

Lecture notes 8 - Minimax Theory - Review

- concerned with best possible estimator $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ of param $\theta \in \Theta$.
- framework uses a loss function and risk function (exp. value of loss).
- loss fn $U(\theta, \hat{\theta})$ - measure of 'quality' of estimator.

examples:

$$U(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2 \quad \text{sq. error}$$

$$U(\theta, \hat{\theta}) = |\theta - \hat{\theta}| \quad \text{abs. error}$$

$$U(\theta, \hat{\theta}) = |\theta - \hat{\theta}|^p \quad L_p\text{-loss}$$

$$U(\theta, \hat{\theta}) = \begin{cases} 0 & \text{if } \theta = \hat{\theta} \\ 1 & \text{if } \theta \neq \hat{\theta} \end{cases} \quad \text{zero-one}$$

$$U(\theta, \hat{\theta}) = I(|\hat{\theta} - \theta| > c) \quad \text{large deviation}$$

$$U(\theta, \hat{\theta}) = \int \log \left(\frac{p(x; \theta)}{p(x; \hat{\theta})} \right) p(x; \theta) \quad -\text{KL}$$

(*) vector losses

$$\theta = (\theta_1, \dots, \theta_K)$$

$$U(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|_2^2 = \sum_{j=1}^K (\hat{\theta}_j - \theta_j)^2 \quad L_2\text{-loss}$$

$$U(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|_p = \left(\sum_{j=1}^K |\hat{\theta}_j - \theta_j|^p \right)^{1/p} \quad L_p\text{-loss}$$

(*) ML-context:

$$\text{- prediction (classifier) } \dots U(Y, h(X)) = I(Y \neq h(X))$$

$$Y \in \{0, 1\}$$

$$\text{class. } h(x)$$

$$\text{- real-valued predic: } U(Y, \hat{Y}) = (Y - \hat{Y})^2$$

- Riskfn:

$$R(\theta, \hat{\theta}) = E_{\theta}[L(\theta, \hat{\theta})] = \int L(\theta, \hat{\theta}(x_1, \dots, x_n)) p(x_1, \dots, x_n; \theta) dx^n$$

Minimax: provides optimal estimator wrt to a given loss

① loss depends on data when we evaluate it:-

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

② Risk does not, as we integrate the data out, by taking expectations of loss wrt the joint distn $p(x^n; \theta)$ that generated the data, at value of true param; risk is a function of the parameter θ .

• L/squared error loss:-

$$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2 \Rightarrow R(\theta, \hat{\theta}) : \text{MSE} = E_{\theta}[(\theta - \hat{\theta})^2] = (E_{\theta}(\hat{\theta}) - \theta)^2 + E[(\hat{\theta} - E_{\theta}(\hat{\theta}))^2] = \sigma^2 + V$$

③ unclear as to why we would write $R(\theta, \hat{\theta})$ if ...

(*) Minimax risk definition

④ feels conceptually tricky.

- let's clarify.
- the definition of minimax risk and minimax estimator involves "two-stage" minimisation/maximisation over both the estimator $\hat{\theta}$ and parameter θ .
- the comparison of two estimators case is best place to start; generalise from there.
- from Bickel & Doksum (1977):- (adapting not.) (clarity on loss, risk functions)
- $\hat{\theta}$ - estimator/procedure /decision rule
- L - loss function
- $X^n = x^n$ - experiment outcome.
- loss: $L(\theta, \hat{\theta}(x_1, \dots, x_n))$
- ⑤ we do not know the value of the loss, in general, as parameter θ is unknown

- require our estimator to have desirable properties not just at one particular x^n , but over a range of possible x^n 's.
- so we use mean/average loss over the sample space
- regard $L(\theta, \hat{\theta}(x_1, \dots, x_n))$ as a random variable and introduce the risk function (not just 'risk') as a measure of the performance of the estimator $\hat{\theta}(x_1, \dots, x_n)$:

$$R(\theta, \hat{\theta}) = E_{\theta} [L(\theta, \hat{\theta}(x_1, \dots, x_n))]$$

- for each $\hat{\theta}$, R maps \mathbb{H} to \mathbb{R}^+ i.e. $R: \mathbb{H} \rightarrow \mathbb{R}^+$ (function of unknown param θ).
- $R(\cdot, \hat{\theta})$ is a prior measure of performance of $\hat{\theta}$.

(i) Should be a lot clearer now; as I was getting confused by instruction of risk function

- (ii) Start intuition from comparison of 2 estimators $\hat{\theta}_i$ and $\hat{\theta}_k$.
- for both estimators; compute their risk functions $R(\theta, \hat{\theta}_i)$, $R(\theta, \hat{\theta}_k)$, which are fractions of the unknown parameter θ .
 - in the examples; these are fractions of the unknown param of our postulated statistical model / distribution family / collection of densities.
 - next bits explain the build-up to minimax risk and minimax est. formulae.

