Lecture Notes 7 - Point Estimation - Review

- Basic setup

$X_1, \ldots, X_n \sim p(x; \theta)$. Estimate $\theta = (\theta_1, \ldots, \theta_K)$ via an estimator from data :-

$$\hat{\theta} = \hat{\theta}_n = w(X_1, \ldots, X_n)$$

(ⓡ): param - fixed, unknown (frequentist)

estimator - r.v. as a deterministic function of random data $(X_1, \ldots, X_n)$

(*) Always found this notational point tricky :-

$$\mathbb{E}_\theta(\hat{\theta}) = \int \ldots \int \hat{\theta}(x_1, \ldots, x_n) \, p(x_1; \theta) p(x_2; \theta) \ldots p(x_n; \theta) \, dx_1, \ldots dx_n$$

- expectation of an estimator under a probability distribution parametrised by value of parameter at its true value $\theta$.

(ⓡ): notion of sampling distribution

- An estimator $\hat{\theta}_n$ is a random variable.

- It has a distribution (like any other r.v.) with mean and variance.

(ⓜ): consistency of an estimator :-

$$\hat{\theta}_n \overset{P}{\to} \theta$$   (convergence in probability

as $n \to \infty$         of an estimator to the parameter value)

i.e. $P(|\hat{\theta}_n - \theta| > \epsilon) \overset{n \to \infty}{\longrightarrow} 0$

2. MOM

- equate $R$ sample moments with $R$ theoretical moments.

- solve to get $\hat{\theta}_{MOM} = (\hat{\theta}_1, \ldots, \hat{\theta}_K)$

- sample moments $m_R = \frac{1}{n} \sum_{i=1}^{n} X_i^R$     · Note $m_R \overset{P}{\longrightarrow} \mu_R(\theta)$

via WLLN.

- theoretical moments     $\mu_R(\theta) = \mathbb{E}[X_i^R]$

- yielding MOM estimator for $R$ moments

- we need to posit a distri family right? (i.e. statistical model)

- Yes, have to posit a distri family with unknown parameters.
- i.e. assume OGP of $X_1, ..., X_n$ is [insert family]
- This is a family of densities/distns indexed by parameters
- We use the data to estimate the parameters via an estimator
- The procedure for generating estimators are covered in the notes

## 3. Maximum likelihood

$$\hat{\theta}_{MLE} = \underset{\theta}{argmax} \; L(\theta)$$

$$L(\theta) = p(X_1, ..., X_n; \theta) \overset{IID}{=} \prod_{i=1}^{n} p(X_i; \theta)$$

$$= \underset{\theta}{argmax} \; \ell(\theta)$$

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

(M) - $U(\mu, \sigma^2)$ was derived earlier.
- use cross-term expansion trick with $\bar{x}$
- constant of prop can be discarded as likelihoods equivalent up to a constant of proport.

## (*) Equivariance and profile likelihood

- Profile likelihood
- standard likelihood: -

$$L(\theta) = p(X_1, ..., X_n; \theta) \overset{IID}{=} \prod_{i=1}^{n} p(X_i; \theta)$$

$$\hat{\theta}_{MLE} \text{ obtained by solving } \frac{\partial \ell(\theta)}{\partial \theta_j} = 0 \quad j = 1, ..., R$$

- profile likelihood: -
- partition parameters: - $\theta = (\eta, \xi)$

$$e.g. \; \theta = (\underbrace{\theta_1, \theta_2, ..., \theta_m}_{\eta}, \underbrace{\theta_{m+1}, ..., \theta_R}_{\xi})$$

Then $L(\theta) = L(\eta, \xi)$

profile likelihood
for $\eta$      is likelihood maximised wrt to the other parameter.

