

LECTURE 3

Let us recall a few key results from Lecture 2 that will be used here:

* $[X_1, X_2, \dots, X_n]^T$ random sample generated by sampling the RV $X \sim (\mu_X, \sigma_X^2)$ (X_1, X_2, \dots, X_n) - i.i.d. \Rightarrow

$$\hat{X} \triangleq \frac{1}{n} \sum X_i$$

$$E_{\hat{X}}(\hat{X}) = \mu_X$$

$$E_{\hat{X}}((\hat{X} - \mu_X)^2) = \sigma_X^2/n$$

* (X_1, X_2, \dots, X_n) converges in distribution to X $\Leftrightarrow \lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x) \quad \forall x \in \mathcal{R}$

where $F_X(\cdot)$ and $F_{X_n}(\cdot)$ are cdfs

* (X_1, X_2, \dots, X_n) converges in probability to a constant c^* $\Leftrightarrow (X_1, X_2, \dots, X_n)$ converges in distribution to a RV X such that: $P(X = c^*) = 1$

* LLN: (X_1, X_2, \dots, X_n) i.i.d. \Rightarrow Denoted with $\hat{X}_i \triangleq \frac{1}{i} \sum_{j=1}^i X_j$, we have that $\{\hat{X}_i\}_i \xrightarrow{P} \mu_X$

* CLT: (X_1, X_2, \dots, X_n) i.i.d. \Rightarrow Denoted with $Z_i \triangleq \frac{\sqrt{i}}{\sigma_X} (\hat{X}_i - \mu_X)$, we have that $\{Z_i\}_i \xrightarrow{D} N(0, 1)$

* Finally, let us recall that, given two RVs X and Y with joint probability function $f(x, y)$ and marginal pdfs $f_X(x)$ and $f_Y(y)$, respectively, we can write:

$$f_{X|Y}(x|y) \triangleq f_X(x|Y=y) = \frac{f(x, y)}{f_Y(y)} \quad \text{- Conditional pdf of } X \text{ given } Y=y$$

A generalization of the Bayes' theorem can be obtained by noticing that:

• $f(x, y) = f_{X|Y}(x|y) f_Y(y) = f_{Y|X}(y|x) f_X(x)$

2

$$f_Y(y) = \int_{-\infty}^{+\infty} f(x, y) dx = \int_{-\infty}^{+\infty} f_{Y|X}(y|x) f_X(x) dx$$

Bayes' Theorem $f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x) f_X(x)}{\int_{-\infty}^{+\infty} f_{Y|X}(y|x) f_X(x) dx}$ (if X is continuous)

$$f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x) f_X(x)}{\sum_x f_{Y|X}(y|x) f_X(x)} \quad (\text{if } X \text{ is discrete})$$

□

A few well-studied probability functions will be considered here:

- Bernoulli RV: $X \sim \text{Bernoulli}(p) \stackrel{\text{DEF}}{\iff} P(X=1) = p; P(X=0) = 1-p$

For this RV, we have: $\mu_X = p; \sigma_X^2 = p(1-p)$

- Binomial RV: $X \sim B(n, p) \stackrel{\text{DEF}}{\iff} P(X=k) = \binom{n}{k} p^k (1-p)^{n-k} \quad k=0, 1, \dots, n$

For this RV, we have: $\mu_X = np; \sigma_X^2 = np(1-p)$

- Poisson RV: $X \sim P(\lambda) \stackrel{\text{DEF}}{\iff} P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!} \quad k=0, 1, 2, \dots$

For this RV, we have: $\mu_X = \lambda; \sigma_X^2 = \lambda$. This RV is a generalization of $B(n, p)$ as $p \rightarrow 0$ and $n \rightarrow \infty$ with $\lambda = np$

Note that in all the RVs considered above, we assume that the binary events they are built upon are INDEPENDENT

□

Let us consider the following scenario: we have a random sample (X_1, X_2, \dots, X_n) generated by a partially-unknown RV X and we would like to estimate X given the random sample.