- we compare the risk functions $R_i(\theta, \hat{\theta}_i)$ and $R_k(\theta, \hat{\theta}_k)$ where we have introduced indexing according to estimator used.
- An estimator $\hat{\theta}_i$ is 'better' than $\hat{\theta}_k$ if

$$\max_{\theta} (R_i(\theta, \hat{\theta}_i)) < \max_{\theta} (R_k(\theta, \hat{\theta}_k)) \quad (\textcircled{R})$$

$$\text{or equivalently } \sup_{\theta} R_i(\theta, \hat{\theta}_i) < \sup_{\theta} R_k(\theta, \hat{\theta}_k) \quad (\textcircled{R})$$

Propose: $\hat{\theta}_i$ is 'better' if maximum of corresponding risk function is smaller than another.

- the issue with comparing risk functions; as examples will show:-
- i) How to deal with a set of candidate estimators $(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_L)$? systematically
- ii) How to deal with issue of inconclusiveness of comparing risk functions (i.e. one does not uniformly dominate another).
- we use maximum risk as we don't know true value of θ . - we assume a worst-case scenario when choosing $\hat{\theta}$ (as higher risk \Rightarrow higher expected loss)
- when choosing an estimator from a set of estimators, we generalise the 2-estimator case; and choose the estimator $\hat{\theta}$ which has the lowest maximum risk over all estimators $\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3, \dots$
- this yields minimax risk:

$$R_n = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta})$$

- for a minimax estimator $\hat{\theta}_{\text{minimax}}$; its maximum risk $(\sup_{\theta} R(\theta, \hat{\theta}_{\text{minimax}}))$ (over parameters θ)

will be the same as the

$$\text{lowest maximum risk } (\inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}) = R_n)$$

(over all
estimators
 $\hat{\theta}$)

↗ (as it is the
minimax
estimator...)

- Hence, an estimator $\hat{\theta}$ that is minimax satisfies: - my defin.)

$$\sup_{\theta} R(\theta, \hat{\theta}) = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}).$$

(ii): conceptual thickness turned out to be more simple than expected.
particular difficulty with minimax estimator condition.

- that's clarity for you.

1.3. Bayes estimators

treat θ as r.v.

- π -prior distn $\pi(\theta)$

- after observing $X^n = (x_1, \dots, x_n)$

- $l(\theta) = p(x^n; \theta)$ (likelihood)

via Bayes, the posterior density:-

$$p(\theta|x^n) = \frac{p(x^n|\theta)\pi(\theta)}{\pi(x^n)} = \frac{p(x^n|\theta)\pi(\theta)}{\int p(x^n|\theta)\pi(\theta)d\theta}$$

note
marginal distn of $x^n = m(x_1, \dots, x_n) = \int p(x^n, \theta)$
 (x_1, \dots, x_n)

A1. Posterior risk

- Recall risk of an estimator $\hat{\theta}$ is the expectation of loss wrt the joint distn that generated the data $p(x^n; \theta)$.
- The posterior risk is integral of loss wrt posterior density/distr $\pi(\theta|x^n)$ or $p(\theta|x^n)$.
- Here density and distn used interchangeably.
- Posterior risk: (of an estimator $\hat{\theta}$)

$$r(\hat{\theta}|x^n) = \int l(\theta, \hat{\theta}(x^n)) p(\theta|x^n) d\theta$$

(*) compared to Bayes risk we are integrating out parameter θ i.e.
taking expectation over posterior rather than prior.

Theorem 6 (Alt/convient form for Bayes risk).

$$B_R(\hat{\theta}) = \int r(\hat{\theta}|x^n) m(x^n) dx^n$$

$B_R(\hat{\theta})$ is Bayes risk

- let $\hat{\theta}(x^n)$ be value of θ that minimises $r(\hat{\theta}|x^n)$; $\hat{\theta}$ is Bayes estimator
i.e. $\hat{\theta}(x^n)$ minimises posterior risk $r(\hat{\theta}|x^n) \rightarrow$ minimises Bayes risk

$$B_R(\hat{\theta})$$

\rightarrow Bayes est.

(6) Bayes risk can be viewed as:-

- i) expectation of usual risk $R(\theta, \hat{\theta})$ wrt prior distri $\pi(\theta)$.
- ii) expectation of posterior risk $r(\hat{\theta}|x^n)$ wrt marginal distri $M(x^n)$.

Theorem 7 (Bayes est with various losses)

- under L_2 loss / sq. error loss

i.e. $U(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$, Bayes est. :-

$$\hat{\theta}(x^n) = \int \theta \pi(\theta|x^n) d\theta = E[\theta|x^n] \quad (\text{i.e. mean of } \theta \text{ wrt posterior density})$$

- under absolute error loss

$U(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$, Bayes est. is median of posterior $\pi(\theta|x^n)$.

- under zero-one loss

$U(\theta, \hat{\theta}) = \begin{cases} 0 & \text{if } \theta = \hat{\theta} \\ 1 & \text{if } \theta \neq \hat{\theta} \end{cases}$, Bayes. est is mode of post. $\pi(\theta|x^n)$.

