

Introduction to Maximum Likelihood Estimation (MLE)

Example with Normal IID variables

Juan Garbayo

Forecasting & Time Series Analysis

CONTENTS :

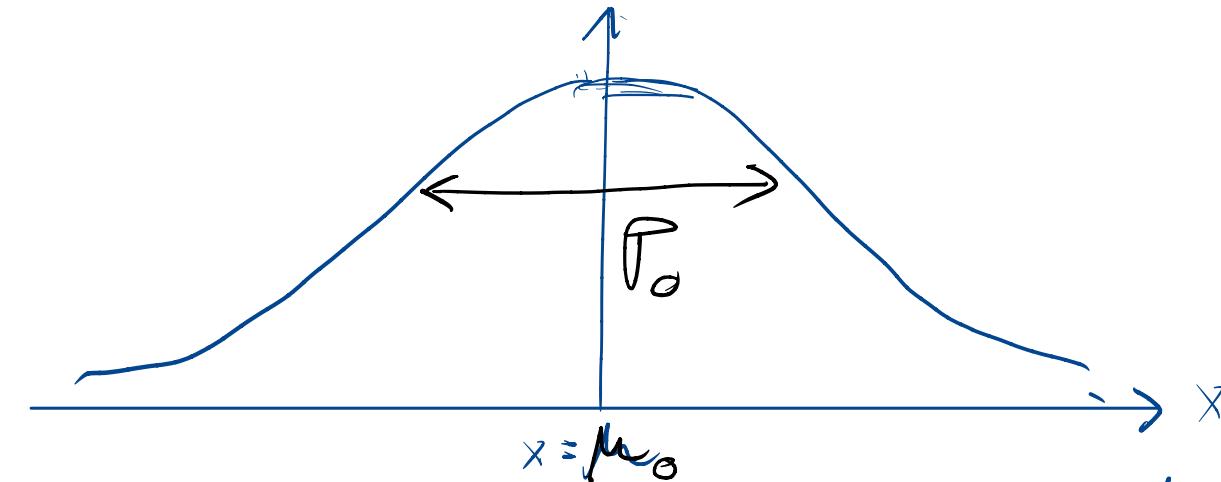
- ① Example with a normal distribution
- ② Generalization from the example
- ③ Relationship L and AIC

MLE → Example Normal Distribution

- f_X → Probability density function of a normal random variable: $X \sim N(\mu_0, \sigma_0^2)$

The sample:

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \cdot e^{-\frac{1}{2} \left[\frac{x - \mu_0}{\sigma_0} \right]^2}$$



- $\underline{X_n}$ → first n elements of an IID sequence of normal random variables with $\underline{\mu_0}$ and $\underline{\sigma_0^2} \rightarrow \{X_1, X_2, \dots, X_n\}$ with X

- Joint Probability Density Function of $\{X_1, \dots, X_n\}$

$$\frac{-1}{2\sigma_0^2} \left[\sum_{i=0}^n (x_i - \mu_0)^2 \right]$$

$$\underline{f_{X_n}(x_1, \dots, x_n)} = \prod_{i=0}^n f_X(x_i) = (2\pi\sigma_0^2)^{-n/2} \cdot e$$

IID \Leftrightarrow product of ind. events!!

PROBLEM : μ_0 and σ^2_0 unknown (parameters of our model) Realization of the RVs
 $\sim N\{x_1, \dots, x_n\}$

Maximum Likelihood Estimation : given a sample $\overbrace{\{x_1, \dots, x_n\}}$ (the sample is now a given, not a variable) obtain the values of μ and σ^2 that maximize the value of the joint probability density function particularized at the given sample:

$$L(\mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \cdot e^{-\frac{1}{2\sigma^2} \left[\sum_{i=1}^n (x_i - \mu)^2 \right]}$$

