

- Thus, Bayes rule is

mean is the best prediction of a r.v. under MSE

reasonable? → a decision function (an estimator)

do it for every  $x$

$$\hat{\theta}(x) = \begin{cases} \int \theta h(\theta|x) d\theta, & \text{in the continuous case} \\ \sum \theta_i h(\theta_i|x), & \text{in the discrete case} \end{cases}$$

- In the case of squared error loss, the Bayes estimator (i.e., Bayes rule) is the mean of the posterior distribution.

5/22

### Example 10.7 (Throw a coin once, Bayes estimator, 2nd Ed., TBp. 584-585)

- A biased coin is thrown once. Estimate  $\theta$  = probability of heads.

- Suppose that we have no idea how biased the coin is  $\Rightarrow$  for  $\theta$ , can use uniform prior:  $g(\theta) = 1, 0 \leq \theta \leq 1$ .

- Let a vague prior  $\rightarrow$  Data  $\rightarrow X = \begin{cases} 1, & \text{if a head appears} \\ 0, & \text{if a tail appears} \end{cases}$

Then the distribution of  $X$  given  $\theta$  is Bernoulli( $\theta$ ):

conditional

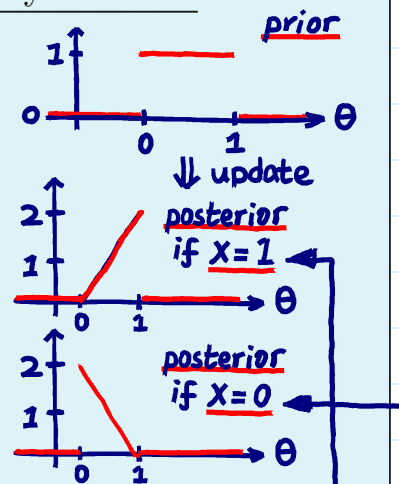
This is the pmf of  $X$  in CH8 & CH9

$$f(x|\theta) = \begin{cases} \theta, & x=1 \\ 1-\theta, & x=0 \end{cases}$$

- The posterior distribution is

conditional  $\theta|x \rightarrow h(\theta|x)$  joint  $(\theta, x) \rightarrow f(x|\theta) \times \frac{1}{g(\theta)}$  marginal  $X \rightarrow \int_0^1 f(x|\theta) \times \frac{1}{g(\theta)} d\theta$

$$h(\theta|x) = \begin{cases} \frac{\theta}{\int_0^1 \theta d\theta} = 2\theta, & x=1 \\ \frac{1-\theta}{\int_0^1 1-\theta d\theta} = 2(1-\theta), & x=0 \end{cases}$$



- The Bayes estimator of  $\theta$  is the posterior mean

posterior median  $\hat{\theta} = \begin{cases} 1/2, & \text{if } x=1 \\ 1/2, & \text{if } x=0 \end{cases}$

$$\int_0^1 \theta \cdot h(\theta|x) d\theta = \begin{cases} \int_0^1 \theta \cdot (2\theta) d\theta = 2/3, & x=1 \\ \int_0^1 \theta \cdot 2(1-\theta) d\theta = 1/3, & x=0 \end{cases}$$

under square error loss ( $L^2$ -norm)

Q: Which one is more reasonable estimator?

- Note. The maximum likelihood estimator is  $X = \begin{cases} 1, & x=1 \\ 0, & x=0 \end{cases}$   $\leftarrow$  Frequentist approach

Note. With one observation, the data carries very little information about  $\theta \Rightarrow$  prior becomes very helpful.

### Theorem 10.3 (Bayes rule for Estimation under Absolute Error Loss)

In the case of absolute error loss, i.e.,

FYI  $l(\theta, d) = |\theta - d| \leftarrow L^1\text{-norm}$

robust estimator

median is the best predictor of a r.v. under absolute error

the Bayes estimator is the median of the posterior distribution.  $\leftarrow$  cf. Thm 10.2 (LNp.15)

50%-quantile

### Definition 10.4 (dominate, strictly dominate, admissible, 2nd Ed., TBp.585)

- Let  $d_1, d_2$  be two decision functions.

Say  $d_1$  dominates  $d_2$ , if  $R(\theta, d_1) \leq R(\theta, d_2)$  for all  $\theta$ .

risk function

$d_1$  is no worse than  $d_2$  But,  $d_1, d_2$  can be equally good.