03/10/16 lecture terminates.

- fill in blanks using notes, Wasserman book.

Example 8

- let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, σ^2 known.

- suppose a $N(a, b^2)$ prior is used for μ and $U(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$ (square error loss)

- the Bayes est. with wrt L_2 loss is posterior mean

$$\hat{\theta}(X_1, \dots, X_n) = \frac{b^2}{b^2 + \sigma^2/n} \bar{X} + \frac{\sigma^2/n}{b^2 + \sigma^2/n} a$$

(*) Again; est. is convex combination of \bar{X} and prior mean.

(*) This is exactly same result as in Example 8 under INF.

(*) However; there we found the posterior mean and variance.

(i.e. compute this to yield point estimator \rightarrow Bayes estimator).

• still an issue of dealing with minimisation of maximum risk over all possible estimators (would be infinite)

(*) Example 4

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

- MLE: $\hat{p}_1 = \bar{x}_n$ (sample mean is unbiased) $\Rightarrow B = E_p[\hat{p}_1] - p = 0$

- As $MSE = B^2 + V$ (or $E_p[(\theta - \hat{\theta})^2] = (E_p[\hat{\theta}] - \theta)^2 + \text{Var}_p(\hat{\theta})$) and $B=0$;

- $MSE = E_p[(p - \hat{p})^2] = \text{Var}_p(\hat{p}) = \text{Var}(\bar{x}_n)$

$$= \frac{p(1-p)}{n} \quad \text{as } \text{Var}(\bar{x}_n) = \frac{p(1-p)}{n}$$

- w.r.t. L_2 loss, i.e. $L(p, \hat{p}) = (p - \hat{p})^2$

- Risk function: - $E_p[(p - \hat{p})^2]$, which is equal to MSE.

Hence $R(p, \hat{p}_1) = \text{MSE} = \frac{p(1-p)}{n}$

- Consider different estimator $\hat{p}_2 = \frac{Y + \alpha}{\alpha + \beta + n}$ where $Y = \sum_{i=1}^N x_i$ $\alpha, \beta > 0$
 (post. mean using Beta (α, β) prior).

(ii): find $R(p, \hat{p}_2)$:

- Unlike \hat{p}_1 (MLE, sample mean); it is unclear whether \hat{p}_2 is unbiased.

- unlike \hat{p}_1 (MLE, sample mean); it is unclear whether \hat{p}_2 is unbiased.

- Hence MSE is not necessarily equal to V .

- However equivalence of risk function $R(p, \hat{p})$ and MSE still holds due to assumption of L_2 loss.

- Evaluate MSE:-

$$\text{MSE} = B^2 + V = (E_p[\hat{p}_2] - p)^2 + \text{Var}_p(\hat{p}_2) =$$

- Have to evaluate each; starting with:-

$$E_p[\hat{p}_2] = E\left[\frac{Y + \alpha}{\alpha + \beta + n}\right] = \frac{E[Y] + \alpha}{\alpha + \beta + n} = \frac{E\left[\sum_{i=1}^N x_i\right] + \alpha}{\alpha + \beta + n} = \frac{\sum_{i=1}^N E[x_i] + \alpha}{\alpha + \beta + n}$$

$$\Rightarrow E_p[\hat{p}_2] = \frac{np+\alpha}{(\alpha+\beta+n)}$$

- AS $X_i \sim \text{Bern}(p)$; $E[X_i] = p$ $\text{Var}(X_i) = p(1-p)$ $E[X_i^2] = p$

$$Y = \sum_{i=1}^n X_i$$

$$\begin{aligned} E[Y] &= \sum_{i=1}^n E[X_i] = np \quad \text{and} \quad \text{Var}[Y] = \text{Var}\left[\sum_{i=1}^n X_i\right] = \sum_{i,j=1}^n \text{Cov}(X_i, X_j) \\ &= \sum_{i=1}^n \text{Var}(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j) \\ &= \sum_{i=1}^n \text{Var}(X_i) = np(1-p) \end{aligned}$$

$$\text{so } E[Y^2] = \text{Var}[Y] + E[Y]^2 = np(1-p) + n^2 p^2$$

$$\text{Var}_p(\hat{p}_2) = E_p[\hat{p}_2^2] - E_p[\hat{p}_2]^2$$

$$= E\left(\frac{(Y+\alpha)^2}{(\alpha+\beta+n)^2}\right) - \frac{(np+\alpha)^2}{(\alpha+\beta+n)^2}$$

$$= \frac{E[Y^2] + 2\alpha E[Y] + \alpha^2}{(\alpha+\beta+n)^2} - \frac{(np+\alpha)^2}{(\alpha+\beta+n)^2}$$

$$= \frac{(n^2 p^2 - np^2 + np + 2\alpha np + \alpha^2)}{(\alpha+\beta+n)^2} - \frac{(np^2 + 2\alpha np + \alpha^2)}{(\alpha+\beta+n)^2}$$

