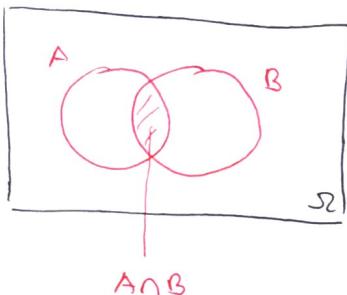


# ASTROSTATISTICS

from distance

## PROBABILITY FUNDAMENTALS



- Sum Rule:  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- Conditional Probability:

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

$$\Rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \begin{matrix} \text{prob of } \\ \text{given } B \end{matrix}$$

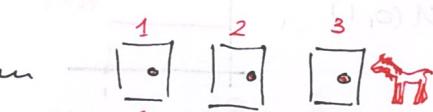
- Law of Total Probability:  $\sum_i P(A|B_i) P(B_i)$  use partitioning  $\cup_i \{B_i\} = \Omega$   
 $B_i \cap B_j = \emptyset \quad i \neq j$

$$\textcircled{a} \text{ Baye's Theorem: } P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{P(A|B) P(B)}{P(A)}$$

- Example: Monty Hall problem

1	2	3	opened door	stay	change
C	G	G	2/3	C	G
G	C	G	3	G	C
G	G	C	2	G	C

better to change!



Monty always opens door w goat  
 $\rightarrow$  2 is ruled out  $\rightarrow$  50-50 chance

- ② Extreme MH:  $N = 1000$  doors, you pick one  
 host opens 998 other doors, all w goat

$$H_i = \text{car is behind } i \quad P(H_i) = 1/1000$$

$d_k$  = "data" observation that  $N-2$  doors have goats

$$P(H_1 | d_k) = \frac{P(d_k | H_1) P(H_1)}{P(d_k)} \quad P(d_k | H_1) = 1/N-1$$

$$P(d_k | H_i) = \begin{cases} 1 & k=i \\ 0 & \text{otherwise} \end{cases}$$

$$\text{LTP: } P(d_k) = \sum_{i=1}^N P(d_k | H_i) P(H_i)$$

$$\bullet P(d_k) = P(d_k | H_1) P(H_1) + P(d_k | H_k) P(H_k) = \frac{1}{N-1} \cdot \frac{N-1}{N} + \frac{1}{N} = \frac{1}{N}$$

$$P(H_1 | d_k) = \frac{\frac{1}{N-1} \cdot \frac{1}{N}}{\frac{1}{N}} = \frac{1}{N} \Rightarrow P(H_{1k} | d_k) = \frac{N-1}{N}$$

- Continuous distributions / random variables

notation:  $P_X(x) \rightarrow \cancel{P(x)} P(x)$

•  $P(x, y) = \text{joint pdf} = P(x|y) P(y) = P(y|x) P(x) = P(x) P(y)$

↑  
if  $x$  and  $y$  are independent

◦ Marginals  $P(x) = \int P(x, y) dy$

use continuous version of LTP  $P(x) = \int P(x|y) P(y) dy$

⇒ Bayes Rule:  $P(y|x) = \frac{P(x|y) P(y)}{P(x)} = \frac{P(x|y) P(y)}{\int P(x|y) P(y) dy}$

- Transformations of RVs

Suppose  $X$  has  $p(x)$ , have  $y = \Phi(x)$  invertible, so  $x = \Phi^{-1}(y)$

$$P(y) = P(x) \left| \frac{dx}{dy} \right| = P(\Phi^{-1}(y)) \left| \frac{d\Phi^{-1}(y)}{dy} \right|$$

→ example:  $X \sim U(0, 1)$

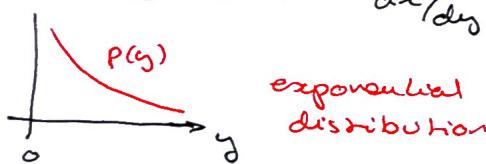
$$P(x) = \begin{cases} 1 & 0 < x < 1 \\ 0 & \text{otherwise} \end{cases} = U(x|0, 1)$$

$$y = -\ln x$$

$$x = e^{-y}$$

$$\frac{dx}{dy} = -e^{-y} \rightarrow P(y) = U(e^{-y}|0, 1) \cdot e^{-y}$$

$$= \begin{cases} e^{-y} & 0 \leq y < \infty \\ 0 & \text{otherwise} \end{cases}$$



