

## LECTURE 23

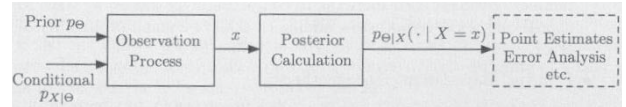
- **Readings:** Section 9.1 through p. 470  
Section 9.2 through p. 482
- Course VI Underground Guide Evaluations  
<https://sixweb.mit.edu/student/evaluate/6.041-s2010>  
until 11:59pm on May 16

## Lecture outline

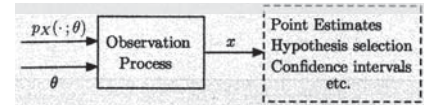
- Classical statistical inference
- Classical parameter estimation
  - bias
  - consistency
- Maximum likelihood estimation
- Confidence intervals

## Bayesian vs. classical inference

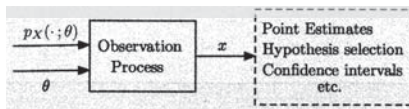
- Want to make inferences about *parameter(s)*  $\theta$
- Bayesian:  $\theta$  is a realization of random variable  $\Theta$



- Classical:  $\theta$  is unknown but not random



## Classical parameter estimation



- $\theta$  is not a random variable
- Estimator  $\hat{\Theta}$  is a (different) random variable for each  $\theta$
- Must care about performance for all possible  $\theta$ 
  - Cannot average over  $\theta$  because  $\theta$  is not random!
- Robustness?
  - No sensitivity to prior because there is no prior
  - Still depends entirely on model for observation generation

## Terminology/properties of estimators

- $\hat{\Theta}_n$ : estimator of  $\theta$  from  $X_1, X_2, \dots, X_n$
- **Estimation error:**  $\tilde{\Theta}_n = \hat{\Theta}_n - \theta$
- **Bias:**  $b_\theta(\hat{\Theta}_n) = E_\theta[\tilde{\Theta}_n] = E_\theta[\hat{\Theta}_n] - \theta$
- **Unbiased:**  $b_\theta(\hat{\Theta}_n) = 0$  or  $E_\theta[\hat{\Theta}_n] = \theta$  (for **all**  $\theta$ )
- **Asymptotically unbiased:**  $\lim_{n \rightarrow \infty} b_\theta(\hat{\Theta}_n) = 0$  for **all**  $\theta$
- **Consistent:**  $\hat{\Theta}_n$  converges in probability to  $\theta$ , for **all**  $\theta$
- **Mean squared error (MSE):** (function of  $\theta$ )
 
$$E_\theta[(\hat{\Theta}_n - \theta)^2] = \text{var}_\theta(\hat{\Theta}_n - \theta) + (E_\theta[\hat{\Theta}_n - \theta])^2$$

$$= \text{var}_\theta(\hat{\Theta}_n) + (b_\theta(\hat{\Theta}_n))^2$$

## Examples

- Parameter  $\theta$  is known to be positive
- $X$  is exponentially distributed with parameter  $\theta$
- $E_\theta[X] = 1/\theta$ , so  $\hat{\Theta} = 1/X$  seems reasonable
  - $b_\theta(\hat{\Theta}) = \int_0^\infty \frac{1}{x} \theta e^{-\theta x} dx = \infty$
- $P_\theta\left(X < \frac{1}{\theta \ln 2}\right) = \frac{1}{2}$ , so  $\hat{\Theta} = \frac{1}{X \ln 2}$  seems reasonable
- $n$  obs:  $\hat{\Theta}_n = n/(X_1 + X_2 + \dots + X_n)$  seems reasonable
  - $(X_1 + X_2 + \dots + X_n)/n$  converges in probability to  $1/\theta$
  - $n/(X_1 + X_2 + \dots + X_n)$  converges in probability to  $\theta$

## Maximum likelihood (ML) estimation

- Pick  $\theta$  that “makes data most likely”:

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} p_X(x; \theta)$$

- Compare to (Bayesian) MAP estimation:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} \frac{p_{X|\Theta}(x|\theta) p_\Theta(\theta)}{p_X(x)}$$

- Advanced properties:
  - **Invariance** to invertible change of parameterization  $\zeta = h(\theta)$ , i.e., ML estimate will be  $\hat{\zeta}_{\text{ML}} = h(\hat{\theta}_{\text{ML}})$
  - **Consistency:** ML estimates  $\hat{\Theta}_n$  from  $(X_1, X_2, \dots, X_n)$  form a consistent sequence
  - **Asymptotic normality:**  $(\hat{\Theta}_n - \theta)/\sigma(\hat{\Theta}_n)$  approaches standard normal

**Example: Exponential distribution parameter**

- $X_1, X_2, \dots, X_n$ : indep. with  $\text{exponential}(\theta)$  distribution
- ML estimate is

$$\begin{aligned}\hat{\theta}_{\text{ML}} &= \arg \max_{\theta} f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n \theta e^{-\theta x_i} \\ &= \arg \max_{\theta} \left( n \log \theta - \theta \sum_{i=1}^n x_i \right)\end{aligned}$$