- We know about X : $X \sim f_x(x) = f_x(x|\theta)$ with f_x belonging to a known class of pdfs
 \uparrow
 parameter vector
- We do NOT know about X : θ^* (i.e., the true value of θ)

- We would like to: $(x_1, x_2, \dots, x_n) \rightarrow$ Estimation $\hat{\theta}$ (*)
 where $X_1 = x_1, X_2 = x_2, \dots$

Condition (*), though, means that the estimation $\hat{\theta}$ may vary when the n -tuple (x_1, x_2, \dots, x_n) changes \Rightarrow We can define a new RV:

Estimator: $T: (x_1, x_2, \dots, x_n) \rightarrow \hat{\theta}$

\Downarrow

$T = T(X_1, X_2, \dots, X_n)$ - is a function of RVs and its own distribution must be determined

The definition of T (and, hence, its distribution) depends on how $\hat{\theta}$ is estimated by using (x_1, x_2, \dots, x_n) , i.e., it depends on the method used to fit $f_x(\cdot)$ on data.

The method of Maximum Likelihood (ML) chooses $\hat{\theta}$ such that:

$$\hat{\theta} = \arg \max_{\theta} f_x(x|\theta)$$

for every observed value x

Ex.: $X \sim B(n, p)$ with p -unknown, i.e., $\theta \triangleq p$. In this case, we have:

$$f_x(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \quad x=0, 1, 2, \dots, n \quad \text{and} \quad f_x(x) \triangleq P(X=x)$$

④

In this case, it is easy to show that the value $\hat{\theta}$ satisfies:

$$\frac{\partial f_x}{\partial \theta} = 0 \Leftrightarrow \frac{\partial}{\partial \theta} \left[\theta^x (1-\theta)^{n-x} \right] = 0$$

We can drop
the constant
 $\binom{n}{x}$

$$\Leftrightarrow x \theta^{x-1} (1-\theta)^{n-x} - (n-x) \theta^x (1-\theta)^{n-x-1} = 0$$

$$\Leftrightarrow \theta^{x-1} (1-\theta)^{n-x-1} \left[x(1-\theta) - (n-x)\theta \right] = 0$$

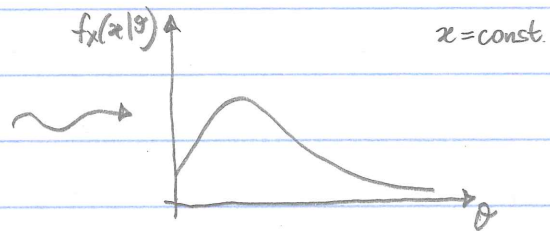
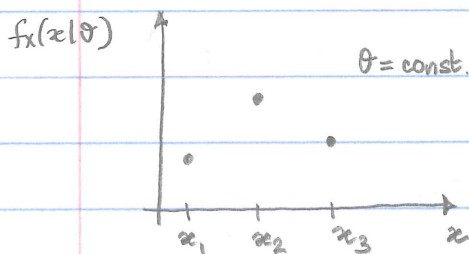
$$\Leftrightarrow \theta^{x-1} (1-\theta)^{n-x-1} \left[x - \theta n \right] = 0$$

$$\Leftrightarrow \hat{\theta} = \frac{x}{n} \quad (**)$$

Formula (***) shows: $T_n: x \rightarrow \hat{\theta} = \frac{x}{n}$ is a RV. In fact: $T_n = \frac{X}{n}$ and, by the LLN, we have: $\frac{X}{n} \xrightarrow{P} \theta^*$, i.e., the true value of the probability that a Bernoulli trial results in 1. $\Rightarrow T_n$ is called the MAXIMUM LIKELIHOOD ESTIMATOR (MLE)

This example highlights a few facts:

- With the ML method, we use $f_x(x|\theta)$ as a function of θ , i.e., we fix x and let θ vary \Rightarrow The maximization problem is now less affected by the amount of data (i.e., the number of samples x_1, x_2, \dots, x_n) we have