$$= \frac{np - np^2}{(\alpha+\beta+n)^2}$$

$$\Rightarrow \text{Var}_p(\hat{p}_2) = \frac{np(1-p)}{(\alpha+\beta+n)^2}$$

$$(*) \text{ Alternatively: } \text{Var}_p(\hat{p}_2) = \text{Var}_p\left(\frac{Y+\alpha}{\alpha+\beta+n}\right) = \text{Var}_p\left(\frac{1}{\alpha+\beta+n}(Y+\alpha)\right)$$

$$= \frac{\text{Var}(Y)}{(\alpha+\beta+n)^2} = \frac{np(1-p)}{(\alpha+\beta+n)^2}$$

- AS $\text{Var}(\alpha X) = \alpha^2 \text{Var}(X)$; $\text{Var}(kX) = k^2 \text{Var}(X)$.

Hence:-

$$R(p, \hat{p}_2) = \left(\frac{np + \alpha}{\alpha + \beta + n} - p \right)^2 + \frac{np(1-p)}{(\alpha + \beta + n)^2}$$

set $\alpha = \beta = \sqrt{n}/4 \Rightarrow$

$$\hat{p}_2 = \frac{\sqrt{n}/4}{n + 2\sqrt{n}/4} = \frac{\sqrt{n}/4}{n + \sqrt{n}} \quad R(p, \hat{p}_2) = \frac{n}{4(n + \sqrt{n})^2}$$

(*)

(6/51): evaluate $R(p, \hat{p}_2)$ - computationally heavy; want to move on.
to get result

(*) Key point: comparing risk functions and even one number summaries of
risk such as maximum risk is imperfect; often neither
risk function uniformly dominates.

(*) summaries of risk fn:-

maximum risk : $\bar{R}(\hat{\theta}) = \sup_{\theta} R(\theta, \hat{\theta})$

- understand Q asked
now; student was
proposing to use
prior + Bayes to
get around
the pointwise
maximisation
of $R(\theta, \hat{\theta})$

Bayes risk (with prior π) : $B_{\pi}(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta$

where $\pi(\theta)$ is a prior on θ .

Q: In context of earlier; Bayes risk amounts to putting a prior $\pi(\theta)$ encoding
what we believe θ should be (note v now); and compute average
of risk under this prior.

example 5

- consider earlier estimators $\hat{p}_1 = \bar{x}_n$ $\hat{p}_2 = \frac{Y + \alpha}{(\alpha + \beta + n)^2}$ $Y = \sum_{i=1}^N x_i$

- Compute maximum risk of both estimators.

$$\bar{R}(\hat{p}_1) = \max_{0 \leq p \leq 1} \frac{p(1-p)}{n} = \frac{1}{4n} \quad (\text{if } p = \frac{1}{2})$$

$$\bar{R}(\hat{p}_2) = \max_p \frac{n}{4(n+5n)^2} = \frac{n}{4(n+5n)^2}$$

As domain of $\bar{R}(\hat{p}_2)$ is larger than that of $\bar{R}(\hat{p}_1)$; $\bar{R}(\hat{p}_2) < \bar{R}(\hat{p}_1)$

(*) now see diagram.

- when n is large; $\bar{R}(\hat{p}_1)$ has smaller risk except for region near $p=1/2$

- so on this basis; could agree for \hat{p}_1 over \hat{p}_2 .

(**) one number summaries e.g. max risk \rightarrow imperfect.

- enduring (from Wasserman):-

- diffuse prior $\pi(p) = 1$

$$B_n(\hat{p}_1) = \int R(p, \hat{p}_1) dp = \int \frac{p(1-p)}{n} dp = \frac{1}{n} \left[\frac{1}{2}p^2 - \frac{1}{3}p^3 \right]_0^1$$

$$= \frac{1}{6n}$$

$$B_n(\hat{p}_2) = \int R(p, \hat{p}_2) dp = \int \frac{n}{4(n+5n)^2} dp = \frac{n}{4(n+5n)^2}$$

$n \gg 20$ $B_n(\hat{p}_2) > B_n(\hat{p}_1) \Rightarrow \hat{p}_2$ is better estimator.

(*) Answer depends on prior choice; illustrates why frequentists prefer maximum risk.

(*) Two risk summaries (maximum risk, Bayes risk) yield two methods for deriving estimators (selecting amongst?).

- select $\hat{\theta}$ to minimise max risk $\bar{R}(\hat{\theta}) \rightarrow$ minimax estimator

- select $\hat{\theta}$ to minimize Bayes risk $B_n(\hat{\theta}) \rightarrow$ Bayes estimator
with respect to prior π

Bayes estimator:

$$\hat{\theta} : B_n(\hat{\theta}) = \inf_{\tilde{\theta}} B_n(\tilde{\theta})$$

- infimum over all estimators.