- "Delta" method  $Y = \Phi(X)$   $\text{Var}(X) = G_x^2$   $E(X) = X_0$

wants to know  $G_Y^2 = \text{Var}(Y)$

- Definitions:  $X \sim P(x)$  random variable distributed by  $P$

$$E[X] = \int x P(x) dx$$

$$E[h(x)] = \int h(x) P(x) dx$$

$$\text{Var}[X] = \int (x - E[X])^2 P(x) dx = E[X^2] - (E[X])^2$$

- if we have  $P(x, y)$

$$E[X] = \int \int x P(x, y) dx dy$$

$$\text{Cov}(X, Y) = \int (x - E[X])(y - E[Y]) P(x, y) dx dy$$

→ if here bilinearity:  $S = \sum_{i=1}^n a_i x_i$      $T = \sum_{j=1}^k b_j y_j$

$$\text{Cov}(S, T) = \text{Cov}\left(\sum a_i x_i, \sum b_j y_j\right) = \sum_i \sum_j a_i b_j \text{Cov}(x_i, y_j)$$

~~Area~~

- Delta method want to know  $\text{Var}(Y)$   $x \sim P(x)$   $y = g(x)$

do a Taylor expansion  $y = g(x) = g(x_0) + \frac{dg}{dx} \Big|_{x=x_0} (x - x_0) + \dots$

$$\Rightarrow \text{Var}(Y) = E(Y - E(Y)) = \left| \frac{dg}{dx} \Big|_{x=x_0} \right| \text{Var}(X)$$

accuracy depends on "wigginess" of  $f(x)$  compared to size of variance. e.g. for  $\log$  not so good

## Limit theorems

- LLN law of large numbers  $X_1, \dots, X_N$  iid rvs w  $M = E[X_i]$

$$\bar{X}_N = \frac{1}{N} \sum_{i=1}^N X_i \quad \begin{matrix} \text{random} \\ \text{sample} \\ \text{mean} \end{matrix} \quad \begin{matrix} \uparrow \\ \text{independent identically} \\ \text{distributed random variables} \end{matrix}$$

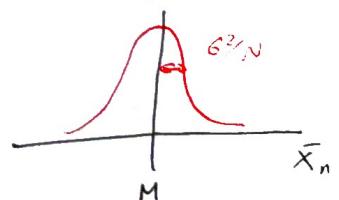
$$\text{as } N \rightarrow \infty \quad \forall \varepsilon > 0 \quad P(|\bar{X}_N - M| > \varepsilon) \rightarrow 0$$

sample mean converges on the true mean

- CLT central limit theorem  $G^2 = \text{Var}(X_i)$   $\frac{1}{\sqrt{2\pi}} e^{-t^2/2G^2}$  gaussian

$$P(\sqrt{n}(\bar{X}_n - M) \leq zG) \rightarrow \Phi(z) = \int_{-\infty}^z N(t|0, 1) dt$$

$$\text{ie } P(\bar{X}_n) \rightarrow N(\bar{X}_n | M, G^2/N)$$



Asymptotic Normality

# STATISTICAL INFERENCE

- suppose  $\underline{D} = \{X_1, \dots, X_N\}$   $X_i \stackrel{iid}{\sim} P(x|G)$

Want to estimate  $G$  from  $\underline{D}$ , different ideas:

① Sample Mean  $\frac{1}{N} \sum_{i=1}^N X_i$

② Sample mean of first  $k < N$ :  $\frac{1}{k} \sum_{i=1}^k X_i$

③  $\frac{1}{N-1} \sum_{i=1}^N X_i$  doesn't have a name...