- The condition $\max_{\theta} f_x(x|\theta)$ is defined up to a positive constant
 \Rightarrow We can choose a function $L(\theta)$ to be maximized such that $L(\theta) \propto f_x(x|\theta)$
 $\Rightarrow L(\theta)$ is the LIKELIHOOD FUNCTION

In our example, $L(\theta) \triangleq \theta^x (1-\theta)^{n-x}$

- The ML method can be generalized to a random sample; i.e., we replace $f_x(\cdot)$ with the joint probability function $f(x_1, x_2, \dots, x_n)$ and choose:

$$L(\theta) \propto f(x_1, x_2, \dots, x_n | \theta)$$

- In many cases, it may be appropriate to maximize the LOG-LIKELIHOOD function $l(\theta) \triangleq \log(L(\theta))$ instead of the likelihood function. In our example:

$$L(\theta) = \theta^x (1-\theta)^{n-x} \Rightarrow l(\theta) = x \log \theta + (n-x) \log(1-\theta)$$

$$\text{Hence: } \frac{\partial L}{\partial \theta} = 0 \Leftrightarrow \frac{\partial l}{\partial \theta} = 0 \Leftrightarrow \frac{x}{\theta} - \frac{n-x}{1-\theta} = 0 \Leftrightarrow \hat{\theta} = \frac{x}{n}$$

Ex.: Let us consider n samples x_1, x_2, \dots, x_n from the RV $X \sim N(\mu, \sigma^2)$ where $\theta \triangleq \mu$ must be estimated and σ^2 is known. \Rightarrow We can think x_1, x_2, \dots, x_n as the outcome of a random sample (X_1, X_2, \dots, X_n) where $X_i \sim N(\mu, \sigma^2) \forall i$ and the RVs are independent \Rightarrow We have:

$$f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n f_{x_i}(x_i | \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x_i - \theta)^2}{\sigma^2}}$$

$$\text{We define: } l(\theta) \triangleq \log \left(\prod_{i=1}^n f_{x_i}(x_i | \theta) \right) = \sum_{i=1}^n \frac{(x_i - \theta)^2}{\sigma^2} =$$

$$= -\frac{1}{\sigma^2} \sum_{i=1}^n (x_i^2 - 2x_i\theta + \theta^2) = -\frac{n}{\sigma^2} \theta^2 + \frac{2}{\sigma^2} n \bar{x} \theta - \frac{1}{\sigma^2} \sum_{i=1}^n x_i^2$$

It is constant

sample mean

6

$$\text{Hence: } \frac{\partial l}{\partial \theta} = 0 \Leftrightarrow -2\theta \frac{n}{\sigma^2} + 2 \frac{n}{\sigma^2} \bar{x} = 0 \Leftrightarrow \hat{\theta} = \bar{x}$$

In this example we have that the MLE $T_n = \hat{X}_n$ (sample mean) and, because of the LLN, we have: $T_n \xrightarrow{P} \mu$.

- The ML estimation is an optimization problem:

$$\begin{aligned} (***) \quad & \max l(\theta) \\ \text{s.t.: } & \theta \in C \quad \text{with } C \text{ set of admissible parameters} \end{aligned}$$

As in any optimization problem, θ can be either a scalar or a vector and the existence (and uniqueness) of the solution depends on the definition of $l(\theta)$ and C . In particular, let us note:

- The class of the function $f_x(\cdot)$ determines whether $l(\theta)$ is concave or not and whether it has one maximum or many local maxima

↓

We will use generalized linear models and Gaussian mixtures

- If $l(\theta)$ is a known mathematical function and concave \Rightarrow The solution exists and can be computed by solving: $\frac{\partial l}{\partial \theta} \Big|_{\theta=\hat{\theta}} = 0$
- If $l(\theta)$ is concave and numerically evaluated \Rightarrow We can apply algorithms for global maximization (e.g. Newton's method, quasi-Newton's methods, secant method, etc.) \Rightarrow These methods are invoked by MLE routines in MATLAB, R, etc.

