

## Lecture Notes 7 - Parametric Point Estimation (continued)

- $X_1, \dots, X_n \sim p(x, \theta)$   $\theta \in \Theta$  -  $\theta$  can be scalar or vector  
- assume vector here
- $\theta = (\theta_1, \dots, \theta_K)$

## 2. Method of Moments

$$\mu_1(\theta) = E[X_i] = \int x p(x, \theta) dx$$

Sample version:  $m_1 = \frac{1}{n} \sum_{i=1}^n X_i$  ; recall  $m_1 \xrightarrow{P} \mu_1(\theta)$  via LLN

- same for 2nd moment

$$\mu_2(\theta) = E[X_i^2] = \int x^2 p(x, \theta) dx$$

Sample version:  $m_2 = \frac{1}{n} \sum_{i=1}^n X_i^2$  ; same reasoning

- for  $k$ th moment

$$\mu_k(\theta) = E[X_i^k] = \int x^k p(x, \theta) dx$$

Sample version:  $m_k = \frac{1}{n} \sum_{i=1}^n X_i^k$  , ——— " ———

- have  $k$  terms that depend on  $\theta$

- via LLN, these terms should be 'close' to sample versions

- equate them to get  $k$  equations

- equate  $\mu_j(\hat{\theta}) = m_j$   $j=1, \dots, k$  and solve

to yield  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3, \dots, \hat{\theta}_k)$

(A): Confused by dm.

- is  $\theta$  a vector in above form

## Example 1

$N(\mu, \sigma^2)$  with  $\theta = (\mu, \sigma^2)$ .

(\*) equate theoretical and sample moments

$$E[X_i] = \beta \quad E[X_i^2] = \text{Var}(X) + (E[X_i])^2 = \sigma^2 + \beta^2$$

- set:-

$$\frac{1}{n} \sum_{i=1}^n X_i = \hat{\beta}$$

$$\frac{1}{n} \sum_{i=1}^n X_i^2 = \hat{\sigma}^2 + \hat{\beta}^2$$

$$\Rightarrow \hat{\beta} = \bar{X}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

W: only works if you can solve the equations

- when you can solve the equations, you get sample estimators
- haven't determined whether they are good estimators
- ~~can~~ currently an algorithm for generating estimators.

### example 2

- Suppose

$$X_1, \dots, X_n \sim \text{Binomial}(K, p)$$

- where  $K$  and  $p$  unknown.

we have:-

$$Kp = \bar{X}_n, \quad \frac{1}{n} \sum_{i=1}^n X_i^2 = Kp(1-p) + Kp^2$$

yielding:- typo

$$\hat{p} = \frac{\bar{X}}{\bar{K}} \quad \hat{K} = \frac{\bar{X}^2}{\bar{X} - \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

(15)

W: note

(1) can be negative  $\Rightarrow$  nonsensical result for estimated coin flips.

- Recall distinction between:-

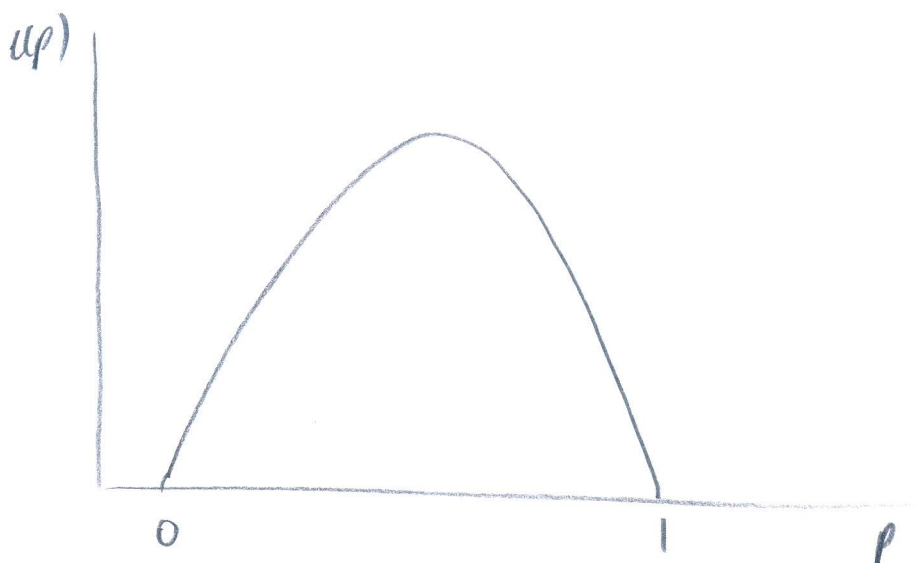
- 1) A procedure for generating estimators
- 2) evaluating their 'quality'

BOTH.

