

Youtube lecture 05/10/16

- * note that you will have to 'splice' this lecture with YouTube lecture 03/10/16; as that lecture had 30 mins missing ~~not yet~~ at the end. So there will be deductions/material in this that might feel unfamiliar.
- ② complete

lecture notes 8 - minimax theory (cont.)

- define $X_1, \dots, X_n \sim N(\theta, 1)$
- use L_2 loss: $l(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$
- Theorem: $\hat{\theta} = \bar{X}$ is minimax
- cannot deal with this as have done previously; finding Bayes estimator $\tilde{\theta}$ does not work.
- W: strategy:- find an upper and lower bound (separately) (which are very close together)
- we have:- $R_n = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta})$ (minimax risk)
- can we upper bound R_n ?
- Recall $R_n = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta})$ - inf'mn' over all estimators of the maximum risk.
- select an arbitrary estimator $\hat{\theta}_*$; and find its maximum risk $\sup_{\theta} R(\theta, \hat{\theta}_*)$ is greater than $\inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta})$
- so we upper bound R_n by taking the maximum risk of some arbitrary estimator $\hat{\theta}_*$
- so Any estimator $\hat{\theta}_*$:-
$$R_n \leq \sup_{\theta} R(\theta, \hat{\theta}_*)$$
- can we lower bound R_n ?
- pick any prior π and get the Bayes estimator $\hat{\theta}_\pi$

(*) compute its Bayes risk $B_{\pi}(\hat{\theta}_{\pi}) \rightarrow$ this is a lower bound on R_n

- let's see why:

- minimax risk: $R_n = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta})$ "min / inf over all estimators of maximum risk"

- (i) note maximum is always greater than the average,

- Hence

$$R_n = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}) \geq \inf_{\hat{\theta}} \int R(\theta, \hat{\theta}) \pi(\theta) d\theta = \inf_{\hat{\theta}} B_{\pi}(\hat{\theta}) = B_{\pi}(\hat{\theta}_{\pi})$$

- note that $\inf_{\hat{\theta}} B_{\pi}(\hat{\theta})$ i.e. 'minimum' / infimum of the Bayes risk $B_{\pi}(\hat{\theta})$ for an est $\hat{\theta}$ over all estimators, is precisely the definition of the Bayes estimator. So, by definition, we have upper and lower bounds on the minimax risk:-

$$B_{\pi}(\hat{\theta}_{\pi}) \leq R_n \leq \sup_{\theta} R(\theta, \hat{\theta}_*)$$

We select an arbitrary Bayes estimator (which is a minimise of the Bayes risk), giving you a lower bound on minimax risk R_n . Select an arbitrary estimator and take its max risk, to give you an upper bound.

- poor choices of $B_{\pi}(\hat{\theta}_{\pi})$ and $\sup_{\theta} R(\theta, \hat{\theta}_*)$ will yield very little info on R_n (i.e. difference is very large).

- clever choices of $B_{\pi}(\hat{\theta}_{\pi})$ and $\sup_{\theta} R(\theta, \hat{\theta}_*)$ i.e. close together can yield info on R_n .

(*) Apply this to the Normal case:-

- as an upper bound; select \bar{x} as an arbitrary estimator.

- Risk of \bar{x} ; note L₂ loss is MSE.

- MSE = Bias² + variance

- Bias = 0 \Rightarrow MSE = variance = $\frac{1}{n}$

- so we have upper bound $R_n \leq \frac{1}{n}$
- to get lower bound on R_n ; select a Normal prior $\pi = N(0, c^2)$
i.e. as an exercise; make sure you are acquainted with conjugacy
(of likelihood, prior, post.)
- under this prior; we have a posterior Normal $p(\theta | x_1, \dots, x_n) \sim N(a, b)$ for some a, b .
- we want Bayes estimator \rightarrow which under L_2 loss is posterior mean
(convex combo
of MLE, prior mean)
- now have to get Bayes risk of $\hat{\theta}_n$ (by getting risk
and then Bayes risk).
- Risk of $\hat{\theta}_n$ under L_2 loss is MSE
- MSE = bias² + variance
- use these to assist:-

$$\text{bias} = E(\hat{\theta}_n) - \theta = \frac{nc^2 E[\bar{X}]}{1+nc^2} - \theta$$

$$= \frac{-\theta}{1+nc^2}$$

- noting: only \bar{X} is an RV.
 $E[\bar{X}] = \theta$.