That is :- $L(\eta) = \sup\limits_{\xi} L(\eta, \xi)$

· $\hat{\eta}_{MLE} = \underset{\eta}{\operatorname{argmax}} L(\eta) = \underset{\eta}{\operatorname{argmax}} \left\{ \sup\limits_{\xi} L(\eta, \xi) \right\}$

· we can therefore find

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

OR

$$\hat{\xi}_{MLE} = \underset{\xi}{\operatorname{argmax}} L(\eta, \xi)$$

$$\hat{\eta}_{MLE} = \underset{\eta}{\operatorname{argmax}} L(\eta)$$

(*) equivariance of MLE

· If $\eta = g(\theta)$      (i.e. an arbitrary function of parameter)

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

· suppose $g$ is invertible so $\eta = g(\theta)$ and $\theta = g^{-1}(\eta)$

· define $L^*(\eta) = L(\theta)$    where $\theta = g^{-1}(\eta)$

· so for any $\eta$ :-

$$L^*(\hat{\eta}) = L(\hat{\theta}) \geqslant L(\theta) = L^*(\eta) \qquad \text{, i.e. value of parameter}$$

· why? Because $\hat{\theta}$ maximises likelihood (it is MLE).

· Hence $\hat{\eta} = g(\hat{\theta})$ maximises $L^*(\eta)$

· For non-invertible functions (?) ; this is still true if we define

$$L^*(\eta) = \sup\limits_{\theta : \tau(\theta) = \eta} L(\theta) \qquad \text{(i.e. profile likelihood)} \quad ⑦$$

- 9.14 (Theorem) - Wasserman

- Let $\tau = g(\theta)$ be a function of $\theta$
- Let $\hat{\theta}_n$ be the MLE of $\theta$.
- Then $\hat{\tau}_n = g(\hat{\theta}_n)$ is the MLE of $\tau$.

PROOF

- Let $h = g^{-1}$ denote the inverse of $g$
- Then $\hat{\theta}_n = h(\hat{\tau}_n)$
- For any $\tau$; $L(\tau) = \prod_{i=1}^{n} f(x_i; h(\tau)) = \prod_{i=1}^{n} f(x_i; \theta) = L(\theta)$

  where $\theta = h(\tau)$
- Hence, for any $\tau$, $L_n(\tau) = L(\theta) \leq L(\hat{\theta}) = L_n(\hat{\tau})$ ∎

4. Bayes Estimator

- move to Bayesian worldview, only for purposes of generating estimator
- treat $\theta$ as r.v.

$$p(\theta | x_1, \ldots, x_n) = \frac{p(x_1, \ldots, x_n | \theta) p(\theta)}{p(x_1, \ldots, x_n)} = \frac{p(x^n | \theta) p(\theta)}{\int p(x^n | \theta) p(\theta) \, d\theta}$$

Bayes est.: $\hat{\theta}_{BAYES} = \mathbb{E}[\theta | x^n] = \int \theta p(\theta | x^n) \, d\theta$

example 7

$X_1, \ldots, X_n \sim Bern(\theta)$   $L(\theta) = \theta^S (1-\theta)^{n-S}$   $S = \sum_{i=1}^{N} x_i$

prior: $\theta \sim Beta(\alpha, \beta)$

$$p(\theta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} \, dt$$

posterior:

$$p(\theta|x^n) \propto p(x^n|\theta)p(\theta)$$

$$\Rightarrow p(\theta|x^n) \propto \theta^S(1-\theta)^{n-S}\left(\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}\theta^{\alpha-1}(1-\theta)^{\beta-1}\right) \quad \text{(omitting norm constant.)}$$

$$\propto \theta^{S+\alpha-1}(1-\theta)^{n-S+\beta-1} \quad \text{(dropping terms with no }\theta\text{).}$$

Here, with appropriate normalisation :-

$$\theta|x^n \sim Beta(S+\alpha, n-S+\beta)$$

__Bayes estimator__: $E[\theta|x^n] = \int \theta\, p(\theta|x^n)\, d\theta$

- in our case, we are looking for mean of the Beta distri.
- For $Beta(\alpha,\beta)$, mean is $\frac{a}{a+\beta}$.

