

Bayesian inference

so far we have tried to minimize  $R(\theta, \delta)$  over all possible  $\delta$ , and hope that will be uniform in  $\theta$ .

Restricted to

unbiased estimation

$$\hat{\delta} = \underset{\delta \text{ unbiased}}{\operatorname{argmin}} R(\theta, \delta)$$

~~IP~~  $\hat{\delta}$  ~~IF~~

$$R(\theta, \hat{\delta}) \leq R(\theta, \delta) \quad \text{if } \delta \text{ unbiased}$$

(1)

we call  $\hat{\delta}$  UMVU. Exponential families.  
equivalent estimators important.

$$\hat{\delta} = \underset{\delta \text{ equivalent}}{\operatorname{argmin}} R(\theta, \delta)$$

Risk does not depend on  $\theta$ , so if

$$R(\theta, \hat{\delta}) \leq R(\theta, \delta) \quad \text{if } \delta \text{ equivalent}$$

we call  $\hat{\delta}$  MRE. Group families important

①

Now change focus, change risk to average risk  
 Treat  $\theta$  as a random variable with a distribution  $A$ , and define the average risk of estimator  $\delta$

$$r(A, \delta) = \int R(\theta, \delta) dA(\theta).$$

Define, given a distribution  $A$ , the Bayes estimator

$$\hat{\delta} = \underset{\delta}{\operatorname{arg\min}} r(A, \delta) = \underset{\delta}{\operatorname{arg\min}} \int R(\theta, \delta) dA(\theta)$$

What does it mean? How do we use it?

### (i) A mathematical tool

Way to solve problems, connection to minimax theory and minimax estimator. Solve problems otherwise intractable.

### (ii) As a way to use past experience

Previous research/experiments/analysts suggest a distribution  $A$  for a parameter. Use this, gather data to update the distribution.

Ex: Coin toss,  $\theta = p(\text{head})$ .

View

\*

(2)  $\theta$  obs of (i) r.v., (ii)  $\sim A$ .

(iii) Add a description of a prior of mind

"Hard-core Bayesianism":

A subjective "feeling about"  $\theta$ . Use data  $x$  obs of  $X$  ( $X|\theta$  has a distribution we model), to update belief or feeling about  $\theta$ .

$$\begin{aligned} \textcircled{1} \quad \theta &\sim \Lambda \quad \text{choose a prior dist. of } \theta \\ X|\theta &\xrightarrow{\Rightarrow} \text{model} \\ \theta|X &\sim \text{posterior dist. of } \theta. \end{aligned}$$

(iv) Add a method to generate estimators:

(i) single prior Bayes  $\theta \sim \Lambda(\lambda)$

(ii) empirical Bayes

$$\theta \sim \Lambda(\lambda)$$

\* estimated from data

(iii) hierarchical Bayes

$$\theta \sim \Lambda(\lambda) \xrightarrow{\text{hyperprior distribution.}}$$

$$\lambda \sim \Gamma(\gamma)$$

\* fixed parameter

(iv) Robust Bayes

make sure the estimator works well for each member of the prior class.

## Some notation:

$$X \sim f(x|\theta) \quad \begin{matrix} \text{density or } X \text{ given } \theta \\ (\text{f.d.f or p.m.f.}) \end{matrix}$$

$$\textcircled{1} \quad \theta \sim \pi, \lambda \quad \begin{matrix} \text{density } \pi(\theta|\lambda), \gamma(\lambda) \\ \text{hyperparameter} \end{matrix}$$

$$\textcircled{2} \quad \frac{\pi(\theta|x, \lambda)}{\lambda|x} \sim \text{posterior distribution} \\ \pi(\theta|x, \lambda), \gamma(\lambda|x) \text{ densities.}$$

Marginal distributions need:

$$\pi(\theta|x, \lambda) = \text{(a posterior distribution of } \theta \text{ given data (sample) and hyperparameter } \lambda)$$