- case where no. of trials and parameter (success prob. are unknown)

- flip coin  $K$  times;  
and we only are given  $X_i$ .  
- no. of heads

- repeat exp., estimate  $K$  and  $p$ .



$$\Rightarrow \hat{p} = \bar{X}$$

in this case MLE and MOM estimators are the same.

#### example 4

$$X_1, \dots, X_n \sim N(\mu, 1)$$

$$L(\mu, \sigma^2) \propto e^{-\frac{n(\bar{X}-\mu)^2}{2\sigma^2}} e^{-\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{2\sigma^2}}$$

$$l(\mu, \sigma^2) = -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2$$

$$\frac{\partial l}{\partial \mu} = 0 \quad \frac{\partial l}{\partial \sigma^2} = 0$$

yields:  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

(A1) - verify this, board unclear.

- messy without applying log-trans.

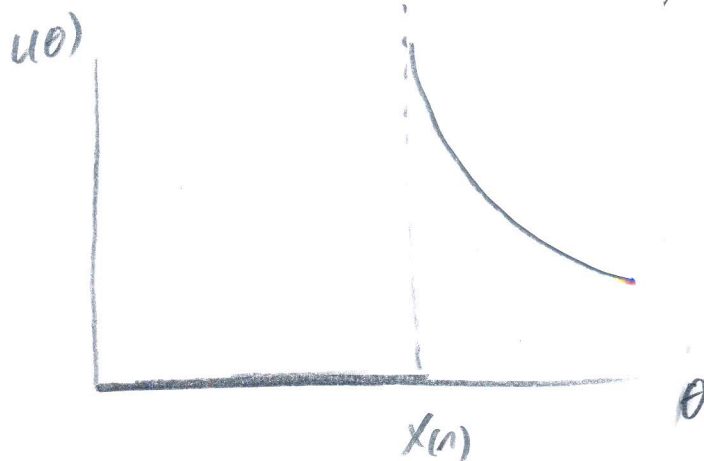
(A2) review calc. ✓

#### example 5

$$X_1, \dots, X_n \sim \text{unif}(0, \theta)$$

$$L(\theta) = \frac{1}{\theta^n} \mathbb{I}(\theta > X_{(n)})$$

$$\hat{\theta} = X_{(n)} \text{ i.e. } \max(X_1, \dots, X_n)$$



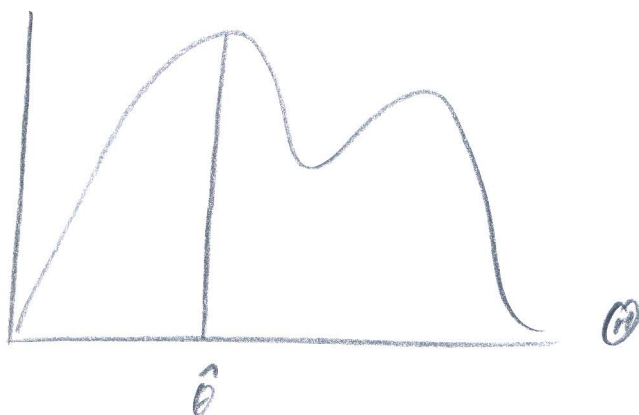
### 3. Maximum likelihood

- Most common in science journals.
- Systematized by Ronald Fisher

$$L(\theta) = p(x_1, \dots, x_n; \theta) = \prod_{i=1}^n p(x_i; \theta)$$

- $\hat{\theta}$  is the point that maximises  $L(\theta)$ .

