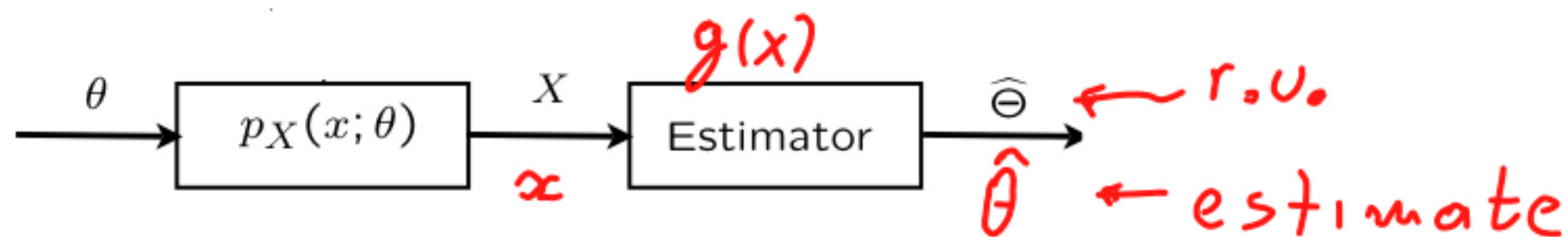


## LECTURE 20: An introduction to classical statistics

- Unknown **constant**  $\theta$  (not a r.v.)
- if  $\theta = \mathbb{E}[X]$ : estimate using the sample mean  $(X_1 + \cdots + X_n)/n$ 
  - terminology and properties
- Confidence intervals (CIs)
  - CIs using the CLT
  - CIs when the variance is unknown
- Other uses of sample means
- Maximum Likelihood estimation

## Classical statistics

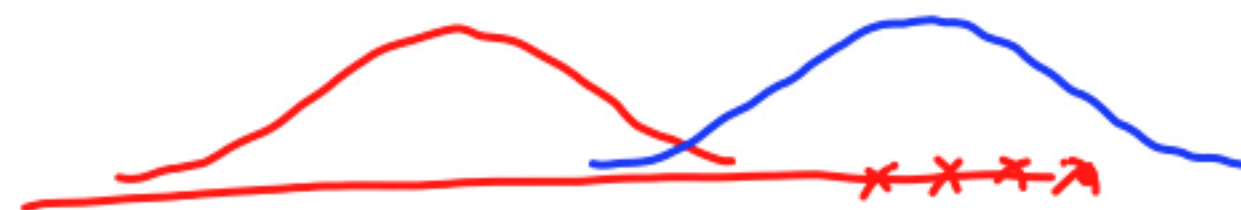
- Inference using the Bayes rule:  
unknown  $\Theta$  and observation  $X$  are both random variables
  - Find  $p_{\Theta|X}$
- Classical statistics: unknown constant  $\theta$



$$P_{\Theta} \quad P_{X|\Theta}$$

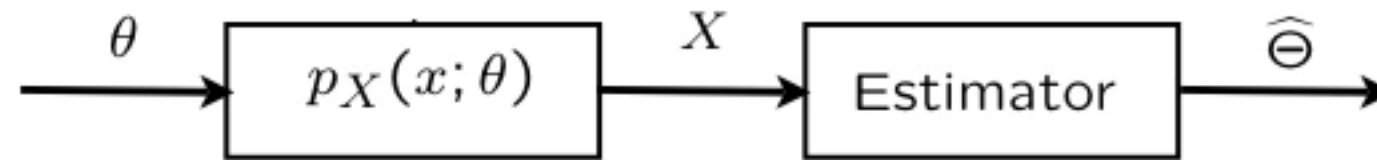
$$P_{X|\Theta}(x|\theta)$$

- also for vectors  $X$  and  $\theta$ :  $p_{X_1, \dots, X_n}(x_1, \dots, x_n; \theta_1, \dots, \theta_m)$
- $p_X(x; \theta)$  are NOT conditional probabilities;  $\theta$  is NOT random
- mathematically: many models, one for each possible value of  $\theta$