Minimax estimator:

$$\hat{\theta} : \sup_{\theta} R(\theta, \hat{\theta}) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \\ = R_n$$

RHS of above, R_n , is the minimax risk

(*) Statistical decision theory involves:-

I) determine minimax risk

II) find an estimator that achieves this risk.

$$R_n = R_n(\Theta) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta})$$

(OTS) \otimes : Not entirely clear: the notation

$\theta \in \Theta$

- But otherwise secure

- think it just means max over entire parameter space Θ

(*) Having found minimax risk ($R_n = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta})$), we find the minimax estimator that achieves

$$\text{this risk. } \hat{\theta} : \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta})$$

Asymptotically $\hat{\theta} : \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \sim \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \quad n \rightarrow \infty$

minimax

estimator

where $a_n \sim b_n$ means $\frac{a_n}{b_n} \rightarrow 1$

(*) May be too difficult, settle for estimator $\hat{\theta}$ that achieves the minimax rate.

minimax rate estimator $\hat{\theta} : \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \asymp \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \quad n \rightarrow \infty$

where $a_n \asymp b_n$ mean that both $\frac{a_n}{b_n}$ and $\frac{b_n}{a_n}$ are both bounded as $n \rightarrow \infty$.

- (*) However; this builds on that result.
- Because we are using T.F to establish a connection between Bayes estimator as posterior mean under l_2 loss
- ### 1.4 Minimax Estimators
- (a) Bayes estimators with a constant risk function are minimax.
- Theorem 9
- let $\hat{\theta}$ be the Bayes estimator for some prior π
 - If / Suppose that $R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta})$ for all θ
- then $\hat{\theta}$ is minimax and π is called a least fav. prior
- Proof
- suppose $\hat{\theta}$ is not minimax
 - then there will be another estimator $\hat{\theta}_0$ such that $\sup_{\theta} R(\theta, \hat{\theta}_0) < \sup_{\theta} R(\theta, \hat{\theta})$
- (*) As the average of a function is always less than or equal to its maximum, we have
- $$B_\pi(\hat{\theta}_0) \leq \sup_{\theta} R(\theta, \hat{\theta}_0)$$
- Hence: $B_\pi(\hat{\theta}_0) \leq \sup_{\theta} R(\theta, \hat{\theta}_0) < \sup_{\theta} R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta})$
- which is a contradiction.
- (*) I am not entirely sure where contradiction is; except it must be in the section I've outlined. NO
- (*) comfortable with supposing $\hat{\theta}$ is not minimax
- (*) Does the contradiction reside in the fact that :-
- $$\sup_{\theta} R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta}) \quad \text{NO}$$
- (152) - see later part.
- see return to this - stackexchange (see wasserman)

Theorem 10

- suppose $\hat{\theta}$ is Bayes est w.r.t prior π ; and $\hat{\theta}$ has constant risk $R(\theta, \hat{\theta}) = c$ for some c
- If the risk is constant, $\hat{\theta}$ is minimax.

Proof:

Bayes risk:- $B_\pi(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta = c$

(?) - Bayes risk or risk constant?

) (?) 6/53

Hence $R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta})$ for all θ

$\Rightarrow \hat{\theta}$ is minimax via T.9.

(*) Inside examples; will help with general's.

(*) Example 11

- see ex 4

- $x_1, \dots, x_n \sim \text{Bern}(p)$ L2 loss i.e. $U(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$

- shared estimator $\hat{p}(x^n) = \frac{\sum_{i=1}^n x_i + \sqrt{n}/4}{n + \sqrt{n}}$ (\hat{p}_i is expt).

- T.3 has a constant risk function

- $R(p, \hat{p}_2) = \frac{n}{4(n+\sqrt{n})^2}$ i.e. not a function of parameter p

- This estimator is posterior mean under Beta prior (α, β) i.e. $\pi(\alpha, \beta) = \text{Beta}(\alpha, \beta)$ (from previous example).

- As we have L2 loss $U(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$, via Theorem 7, this estimator is the Bayes estimator (i.e. posterior mean w.r.t L2 loss is Bayes est.)

- That is, minimises the Bayes risk under a beta (α, β) prior

with $\alpha = \beta = \sqrt{n}/4$.

- via Theorem 10; as risk function $R(p, \hat{p}_2)$ is constant and \hat{p}_2 is Bayes est, \hat{p}_2 is minimax

Example 12

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

- define following loss:

$$l(p, \hat{p}) = \frac{(p - \hat{p})^2}{p(1-p)}$$

- define following estimator:-

$$\hat{p}(X^n) = \hat{p} = \frac{\sum_{i=1}^n X_i}{n}$$

recall r.v. is the estimator
expectation not distri param
by p

The risk is:-

$$R(p, \hat{p}) = E[l(p, \hat{p})] = E\left(\frac{(p - \hat{p})^2}{p(1-p)}\right) = \frac{1}{p(1-p)} \underbrace{E_p[(p - \hat{p})^2]}_{\text{MSE}}$$

$$= \frac{1}{p(1-p)} \text{Var}_p(\hat{p})$$

$$= \frac{1}{p(1-p)} \frac{p(1-p)}{n}$$

$$= \frac{1}{n}$$

- As \hat{p} is unbiased
 $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$

Hence $R(p, \hat{p})$, as fun of p, is constant
(*) can be shown that for this loss function $\hat{p}(X^n)$ is Bayes estimator

under prior $\pi(p) = 1$.