- for full ref: bias in terms
of sample size is $O(\frac{1}{n})$

$$\text{var}(\hat{\theta}_n) = \frac{n^2 c^4}{(1+nc^2)^2} \text{Var}[\bar{X}] = \frac{n^2 c^4}{(1+nc^2)^2} \cdot \frac{1}{n} = \frac{nc^4}{(1+nc^2)^2}$$

- the risk in this case is MSE = bias² + variance (L_2 loss)
- if not L_2 loss; will have to take expected value of loss to get risk.

$$R(\theta, \hat{\theta}_n) = \text{bias}^2 + \text{variance} = \frac{\theta^2 + nc^4}{(1+nc^2)^2}$$

- so we now have the risk; but we require Bayes risk

① exercise
review

(prior, post,
like).

② - check
discrep
against
notes

- Bayes risk of the Bayes estimator :-

$$B\pi(\hat{\theta}_n) = \int \frac{\theta^2 + nc^4}{(1+nc^2)^2} \pi(\theta) d\theta \sim \text{?} \text{ cm}^4 \text{ see.}$$

• $\pi(\theta)$ is just a normal prior with mean 0, variance c^2

$$= \frac{\mathbb{E}\pi[\theta^2] + nc^4}{(1+nc^2)^2}$$

• we are just taking expected value under prior $\pi(\theta)$

$\frac{\theta^2 + nc^4}{(1+nc^2)^2}$ and treating θ as a r.v. (that is normal).

$$= \frac{c^2 + nc^4}{(1+nc^2)^2}$$

$$= \frac{c^2(1+nc^2)}{(1+nc^2)^2} = \frac{c^2}{(1+nc^2)}$$

$$\mathbb{E}\pi[\theta^2] = c^2$$

- 2nd moment of θ , with mean 0 is just variance $\pi \sim N(0, c^2)$

- Hence, we have an upper and lower bound on R_n :-

$$\frac{c^2}{(1+nc^2)} \leq R_n \leq \frac{1}{n}$$

- As we selected arbitrary value of c (sd. of prior on θ); above holds for $\forall c$

- so we can take limit of LHS:- $\lim_{c \rightarrow \infty} \left(\frac{c^2}{1+nc^2} \right) = \frac{1}{n}$ (a3) (?)
 - cm⁴ see this limit intuitively

$$\text{- Hence } \frac{1}{n} \leq R_n \leq \frac{1}{n} \Rightarrow R_n = \frac{1}{n}$$

(*) finding Bayes estimator doesn't work; as we would need to find a prior that is uniform over the whole real line

(*) hence we know minimax risk (in context) is $\frac{1}{n}$.

- To find minimax estimator that achieves the minimax risk; we have to find an estimator whose maximum risk ($\sup_{\theta} R(\theta, \hat{\theta})$)