• From a numerical standpoint, the optimization problem is often formulated as a minimization problem \Rightarrow We solve for $\hat{\theta} = \arg \min_{\theta \in G} (-l(\theta))$

Because the MLE T_n is a RV, we must characterize the probability distribution of $T_n \Rightarrow$ The features of T_n as RV will justify its wide use as an estimation technique:

* ML Estimators are consistent and asymptotically normal

T_n - consistent $\stackrel{\text{DEF}}{\iff} T_n \xrightarrow{P} \theta^*$ (true set of parameters), i.e.,
 $\forall \epsilon > 0 \quad P(|T_n - \theta^*| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow +\infty$

T_n - asymptotically normal $\stackrel{\text{DEF}}{\iff} \frac{T_n - \theta^*}{\sigma_{T_n}} \xrightarrow{D} N(0, 1)$ where $\sigma_{T_n} \triangleq$ standard dev. of T_n

The consistency of the MLE can be shown by noticing that MLE minimizes the mean square error (MSE). In fact, denoted with $MSE(T_n) \triangleq E_{T_n}((T_n - \theta^*)^2)$, the following inequality holds:

$$\forall \epsilon > 0 \quad P(|T_n - \theta^*| > \epsilon) < \frac{E_{T_n}((T_n - \theta^*)^2)}{\epsilon^2}$$

Hence, the condition $MSE(T_n) \rightarrow 0$ shall imply that $T_n \xrightarrow{P} \theta^*$. It is interesting that MLE minimizes MSE because the following occurs:

$$(T_n - \theta^*)^2 = \underbrace{(T_n - \mu_n)}_{\substack{\uparrow \\ \cong E_{T_n}(T_n) \\ \text{(expected value)}}}^2 + (\mu_n - \theta^*)^2 + 2(T_n - \mu_n)(\mu_n - \theta^*) \Rightarrow \text{By applying the expected value operator and exploiting the linearity:}$$

8

$$\text{MSE}(T_n) = \underbrace{\mathbb{E}_{T_n} \left((T_n - \mu_n)^2 \right)}_{\substack{\uparrow \\ \text{variance} \\ \text{of } T_n}} + \underbrace{(\mu_n - \theta^*)^2}_{\substack{\uparrow \\ \text{a constant}}} + 2(\mu_n - \theta^*) \underbrace{\mathbb{E}_{T_n} (T_n - \mu_n)}_0 \Rightarrow$$

$$\text{MSE}(T_n) = \mathbb{E}_{T_n} \left((T_n - \mu_n)^2 \right) + \left(\mathbb{E}_{T_n} (T_n - \theta^*) \right)^2 \Rightarrow \text{To minimize MSE is equivalent to minimize the bias (squared) plus the variance of } T_n$$

Ex.: Let us suppose that a train of neural spikes follow a Poisson distribution with unknown parameter λ (i.e., $\theta^* = \lambda$) and let us assume that we have n independent spike counts:

$$X \sim P(\lambda) \Rightarrow f_x(x) \triangleq P(X=x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

$$f(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_x(x_i) = e^{-n\lambda} \frac{\lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}$$

$$\text{Let us define: } \ell(\theta) = \log \left(e^{-n\theta} \theta^{\sum_{i=1}^n x_i} \right) = -n\theta + \log \theta \sum_{i=1}^n x_i$$

$$\text{By applying the ML method, we have: } \frac{\partial \ell}{\partial \theta} = 0 \Leftrightarrow -n + \frac{1}{\theta} \sum_{i=1}^n x_i = 0$$

$$\Leftrightarrow \hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i, \text{ i.e., the MLE is the sample mean } \bar{X}$$

In this case, we already know that $\mathbb{E}_{T_n} (T_n) = \mathbb{E}_{\bar{X}} (\bar{X}) = \lambda \Rightarrow \text{MSE}(T_n) = \mathbb{E}_{T_n} \left((T_n - \lambda)^2 \right)$
i.e., the MLE is unbiased