## Problem types in classical statistics

- Classical statistics: unknown constant  $\theta$



- Hypothesis testing:  $H_0 : \theta = 1/2$  versus  $H_1 : \theta = 3/4$
- Composite hypotheses:  $H_0 : \theta = 1/2$  versus  $H_1 : \theta \neq 1/2$
- Estimation: design an **estimator**  $\hat{\Theta}$ , to “keep estimation **error**  $\hat{\Theta} - \theta$  small”

*Art!* •

## Estimating a mean

- $X_1, \dots, X_n$ : i.i.d., mean  $\theta$ , variance  $\sigma^2$

$$\hat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \dots + X_n}{n}$$

$\hat{\Theta}_n$ : estimator (a random variable)

### Properties and terminology:

- $E[\hat{\Theta}_n] = \theta$  (unbiased)

for all  $\theta$

这个等式成立

- WLLN:  $\hat{\Theta}_n \xrightarrow{i.p.} \theta$  (consistency)  
for all  $\theta$

$$\hat{\Theta} = g(x)$$
$$E[\hat{\Theta}] = \sum_x g(x) P_x(x; \theta)$$

- mean squared error (MSE):  $E[(\hat{\Theta}_n - \theta)^2] = \text{var}(\hat{\Theta}_n) = \frac{\sigma^2}{n}$  .

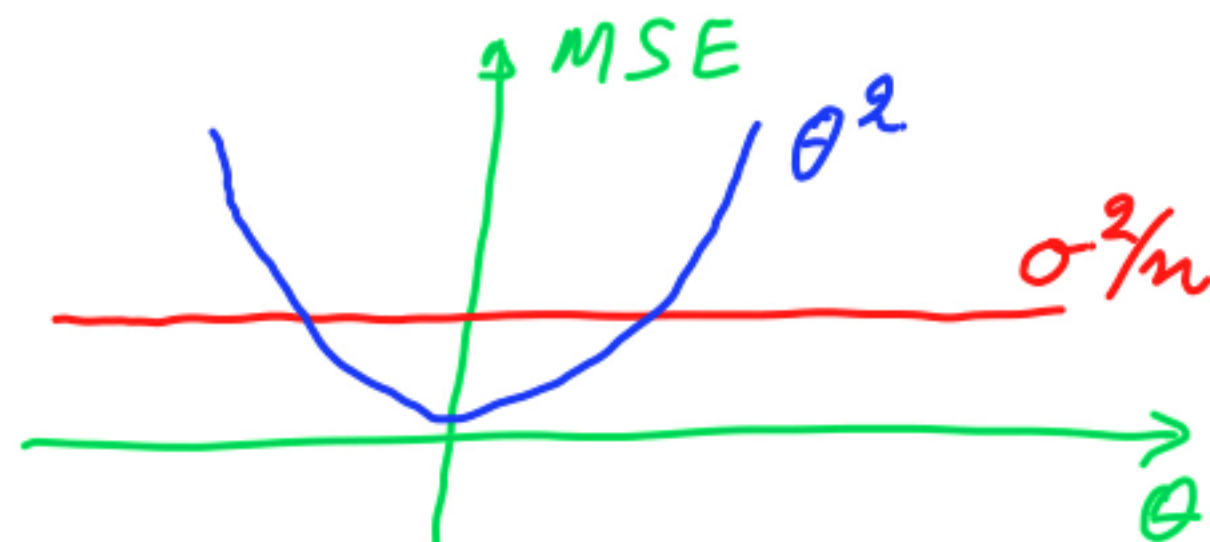
## On the mean squared error of an estimator

