

Introduction to Maximum Likelihood Estimation (MLE)

Example with Normal IID variables

Juan Garbayo

Forecasting & Time Series Analysis

CONTENTS :

- ① Example with a normal distribution $N(\mu, \sigma^2)$
- ② Generalization from the example to TS Models.
- ③ Relationship L and AIC

MLE: EXAMPLE WITH NORMAL DISTRIBUTION

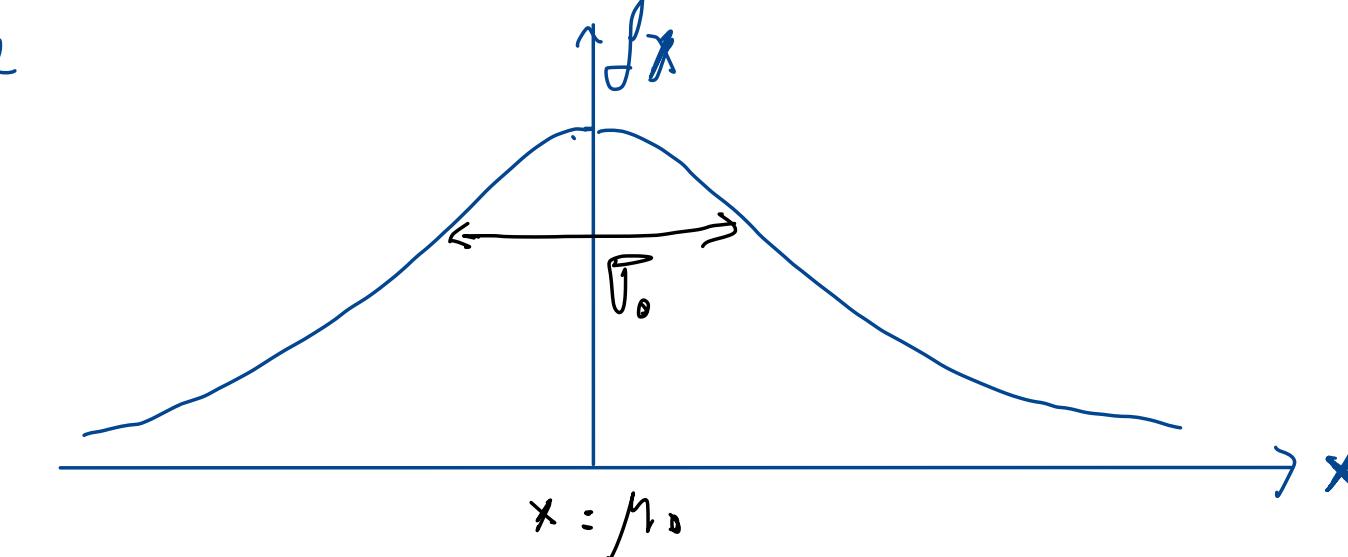
f_X : probability density function of a normal random variable $X \sim N(\mu, \sigma^2)$

$X \rightarrow$ random variable

$x \rightarrow$ variable of the pdf. Its

domain is all the possible outcomes of X .

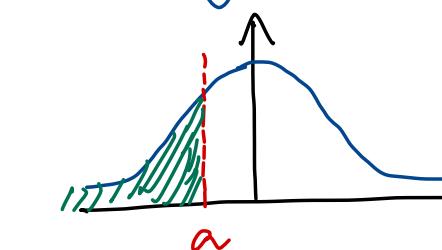
$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \left[\frac{(x-\mu)}{\sigma} \right]^2}$$



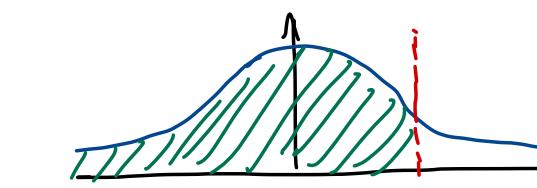
Reminder: $f_X(x)$ is a probability DENSITY

Because X is a continuous RANDOM VARIABLE, we do not compute the probability of X taking a specific value, but rather the probability of X being within a specific interval.

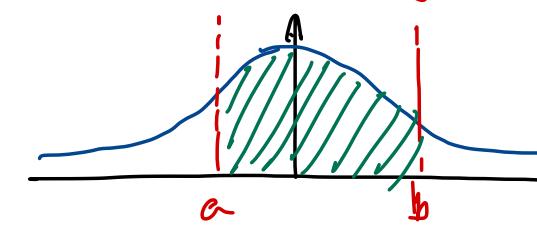
Example ①: $P(X \leq a) = \int_{-\infty}^a f_X(x) dx$



Example ②: $P(X \leq b) = \int_{-\infty}^b f_X(x) dx$



Example ③: $P(a \leq X \leq b) = \int_a^b f_X(x) dx$



Of course, because of the axioms of probability:

$$\int_{-\infty}^{+\infty} f_X(x) dx = 1$$

↳ NOTE: $f_X(x)$ may be > 1 in a subset.