④  $\frac{1}{2} [\min(x) + \max(x)]$  Midrange

⑤ median

⑥ bin data  $\rightarrow$  take mode

ESTIMATOR  $\hat{G} = f(\underline{\hat{X}})$  for ESTIMAND  $M_G = \int x p_G(x) dx$

$\hat{G} = f(\underline{\hat{X}})$  is an estimator for  $G$

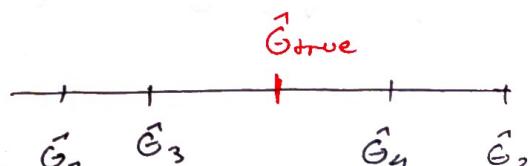
Frequency properties

$\hat{G}$  is unbiased if  $E[\hat{G}] = G$

- Imagine  $k$  experiments:  $j = 1, \dots, k$   $\underline{X}_j = (X_{1,j}, \dots, X_{n,j})$

get  $\underline{X}_1, \dots, \underline{X}_j, \dots, \underline{X}_k$

$$\hat{G}_1 = \hat{G}(\underline{X}_1) \quad \hat{G}_j = \hat{G}(\underline{X}_j)$$



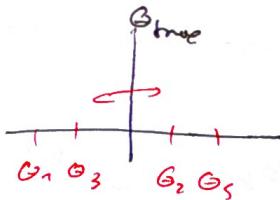
$$\frac{1}{k} \sum_{j=1}^k \hat{G}_j \rightarrow G \text{ as } k \rightarrow \infty$$

BIAS  $B(\hat{G}) = E(\hat{G}) - G$

consistency as you gather more data  $N \rightarrow \infty$

$$\hat{G} \rightarrow G, \forall \varepsilon > 0 \quad P(|\hat{G} - G| \geq \varepsilon) \rightarrow 0$$

Efficiency:  $\text{Var}(\hat{G})$



Mean Square Error  $MSE(\hat{G}) = E[(\hat{G}(x) - G)^2]$

$$= E[\underbrace{(\hat{G} - E(\hat{G}))^2}_{\text{bias as variables}} + \underbrace{E(\hat{G}) - G}_{\text{variance}}^2]$$

$$\dots = \text{Var}(\hat{G}) + B(\hat{G})^2$$

so for unbiased estimators smallest MSE for most efficient

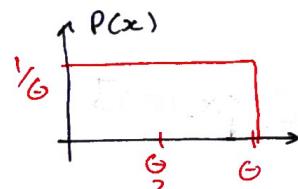
- MVUE among all unbiased estimators  $\rightarrow$  smallest variance

$\Rightarrow$  Bias-Variance Tradeoff not always unbiased best

	$G=1, N=10$	$G^2/N$	$\frac{(b-a)^2}{12}$
sample mean	0.10	0.10	0.10
unbiased	$\frac{\pi^2 G^2}{24 \ln(G)}$	0.18	$\frac{6G^2}{(N+2)(N+1)}$
Variance			0.06

for uniform distribution  
unbiased better!

(2) Uniform distribution



$$x_i \sim U(0, G)$$

$\rightarrow$  estimator for  $G$   $\hat{G} = \max(\underline{x}) < G$  biased

$$\hat{\phi} = 2\underline{x} \rightarrow E[\underline{x}] = \frac{G}{2} \rightarrow 2E[2\underline{x}] = G$$

$$\text{true: } \underline{x} = (0.32, 0.16, 0.97, 2.77, 7.06)$$

$$\hat{\phi}(\underline{x}) = 4.63 \quad \hat{G}(\underline{x}) = 7.06 \quad \hat{\phi} \text{ obviously wrong!}$$

maximum likelihood actually  $\hat{G}$ !