Now, in order to show that $\text{MSE}(T_n)$ is less than the MSE obtained with other estimators, let us consider the estimator:

$$S_n^2 \triangleq \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad (\text{sample variance})$$

$$\text{We have: } E(S_n^2) = \frac{1}{n-1} \sum_{i=1}^n E((X_i - \bar{X})^2) =$$

$$= \frac{1}{n-1} \sum_{i=1}^n E\left((X_i - \lambda)^2 + (\lambda - \bar{X})^2 - 2(X_i - \lambda)(\lambda - \bar{X})\right)$$

$$= \frac{1}{n-1} \left[n\sigma_X^2 + n \cdot \frac{\sigma_X^2}{n} - 2n \frac{\sigma_X^2}{n} \right] = \sigma_X^2 = \lambda$$

↑
Poisson

$$\text{Hence: } \text{MSE}(S_n^2) = E_{S_n^2}((S_n^2 - \lambda)^2) \text{ - unbiased}$$

Now, by computing the variances, one obtains:

$$\text{MSE}(T_n) = \frac{\lambda}{n} \quad (\text{see lecture 2, page 13})$$

$$\text{MSE}(S_n^2) = \frac{\lambda}{n} + \frac{2\lambda^2}{n-1}$$

$$\Rightarrow \text{MSE}(T_n) < \text{MSE}(S_n^2) \quad \forall n > 0$$

□

With regard to the asymptotic normality, one can easily prove it in this case:

$$\max_{\theta \in G} \ell(\theta) \Leftrightarrow \left. \frac{\partial \ell}{\partial \theta} \right|_{\theta = \hat{\theta}} = 0$$

i.i.d. assumption

$$\text{Let us assume: } \ell(\theta) = \frac{1}{n} \log f(x_1, x_2, \dots, x_n | \theta) \stackrel{\uparrow}{=} \frac{1}{n} \sum_{i=1}^n \log f_X(x_i | \theta)$$

$$\text{So we can write: } \left. \frac{\partial \ell}{\partial \theta}(\hat{\theta}) \right|_{\theta = \hat{\theta}} = 0 \Leftrightarrow \frac{1}{n} \sum_{i=1}^n \left. \frac{\partial}{\partial \theta} (\log f_X(x_i | \theta)) \right|_{\theta = \hat{\theta}} = 0$$

$$\text{Let us remind: } \frac{\partial}{\partial \theta} (\log f_X(x | \theta)) = \frac{f'_X(x | \theta)}{f_X(x | \theta)} \quad \text{where } f'_X(x | \theta) \triangleq \frac{\partial f_X(x | \theta)}{\partial \theta}$$

$$\frac{\partial^2}{\partial \theta^2} (\log f_X(x | \theta)) = \frac{f''_X(x | \theta)}{f_X(x | \theta)} - \frac{[f'_X(x | \theta)]^2}{f_X^2(x | \theta)} \quad \text{where } f''_X(x | \theta) \triangleq \frac{\partial^2 f_X(x | \theta)}{\partial \theta^2}$$

(10)

$$\text{Therefore: } E_X \left(\frac{\partial}{\partial \theta} \log f_X(X|\theta) \right) = \int_{-\infty}^{+\infty} f'_X(x|\theta) dx = \frac{\partial}{\partial \theta} \int_{-\infty}^{+\infty} f_X(x|\theta) dx = 0$$

$$E_X \left(\frac{\partial^2}{\partial \theta^2} \log f_X(X|\theta) \right) = \underbrace{\int_{-\infty}^{+\infty} f''_X(x|\theta) dx}_{=0} - \int_{-\infty}^{+\infty} \frac{[f'_X(x|\theta)]^2}{f_X(x|\theta)} f_X(x|\theta) dx$$

$$= - E_X \left(\left(\frac{\partial}{\partial \theta} \log f_X(X|\theta) \right)^2 \right)$$