Hence $\widetilde{\theta} = E[\theta|x^n] = \dfrac{S+\alpha}{(S+\alpha)+(n-S+\beta)} = \dfrac{S+\alpha}{\alpha+\beta+n}$

· Note $\widetilde{\theta} = \dfrac{S+\alpha}{\alpha+\beta+n} = \lambda\bar{\theta} + (1-\lambda)\hat{\theta}_{MLE}$

$$\lambda = \frac{\alpha+\beta}{\alpha+\beta+n} \qquad \bar{\theta} = \frac{n}{\alpha+\beta}$$

⊛ : Properties of MOM, MLE → lore in LN9     (O/S 2)
    (from wasserman)

__Example 8__

(A1) $x_1,\ldots,x_n \sim N(\mu,\sigma^2)$ $\sigma^2$ known

Assume conjugate prior on $\mu$
i.e. $p(\mu) \sim N(m, \tau^2)$    $m, \tau^2$ fixed with

- <u>posterior</u>:
$$p(\mu|x_1,\dots,x_n) = \frac{p(x_1,\dots,x_n|\mu)p(\mu)}{\int p(x_1,\dots,x_n|\mu)p(\mu)\,d\mu}$$

-drop norm. constant.

$$\propto p(x_1,\dots,x_n|\mu)p(\mu)$$

$$= \prod_{i=1}^{n} \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(x_i-\mu)^2\right\} \frac{1}{\tau\sqrt{2\pi}} \exp\left\{-\frac{1}{2\tau^2}(\mu-m)^2\right\}$$

$$= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left\{-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(x_i-\mu)^2\right\} \frac{1}{\tau\sqrt{2\pi}} \exp\left\{-\frac{1}{2\tau^2}(\mu-m)^2\right\}$$

$$= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \left(\frac{1}{\tau\sqrt{2\pi}}\right) \exp\left\{-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(x_i-\bar{x})^2\right\} \cdot$$

$$\exp\left\{-\frac{n}{2\sigma^2}(\bar{x}-\mu)^2\right\} \exp\left\{-\frac{1}{2\tau^2}(\mu-m)^2\right\}$$

· we drop the normalisation
constants, i.e. any terms not containing $\mu$.

$$\Rightarrow p(\mu|x_1,\dots,x_n) \propto \exp\left\{-\frac{n}{2\sigma^2}(\bar{x}-\mu)^2\right\} \exp\left\{-\frac{1}{2\tau^2}(\mu-m)^2\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\frac{n(\bar{x}-\mu)^2}{\sigma^2} + \frac{(\mu-m)^2}{\tau^2}\right)\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\frac{\tau^2 n(\bar{x}-\mu)^2 + \sigma^2(\mu-m)^2}{\sigma^2\tau^2}\right)\right\}$$

$$= \exp\left\{-\frac{1}{2\sigma^2\tau^2}\left(\tau^2 n(\bar{x}-\mu)^2 + \sigma^2(\mu-m)^2\right)\right\}$$

$$= \exp\left\{-\frac{1}{2\sigma^2\tau^2}\left(\tau^2 n\bar{x}^2 - 2\tau^2 n\bar{x}\mu + \tau^2 n\mu^2 + \sigma^2\mu^2 - 2\sigma^2 m\mu + \sigma^2 m^2\right)\right\}$$

$$= \exp\left\{-\frac{1}{2\sigma^2\tau^2}\left((\tau^2 n + \sigma^2)\mu^2 - 2(\tau^2 n\bar{x} + \sigma^2 m)\mu + (\tau^2 n\bar{x}^2 + \sigma^2 m^2)\right)\right\}$$

- Completing the square formula for an arbitrary quadratic:-

$$ax^2 + bx + c = a\left(x + \frac{b}{2a}\right)^2 + \left(c - \frac{b^2}{4a}\right)$$

- We complete square in $\mu$, setting:-

$$a = (\tau^2 n + \sigma^2) \quad b = -2(\tau^2 n \bar{x} + \sigma^2 m) \quad c = (\tau^2 n \bar{x}^2 + \sigma^2 m^2)$$