~~all  $\lambda$  as a fixed parameter and Bayes theorem to get this~~

$$= \frac{f(x|\theta) \pi(\theta|\lambda)}{m(x|\lambda)}$$

(  $\int f(x|\theta) \pi(\theta|\lambda) d\theta$  marginal distribution of  $x$  given  $\lambda$ , integrating out  $\theta$ . )

This means, if we forget about the hyperparameter  $\lambda$

$$\pi(\theta|x) = \frac{f(x|\theta) \pi(\theta)}{\int f(x|\theta) \pi(\theta) d\theta}$$

(i) but everything depends on hyperparameter  $\lambda$ , make explicit above

How do we calculate the Bayes estimator?  
Given a loss function  $L(\theta, \hat{\theta})$ .

Before data is sampled:

$$\hat{\theta} \sim \lambda$$

Risk is

$$E L(\hat{\theta}, \theta)$$

for estimating

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

Then

$$\hat{\theta} = \underset{\hat{\theta}}{\operatorname{arg\min}} E(L(\hat{\theta}, \theta))$$

is Bayes estimator.

After data is sampled:  $x$  observed of  $X | \hat{\theta}$

Prior is replaced by posterior distribution of  
 $\hat{\theta} | x$

The Bayes estimator is

$$\hat{\theta}(x) = \underset{\hat{\theta}}{\operatorname{arg\min}} E(L(\hat{\theta}, g(x)) | x)$$

### Theorem

Assume  $\Theta \sim \Lambda$ ,  $X|\Theta = \theta \sim P_\theta$ . Loss function

$L(\theta, \delta)$  non-negative for estimating  $g(\Theta)$ . Assume

a) There is an estimator with finite risk.

b) For almost all  $x$  ( $P_\theta$ -a.s.), ~~exists~~

$$\delta_\lambda(x) = \underset{\delta}{\operatorname{argmin}} E(L(\Theta, \delta(x)) | X=x)$$

exists.

Then  $\delta_\lambda(X)$  is a Bayes estimator.

### Proof.

Assume  $\delta$  is an arbitrary estimator with

$$\int R(\theta, \delta) d\Lambda(\theta) = \int E_L(L(\Theta, \delta(\Theta))) d\Lambda(\Theta)$$

(there is at least one such by a). Then

$$= E_\Lambda \int L(\Theta, \delta(\Theta)) dP_\Theta(x)$$

$$= \int E_\Lambda(L(\Theta), \delta(x)) | X=x dP_\Theta(x) < \infty$$

(there is at least one such by a). Since  $L \geq 0$  and  $P_\Theta$  is a prob. measure, that means that

$E_{\lambda}(\mathcal{L}(\theta, \delta(x)) | \bar{x}=x) < \infty$

and thus by b)

Ute Bayes  
Theore-  
und Th.  
→ posterior  
distribution

$$E(\mathcal{L}(\theta, \delta(x)) | \bar{x}=x) \geq E(\mathcal{L}(\theta, \delta_{\lambda}(x)) | \bar{x}=x)$$

$P_{\theta}$ -a.s. Do  $\int \dots dP_{\theta}$  of both sides to obtain

$$\int R_{\lambda}(\theta, \delta) d\lambda(\theta) \geq \int R(\theta, \delta_{\lambda}) d\lambda(\theta)$$

for all  $\delta$ , i.e.  $\delta_{\lambda}$  is a Bayes estimator #

By taking a few particular loss functions we get the following special cases

Corollary

Assume the assumptions of the theorem hold.

(i) If  $L(\theta, d) = (d - g(\theta))^2$  is quadratic loss then

$$\delta_{\lambda}(x) = \frac{E(g(\theta)|x)}{\lambda}$$

(If  $L(\theta, d) = (d - g(\theta))^2 \omega(\theta)$  quadratic loss with weights then

$$\delta_{\lambda}(x) = \frac{E(g(\theta)\omega(\theta)|x)}{E(\omega(\theta)|x)}$$

(ii) If  $L(\theta, d) = |d - g(\theta)|$  is absolute value loss then  $\hat{\delta}_L(x)$  is (any) median of the conditional distribution of  $\mathbb{U} | \bar{x}=x$ .