(*) under Theorem 10, \hat{p} is minimax

$$\hookrightarrow \text{Bayes risk: } B_{\pi}(\hat{p}) = \int_0^1 R(p, \hat{p}) \pi(p) dp = \int_0^1 \frac{1}{n} \cdot 1 dp = \frac{1}{n} [p]_0^1 = \frac{1}{n}$$

(O/S 4) - show this. Don't understand why \hat{p} is a minimax estimator,
as I can't see how it is the Bayes estimator i.e. \hat{p} minimises
Bayes risk.
- surely my \hat{p} is Bayes est. in this case?

Theorem 9 (redo \rightarrow Wasserman book helps) (did not understand proof by contra)

- let $\hat{\theta}$ be the Bayes estimator for some prior π ; then $\hat{\theta}$ satisfies
- $B_\pi(\hat{\theta}) = \inf_{\hat{\theta}} B_\pi(\hat{\theta})$ (via definition of Bayes est.)
- suppose that $R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta})$ for all θ
- then $\hat{\theta}$ is minimax and π is a least favourable prior
- Proof (Focus on structure)
 - assume that $\hat{\theta}$ (Bayes estimator) is not minimax.
not prior π
 - contradiction is that we are able to find another estimator $\hat{\theta}_0$ such that $B_\pi(\hat{\theta}_0) < B_\pi(\hat{\theta})$; which contradicts our assumption that $\hat{\theta}$ is the Bayes estimator; and minimises the Bayes risk - we have found an estimator with lower Bayes risk $B_\pi(\hat{\theta}_0)$ as consequence of assuming $\hat{\theta}$ is not minimax.
 - we therefore conclude that $\hat{\theta}$ is minimax
- on the issue of $\sup_{\theta} R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta})$ (which you thought was contradict.)
 - we assume that $R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta})$
i.e. for our Bayes estimator $\hat{\theta}$, the risk function $R(\theta, \hat{\theta})$ is smaller than the Bayes risk over all values of the parameter θ .
 - (*) this means that $\sup_{\theta} R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta})$ by assumption
 $\sup_{\theta} R(\theta, \hat{\theta})$ is just maximal value of $R(\theta, \hat{\theta})$ evaluated at a particular value of θ . (param)
(*) we assume $R(\theta, \hat{\theta}) \leq B_\pi(\hat{\theta}) + \theta$ excellent

- minimax estimators - normal model
- A function ℓ is bowl-shaped if
 - sets $\{x : \ell(x) \leq c\}$ are convex and symmetric about origin
- A loss function L is bowl-shaped if $L(\theta, \hat{\theta}) = \ell(\theta - \hat{\theta})$ for some bowl-shaped function ℓ .

Theorem 13

- suppose i.i.d. x has a Normal distn with mean vector θ and cov matrix Σ .
- if the loss-function is bowl-shaped then \bar{x} is the unique ^{*} minimax estimator of θ .
- up to sets of measure zero.

(*) param space restricted; above does not apply.

① Wasserman states in T.12.14 for $\bar{x} \rightarrow$ is there a typo?

Example 14

- suppose $x \sim N(0, 1)$ and $\theta \in [-m, m]$ where $m > 0$ (i.e. restricted param space)
 - under L_2 loss / squared error loss; unique minimax est is:-
- $$\hat{\theta}(x) = m \tanh(mx) = \frac{e^{mx} - e^{-mx}}{e^{mx} + e^{-mx}}$$

- this is Bayes est. not prior that puts mass $1/2$ at m and mass $1/2$ at $-m$.

- risk is not constant but it does satisfy $R(\theta, \hat{\theta}) \leq C_R(\hat{\theta}) \ \forall \theta$.

$$R(\theta, \hat{\theta})$$

- via 1.9. $\hat{\theta}$ is minimax

② Modern minimax theory \rightarrow now is param space restricted?

(iii) clarity on least favorable priors

WIKI:

- estimator is minimax when it does best in the worst case.
- A minimax estimator should be a Bayes estimator w.r.t. least favorable prior distribution of θ .
- Bayes risk of a Bayes estimator $\hat{\theta}$ w.r.t. prior $\pi(\theta)$:-
$$B_\pi(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta$$

- A prior distri $\pi(\theta)$ is least favourable if for every other prior distri π' the Bayes risk satisfies:-

$$B_{\pi'}(\hat{\theta}) \geq B_{\pi}(\hat{\theta})$$

(Bayes risk of Bayes est.
under least favourable
prior is the highest)
(worse case)

- If $B_\pi(\hat{\theta}) = \sup_{\theta} R(\theta, \hat{\theta})$

then i) $\hat{\theta}$ is minimax

ii) $\hat{\theta}$ is unique Bayes estimator and unique minimax est.

iii) π is least favourable

6/5/5: notice how to reconcile with UW notes.