Yielding:

$$(\tau^2 n + \sigma^2)\mu^2 - 2(\tau^2 n \bar{x} + \sigma^2 m)\mu + (\tau^2 n \bar{x}^2 + \sigma^2 m^2)$$

$$= (\tau^2 n + \sigma^2)\left(\mu - \frac{2(\tau^2 n \bar{x} + \sigma^2 m)}{2(\tau^2 n + \sigma^2)}\right)^2 + \left(\tau^2 n \bar{x}^2 + \sigma^2 m - \frac{4(\tau^2 n \bar{x} + \sigma^2 m)^2}{4(\tau^2 n + \sigma^2)}\right)$$

$$p(\mu \mid x_1, \ldots, x_n) \propto \exp\left\{ -\frac{1}{2\sigma^2\tau^2}\left[(\tau^2 n + \sigma^2)\left(\mu - \frac{\tau^2 n \bar{x} + \sigma^2 m}{\tau^2 n + \sigma^2}\right)^2 + \left(\tau^2 n \bar{x}^2 + \sigma^2 m - \frac{(\tau^2 n \bar{x} + \sigma^2 m)^2}{(\tau^2 n + \sigma^2)}\right)\right] \right\}$$

$$\propto \exp\left\{ -\frac{1}{2\left(\frac{\sigma^2\tau^2}{\tau^2 n + \sigma^2}\right)}\left(\mu - \frac{\tau^2 n \bar{x} + \sigma^2 m}{\tau^2 n + \sigma^2}\right)^2 \right\} \exp\left\{ -\frac{1}{2\sigma^2\tau^2}\left(\tau^2 n \bar{x}^2 + \sigma^2 m - \frac{(\tau^2 n \bar{x} + \sigma^2 m)^2}{\tau^2 n + \sigma^2}\right) \right\}$$

- Discard right hand exp term (contains no $\mu$ terms).

Hence:-

$$p(\mu \mid x_1, \ldots, x_n) \propto \exp\left\{ -\frac{1}{2\left(\frac{\sigma^2\tau^2/n}{\tau^2 + \sigma^2/n}\right)}\left(\mu - \frac{\tau^2 n \bar{x} + \sigma^2 m}{\tau^2 n + \sigma^2}\right)^2 \right\}$$

- Hence the posterior of the mean parameter is Normal, with

$$\mathbb{E}[\mu \mid x_1, \ldots, x_n] = \frac{\tau^2 n}{\tau^2 n + \sigma^2}\bar{x} + \frac{\sigma^2}{\tau^2 n + \sigma^2}m \qquad \text{(convex comb of MLE sample mean / prior mean)}$$

$$\text{var}[\mu \mid x_1, \ldots, x_n] = \frac{\sigma^2\tau^2/n}{\tau^2 + \sigma^2/n}$$

## 5. MSE

(clue is in name).

- mean squared error

$$\mathbb{E}_\theta[(\hat\theta - \theta)^2] = \int \cdots \int (\hat\theta(x_1, \ldots, x_n) - \theta)^2 \, p(x_1;\theta) \ldots p(x_n;\theta) \, dx_1 \ldots dx_n$$

- BiAS $= \mathbb{E}_\theta(\hat\theta) - \theta$

- variance $V = Var_\theta(\hat\theta) = \mathbb{E}_\theta[(\hat\theta - \mathbb{E}_\theta(\hat\theta))^2]$

⓶ · expectation wrt joint distri (that generated the data), not over a distri for $\theta$ !

$$\mathbb{E}_\theta[(\hat\theta - \theta)^2] = \int (\hat\theta(x_1, \ldots, x_n) - \theta)^2 \, p(x^n, \theta) \, dx^n$$

- IID decomposes joint $p(x^n; \theta)$ into $p(x_1;\theta) p(x_2;\theta) \ldots p(x_n;\theta)$

- MSE $= B^2 + V$

- MSE is a 'metric' (preliminary for evaluating estimators)
- unbiasedness → bias $B = \mathbb{E}_\theta[\hat\theta] - \theta = 0 \Rightarrow \mathbb{E}_\theta[\hat\theta] = \theta$
- when this occurs then MSE = variance.
- ⓞ/51 : integrate this with presentation in Bishop and with the various ways of understanding this to get holistic understanding
    - Bishop + interpret. of bias, variance

    - diagram.