Joint Probability Density Function of sequence of IID $N(\mu_0, \sigma_0^2)$ (1/2)

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

$\rightarrow X \sim N(\mu_0, \sigma_0^2)$ Known Parameters

Single R.V called X

Experiment: first n elements of an IID sequence of normal random variables

$$\bar{X} = [X_1, X_2, \dots, X_n] \text{ where } X_i \sim N(\mu_0, \sigma_0^2); \quad \bar{X} \text{ is a RANDOM VECTOR}$$

The random vector \bar{X} is the process (experiment) of extracting the first n elements of a sequence of independent random variables

The joint probability density function of \bar{X} is the mathematical characterization of that process.

↳ Next page

Joint Probability Density Function of sequence of IID $N(\mu_0, \sigma_0)$ (2/2)

Probability density function of normal R.V $X_i \sim N(\mu, \sigma)$

$$f_{X_i}(x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \left[\frac{(x_i - \mu_0)}{\sigma_0} \right]^2}$$

Joint Probability Density Function of $\bar{X} = \{X_1, X_2, \dots, X_n\}$

$$f_{\bar{X}}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i) = (2\pi\sigma^2)^{-n/2} \cdot e^{-\frac{1}{2\sigma^2} \left[\sum_{i=1}^n (x_i - \mu_0)^2 \right]}$$

IID \rightarrow product of independent events !!

Reminder $e^a \cdot e^b = e^{a+b}$

Example of Joint Probability Density Function with n=2

$$\cdot \underline{\int X_1(x_1)} = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(x_1 - \mu_0)^2}{2\sigma^2}}$$
$$\cdot \underline{\int X_2(x_2)} = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(x_2 - \mu_0)^2}{2\sigma^2}}$$

$$\bar{X} = \{X_1, X_2\}$$

$$\underline{\int \bar{X}(x_1, x_2)} = \int X_1(x_1) \cdot \int X_2(x_2) = \frac{1}{2\pi\sigma^2} \cdot e^{-\frac{[(x_1 - \mu_0)^2 + (x_2 - \mu_0)^2]}{2\sigma^2}}$$

Example of computation of Joint Probabilities

$$P(X_1 \in [a_1, b_1], X_2 \in [a_2, b_2]) = \int_{x_1=a_1}^{x_1=b_1} \int_{x_2=a_2}^{x_2=b_2} \int \bar{X}(x_1, x_2) dx_1 dx_2$$

Estimate μ_0 and σ_0^2 to best fit a sample

- . Assumption: data follows a normal distribution $N(\mu, \sigma^2) \rightarrow \mu, \sigma^2$ UNKNOWN
- . Problem: find μ_0 and σ_0^2 corresponding to the normal distribution that best fits the sample

Data given: $\underbrace{\{x_1, \dots, x_n\}}_{\text{SAMPLE}} \rightarrow$ specific outcome of the random vector $\bar{X} = \{X_1, X_2, \dots, X_n\}$

- . Approach: find the values of μ and σ^2 that maximize the probability of the random vector

$\bar{X} = \{X_1, X_2, \dots, X_n\}$ taking the values of the sample $\{x_1, x_2, \dots, x_n\}$.

→ Likelihood

$$l(\mu, \sigma^2, \underbrace{x_1, x_2, \dots, x_n}_{\text{sample given}}) =$$

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

Joint probability density function
of the process \bar{X} particularized at
the specific outcome (the sample
 $\{x_1, x_2, \dots, x_n\}$)