(iv) proof that

\bar{x}_n is minimax mle squared error loss (fairly involved)

- strategy consists of

- i) upper bound minimax risk R_n using maximum risk of arbitrary est.

- ii) lower bound minimax risk R_n using Bayes risk eval. at Bayes estimator (minimise of Bayes risk).

- computation.

• see notes on justifying how we bound minimax risk R_n

• (2) idea:-

$$B_n(\hat{\theta}_n) \leq \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}) \leq \sup_{\theta} R(\theta, \hat{\theta}_n)$$

$\underbrace{\phantom{B_n(\hat{\theta}_n) \leq \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}) \leq \sup_{\theta} R(\theta, \hat{\theta}_n)}}$
 R_n arb.

• $B_n(\hat{\theta}_n)$ - Bayes risk eval. at Bayes est. not prior π (i.e. minimiser of Bayes risk)

• $\sup_{\theta} R(\theta, \hat{\theta}_n)$ - maximum risk for an arbitrary est. $\hat{\theta}_n$

• Computation of upper and lower bounds.

- select sample mean \bar{x} as arbitrary est.

- under L_2 loss, risk of estimator $\hat{\theta} = \bar{x}$ and MSE are the same.

Risk of $\hat{\theta}$: $R(\theta, \hat{\theta}) = E_{\theta}[L(\theta, \hat{\theta}(x_1, \dots, x_n))] = E_{\theta}[(\theta - \hat{\theta})^2]$

L_2 loss - $L(\theta, \hat{\theta}(x_1, \dots, x_n)) = (\theta - \hat{\theta})^2$

MSE of $\hat{\theta}$: $MSE = E_{\theta}[(\hat{\theta} - \theta)^2] = \sigma^2 + V$ P10
corr. \rightarrow

(0/5) 6: - got confused by L.W. claiming that risk and MSE under L_2 loss are the same; in particular due to sign differences

i.e. $E_{\theta}[(\theta - \hat{\theta})^2]$ vs $E_{\theta}[(\hat{\theta} - \theta)^2]$

- guessing this is immaterial.

- this is elementary, but needed to be cleared intuitively and formally.

- intuit: squared differences \rightarrow order $(\hat{\theta} - \theta)^2$ vs $(\theta - \hat{\theta})^2$ doesn't matter.

- formally: $(\hat{\theta} - \theta)^2 = \hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2 = (\theta - \hat{\theta})^2$!

- well done for raising this and clearing. \rightarrow prevents confusion

$$MSE = \hat{\theta}^2 + V$$

$= V$ as \bar{X} is unbiased.

$$\text{var}_\theta(\bar{X}) = \frac{1}{n}$$

Q15.7: This is elementary; but I can't quite see why

$$\text{var}_\theta(\bar{X}) = \frac{1}{n}. \text{ return to this.}$$

- easy. we assumed that $X_1, \dots, X_n \sim N(\theta, 1)$!

$$\text{so } \text{var}_\theta(\bar{X}) = \text{var}_\theta\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{var}_\theta\left(\sum_{i=1}^n X_i\right)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \text{var}_\theta(X_i) \text{ as } X_1, \dots, X_n \text{ IID}$$

$$= \frac{n}{n^2} = \frac{1}{n}$$

$$\text{hence } R(\theta, \bar{X}_n) = \frac{1}{n} + f(\theta)$$

$$\text{so } \sup_\theta R(\theta, \bar{X}_n) = \frac{1}{n} \Rightarrow R_n \leq \frac{1}{n} \quad (\text{upper bound})$$

recall $X_1, \dots, X_n \sim N(\theta, 1)$

• lower bounding R_n .

- assume conjugate prior $\pi(\theta) = N(\theta, c^2)$

- posterior $p(\theta | x_1, \dots, x_n) \sim N(a, b)$ (due to conjugacy).

- evaluate using same stat as eg. 8 in INT
(via trick: square completion).