- For any estimator, using  $E[Z^2] = \text{var}(Z) + (E[Z])^2$ :  $Z = \hat{\Theta} - \theta$

$$E[(\hat{\Theta} - \theta)^2] = \text{var}(\hat{\Theta} - \theta) + \underbrace{(E[\hat{\Theta} - \theta])^2}_{\text{bias}^2} = \text{var}(\hat{\Theta}) + (\text{bias})^2$$

$$\hat{\Theta}_n = M_n : \text{MSE} = \sigma^2/n + 0$$

$$\hat{\Theta} = 0 : \text{MSE} = 0 + \theta^2$$



- $\sqrt{\text{var}(\hat{\Theta})}$  is called the **standard error**





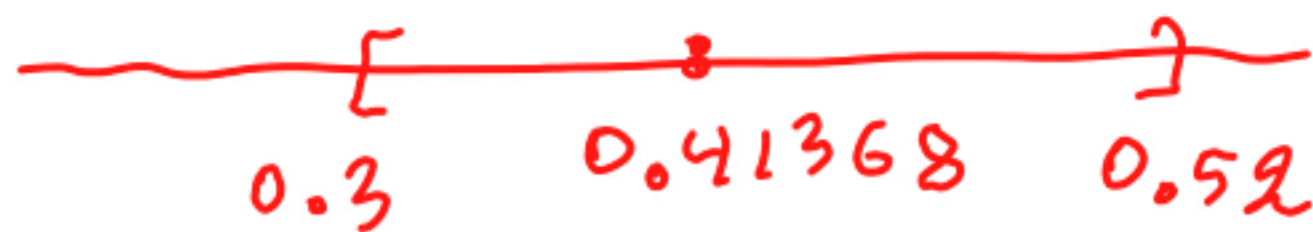
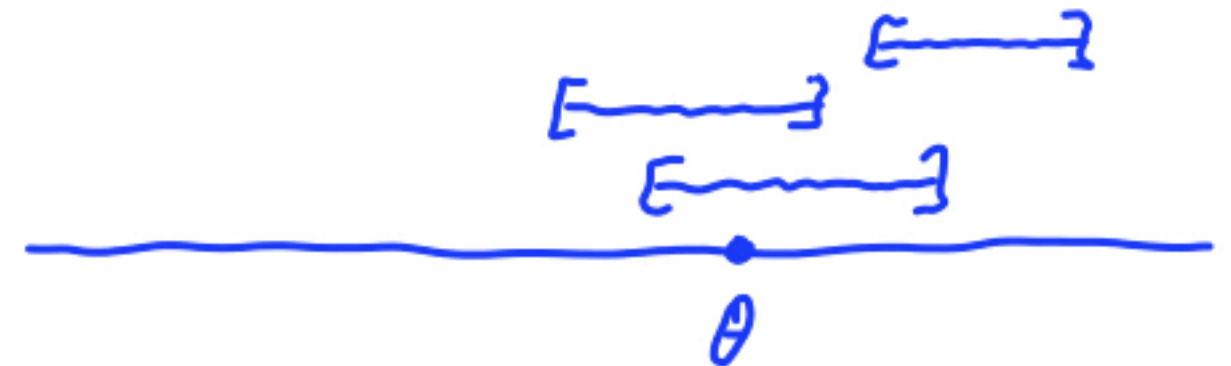
## Confidence intervals (CIs)

- The value of an estimator  $\hat{\Theta}$  may not be informative enough

- An <sup>95%</sup>  $1 - \alpha$  **confidence interval** is an interval  $[\hat{\Theta}^-, \hat{\Theta}^+]$ ,

s.t.  $P(\hat{\Theta}^- \leq \theta \leq \hat{\Theta}^+) \geq 1 - \alpha$ , for all  $\theta$

- often  $\alpha = 0.05$ , or 0.025, or 0.01
- interpretation is subtle



$$P(0.3 < \theta < 0.52) \geq 0.95$$

not in one specific sample

## CI for the estimation of the mean

$$\hat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \dots + X_n}{n}$$

95%

normal tables:  $\Phi(1.96) = 0.975 = 1 - 0.025$

90%

$$\Phi(1.645) = 0.95$$

$$P\left(\frac{|\hat{\Theta}_n - \theta|}{\sigma/\sqrt{n}} \leq 1.96\right) \approx 0.95 \quad (\text{CLT})$$