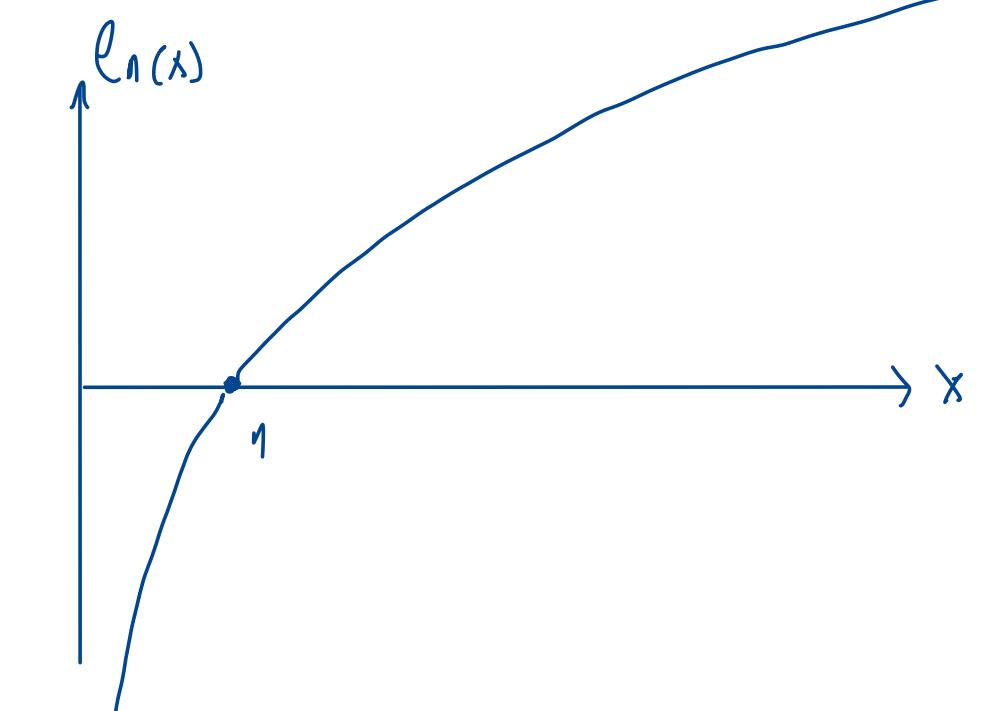
↓
Model parameters
to be estimated

In other words: assuming your data was generated by
a $N(\mu, \sigma^2)$, find μ_0 and σ_0^2 that maximize the probability
of occurrence of the observed data.

FORMAL WAY TO WRITE THE OPTIMIZATION PROBLEM:

$$\{M_0, \Sigma_0\} = \underset{\{M, \Sigma\} \in \mathbb{R} \times \mathbb{R}^+}{\operatorname{argmax}} L(M, \Sigma)$$

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



$$\underset{\{M, \Sigma\} \in \mathbb{R} \times \mathbb{R}^+}{\operatorname{argmax}} L(M, \Sigma)$$

Same optimization problem

$$= \underset{\{M, \Sigma\} \in \mathbb{R} \times \mathbb{R}^+}{\operatorname{argmax}} \ln(L(M, \Sigma))$$

↳ log-likelihood.

Natural logarithm (base e).

PROBLEM ①

PROBLEM ②

LOG-LIKELIHOOD FOR NORMAL DISTRIBUTION CASE

Likelihood

$$L(\mu, \sigma^2; x_1, \dots, x_n) \stackrel{\text{IID}}{=} \prod_{i=1}^n f(x_i | \mu, \sigma^2) = \underbrace{(2\pi\sigma^2)^{-n/2}}_{\cdot e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2}}$$

Reminder:

$$\begin{cases} \ln(a \cdot b) = \ln(a) + \ln(b) \\ \ln(a^b) = b \ln(a) \\ e^{\ln(a)} = a \end{cases}$$

Log-Likelihood

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

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

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

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

Mean sample.

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

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

↳ 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 $\rightarrow \{x_1, \dots, x_n\}$; Process $\rightarrow \{X_1, \dots, X_n\}$
 - 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 of the model to the given sample.

MAXIMUM LIKELIHOOD ESTIMATION - More GENERAL NOTATION

PROBLEM: given a sample of n elements $\{x_1, \dots, x_n\} = \bar{x}$ and a statistical model dependant on p parameters $\{\theta_1, \dots, \theta_p\} = \bar{\theta}$ find the values of the parameters that maximize the joint probability of occurrence of the sample under the assumption that the model generated the sample.

$$L(\bar{\theta}, \bar{x}) = f_{\bar{\theta}}(\bar{\theta}, \bar{x}) \mid \bar{x} = \bar{x}$$

↳ Likelihood
Model Sample
params $\{x_1, x_2, \dots, x_n\}$
 $\bar{\theta} = \{\theta_1, \theta_2, \dots, \theta_p\}$

↳ joint probability
density function
Random vector $\{X_1, X_2, \dots, X_n\}$

④ set of possible values for
the parameters $\bar{\theta} = \{\theta_1, \dots, \theta_p\}$

$$\hat{\theta} = \underset{\theta \in \mathbb{H}}{\operatorname{argmax}} L(\bar{\theta}, \bar{x})$$

→ Optimization problem to be solved.

→ set of possible values taken by the parameters $\{\theta_1, \dots, \theta_p\}$
→ Values of the params that maximize the likelihood

MLE on Time Series Models

- model: ETS, ARIMA ...
- process / model $\bar{Y}_T = \{Y_1, Y_2, \dots, Y_T\} \rightarrow$ the collection of R.V that describe the stochastic process
- sample: $\bar{y}_T = \{y_1, y_2, \dots, y_T\} \rightarrow$ the TS data given \rightarrow the specific realization of the TS.
- parameters $\bar{\Theta} = \{\theta_1, \theta_2, \dots, \theta_p\}$
- Likelihood function $L(\bar{\Theta}, \bar{y}_T) = \int_{\bar{\Theta}} \phi(\bar{\Theta}, \bar{Y}_T) | \bar{Y}_T = \bar{y}_T$
 - Joint probability density function of the TS (of the collection of R. Variables) assuming the TS process obeys the equations of the model.
 - particularized at the sample.
- Optimization problem:
$$\hat{\Theta} = \underset{\bar{\Theta} \in \Theta}{\operatorname{argmax}} L(\bar{\Theta}, \bar{y}_T)$$
 - values of parameters that maximize L
 - Set of possible values for the vector $\bar{\Theta}$ of parameters

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 dist.

Check behavior

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

L
↓
 > 0

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