The quantity $I_F(\theta) \triangleq - E_X \left(\frac{\partial^2}{\partial \theta^2} \log f_X(X|\theta) \right) = E_X \left(\left(\frac{\partial}{\partial \theta} \log f_X(X|\theta) \right)^2 \right)$ is called

"Fisher Information" of the RV X

We can apply the mean value theorem (MVT) to $\frac{\partial \ell}{\partial \theta}$ between θ^* (i.e., the true parameter) and $\hat{\theta}$ (i.e., the ML estimation):

$$\text{MVT: } \frac{f(a) - f(b)}{a - b} = f'(c) \quad c \in [a, b] \Rightarrow a \triangleq \hat{\theta} \quad b \triangleq \theta^* \quad f(\theta) \triangleq \frac{\partial \ell}{\partial \theta}$$

$$\frac{\partial \ell}{\partial \theta}(\hat{\theta}) = \frac{\partial \ell}{\partial \theta}(\theta^*) + \frac{\partial^2 \ell}{\partial \theta^2}(\tilde{\theta})(\hat{\theta} - \theta^*) \quad \text{with } \tilde{\theta} \in [\hat{\theta}, \theta^*] \text{ (assuming } \hat{\theta} < \theta^*)$$

$$\Rightarrow \sqrt{n}(\hat{\theta} - \theta^*) = - \frac{\sqrt{n} \frac{\partial \ell}{\partial \theta}(\theta^*)}{\frac{\partial^2 \ell}{\partial \theta^2}(\tilde{\theta})}$$

Let us note: $\frac{\partial \ell}{\partial \theta}(\theta^*) \triangleq \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \left(\log f_X(X_i|\theta) \right) \Big|_{\theta=\theta^*}$ - If we define:

$$Y_i \triangleq \frac{\partial}{\partial \theta} \left(\log f_X(X_i|\theta) \right) \Big|_{\theta=\theta^*} \Rightarrow \frac{\partial \ell}{\partial \theta}(\theta^*) = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{- sample mean}$$

Hence: - By LLN: $E \left(\frac{\partial \ell}{\partial \theta}(\theta^*) \right) = E_X \left(\frac{\partial}{\partial \theta} \log f_X(X|\theta) \Big|_{\theta=\theta^*} \right) = 0$

- By CLT: $\sqrt{n} \frac{\partial \ell}{\partial \theta}(\theta^*) \xrightarrow{D} N\left(0, \text{var}\left(\frac{\partial}{\partial \theta} \log f_X(X|\theta) \Big|_{\theta=\theta^*}\right)\right) =$

$$= N(0, I_F(\theta^*)) \quad (1)$$

Let us also note: $\frac{\partial^2 \ell}{\partial \theta^2}(\tilde{\theta}) = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} (\log f_X(x_i|\theta)) \Big|_{\theta=\tilde{\theta}}$ - If we define:

$$Z_i \triangleq \frac{\partial^2}{\partial \theta^2} (\log f_X(x_i|\theta)) \Big|_{\theta=\tilde{\theta}} \Rightarrow \frac{\partial^2 \ell}{\partial \theta^2}(\tilde{\theta}) = \frac{1}{n} \sum_{i=1}^n Z_i \text{ - sample mean}$$

Hence, by LLN: $E\left(\frac{\partial^2 \ell}{\partial \theta^2}(\tilde{\theta})\right) = -I_F(\tilde{\theta})$

Also, since $T_n = \hat{\theta}$ is a consistent estimator, we have:

$$T_n \xrightarrow{P} \theta^* \Rightarrow I_F(\tilde{\theta}) \xrightarrow{P} I_F(\theta^*) \quad (2)$$

By combining (1) and (2) we have: $\sqrt{n} (T_n - \theta^*) \xrightarrow{D} N\left(0, \frac{1}{I_F(\theta^*)}\right)$