is $\frac{1}{n}$. That is

- (*) If we take $\hat{\theta} = \bar{X}$; $\sup_{\theta} R(\theta, \hat{\theta}) = \frac{1}{n} = R_n$
- ~~~~~(a)
- (*) max risk does not depend on θ here, but we write \sup anyway.
- so we have:-
- 1) Found minimax risk ($R_n = \frac{1}{n}$)
 - 2) Found an estimator that achieves minimax risk ($\hat{\theta} = \bar{X}$)
- Hence \bar{X} in this context is the minimax estimator ✓ test prep for gen. case
- uw: details here make proof strategy more explicit than notes.
- uw: Above strategy is contingent on artful choices of $B_n(\hat{\theta}_n)$ and $\sup_{\theta} R(\theta, \hat{\theta}_n)$.

1.5 maximum likelihood

(*) parametric models that satisfy certain conditions (weak regularity);

$\hat{\theta}_{MLE}$ is approximately minimax.

Making this precise \rightarrow tricky (pioneered by Lucien LeCam).
- beyond scope of this course.

- You will get a feeling for assumptions in parametric models for above to be true.

- This assumes that the dimension of θ is fixed; n is increasing

- fixed, finite dimensional parametric models with smoothness condition;

$\hat{\theta}$ is approximately minimax.

(*) See "Asymptotic Statistics", van de Waart for further exposition

1.6 Hodges example

- skip; give it a read.

- Related to lecture notes 9

YouTube Lecture 05/10/16

Lecture Notes 9 - Asymptotic Theory- Notation: - Recall:-

$$\mathbb{E}_\theta[g(x)] = \int g(x) p(x; \theta) dx$$

- DOES NOT MEAN θ IS random
(in this course)

- A reminder that the true distn is parenthesised by θ .

- θ is fixed, not random
(untrue)

- 1. Review of θ and $\hat{\theta}$

- These are included

in notes for reference;

should all be familiar and intuitive.

- 2. Distances between probability distributions

- Included for reference; not expected to memorise.

- Common distances:-

K-L distance: - $K(P, Q) = \int p \log(p/q) \quad \text{mean}$

Hellinger: - $h(P, Q) = \sqrt{\int (p - q)^2}$

we are typically interested in behaviour/properties of an estimator as $n \rightarrow \infty$ (i.e. we have more and more samples).

- We are concerned with:- for an est. $\hat{\theta}_n = g(x_1, \dots, x_n)$ 1) Consistency : $\hat{\theta}_n \xrightarrow{P} \theta$ as $n \rightarrow \infty$ i.e. $\hat{\theta}_n - \theta = o_p(1)$
of an est.m.

i.e. estimator
converges in probability to the true value of the parameter.
(in particular ML estimator)

2) Asymptotic normality
under some conditions; a CLT for in particular ML estimators.

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \sigma^2)$$

3. Consistency

- 2 modes of showing consistency; directly show via definition:-

Method 1 :- Show that for all $\epsilon > 0$ $\hat{\theta}_n \xrightarrow{P} \theta$ (as $n \rightarrow \infty$)

$$P(|\hat{\theta}_n - \theta| \geq \epsilon) \rightarrow 0$$

Method 2 : Prove convergence in quadratic mean :- $\hat{\theta}_n \xrightarrow{Q.M} \theta$

$$E[(\hat{\theta}_n - \theta)^2] \rightarrow 0$$

- But in more updated vocab:-

$$MSE(\hat{\theta}_n) = \text{Bias}^2(\hat{\theta}_n) + \text{Var}(\hat{\theta}_n) \rightarrow 0$$

