**Last time:** ① Examples of Finding MLE:

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\arg\max} \; L(\theta; x_1, \cdots, x_N) \text{ or } l(\theta; x_1, \cdots, x_N)$$

where $\bar{X}_1, \cdots, \bar{X}_N \overset{iid}{\sim} f_{\bar{X}}(x; \theta_0)$ for some $\theta_0 \in \Theta$

② $$\theta_0 = \underset{\theta \in \Theta}{\arg\max} \; E\Big[l(\theta; \bar{X})\Big] \quad \text{where } \bar{X} \sim f_{\bar{X}}(x; \theta_0)$$

where $l(\theta; x) = \ln\Big(f_{\bar{X}}(x; \theta)\Big)$

By the WLLN:

$$\frac{1}{N}\sum_{i=1}^{N} l(\theta; x_i) \xrightarrow{P} E\Big[l(\theta; \bar{X})\Big] \quad \text{for all } \theta \in \Theta$$

③ Bias, Variance and MSE:

$$MSE(\hat{\theta}) = Var(\hat{\theta}) + Bias(\hat{\theta})^2$$

where $Bias(\hat{\theta}) = E[\hat{\theta}] - \theta_0 = E[\hat{\theta} - \theta_0]$

$\hat{\theta}$ unbiased $\iff Bias(\hat{\theta}) = 0$ for all $\theta_0 \in \Theta$

$$\searrow MSE(\hat{\theta}) = Var(\hat{\theta})$$

There is a trade-off between bias and variance.

Today : ① Efficiency + Cramer-Rao Lower Bound

② Consistency

③ Bayesian Estimation

① Efficiency : See previous lecture notes

"Lecture-Nov09-Sheng.pdf"

② Consistency :

An estimator $\hat{\theta}_N$ based on $\bar{x}_1, \cdots, \bar{x}_N$ is consistent if $\hat{\theta}_N \xrightarrow{P} \theta_0$ :

For every $\varepsilon > 0$, $\lim\limits_{N \to \infty} P\left(|\hat{\theta}_N - \theta_0| > \varepsilon\right) = 0$

Example 1: $\bar{x}_1, \bar{x}_2, \cdots \overset{iid}{\sim} Ber(\theta_0)$

let $\hat{\theta}_N = \frac{1}{N} \sum\limits_{i=1}^{N} \bar{x}_i$

$$P\left(|\hat{\theta}_N - \theta_0| > \varepsilon\right) = P\left(|\hat{\theta}_N - \theta_0|^2 > \varepsilon^2\right)$$

Markov's Inequality $\longrightarrow \leq \dfrac{E\left[|\hat{\theta}_N - \theta_0|^2\right]}{\varepsilon^2}$

$E[\hat{\theta}_N] = \theta_0 \longrightarrow = \dfrac{Var(\hat{\theta}_N)}{\varepsilon^2}$

$$= \frac{\theta_0(1-\theta_0)}{N\varepsilon^2} \longrightarrow 0 \quad as \ N \to \infty$$

So $\hat{\theta}_N$ is consistent.

Example 2: $\bar{x}_1, \bar{x}_2, \cdots \overset{iid}{\sim} N(\mu, \sigma^2)$

let $\hat{\mu}_N = \frac{1}{N} \sum\limits_{i=1}^{N} \bar{x}_i \sim N\left(\mu, \frac{\sigma^2}{N}\right)$

$$P\left(|\hat{\mu}_N - \mu| > \varepsilon\right)$$

$$= P\left(\frac{|\hat{\mu}_N - \mu|}{\frac{\sigma}{\sqrt{N}}} > \frac{\varepsilon}{\frac{\sigma}{\sqrt{N}}}\right)$$

$$= P\left(|Z_N| > \frac{\sqrt{N}\varepsilon}{\sigma}\right) \qquad \left[Z_N = \frac{\hat{\mu}_N - \mu}{\frac{\sigma}{\sqrt{N}}} \sim N(0, 1)\right]$$

$$= 2\left(1 - \Phi\left(\frac{\sqrt{N}\,\varepsilon}{\sigma}\right)\right) \longrightarrow 0 \qquad \left[\begin{array}{l}\Phi(\cdot) = \text{standard normal CDF} \\ \lim_{z \to \infty} \Phi(z) = 1\end{array}\right]$$

Let $\hat{\sigma}_N^2 = \frac{1}{N}\sum_{i=1}^{N}\left(x_i - \bar{\bar{x}}\right)^2$

$$E\left[\hat{\sigma}_N^2\right] = \frac{N-1}{N}\sigma^2 \quad, \quad Var\left(\hat{\sigma}_N^2\right) = \frac{2(N-1)}{N^2}\sigma^4$$

$$\Rightarrow MSE\left(\hat{\sigma}_N^2\right) = E\left[|\hat{\sigma}_N^2 - \sigma^2|^2\right] = \frac{2N-1}{N^2}\sigma^4$$

So $P\left(|\hat{\sigma}_N^2 - \sigma^2| > \varepsilon\right) \leq \dfrac{E\left[|\hat{\sigma}_N^2 - \sigma^2|^2\right]}{\varepsilon^2}$

$$= \frac{2N-1}{N^2}\frac{\sigma^4}{\varepsilon^2} \longrightarrow 0 \qquad \text{as } N \to \infty$$

<span style="color:red">Biased but consistent example!</span>

<span style="color:red">How about unbiased but not consistent?</span>

<span style="color:red">For example: $\bar{x}_1, \bar{x}_2, \dots \overset{iid}{\sim}$ with mean $\mu$, Let $\hat{\mu}_N = \bar{x}_N$
then $E[\hat{\mu}_N] = \mu$ but $\hat{\mu}_N \nrightarrow \mu$ as $N \to \infty$.</span>

# ③ Bayesian Estimation

When estimating a parameter $\theta_0 \in \Theta$, the MLE framework makes almost no assumptions.