- result from ex 8:

- for $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ σ^2 known

conjugate prior $p(\mu) \sim N(m, c^2)$

then mean and variance of posterior $p(\mu | x_1, \dots, x_n)$ is given by:-

$$E[\mu | x_1, \dots, x_n] = \frac{c^2 n}{c^2 n + \sigma^2} \bar{x} + \frac{\sigma^2}{c^2 n + \sigma^2} m$$

$$\text{var}[\mu | x_1, \dots, x_n] = \frac{\sigma^2 c^2 / n}{c^2 + \sigma^2 / n}$$

- now case; $\mu=0$ $\sigma^2=1$ $n=0$ $c^2=c^2$

post.

$$\rho(\theta | X_1=x_1, X_2=x_2, \dots, X_n=x_n)$$

$$\text{has } E[\theta | X_1, \dots, X_n] = \frac{c^2 n}{c^2 n + 1} \bar{x} + \frac{1}{c^2 n + 1} (0)$$

$$= \frac{n c^2 \bar{x}}{n c^2 + 1}$$

$$\text{Var}[\theta | X_1, \dots, X_n] = \frac{c^2/n}{c^2 + 1/n} = \frac{c^2}{n c^2 + 1}$$

$$\Rightarrow \theta | X^n = x_n \sim N\left(\frac{n c^2 \bar{x}}{n c^2 + 1}, \frac{c^2}{n c^2 + 1}\right)$$

- Under l_2 loss, Bayes est. is post mean

$$\Rightarrow \hat{\theta}_{\pi} = \frac{n c^2 \bar{x}}{n c^2 + 1} \quad (\text{Bayes est.})$$

(*) again using fact that risk and MSE are same under l_2 loss.

$$R(\theta, \hat{\theta}_{\pi}) = \text{MSE} = B^2 + V$$

$$\text{MSE} = E_{\theta}[(\hat{\theta}_{\pi} - \theta)^2] = (E_{\theta}[\hat{\theta}_{\pi}] - \theta)^2 + \text{Var}_{\theta}(\hat{\theta}_{\pi})$$

0158

$$\text{Bias: } E_{\theta}(\hat{\theta}_{\pi}) - \theta = \frac{n c^2 E[\bar{x}]}{n c^2 + 1} - \theta = \frac{-\theta}{1+n c^2} \quad \text{i.e. } O\left(\frac{1}{n}\right)$$

$$\text{variance: } \text{Var}_{\theta}(\hat{\theta}_{\pi}) = \text{Var}_{\theta}\left(\frac{n c^2 \bar{x}}{n c^2 + 1}\right) = \text{Var}_{\theta}\left(\frac{n c^2}{n c^2 + 1} \bar{x}\right)$$

$$= \left(\frac{n c^2}{n c^2 + 1}\right)^2 \text{Var}_{\theta}(\bar{x}) = \frac{n^2 c^4}{(n c^2 + 1)^2} \cdot \frac{1}{n} = \frac{n c^4}{(1+n c^2)^2}$$

- Not L_2 loss \rightarrow take exp. of loss to get risk.

$$R(\theta, \hat{\theta}_n) = \text{MSE} = \theta^2 + V = \left(\frac{-\theta}{nc^2+1} \right)^2 + \frac{nc^4}{(nc^2+1)^2}$$

$\stackrel{\text{not } L_2 \text{ loss}}{=}$

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

- This is the IBK function $R(\theta, \hat{\theta})$.

- we want the Bayes risk of the Bayes est. :-

$$B_n(\hat{\theta}_n) = \int R(\theta, \hat{\theta}_n) \pi(\theta) d\theta \quad \pi(\theta) \sim N(0, c^2)$$

$$= \int \frac{\theta^2 + nc^4}{(1+nc^2)^2} \pi(\theta) d\theta$$

$$= E_{\pi} \left[\frac{\theta^2 + nc^4}{(1+nc^2)^2} \right]$$

$$E_{\pi}[\theta^2] = c^2 \text{ as } E_{\pi}[\theta] = 0$$

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

$$\Rightarrow R_n \geq B_n(\hat{\theta}_n) = \frac{c^2}{(1+nc^2)}$$

$$\Rightarrow \frac{c^2}{(1+nc^2)} \leq R_n \leq \frac{1}{n} \quad \forall c \quad \text{arbitrary}$$

$$\text{Hence } \lim_{n \rightarrow \infty} \left(\frac{c^2}{1+nc^2} \right) = \frac{1}{n}$$

$$\Rightarrow \frac{1}{n} \leq R_n \leq \frac{1}{n} \Rightarrow R_n = \frac{1}{n}$$

- Hence minimax risk R_n is $\frac{1}{n}$

- we now need to find estimator $\hat{\theta}$ that achieves the minimax risk. R_n

use $\hat{\theta} = \bar{X}$; $\sup_{\theta} R(\theta, \hat{\theta}) = \frac{1}{n} = R_n$

Hence \bar{X}_n is minimax

above

(W) A little confused by what context ~~this~~ is w.r.t.

- notes say $n=1$ and $\sigma^2=1$, but I think this means $d=1$ and $\sigma^2=1$

- i.e. univariate Gaussian, $x_1, \dots, x_n \sim N(\theta, 1)$

$\sigma^2=1$ variance

- more generally, we can be d -dim multivariate Gaussian

and σ^2 i.e. $x_1, \dots, x_n \sim N(\theta, \sigma^2 \mathbb{I}_d)$

where $x_i \in \mathbb{R}^d$, $\theta \in \mathbb{R}^d$, $\sigma^2 \mathbb{I}_d \in \mathbb{R}^{d \times d}$