Finally one can note:

$$\left. \begin{aligned} I_n(T_n) &\triangleq - \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} (\log f_X(x_i|\theta)) \Big|_{\theta=\hat{\theta}} = n I_F(\hat{\theta}) \\ T_n &\xrightarrow{P} \theta^* \Rightarrow I_F(\hat{\theta}) \xrightarrow{P} I_F(\theta^*) \end{aligned} \right\} \begin{array}{l} \text{Slutsky's} \\ \text{Theorem} \\ \uparrow \\ \Rightarrow \sqrt{I_n(T_n)} (T_n - \theta^*) \xrightarrow{D} N(0, 1) \end{array}$$

In this derivation, we have $\sigma_{T_n}^2 = \frac{1}{I_n(T_n)}$. It can be proved that the variance of

MLE is actually $1/I_n(T_n)$. To give a sense of this, let us consider $f_X(x|\theta) \triangleq \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\theta)^2}{2\sigma^2}}$ where σ is known. In this case, we have:

(12)

$$l(\theta) = \log e^{-\frac{(x-\theta)^2}{2\sigma^2}} = -\frac{(x-\theta)^2}{2\sigma^2} \Rightarrow \frac{\partial l}{\partial \theta} = \frac{x-\theta}{\sigma^2} = 0 \Rightarrow \hat{\theta} = x$$

$$\frac{\partial^2 l}{\partial \theta^2} = -\frac{1}{\sigma^2} \Rightarrow \sigma^2 = \frac{1}{I_T(\hat{\theta})}$$

More in general, the Fisher Information can be defined for any estimator T as:

$$I^T(\theta) \triangleq E \left(-\frac{\partial^2}{\partial \theta^2} \log f_T(X|\theta) \right)$$

↑
pdf chosen by
the estimator T

It can be shown that: $I^T(\theta) \leq I_n(\theta)$, i.e., a generic estimator T will not bear more information than the random sample \Rightarrow MLE maximizes the information carried by the estimator (i.e., it reaches the upper bound).

* ML Estimators are efficient

T_n -efficient $\stackrel{\text{DEF}}{\iff}$ T_n -consistent and asymptotically normal AND

$$\frac{I_n(T_n)}{I_n(\theta^*)} \xrightarrow{P} 1 \text{ as } n \rightarrow \infty \Rightarrow \sqrt{I_n(\theta^*)} (T_n - \theta^*) \xrightarrow{D} N(0,1)$$

"Efficiency" means that the estimator contains the maximal amount of information supplied by the random sample about the value of the parameter θ^*

From a practical standpoint, we exploit the fact that $T_n \xrightarrow{P} \theta^* \Rightarrow \frac{I_n(T_n)}{I_n(\theta^*)} \xrightarrow{P} 1$

and consider the form: $\sqrt{I_n(T_n)} (T_n - \theta^*) \xrightarrow{D} (0,1)$

In this case, an approximation of the Fisher Information $I_n(T_n)$ is given

by the sampled (i.e., observed) numerical derivative $-\frac{\partial^2 l}{\partial \theta^2}(\bar{I}_n)$, i.e., we replace $I_n(\hat{\theta})$ with:

$$I_{\text{obs}}(\hat{\theta}) \triangleq -\frac{\partial^2 l}{\partial \theta^2}(\hat{\theta})$$

It can be shown that, as $n \rightarrow \infty$, the standard error of the MLE is:

$$SE = \frac{1}{\sqrt{I_{\text{obs}}(\hat{\theta})}}$$

and an approximate 95% confidence interval is $(\hat{\theta} - 2SE, \hat{\theta} + 2SE)$.