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$$



- (\*)  $\hat{\theta}$  depends on likelihood function  $L(\theta)$ , which depends on data  $x_1, \dots, x_n$ .

- (\*) Numerically easier to maximise:-

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta) = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

$\log(\cdot)$  - monotone fn

$$\ell(\theta) = \log L(\theta)$$

- (\*) Certain very desirable properties that ML yields an estimator.

- often solved via:-

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = 0 \quad j=1, \dots, k$$

(for sufficiently 'nice' functions)

- classroom settings  $\rightarrow$  analytic sol

- real world  $\rightarrow$  numerical methods.  
(10.725 convex opt.)

Example - Bernoulli

$$x_1, \dots, x_n \sim \text{Ber}(p)$$

$$L(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^S (1-p)^{n-S}$$

$$S = \sum_{i=1}^n x_i$$

$$\ell(p) = S \log(p) + (n-S) \log(1-p)$$



- MOM estimator (for  $\text{Unif}(0, \theta)$ )

$$E[X_i] = \frac{\theta}{2} \quad \text{set} \quad \frac{\theta}{2} = \bar{X} \Rightarrow \hat{\theta} = 2\bar{X}$$

(\*) MLE and MOM yield distinct estimators  $\Rightarrow$  how to assess?

(\*) An important property of MLE.

- introduce some terminology to disclose property.

- suppose  $\theta = (\eta, \xi)$

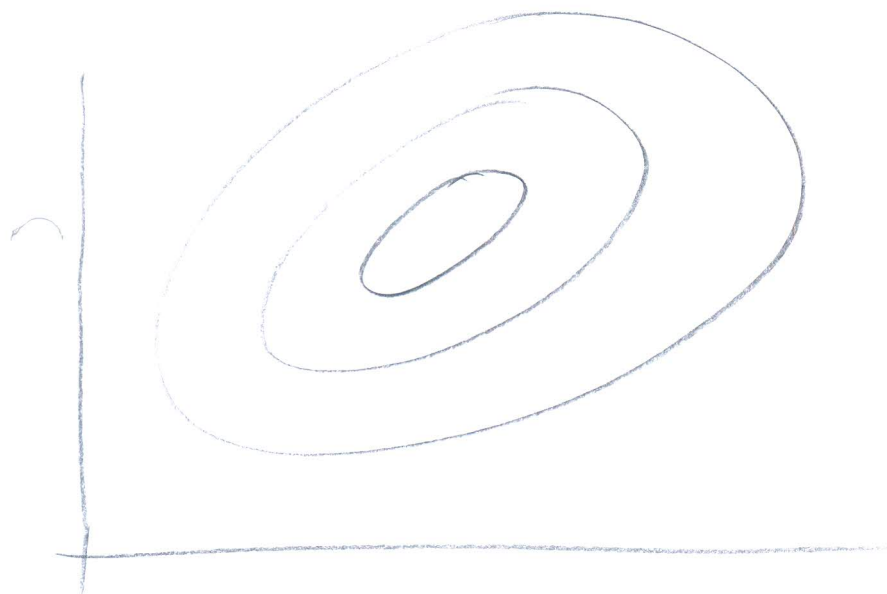
-  $L(\theta)$ .  $\hat{\theta} = (\hat{\eta}, \hat{\xi})$  (analogous to  $\hat{\mu}, \hat{\sigma}$  for normal case)

- focus on likelihood function as function over one parameter 'at a time'.

- Profile likelihood:-  $L(\eta) = \sup_{\xi} L(\eta, \xi)$   
(for  $\eta$ )

- note:-  $\hat{\eta} = \underset{\eta}{\operatorname{argmax}} L(\eta) = \underset{\eta}{\operatorname{argmax}} \left( \sup_{\xi} L(\eta, \xi) \right)$

(A3) - fill in details  
of geometric  
interpretation  
of maximisation



(\*) Overall maximiser of  $L(\theta)$  and maximiser of profile likelihood are equivalent