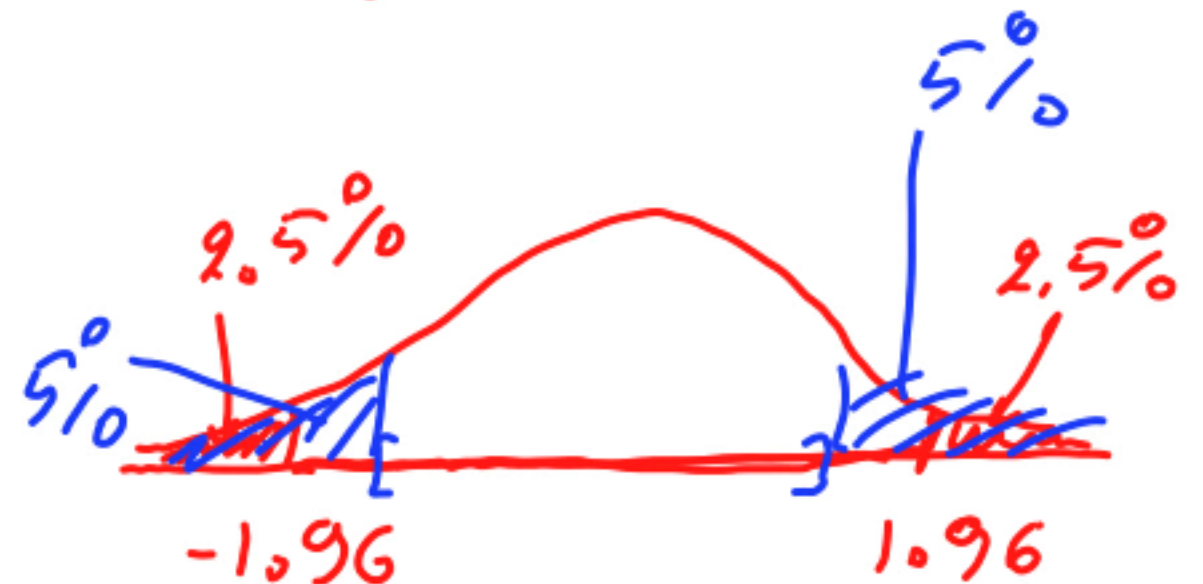
std of the sample mean(theta\_n)  $Z_n = \frac{(S_n - n\mu)/n}{\sqrt{n}\sigma/n}$

$$P\left(\hat{\Theta}_n - \frac{1.96\sigma}{\sqrt{n}} \leq \theta \leq \hat{\Theta}_n + \frac{1.96\sigma}{\sqrt{n}}\right) \approx 0.95$$

$\Theta^-$

$\Theta^+$

iid  $\theta$   $\sigma^2$



## Confidence intervals for the mean when $\sigma$ is unknown

$$\hat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \cdots + X_n}{n}$$

$$P\left(\hat{\Theta}_n - \frac{1.96\sigma}{\sqrt{n}} \leq \theta \leq \hat{\Theta}_n + \frac{1.96\sigma}{\sqrt{n}}\right) \approx 0.95$$

- **Option 1:** use **upper bound** on  $\sigma$ 
  - if  $X_i$  Bernoulli:  $\sigma \leq 1/2$
- **Option 2:** use **ad hoc estimate** of  $\sigma$ 
  - if  $X_i$  Bernoulli:  $\hat{\sigma} = \sqrt{\hat{\Theta}_n(1 - \hat{\Theta}_n)}$

$$\sigma = \sqrt{\theta(1 - \theta)}$$



## Confidence intervals for the mean when $\sigma$ is unknown

$$\mathbf{P}\left(\widehat{\Theta}_n - \frac{1.96 \sigma}{\sqrt{n}} \leq \theta \leq \widehat{\Theta}_n + \frac{1.96 \sigma}{\sqrt{n}}\right) \approx 0.95$$

- **Option 3:** Use **sample mean estimate** of the variance

- Two approximations involved here:
  - CLT: approximately normal
  - using estimate of  $\sigma$
- correction for second approximation ( $t$ -tables) used when  $n$  is small

Start from  $\sigma^2 = \mathbf{E}[(X_i - \theta)^2]$

$$\frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 \rightarrow \sigma^2$$

(but do not know  $\theta$ )

$$\frac{1}{n-1} \sum_{i=1}^n (X_i - \widehat{\Theta}_n)^2 \rightarrow \sigma^2$$

## Other natural estimators

- $\theta_X = \mathbf{E}[X]$        $\widehat{\Theta}_X = \frac{1}{n} \sum_{i=1}^n X_i$