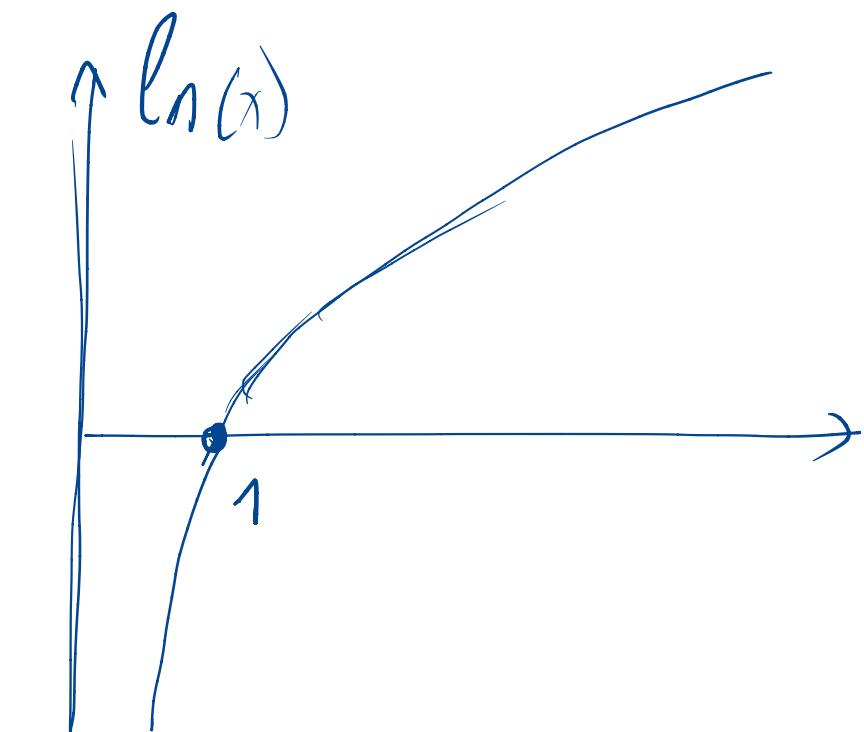
In other words : assuming the sample was generated by a $N(\mu, \sigma^2)$,

find the values of μ and σ^2 that maximize the probability of occurrence of the sample.

FORMALLY

$$\{M_0, T_0\} = \underset{\{M, T\} \in \mathbb{R} \times \mathbb{R}^+}{\arg \max} L(M, T)$$

Because the function $\ln(x)$ is monotonically increasing:



$$\underset{\{M, T\} \in \mathbb{R} \times \mathbb{R}^+}{\arg \max} L(M, T) = \underset{\{M, T\} \in \mathbb{R} \times \mathbb{R}^+}{\arg \max} \ln(L(M, T))$$

$\hookrightarrow L(M, T) \rightarrow \text{log-likelihood}$

MLE \rightarrow Normal Distr. (Cont'd)

$$\mathcal{L}(\mu, \sigma^2, x_1, \dots, x_n) = \prod_{i=1}^n f(x_i) = (2\pi\sigma^2)^{-n/2} \cdot e^{-\frac{1}{2\sigma^2} \left[\sum_{i=1}^n (x_i - \mu)^2 \right]}$$

Likelihood function

IID

Reminder:

- $\ln(a \cdot b) = \ln(a) + \ln(b)$
- $\ln(a^b) = b \cdot \ln(a)$
- $e^{\ln(a)} = a$

Take logarithms:

$$l = \ln(L)$$

Log-Likelihood
function

$$l(\mu, \sigma^2, x_1, \dots, x_n) = \underbrace{-\frac{1}{2} \ln(2\pi)}_{\text{constant!!}} - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

MLE - Normal Distribution $\rightarrow \mu_0, \sigma_0$

$$\bullet \frac{\partial}{\partial \mu} \ell(\mu, \sigma, x_1, \dots, x_n) = 0 \Rightarrow \frac{-1}{\sigma^2} \sum_{i=0}^n (x_i - \hat{\mu}) \cdot (-1) = 0 \Rightarrow \sum_{i=0}^n (x_i - \hat{\mu}) = 0$$

$$\Rightarrow \sum_{i=0}^n x_i - n\hat{\mu} = 0 \Rightarrow \boxed{\hat{\mu} = \frac{1}{n} \sum_{i=0}^n x_i}$$

Mean sample.

$$\bullet \frac{\partial}{\partial \sigma} \ell(\mu, \sigma, x_1, \dots, x_n) = 0 \Rightarrow -\frac{n}{2} \frac{1}{\sigma^2} \cdot 2\hat{\sigma} - \frac{1}{2} \cdot (-2) \hat{\sigma}^3 \sum_{i=0}^n (x_i - \hat{\mu})^2 = 0$$