## Deriving Estimators

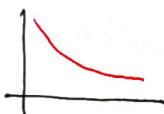
① Method of Moments: match theoretical moments of distribution w/ sample moment

$$X_1, \dots, X_N \stackrel{iid}{\sim} P_G(x) \quad \text{Theoretical} \quad M_k = \mathbb{E}[X^k] = \int x^k P_G(x) dx = \text{function}(G)$$

$$\hat{\mu}_k = \frac{1}{N} \sum_{i=1}^N x_i^k \quad \text{average of } x^k$$

$$\hat{\mu}_k = \mu_k(G) \rightarrow \text{solve for } G$$

e.g.:  $X_i \sim \exp(\lambda)$



$$p(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{E}[X] = 1/\lambda, \bar{x} = \frac{1}{N} \sum x_i \Rightarrow \hat{\lambda} = 1/\bar{x}$$

② Method of Least Squares:  $\mu = \int x P_G(x) dx$

$$\hat{\mu}_{LS} = \underset{\mu}{\operatorname{argmin}} \left[ \sum_{i=1}^N (x_i - \mu)^2 \right] \quad \text{ss} \quad \rightarrow \hat{\mu}_{LS} = \frac{1}{N} \sum_{i=1}^N x_i$$

③ Maximal Likelihood:

let  $X_i$  has pdf  $f(x_i | \theta)$ ,  $X_i \sim f(x_i | \theta)$

$\theta$  = true value, view as function of  $X$   
sampling distribution of  $X$

experiment  $\rightarrow \bar{x} = (x_1, \dots, x_N)$

$$P(\bar{x} | \theta) = \prod_{i=1}^N f(x_i | \theta) = L(\theta)$$

Likelihood Fnc. of  $\theta$  ↑

Sometimes useful to calculate:

$$e(\theta) = \ln(L(\theta)) = \sum_{i=1}^N \log[f(x_i | \theta)] \quad \text{LOG LIKELIHOOD}$$

$$\Rightarrow \text{MLE} \quad \hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} [L(\theta)] = \underset{\theta}{\operatorname{argmax}} [e(\theta)]$$

need to know how to model the problem

data  $\bar{x}$   
fixed

Properties of MLE if model assumptions are true, i.e.  $L(\theta)$  true

- consistent  $\hat{\theta}_{MLE} \rightarrow \theta_{true}$  as  $N \rightarrow \infty$
- asymptotically unbiased  $E[\hat{\theta}_{MLE}] \rightarrow \theta_{true}$  as  $N \rightarrow \infty$
- asymptotically normal  $(\hat{\theta}_{MLE} - \theta) \xrightarrow{d} N(0, I^{-1})$  if  $\theta$  vector  
then  $I$  matrix

$$\Rightarrow I(\theta) = E \left[ \left( \frac{\partial \log[L(\theta)]}{\partial \theta} \right)^2 \right] = -E \left[ \frac{\partial^2 \ell(\theta)}{\partial \theta^2} \right]$$

EXPECTED  
FISHER  
INFORMATION

curvature of likelihood function, 2nd derivative

$$\hat{I} = - \left. \frac{\partial^2 \ell(\theta)}{\partial \theta^2} \right|_{\theta = \hat{\theta}_{MLE}}$$

OBSERVED  
FISHER  
INFORMATION

$$\hat{I} \rightarrow I(\theta) \text{ as } N \rightarrow \infty$$

RLB

- Efficiency: Cramér-Rao Lower Bound  $\text{Var}(\hat{\theta}_{MLE}) \geq I^{-1}(\theta)$   
asymptotically  $N \rightarrow \infty$  achieves the RLB so efficiency  $I^{-1}$
- Functionally invariant  $\alpha = g(\theta) \Rightarrow \hat{\alpha}_{MLE} = g(\hat{\theta}_{MLE})$

### Multiparameter MLE

$$X_i \sim_{iid} f(x_i | \vec{\theta}) \quad L(\vec{\theta}) = P(\vec{x} | \vec{\theta}) = \prod_{i=1}^N f(x_i | \vec{\theta})$$

$\theta$  always vector here, all dofs the same except Fisher info a matrix

$$I_{jk}(\theta) = E \left[ - \frac{\partial^2 \log \ell(\theta)}{\partial \theta_j \partial \theta_k} \right]$$