## example 10

- note $\hat{S}_n^2 = \frac{n}{n-1} \hat\sigma_{MLE}^2$   (correction for biasedness]

- why is

$\mathbb{E}[S_n^2] = \sigma^2 \longrightarrow$ see 1.3.17. (and HW question for proof)

$X_1, \ldots, X_n \sim N(\mu, \sigma^2)$

- consider
  the following MLE :- $\quad\quad$ ① Don't forget; this assumes
  $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ normal Distri.

  $$\hat{\mu}_{MLE, MOM} = \bar{X}_n \quad\quad \hat{\sigma}^2_{MLE} = \frac{1}{n}\sum_{i=1}^{n}(X_i - \bar{X})^2$$

- $\mathbb{E}[\hat{\mu}_{MLE, MOM}] = \mathbb{E}[\bar{X}_n] = \mu$
- instead of using $\hat{\sigma}^2_{MLE}$ (without bias adjustment); use <u>unbiased</u>
  <u>sample variance $\hat{S}_n^2$</u>

- $\mathbb{E}[\hat{S}_n^2] = \sigma^2$
- $MSE(\hat{\mu}_{MLE}) = MSE(\bar{X}_n) = B^2 + V = V \quad$ as $B = 0 \quad$ (unbiased)

  $$= Var_\mu(\bar{X}_n) = \mathbb{E}_\mu[(\bar{X}_n - \mathbb{E}_\mu[\bar{X}_n])^2] = \mathbb{E}[(\bar{X}_n - \mu)^2]$$

  $$= \frac{\sigma^2}{n}$$

———— A

- $MSE(\hat{S}_n^2) = B^2 + V = V \quad$ as $B = 0 \quad\quad\quad$ (unbiased)

  $$= Var_{\sigma^2}(\hat{S}_n^2)$$

  $$= \mathbb{E}_{\sigma^2}[(\hat{S}_n^2 - \mathbb{E}_{\sigma^2}[\hat{S}_n])^2] = \mathbb{E}[(\hat{S}_n^2 - \sigma^2)^2]$$

  $$= \mathbb{E}[\hat{S}_n^4 - 2\hat{S}_n^2\sigma^2 + \sigma^4]$$

  $$= \mathbb{E}[\hat{S}_n^4] - 2\sigma^2\mathbb{E}[\hat{S}_n^2] + \mathbb{E}[\sigma^4]$$

  $$= \mathbb{E}[\hat{S}_n^4] - 2\sigma^2(\sigma^2) + \sigma^4$$

  $$= \mathbb{E}[\hat{S}_n^4] - \sigma^4 \quad\quad\quad\quad\quad\text{(also obt. via } Var(\hat{S}_n^2) = \mathbb{E}[\hat{S}_n^4] - (\mathbb{E}[\hat{S}_n^2])^2$$

- Not sure how to <u>compute</u> $4^{th}$ <u>moment of $\hat{S}_n$.</u>
- seems like an appeal to $\chi^2$ distri is used. (or a lot of tedious
  $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ derivation).
- Going to put this here :-

WIKI:
  - Distn of sample variance (itself an r.v.)

WIKI: (also from LN1):-

· For $X_1, ..., X_n \sim N(\mu, \sigma^2)$

$$\frac{(n-1)S_n^2}{\sigma^2} \sim \chi^2_{n-1} \qquad \text{NOTE RHS IS:-} \quad \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}$$

· $\frac{(n-1)S_n^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}$ is an r.v. with $\chi^2$ distri $(n-1)$ degrees of free.

(*) For an r.v. Y with $\chi^2$ distri and $M$ degrees of freedom:-

$$\mathbb{E}[Y] = M \qquad Var[Y] = 2M$$

Hence:

$$Var\left[ \frac{n-1}{\sigma^2} S_n^2 \right] = 2(n-1)$$

$$\cdot Var\left[ \frac{n-1}{\sigma^2} S_n^2 \right] = \left( \frac{n-1}{\sigma^2} \right)^2 Var[S_n^2] = 2(n-1)$$

$$\Rightarrow Var_{\sigma^2}(S_n^2) = Var(S_n^2) = \frac{\sigma^4}{(n-1)^2} 2(n-1) = \frac{2\sigma^4}{(n-1)}$$

- As required.