W. MLE has a 'good' property called equivariance

(\*) suppose I have an arbitrary function of  $\theta$ :- AB

$$\eta = g(\theta)$$

$$\hat{\eta} = g(\hat{\theta})$$

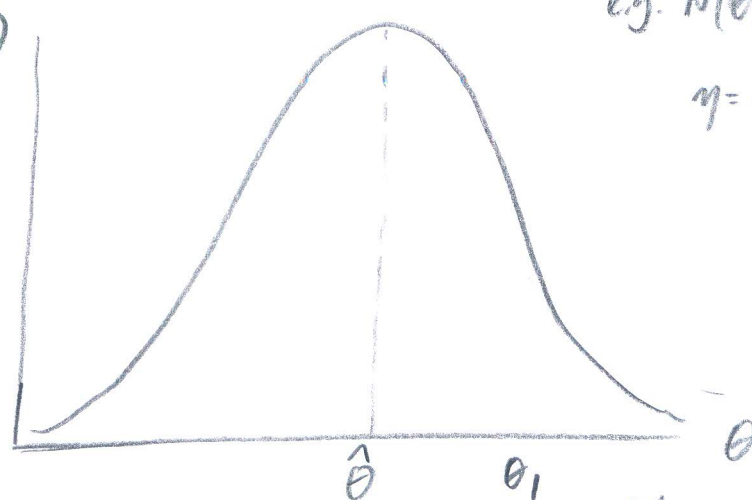
eg.  $N(\theta, 1)$

$$\eta = e^{\theta}$$

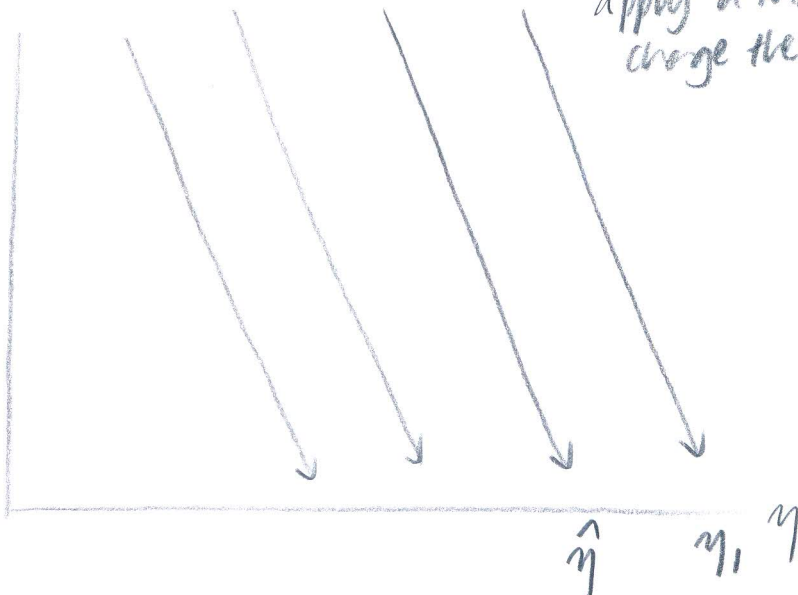
nothing stopping us from modeling distri with a suitable transformation

likelihood is just density fn of  $\eta$  (or  $\hat{\eta}$ ?)

Does not have to be a monotonic transformation/one-to-one



apply a mapping/  
change the scale.



(\*) need to review for clarity if example doesn't clear it up.

example - Bernoulli

$$\hat{p} = \bar{X}$$

$$\psi = \log\left(\frac{p}{1-p}\right)$$

noting  $p = \frac{e^{\psi}}{1+e^{\psi}}$

(\*) for MLE; via invariance,  $\hat{\psi} = \log\left(\frac{\hat{p}}{1-\hat{p}}\right)$

(\*) If you took original likelihood, reexpressed it in terms of  $\psi$ ; we can effect a substitution

(\*) useful if we are interested in a function of the parameter we are trying to estimate.