HESSIAN MATRIX

in astrophysics  $\hat{I}_{jk} = - \left. \frac{\partial^2 \log \ell(\theta)}{\partial \theta_j \partial \theta_k} \right|_{\theta = \hat{\theta}_{MLE}} \approx I_{jk}(\theta) \quad N \rightarrow \infty$

large sample limit  
in freq this only care about  $\theta$

-  $\text{Var}(\hat{\theta}_{MLE,i}) \approx (I^{-1})_{ii}$  by CRLB

$\Rightarrow \text{Cov}(\hat{\theta}_{MLE}) \rightarrow I^{-1}$  as  $N \rightarrow \infty$

- For  $X_i \sim N(\mu, G^2)$   $i = 1, \dots, N$  Example for Normal RVs

$$L(\mu, G^2) = \prod_{i=1}^N N(X_i | \mu, G^2) = \prod_{i=1}^N \frac{1}{G\sqrt{2\pi}} e^{-\frac{(X_i - \mu)^2}{2G^2}}$$

$$\ell(G) = \sum_{i=1}^N -\frac{1}{2} \log(2\pi G^2) - \frac{1}{2} \frac{(X_i - \mu)^2}{G^2}$$

overbar,  
not vector.  
↓

①

suppose  $G^2$  is unknown  $\frac{\partial \ell}{\partial \mu} = \sum_{i=1}^N \frac{(X_i - \mu)}{G^2} = 0 \rightarrow \hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N X_i = \bar{X}$

$$\Rightarrow \frac{\partial^2 \ell}{\partial \mu^2} = -\frac{N}{G^2} \Rightarrow I = -\mathbb{E}\left[\frac{\partial^2 \ell}{\partial \mu^2}\right] = \frac{N}{G^2} \Rightarrow I^{-1} = G^2/N$$

so by CRLB  $\text{Var}(\hat{\mu}_{MLE}) \geq G^2/N$

could have calculated more directly

$$\begin{aligned} \text{Var}(\hat{\mu}_{MLE}) &= \text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{N} \sum X_i\right) = \text{Cov}\left(\frac{1}{N} \sum X_i, \frac{1}{N} \sum X_j\right) \\ &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \underbrace{\text{Cov}(X_i, X_j)}_{\text{S.i.i. } G^2 \text{ bc of iid assumption}} = \frac{G^2}{N} \quad \text{saturates the CRLB} \end{aligned}$$

② suppose  $G^2$  is also unknown  $\boldsymbol{\theta} = \begin{pmatrix} \mu \\ G^2 \end{pmatrix}$  multiparameter now!

$$\frac{\partial \ell}{\partial G^2} = \sum_{i=1}^N -\frac{1}{2} \frac{1}{G^2} + \frac{1}{2} \frac{(X_i - \mu)^2}{(G^2)^2} = 0$$

$$\Rightarrow \sum_{i=1}^N \frac{1}{G^2} = \sum_{i=1}^N \frac{(X_i - \mu)^2}{(G^2)^2} \Rightarrow \hat{G^2}_{MLE} = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$$

Want to calculate the Fisher matrix:

$$\frac{\partial^2 \ell}{\partial (G^2) \partial \mu} = \sum_{i=1}^N -\frac{(X_i - \mu)}{(G^2)^2} \Rightarrow \mathbb{E}\left[-\frac{\partial^2 \ell}{\partial (G^2) \partial \mu}\right] = \sum_{i=1}^N \frac{\mathbb{E}(X_i - \mu)}{G^4} = 0$$

$$\frac{\partial^2 \ell}{\partial (G^2)^2} = \sum_{i=1}^N \frac{1}{2} \frac{1}{G^4} - \frac{(X_i - \mu)^2}{G^6} \Rightarrow \mathbb{E}\left[\frac{\partial^2 \ell}{\partial (G^2)^2}\right] = -\frac{N}{2G^4}$$

$$\Rightarrow I = \begin{pmatrix} N/G^2 & 0 \\ 0 & N/2G^4 \end{pmatrix} \Rightarrow I^{-1} = \begin{pmatrix} G^2/N & 0 \\ 0 & 2G^4/N \end{pmatrix}$$

Again,  $\text{Var}(\hat{\mu}_{MLE}) = G^2/N$ ,  $\text{Var}(\hat{G^2}_{MLE}) = 2G^4/N$  both saturate the CRLB