Find critical point by differentiating w.r.t.  $\theta$ :

$$\begin{aligned}\hat{\theta}_{\text{ML}} &= \frac{n}{x_1 + x_2 + \dots + x_n} \\ \hat{\Theta}_n &= \frac{n}{X_1 + X_2 + \dots + X_n}\end{aligned}$$

**Example: Normal distribution mean**

- $X_1, X_2, \dots, X_n$ : indep. with mean  $\theta$  and variance 1
- ML estimate is

$$\begin{aligned}\hat{\theta}_{\text{ML}} &= \arg \max_{\theta} f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-(x_i - \theta)^2/2} \\ &= \arg \min_{\theta} \sum_{i=1}^n (x_i - \theta)^2\end{aligned}$$

Find critical point by differentiating w.r.t.  $\theta$ :

$$\begin{aligned}\hat{\theta}_{\text{ML}} &= \frac{x_1 + x_2 + \dots + x_n}{n} \\ \hat{\Theta}_n &= \frac{X_1 + X_2 + \dots + X_n}{n}\end{aligned}$$

**Confidence interval (CI)**

- An estimate  $\hat{\Theta}_n$  is only a number with no reliability
- $[\hat{\Theta}_n^-, \hat{\Theta}_n^+]$  is a  $1 - \alpha$  **confidence interval** for  $\theta$  when

$$\mathbf{P}_{\theta}(\hat{\Theta}_n^- \leq \theta \leq \hat{\Theta}_n^+) \geq 1 - \alpha \quad \text{for all } \theta$$

- often  $\alpha = 0.01$  or  $0.05$
- interpretation is subtle
- sometimes want one-sided confidence interval

$$(-\infty, \hat{\Theta}_n^+] \quad \text{or} \quad [\hat{\Theta}_n^-, \infty)$$

**Example: Exponential distribution parameter**

- $X$ : exponential with parameter  $\theta$
- Analyze  $[a/X, b/X]$  as a confidence interval for  $\theta$ :

$$\begin{aligned}\mathbf{P}_{\theta}\left(\frac{a}{X} \leq \theta \leq \frac{b}{X}\right) &= \mathbf{P}_{\theta}\left(\frac{a}{\theta} \leq X \leq \frac{b}{\theta}\right) \\ &= \int_{a/\theta}^{b/\theta} \theta e^{-\theta x} dx \\ &= -e^{-\theta x} \Big|_{x=a/\theta}^{x=b/\theta} = e^{-a} - e^{-b}\end{aligned}$$

No dependence on  $\theta$ , so have a confidence interval

- Example:  $[\frac{1}{4X}, \frac{4}{X}]$  is a 0.76 confidence interval ("76% CI")

**Example: Normal distribution mean**

- $X_1, X_2, \dots, X_n$ : indep. with mean  $\theta$  and variance 1
- $\hat{\Theta}_n = (X_1 + \dots + X_n)/n$  has mean  $\theta$ , variance  $1/n$
- Analyze  $[\hat{\Theta}_n - \delta, \hat{\Theta}_n + \delta]$  as a confidence interval for  $\theta$ :

$$\begin{aligned}\mathbf{P}_{\theta}(\hat{\Theta}_n - \delta \leq \theta \leq \hat{\Theta}_n + \delta) &= \mathbf{P}_{\theta}(-\delta \leq -\hat{\Theta}_n + \theta \leq \delta) \\ &= \mathbf{P}_{\theta}(-\delta \leq \hat{\Theta}_n - \theta \leq \delta) \\ &= \mathbf{P}_{\theta}\left(-\delta\sqrt{n} \leq \underbrace{\sqrt{n}(\hat{\Theta}_n - \theta)}_{\text{standard normal}} \leq \delta\sqrt{n}\right)\end{aligned}$$

- For 0.95 CI, want  $\Phi(-\delta\sqrt{n}) = 0.025$ , or  $\delta\sqrt{n} \approx 1.96$

$$\left[ \hat{\Theta}_n - \frac{1.96}{\sqrt{n}}, \hat{\Theta}_n + \frac{1.96}{\sqrt{n}} \right]$$

**Mean and variance estimation with unknown distribution**

- $X_1, X_2, \dots, X_n$ : i.i.d., mean  $\mu$ , variance  $\sigma^2$  (both unknown)

- Sample mean  $M_n = \frac{X_1 + \dots + X_n}{n}$  well studied

- $\mathbf{E}[M_n] = \mu$  (unbiased)
- $M_n$  converges in probability to  $\theta$  (WLLN, consistency)

- Natural variance estimate  $\bar{S}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - M_n)^2$

- $\mathbf{E}[\bar{S}_n^2] = \frac{n-1}{n} \sigma^2$  (biased, asymptotically unbiased)
- $\hat{S}_n^2 = \frac{n}{n-1} \bar{S}_n^2$  unbiased
- $\bar{S}_n^2$  and  $\hat{S}_n^2$  both consistent