Ex.: $X_i \sim \text{Exp}(\lambda)$ $i=1, 2, \dots, n$

$$\theta \triangleq \lambda \quad l(\theta) = \log\left(\prod_{i=1}^n \theta e^{-\theta x_i}\right) = n \log \theta - n\theta \bar{x}$$

$$\text{with } \bar{x} \triangleq \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{We have: } \frac{\partial l}{\partial \theta} = \frac{n}{\theta} - n\bar{x} = 0 \Leftrightarrow \hat{\theta} = \frac{1}{\bar{x}}$$

$$\frac{\partial^2 l}{\partial \theta^2} = -\frac{n}{\theta^2} \Rightarrow I_n(\hat{\theta}) = n\bar{x}^2 \Rightarrow SE = \frac{1}{\bar{x}\sqrt{n}}$$

Ex.: $X \sim B(n, p)$

$$\theta \triangleq p \quad l(\theta) \triangleq \log(\theta^x (1-\theta)^{n-x}) = x \log \theta + (n-x) \log(1-\theta)$$

$$\frac{\partial l}{\partial \theta} = \frac{x}{\theta} - \frac{n-x}{1-\theta} = 0 \Leftrightarrow \hat{\theta} = \frac{x}{n}$$

$$\frac{\partial^2 l}{\partial \theta^2} = -\frac{x}{\theta^2} - \frac{n-x}{(1-\theta)^2} \Rightarrow I_n(\hat{\theta}) = \frac{n^2}{x} + \frac{n^2(n-x)}{(n-x)^2} = \frac{n^3}{x(n-x)}$$

$$= \frac{n}{\hat{\theta}(1-\hat{\theta})}$$

$$\Rightarrow SE = \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

(14)

* ML Estimators are approximately Bayesian

Let us consider the general form of the Bayes' theorem:

$$f_{Y|X}(y|x) = \frac{f_{X|Y}(x|y) f_Y(y)}{\int_{-\infty}^{+\infty} f_{X|Y}(x|y) f_Y(y) dy}$$

Let us assume: $Y = T_n = \theta \Rightarrow f_{\theta|X}(\theta|x) = \frac{f_X(x|\theta) f_{\theta}(\theta)}{\int_{-\infty}^{+\infty} f_X(x|\theta) f_{\theta}(\theta) d\theta}$

Let us call: $\pi(\theta) \triangleq f_{\theta}(\theta) \triangleq$ "prior distribution" \Rightarrow It represents our knowledge before seeing the data

$f_{\theta|X}(\theta|x) \triangleq$ "posterior distribution"

Because the likelihood function $L(\theta) \propto f_X(x|\theta)$, we can write:

$$f_{\theta|X}(\theta|x) = \frac{L(\theta) \pi(\theta)}{\int_{-\infty}^{+\infty} L(\theta) \pi(\theta) d\theta}$$

It can be proved that, for large samples, $f_{\theta|X}(\theta|x) \sim N\left(\hat{\theta}, \frac{1}{-l''(\hat{\theta})}\right)$

where $l''(\theta) \triangleq \frac{\partial^2 l}{\partial \theta^2}$

Finally, one can observe the following property:

$$X_1 \sim N(\mu, \sigma_1^2) \quad X_2 \sim N(\mu, \sigma_2^2)$$

In this case, we define:

$$\theta \triangleq \mu \quad \ell(\theta) \triangleq -\frac{(x_1 - \theta)^2}{2\sigma_1^2} - \frac{(x_2 - \theta)^2}{2\sigma_2^2}$$

The MLE is obtained by imposing:

$$\frac{\partial \ell}{\partial \theta} = \frac{x_1 - \theta}{\sigma_1^2} + \frac{x_2 - \theta}{\sigma_2^2} = 0 \Leftrightarrow \hat{\theta} = \frac{\frac{1}{\sigma_1^2} x_1 + \frac{1}{\sigma_2^2} x_2}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}}$$

\Rightarrow Denoted with $w_1 \triangleq \frac{1/\sigma_1^2}{1/\sigma_1^2 + 1/\sigma_2^2}$ $w_2 \triangleq \frac{1/\sigma_2^2}{1/\sigma_1^2 + 1/\sigma_2^2}$, we have:

$T_w = w_1 X_1 + w_2 X_2 \Rightarrow T_w$ is the weighted mean of the RVs X_1 and X_2 \square

References:

Textbook: ch.5 (section 5.1 and 5.2)

ch.7 (section 7.2 and 7.3)

ch.8 (section 8.1, 8.2, 8.3, and 8.4)