It is often the case that some values of $\theta \in \Theta$ are a priori more likely than others.

Bayes rule allows us to incorporate information like this.

Let $A, B$ be two events, since

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

then we have

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

More generally, in Bayesian estimation, the unknown parameter $\theta$ is treated as a random variable, and prior information for $\theta$ is encoded in a distribution $f_\Theta(\theta)$:

$$\Theta \sim f_\Theta(\theta)$$

The observable random variable $\underline{x}$ is related to $\theta$ through the conditional distribution:

$$\underline{x} \sim f_{\underline{x}}(x | \Theta = \theta)$$

Given an observation $\overline{X}=x$, we update our model for $\Theta$ according to the Bayes rule:

$$f_{\Theta}(\theta \mid \overline{X}=x) = \frac{f_{\overline{X}}(x \mid \Theta=\theta)\, f_{\Theta}(\theta)}{f_{\overline{X}}(x)}$$

likelihood → $f_{\overline{X}}(x \mid \Theta=\theta)$

prior ← $f_{\Theta}(\theta)$

posterior ↑ $f_{\Theta}(\theta \mid \overline{X}=x)$

marginal likelihood:

$$f_{\overline{X}}(x) = \int_{\theta \in \Theta} f_{\overline{X}}(x \mid \Theta=\theta)\, f_{\Theta}(\theta)\, d\theta$$

How do we turn the posterior $f_{\Theta}(\theta \mid \overline{X}=x)$ into an estimate of $\theta$? There are two popular approaches:

① Posterior mean or MMSE:

$$\hat{\Theta}_{MMSE} = E[\Theta \mid \overline{X}=x] = \underset{\theta}{\arg\min}\ E[|\Theta - \theta|^2 \mid \overline{X}=x]$$

② Posterior mode or MAP:

$$\hat{\Theta}_{MAP} = \underset{\theta \in \Theta}{\arg\max}\ f_{\Theta}(\theta \mid \overline{X}=x)$$

$$= \underset{\theta \in \Theta}{\arg\max}\ f_{\overline{X}}(x \mid \Theta=\theta)\, f_{\Theta}(\theta)$$

If $f_{\Theta}(\theta) = \dfrac{1}{|\Theta|}$ for all $\theta \in \Theta$, then

$$\hat{\Theta}_{MAP} = \underset{\theta \in \Theta}{\arg\max}\ f_{\overline{X}}(x \mid \Theta=\theta) = \hat{\Theta}_{MLE}$$

Example: $X_1, \ldots, X_N \overset{iid}{\sim} Ber(\theta)$, $\theta \in Uniform(0,1)$

observations: $X_1 = x_1, \ldots, X_N = x_N$

$$f_{\underline{X}}(x_1, x_2, \ldots, x_N \mid \Theta = \theta) f_{\Theta}(\theta) = \prod_{i=1}^{N} \theta^{x_i} (1-\theta)^{1-x_i} \cdot 1$$

$$= \theta^{\sum_{i=1}^{N} x_i} (1-\theta)^{N - \sum_{i=1}^{N} x_i}$$

$$= \theta^{S_N} (1-\theta)^{N-S_N} \quad \text{where } S_N \triangleq \sum_{i=1}^{N} x_i$$

$$f_{\underline{X}}(x_1, \ldots, x_N) = \int_0^1 f_{\underline{X}}(x_1, \ldots, x_N \mid \Theta = \theta) f_{\Theta}(\theta) \, d\theta$$

$$= \int_0^1 \theta^{S_N} (1-\theta)^{N-S_N} d\theta$$

$$= B(S_N + 1, N - S_N + 1) \quad \text{where } B(x, y) = \frac{\Gamma(x) \Gamma(y)}{\Gamma(x+y)} \text{ is}$$

the Beta function

therefore

$$f_{\Theta}(\theta \mid x_1, \ldots, x_N) = \frac{f_{\underline{X}}(x_1, \ldots, x_N \mid \Theta = \theta) f_{\Theta}(\theta)}{f_{\underline{X}}(x_1, \ldots, x_N)}$$

$$= \frac{\theta^{S_N} (1-\theta)^{N-S_N}}{B(S_N + 1, N - S_N + 1)}$$

This is the PDF of the $Beta(S_N + 1, N - S_N + 1)$ distribution. We compute explicitly the marginal $f_{\underline{X}}(x_1, \ldots, x_N)$ above, but this was not necessary to find the posterior. Indeed,

$$f_{\underline{X}}(x_1, \ldots, x_N | \Theta = \theta) f_{\Theta}(\theta) = \theta^{S_N}(1-\theta)^{N-S_N}$$

tells us the PDF of the posterior distribution of $\Theta$.

Since $\Theta | x_1, \ldots, x_N \sim \text{Beta}(S_N + 1, N - S_N + 1)$ has

mean $E[\Theta | x_1, \ldots, x_N] = \dfrac{S_N + 1}{(S_N + 1) + (N - S_N + 1)} = \dfrac{S_N + 1}{N + 2}$

mode $\dfrac{(S_N + 1) - 1}{(S_N + 1) + (N - S_N + 1) - 2} = \dfrac{S_N}{N}$

then $\hat{\Theta}_{MMSE} = \dfrac{S_N + 1}{N + 2}$ and $\hat{\Theta}_{MAP} = \dfrac{S_N}{N}$