	(0.489)	(0.4908)	(0.492)	(0.481)
τ 1.0	0.994	0.997	0.990	0.981	0.995
	(0.066)	(0.048)	(0.054)	(0.068)	(0.046)
α 0.95	0.941	0.941	0.939	0.937	0.941
	(0.033)	(0.020)	(0.022)	(0.024)	(0.020)
$\sigma$ 0.55	0.535	0.551	0.560	0.571	0.554
	(0.160)	(0.064)	(0.072)	(0.087)	(0.061)

- be able to compute MLE estimate for rather complex models,
- be able to compute the asymptotic variance of the MLE estimate,
- be able to derive the expression of the asymptotic variance of simple estimates different from MLE.