(iii) If

$$L(\theta, d) = \begin{cases} 0 & |d - \theta| \leq c \\ 1 & |d - \theta| > c \end{cases} = 1\{|d - \theta| > c\}$$

$\downarrow g(\theta)$

then  $\hat{\delta}_L(x)$  is the midpoint of interval of length  $2c$  which maximizes  $P(\mathbb{U} \in I | x)$

#

Proof.

(i) We want to find  $\hat{\delta}_L(x)$  that minimizes  $E((g(\mathbb{U}) - \hat{\delta}_L(x))^2 | x)$ .

That is possible by the theorem and the minimizing value is (conditional  $L^2$ -projection...)

$$E(g(\mathbb{U}) | x)$$

(ii) Yes, and conditional median solve it

(iii) Yes, just use  $E(1\{|d - \mathbb{U}| > c\} | x) = P(|d - \mathbb{U}| > c | x)$  and change minimize this to maximize  $P(|d - \mathbb{U}| \leq c | x)$ .

(8)

#

We can next state a result on uniqueness of the Bayes estimator.

Corollary.

Assume the loss function  $\lambda(\theta, \delta)$  is strictly convex in  $\delta$ . Assume

- $\int L(\theta, \delta_\lambda) d\Lambda(\theta)$  is finite
- Let  $Q$  be the marginal distribution of  $X$

$$Q(A) = \int P_\theta(X \in A) d\Lambda(\theta)$$

Then ~~either  $\delta_\lambda$  or  $\delta$  is Q-a.s.~~ implies  $P_{\theta \sim Q}(\delta_\lambda)$  for all  $\theta$ .

~~Proof (read on your own, carefully).~~ #

Then the Bayes estimator  $\delta_\lambda$  is unique  $P_{\theta \sim Q}$ -a.s. for all  $\theta$ .

Proof (read on your own, carefully) #

An example on when condition (ii) is not satisfied is

Assume

Ex:  ~~$B(X \sim \text{Bin}(n, p))$~~  distribution and  $\Lambda$  is two-point distribution:  $\Lambda$  puts mass  $1/2$  on  $p=0$  and mass  $1/2$  on  $p=1$ . Let  $\delta(X)$  be an arbitrary estimator such that  $\delta(0)=0$ ,  $\delta(n)=1$  ( $\delta(k)$  for  $k=1, \dots, n-1$  does not matter). Assume we have  $n$  general

(9)

error loss. Then the Bayes risk is

$$\begin{aligned} & \int R(\theta, \delta) d\lambda(\theta) \\ &= \frac{1}{2} R(0, \delta) + \frac{1}{2} R(1, \delta) \\ &= \frac{1}{2} E_0 ((0 - \delta(x))^2) + \frac{1}{2} E_1 ((1 - \delta(x))^2) \end{aligned}$$

But if  $p=0$ ,  $X=0$  with prob 1 and if  $p=1$ ,  $X=n$  with prob 1 so the above is

$$\begin{aligned} & \cancel{\frac{1}{2} E_0 ((0 - \delta(0))^2) + \frac{1}{2} E_1} \\ &= \frac{1}{2} (0 - \delta(0))^2 + \frac{1}{2} (1 - \delta(n))^2 \end{aligned}$$

which is  $= 0$  if  $\delta(x)$  is as above, and then is the smallest possible value. So then the Bayes estimator is not unique.