• If $\text{Bias} \rightarrow 0$ and $\text{Var} \rightarrow 0$ then $\hat{\theta}_n \xrightarrow{Q.M} \theta$ (^{estimator} converges in q.m to param.)

• $\hat{\theta}_n \xrightarrow{Q.M} \theta \Rightarrow \hat{\theta}_n \xrightarrow{P} \theta \Rightarrow$ estimator $\hat{\theta}_n$ is consistent

W: Run through examples of consistent estimators and estimators which are not.

Example 1

- $X_1, \dots, X_n \sim \text{Bern}(p)$

$$\hat{p}_{MLE} = \bar{X} \quad \hat{p}_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\text{Bias} = E[\hat{p}_{MLE}] - p = 0$$

$$\text{variance} = \frac{p(1-p)}{n} \rightarrow 0 \text{ as } n \rightarrow \infty$$

- Hence $MSE(\hat{p}_{MLE}) \rightarrow 0$; and $\hat{p}_{MLE} \xrightarrow{Q.M} p \Rightarrow \hat{p}_{MLE} \xrightarrow{P} p$

Suppose we want to compute log-odds $\phi = \log\left(\frac{p}{(1-p)}\right)$

$$\text{via equivariance: } \hat{\phi} = \log\left(\frac{\hat{p}_{MLE}}{(1-\hat{p}_{MLE})}\right)$$

④ Is $\hat{\phi}$ a consistent estimator of ϕ ? YES; via continuous mapping theorem.

Example 2

(AU) - review proof

$$x_1, \dots, x_n \sim \text{unif}(0, \theta)$$

$$\hat{\theta}_{n, \text{mle}} = x_{(n)} = \max\{x_1, \dots, x_n\}$$

$$\hat{\theta}_n \xrightarrow{P} \theta \quad (\text{proof directly on previous lectures})$$

(*) To refresh:-

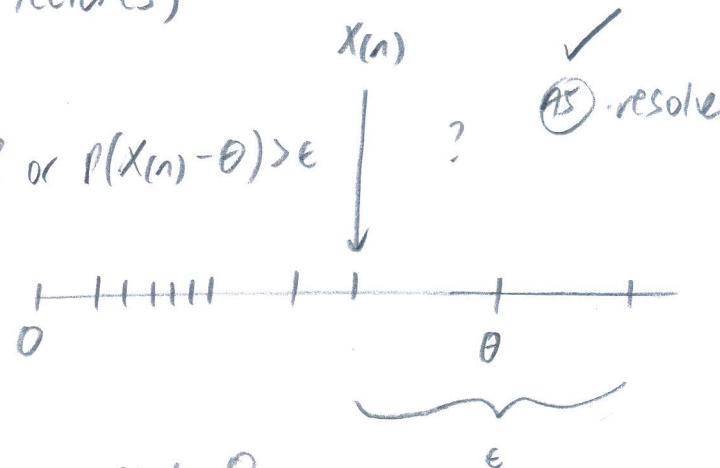
$$\text{what is } P(\max\{x_1, \dots, x_n\} - \theta > \epsilon) ? \text{ or } P(x_{(n)} - \theta) > \epsilon$$

- probability of being outside

$$\text{interval size } \epsilon \Rightarrow x_{(n)} - \theta$$

all of observations outside
the interval or less than $(\theta - \epsilon)$.

- can compute that probability; that it converges to 0.



- not hugely important
- move on.

Then: $P(|\hat{\theta}_n - \theta| > \epsilon) \rightarrow 0 \quad \forall \epsilon > 0$

and: $\hat{\theta}_{n, \text{mle}} \xrightarrow{P} \theta \quad \text{i.e. } \hat{\theta}_{n, \text{mle}} \text{ is consistent}$

Example 5

- Suppose I decide to measure everyone's blood pressure

- blood pressure is not i.v.

- we each have a mean blood pressure μ_i

- there is measurement noise from the machine σ^2

- take 2 measurements per person

$$Y_{11}, Y_{12} \sim N(\mu_1, \sigma^2)$$

$$Y_{21}, Y_{22} \sim N(\mu_2, \sigma^2)$$

⋮

$$Y_{n1}, Y_{n2} \sim N(\mu_n, \sigma^2)$$

(*) find likelihood

$L(\theta)$ which is

just a large product

μ - assume blood pressure
is fixed, unknown no.
(as opposed to varying)

$$(*) \hat{\sigma}^2 = \sum_{i=1}^n \sum_{j=1}^2 \frac{(Y_{ij} - \bar{Y}_i)^2}{2n}$$

\bar{Y}_i - average of
2 measurements
for each individual i .

(*) - check calc
of $\hat{\sigma}^2$ and
of $\hat{\sigma}^2$
(test prep).

- we can show:-

$$\hat{\sigma}^2 \xrightarrow{P} \frac{\sigma^2}{2}$$

(also $\hat{\sigma}^2$ is not consistent), but $2\hat{\sigma}^2$ is consistent

- some of issue:- no. of params increases with no. of observations
- an issue for maximum likelihood (not psh?)
- dimension of parameter space increasing with n .

QW: when is the maximum likelihood estimator consistent?

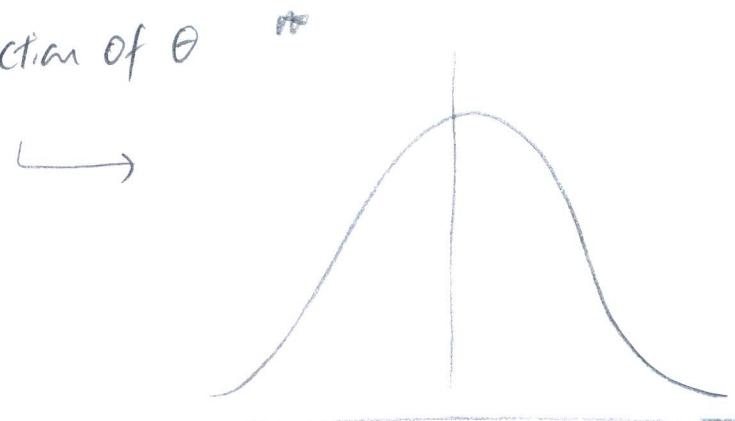
- full answer → see van de wael