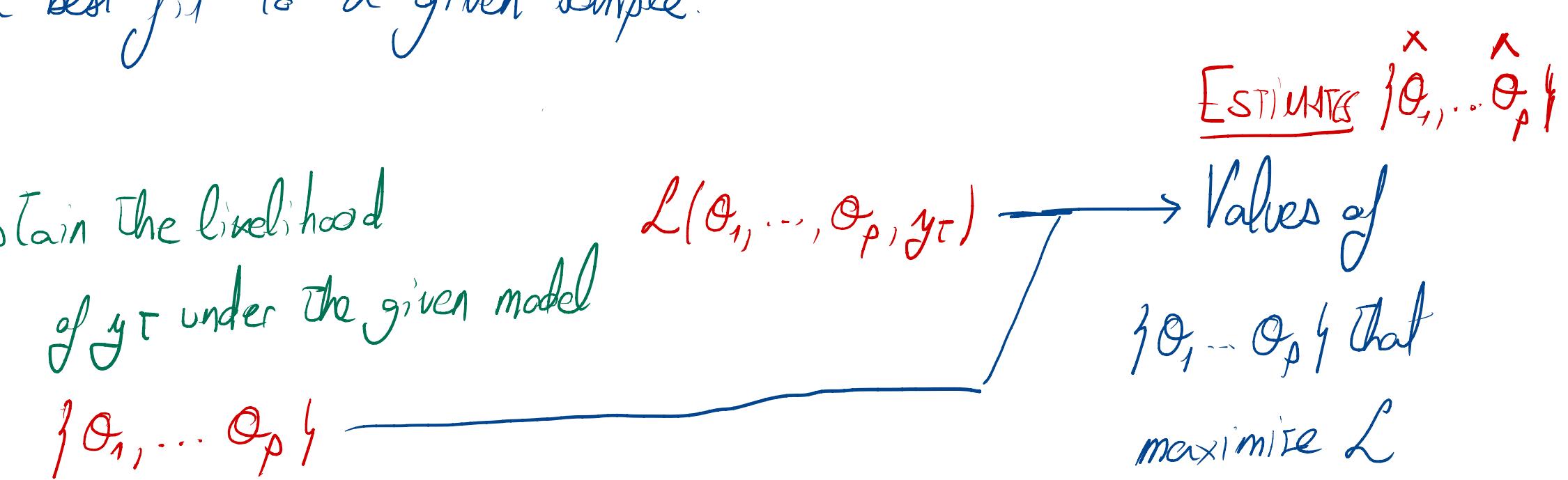
$$\begin{aligned} & \rightarrow \frac{1}{\hat{\sigma}^2} \left[-n + \frac{1}{\hat{\sigma}^2} \sum_{i=0}^n (x_i - \hat{\mu})^2 \right] = 0 \\ & \neq 0 (\hat{\sigma} \in \mathbb{R}^+) \quad = 0 \rightarrow n \hat{\sigma}^2 = \sum_{i=0}^n (x_i - \mu)^2 \rightarrow \boxed{\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=0}^n (x_i - \hat{\mu})^2} \end{aligned}$$

↳ Unadjusted sample variance.

Conclusions MLE

- General method to estimate the parameters of a statistical model to maximize its fit to a given sample.
- In the example we have:
 - model $\rightarrow N(\mu, \sigma^2)$
 - sample $\rightarrow \{x_1, \dots, x_n\}$
 - parameters $\rightarrow \mu, \sigma^2$
- The L can be computed for other statistical models and maximized to estimate the parameters that result in a best fit to a given sample.
- Example:

model \rightarrow ETS, ARIMA ...
sample $\rightarrow y_T$
parameters of ETS or ARIMA $\{\theta_1, \dots, \theta_p\}$



Relationship with AIC \rightarrow Rule: the smaller the AIC, the better (defined to work in this manner).

- The likelihood (and log likelihood) L is a metric of how well the model fits the sample. $\rightarrow L \uparrow \rightarrow$ fit improves.

$$\text{AIC} = f_1(L) + f_2(p)$$

$f_1(L)$ ↗ Likelihood
 $f_2(p)$ ↗ Number of parameters (complexity)

- $f_1(L)$ such that if $L \uparrow \rightarrow f_1(L) \downarrow \rightarrow \text{AIC} \downarrow$
- $f_2(p)$ such that if $p \uparrow \rightarrow f_2(p) \uparrow \rightarrow \text{AIC} \uparrow$
- Case ①: L increases without increasing p too much $\Rightarrow \underline{L \uparrow \uparrow \uparrow}$ and $\underline{p \uparrow} \rightarrow \underline{\text{AIC} \downarrow}$ ✓
- Case ②: L increases at the cost of increasing p too much $\rightarrow \underline{L \uparrow}$ and $\underline{p \uparrow \uparrow \uparrow} \rightarrow \underline{\text{AIC} \uparrow}$ ✗
 - Prevents overfitting

EXAMPLE OF AIC

$$AIC_{ETS \text{ Models}} = -2 \log(L) + 2P$$

$\sim f_1(L)$

$\sim f_2(p)$

→ Number of parameters in the model

Likelihood of a TS

assuming it was generated

by an ETS model

↙ More complex than the example

of the normal dists

Check behavior

$$\cdot \frac{df_1}{dL} = -\frac{2}{L} < 0 \Rightarrow L^+ \Rightarrow AIC^+$$

L
↗
0

$$\cdot \frac{df_2}{dp} : 2 > 0 \Rightarrow p^+ \Rightarrow AIC^+$$