- $\theta = \mathbf{E}[g(X)]$        $\widehat{\Theta} = \frac{1}{n} \sum_{i=1}^n g(X_i)$

- $v_X = \text{var}(X) = \mathbf{E}[(X - \theta_X)^2]$

$$\widehat{v}_X = \frac{1}{n} \sum_{i=1}^n (X_i - \widehat{\Theta}_X)^2$$

- $\text{cov}(X, Y) = \mathbf{E}[(X - \theta_X)(Y - \theta_Y)]$   
 $(x_i, y_i)$

$$\widehat{\text{cov}}(X, Y) = \frac{1}{n} \sum_{i=1}^n (X_i - \widehat{\Theta}_X) (Y_i - \widehat{\Theta}_Y)$$

- $\rho = \frac{\text{cov}(X, Y)}{\sqrt{v_X} \cdot \sqrt{v_Y}}$

$$\widehat{\rho} = \frac{\widehat{\text{cov}}(X, Y)}{\sqrt{\widehat{v}_X} \cdot \sqrt{\widehat{v}_Y}}$$

- next steps: find the distribution of  $\widehat{\Theta}$ , MSE, confidence intervals,...

## Maximum Likelihood (ML) estimation

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

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

- also applies when  $x$ ,  $\theta$  are vectors or  $x$  is continuous

- compare to Bayesian posterior:  $p_{\Theta|X}(\theta | x) = \frac{p_{X|\Theta}(x | \theta) \overbrace{p_{\Theta}(\theta)}^{\text{constant}}}{\cancel{p_X(x)}}$

- interpretation is very different

## Comments on ML

- maximize  $p_X(x; \theta)$
- maximization is usually done numerically
- if have  $n$  i.i.d. data drawn from model  $p_X(x; \theta)$ , then, under mild assumptions:
  - consistent:  $\hat{\Theta}_n \rightarrow \theta$
  - asymptotically normal:  $\frac{\hat{\Theta}_n - \theta}{\sigma(\hat{\Theta}_n)} \rightarrow N(0, 1)$  (CDF convergence)
- analytical and simulation methods for calculating  $\hat{\sigma} \approx \sigma(\hat{\Theta}_n)$ 
  - hence confidence intervals  $\mathbf{P}\left(\hat{\Theta}_n - 1.96 \hat{\sigma} \leq \theta \leq \hat{\Theta}_n + 1.96 \hat{\sigma}\right) \approx 0.95$
  - asymptotically “efficient” (“best”)



## ML estimation example: parameter of binomial

- $K$ : binomial with parameters  $n$  (known), and  $\theta$  (unknown)

$k$

$$p_K(k; \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

$$\log \left[ \binom{n}{k} \right] + k \log \theta + (n-k) \log (1-\theta)$$

$$0 + \frac{k}{\theta} - \frac{n-k}{1-\theta} = 0 \Rightarrow k - \cancel{k\theta} = n\theta - \cancel{k\theta}$$

$$\hat{\theta}_{\text{ML}} = \frac{k}{n} \quad \hat{\Theta}_{\text{ML}} = \frac{K}{n}$$

- same as MAP estimator with uniform prior on  $\theta$



## ML estimation example — normal mean and variance

- $X_1, \dots, X_n$ : i.i.d.,  $N(\mu, v)$   $f_X(x; \mu, v) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi v}} \exp \left\{ -\frac{(x_i - \mu)^2}{2v} \right\}$

minimize  $\frac{n}{2} \log v + \sum_{i=1}^n \frac{(x_i - \mu)^2}{2v}$

– minimize w.r.t.  $\mu$ :  $\hat{\mu} = \frac{x_1 + \dots + x_n}{n}$

$$\frac{1}{v} \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \sum x_i = n\mu$$

– minimize w.r.t.  $v$ :  $\hat{v} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$

$$\frac{n}{2} \cdot \frac{1}{v} \Rightarrow \sum_{i=1}^n \frac{(x_i - \mu)^2}{2v} = 0$$