- after m notes to give intuition on when ML estimator is consistent
and what makes it consistent

QW: we need two regularity conditions on model

1) $\dim(\Theta) = d$ is fixed

2) $p(x; \theta)$ is a smooth function of θ



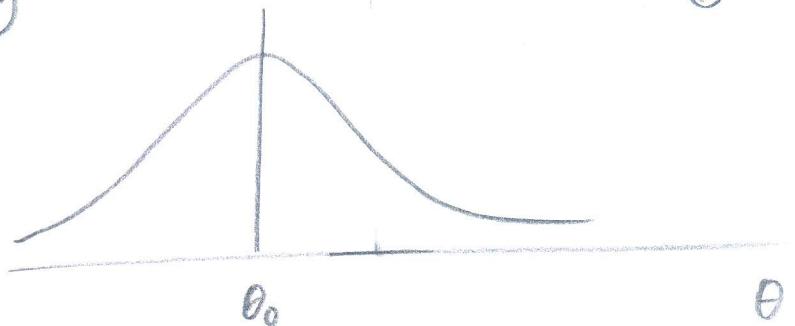
- intuitively:- perturbing θ only
informally shifts $p(x; \theta)$ a small amount

Theorem 6

- under regularity conditions on the model $\{p(x; \theta) : \theta \in \Theta\}$, the maximum likelihood estimator is consistent
- proof: needs probability tools not covered. But show some intuition with other theorems.
- ↳: theorem 3 is optional. Focus on theorem 4: also reveals an affinity between maximum likelihood and KL-divergence

Theorem 4

- assume we have a true value of the parameter θ_0 $x_1, \dots, x_n \sim$?
- generating data from distribution
- true param θ_0
- intuitively, would like to show that likelihood at θ_0 is high; likelihood at all other points is small.
- select another arbitrary value of θ ; compare the likelihoods:-



$$\frac{L(\theta_0)}{L(\theta)} > 1 \quad \leftarrow \text{this encodes above intuition.}$$

- switching to log-likelihood form:-

$$\ell(\theta_0) - \ell(\theta) > 0$$

$$\Rightarrow \frac{1}{n}(\ell(\theta_0) - \ell(\theta)) > 0$$

- we now assess if this is positive:

$$\begin{aligned} \frac{1}{n}(\ell(\theta_0) - \ell(\theta)) &= \frac{1}{n} \sum_{i=1}^n \log p(x_i; \theta_0) - \frac{1}{n} \sum_{i=1}^n \log p(x_i; \theta) \\ &= \frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i; \theta_0)}{p(x_i; \theta)} \end{aligned}$$

(4) Doing some variable renaming:-

$$\log \frac{p(x_i; \theta_0)}{p(x_i; \theta)} = y_i$$

- we can view $\frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i; \theta_0)}{p(x_i; \theta)}$ = $\frac{1}{n} \sum_{i=1}^n y_i$

- i.e.
- Knowledge of IITD C.V.S.

- Analogue of IID r.v.s.
- We know via WLLN that $\frac{1}{n} \sum_{i=1}^n y_i \xrightarrow{P} \mu_Y$

- And that in this case:-

$$H_{\mathcal{Y}} = \int \log \frac{p(x; \theta_0)}{p(x; \theta)} p(x; \theta_0) dx$$

i.e.

Q6 - why using $p(x; \theta_0)$ as density?

-(* 50 :-

$$\frac{1}{n}(\ell(\theta_0) - \ell(\theta)) \xrightarrow{\rho} \text{KL}\left(p(x; \theta_0) \parallel p(x; \theta)\right) > 0$$

- this shows we need to assume identifiability (under regularity conditions)

- Two different values of ρ or Θ_1, Θ_2 ; $\Theta_1 \neq \Theta_2$

then $KL(\theta_1, \theta_2) > 0$

then $KL(\theta_1, \theta_2) > 0$
 i.e. we require that there be different densities/distinct.

(*) So $\frac{1}{n}(\ell(\theta_0) - \ell(\theta))$ is converging in probability to a strictly positive number

Now this informs us that $L(0_0)/L(0) \xrightarrow{P} R$ where $R > 1$

- (x) so likelihood of the value of parameter θ_0 as $n \rightarrow \infty$, is always larger than the likelihood of another arbitrary value θ .

(*) suggests that maximum should somehow be close to 0.0.
likelihood? $\hat{\theta}_{n, \text{MLE}}$?
est.

- making the deduction that

$\hat{\theta}_{n,MLE}$ is consistent: $\hat{\theta}_{n,MLE} \rightarrow \theta_0$ requires more work on top of what we have shown.