Furthermore the marginal distribution of  $X$  is

$$\begin{aligned} Q(A) &= \cancel{P_{X \in A}} = \int P_\theta(X \in A) d\lambda(\theta) \\ &= P_0(X \in A) \frac{1}{2} + P_1(X \in A) \frac{1}{2} \end{aligned}$$

and the set  $A=\{0, n\}$  has  $Q$  more 1, last for instance

$$\cancel{Q(A) =}$$

if  $p=\frac{1}{2}$

$$\begin{aligned} P_\theta(X \in A) &= P_\theta(0) + P_\theta(n) \\ &= \left(\frac{1}{2}\right)^n + \left(\frac{1}{2}\right)^n < 1 \end{aligned}$$

so the condition (ii) is not satisfied

How does one evaluate a Bayes estimator?

Different (several) risk functions:

(i) The Bayes risk

$$\int R(\theta, \delta) d\lambda(\theta)$$

(ii) The conditioned risk given  $\Theta = \theta$ , the previously  
studied risk

$$R(\theta, \delta)$$

(iii) The a posteriori risk

$$E(L(\theta, \delta(x)) \mid X=x)$$

The Bayes risk is the criterion for the Bayes estimator.

(ii) can be used to make sure that one does not get "very poor" results simply by choosing a "proper" prior  $\lambda$ . (iii) is of use to the hard core Bayesian.

Ex: Assume  $X \in \text{Bin}(n, p)$ . Choose as prior distribution for  $p$  a Beta distribution  $B(a, b)$  with density

$$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}; \quad 0 < x < 1$$

and

$$(11) \quad E(p) = \frac{a}{a+b} \quad \text{Var}(p) = \frac{ab}{(a+b)^2 (a+b+1)}$$

We want to find an estimator (the Bayes estimator) for  $g(p)$ . We need to get conditional distribution of  $p$  given  $x$ . Joint density of  $x$  and  $p$  is (w.r.t. product of Lebesgue measure and counting measure)

$$(n) \quad \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{x+a-1} (1-p)^{n-x+b-1}$$

$x = 0, 1, 2, \dots, n$ ,  $0 < p < 1$ . Integrate out  $p$  to get the marginal density (pmf) of  $x$ ; this will be a function of only  $x$ . Divide by this to get conditional density of  $p$  given  $x$

$$c(a, b, x) p^{x+a-1} (1-p)^{n-x+b-1}$$

Recognise as Beta distribution with parameters

$$a' = a+x$$

$$b' = b+n-x$$

Assume we have a quadratic loss and  $g(p) = p$ . Then we know

Beta distri.

$$\delta_1(x) = E(p|x) \stackrel{?}{=} \frac{a'}{a'+b'} = \frac{a+x}{a+b+n}$$

$$= \left( \frac{a+b}{a+b+n} \right) \frac{a}{a+b} + \left( \frac{n}{a+b+n} \right) \frac{x}{n}$$

Bayes estimator  
before observations  
taken

Standard  
estimator  
(UMVU etc.)

and the Bayes estimator is a weighted average of the two!

(i)  $a \rightarrow \infty, b \rightarrow \infty, \frac{a}{b}$  fixed  $\Rightarrow \text{Var}(p) \rightarrow 0$   
 all mass of  $p$  concentrates on  $\frac{a}{a+b}$   
 so  $\lambda_{a,b} \rightarrow \underline{\delta}_{\frac{a}{a+b}}$  nice measure

Then

$$\underline{\delta}_\lambda \rightarrow \underline{\delta}_{\frac{a}{a+b}} \quad \begin{array}{l} \text{Bayes before data gathered} \\ \text{so data has no effect.} \end{array}$$

(ii)  $a, b$  fixed,  $n \rightarrow \infty$

$$\Rightarrow \underline{\delta}_\lambda \xrightarrow{n \rightarrow \infty} \frac{x}{n} \quad (\text{rather } \frac{\underline{\delta}_\lambda(x)}{x/n} \rightarrow 1)$$

and the a priori distribution has no effect on the Bayes estimator

#

Finally note that  $a=b=0$  gives

$$\underline{\delta}_\lambda(x) = \frac{x}{n} \quad (\text{the ordinary estimator}).$$

But  $B(0, v)$  is not a proper distribution since

$$\text{then } \int_0^1 f_p(p) dp = \infty$$