## equivariance example

example - normal

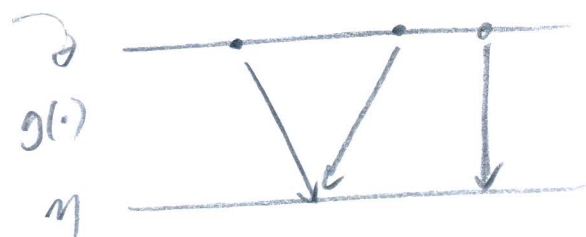
$$N(\mu, \sigma^2) \quad \hat{\mu} = \bar{X} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

via equivariance

$$\hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2} \quad \frac{\hat{\mu}}{\hat{\sigma}} = \frac{\hat{\mu}}{\hat{\sigma}} \quad (\text{plug in})$$

(\*) Argument for equivariance when the transformation is not one-to-one

(\*\*): logic of equivariance, MLE - review



## 4. Bayes Estimator

- this is not Bayesian inference

- want to view this as another algorithm for generating an estimator

- against frequentist ass. of course so far and as

(\*) Assume that the parameter  $\theta$  is a random variable

- define a prior distribution  $p(\theta)$

-  $p(x_1, \dots, x_n | \theta)$  - joint probability of  $x_1, \dots, x_n$ , parametrized by  $\theta$ .

(\*) treating  $\theta$  as a random variable  $\Rightarrow$  we can condition on it (as it assumed to be i.v.)

$p(x_1, \dots, x_n | \theta)$  - joint prob, condit. on  $\theta$ .

- consider:-  $p(x_1, \dots, x_n | \theta) p(\theta) = p(x_1, \dots, x_n, \theta)$



compute the posterior via Bayes theorem

$$p(\theta | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | \theta) p(\theta)}{p(x_1, \dots, x_n)} = \frac{p(x_1, \dots, x_n | \theta) p(\theta)}{\int p(x_1, \dots, x_n | \theta) p(\theta) d\theta}$$

- Recall standard Bayes theorem over r.v. param semantics.

posterior  $\propto$  likelihood  $\times$  prior

- Have a density function for  $\theta$ , can extract a point estimator (but could go further).

Bayes estimator  $\hat{\theta} = E(\theta | x_1, \dots, x_n) = \int \theta p(\theta | x_1, \dots, x_n) d\theta$

(\*) ultimately end up with  $\hat{\theta}$ , an estimator i.e. function of the data.

W: Prefer to see this as purely an algorithm to get an estimator; and not to discuss/endow the semantics of Bayesian formalism with a deeper meaning worthy.

example 7

$x_1, \dots, x_n \sim \text{Bern}(\theta)$   $U(\theta) = \theta^S (1-\theta)^{n-S}$   $S = \sum_{i=1}^n x_i$

$\rightarrow$  (PS): can't see board.

- select prior:-

$$p(\theta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \quad \text{where} \quad \Gamma(\alpha) = \int_0^1 t^{\alpha-1} e^{-t} dt$$

-  $\alpha, \beta$  hyperparam; Bayes est in this case is an infinity of possible est.  
(\*) Beta distribution family

- posterior:

$$p(\theta | x_1, \dots, x_n) \propto \theta^{S+\alpha-1} \cdot (1-\theta)^{n-S+\beta-1}$$

note functional form similarities

(\*) note functional form similarity of prior and posterior

(\*) Hence posterior is a beta distribution  $\text{Beta}(\alpha+S, n-S+\beta)$



AS Review

$$\hat{\theta}_{\text{BAYES}} = E[\theta | x_1, \dots, x_n]$$

$\downarrow \bar{\theta}$

$$= \frac{S + \alpha}{\alpha + \beta + n}$$

$$= (1 - \lambda) \hat{\theta}_{\text{MLE}} + \lambda \bar{\theta}$$

where  $\bar{\theta} = \frac{\alpha}{\alpha + \beta}$

$$\lambda = \frac{\alpha + \beta}{\alpha + \beta + n}$$

$$\hat{\theta}_{\text{MLE}} = \frac{S}{n}$$

(mean of  $\theta$  w.r.t prior distn) in this context.

(\*) Bayes estimator is a convex combination of MLE estimator and the mean of prior

(\*) note collapsement of  $\hat{\theta}_{\text{BAYES}}$  to  $\hat{\theta}_{\text{MLE}}$  or  $\bar{\theta}$

iw: (or) note that at the end of all 3 algorithms, we end up with an estimator:-

$$\hat{\theta} = g(x_1, \dots, x_n)$$