Say  $d_1$  strictly dominates  $d_2$ , if  $R(\theta, d_1) < R(\theta, d_2)$  for all  $\theta$ , and the inequality is strict for some  $\theta$ .  $\leftarrow$   $d_1$ : no worse than  $d_2$ , but better for some  $\theta$

- A decision function  $d$  is called admissible if  $d$  is not strictly dominated by any other decision function.

cf.  $\rightarrow$  inadmissible (not admissible)

- For estimation, an estimator is called admissible if it is not strictly dominated by any other estimator. (Note: Admissibility is a rather weak property. There could be many admissible estimators.)  $\leftarrow$  cf.  $\rightarrow$  UMVUE  $\Rightarrow$  only one admissible est'or

**Theorem 10.4 (2nd Ed., TBp.586)**← true for any decision problem, not only estimation.Suppose that one of the following two assumptions holds:

1. <sup>(a)</sup>  $\Omega$  is discrete, <sup>(b)</sup>  $d^*$  is a Bayes rule w.r.t. a prior pmf  $g(\theta)$  such that <sup>(c)</sup>  $g(\theta) > 0$  for all  $\theta \in \Omega$ .

countable

2. <sup>(a)</sup>  $\Omega$  is an interval (perhaps infinite) and <sup>(b)</sup>  $d^*$  is a Bayes rule w.r.t a prior pdf  $g(\theta)$  such that <sup>(c)</sup>  $g(\theta) > 0$  for all  $\theta \in \Omega$  and <sup>(d)</sup>  $R(\theta, d)$  is a continuous function of  $\theta$  for all  $d$ .

uncountable

← risk function

Then  $d^*$  is admissible. ← Intuition: If  $d'$  strictly dominates  $d^*$ ,  $d'$  has smallerBayes risk than  $d^*$ .Suppose that  $d^*$  is inadmissible. Then, there is another estimator  $d'$  such that $d'$  strictly dominates  $d^*$ 

$$R(\theta, d^*) \geq R(\theta, d'), \quad \text{for all } \theta$$

and with strict inequality for some  $\theta$ , say  $\theta_0$ .Since  $R(\theta, d^*) - R(\theta, d')$  is a continuous function of  $\theta$ , there exist an  $\epsilon > 0$  and an  $h > 0$  such that

$$R(\theta, d^*) - R(\theta, d') > \epsilon > 0 \quad \text{for } \theta_0 - h \leq \theta \leq \theta_0 + h$$

Then,

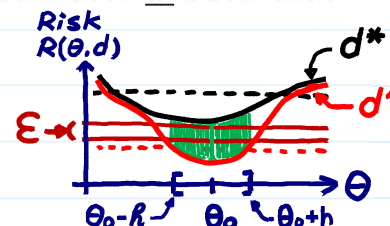
$$\int_{-\infty}^{\infty} [R(\theta, d^*) - R(\theta, d')] g(\theta) d\theta \geq \int_{\theta_0-h}^{\theta_0+h} [R(\theta, d^*) - R(\theta, d')] g(\theta) d\theta$$

$$B(d^*) - B(d')$$

Bayes risk

$$\geq \epsilon \int_{\theta_0-h}^{\theta_0+h} g(\theta) d\theta > 0$$

∴ (c)

But this contradict the fact that  $d^*$  is Bayes rule with respect to  $g(\theta)$ , that is

$$B(d^*) - B(d') = \int_{-\infty}^{\infty} [R(\theta, d^*) - R(\theta, d')] g(\theta) d\theta \leq 0.$$

The proof is complete.

**Notes.**The theorem can be regarded as both a positive and a negative result.positive part: It identifies a certain class of estimators (i.e., Bayes estimators) as being admissible.← for estimation  $d$ (Note: In general, it could be difficult to examine whether a particular estimator is admissible, since by definition, one have to check that the estimator was not strictly dominated by any other estimators.)negative part: There are so many admissible estimators — one for every prior distribution satisfying the hypotheses of the theorem.⇒ admissibility is a rather weak property ⇒ can put more restriction.

❖ Reading: textbook (2nd ed.), 15.2.4

e.g., unbiased → UMVUE ← unique**• The Subjectivist Point of View — where the prior distributions come from?**

- Different viewpoints on probability.

Question. how to interpret the statement:objective (客)  
subjective (主)“the probability that a coin will land heads up is  $1/2$ ”?a distribution assigned on  $\Omega$   
parameter: fixed value  
→ random variable←  $\theta \in \Omega$



**Frequentists' interpretation:** the long run relative frequency of heads approaches  $1/2$ .  $\leftarrow$  long-run average  $\leftarrow$  LLN  $\leftarrow$  objective

**Bayesian:** the prior opinion ( $1/2$ ) is such that they would as soon guess heads or tails if the rewards are equal.  $\leftarrow$  subjective

a fixed world ( $\theta$ )  
probability that Shakespeare wrote Hamlet

### Bayesian View of Probability (Personal Opinion) (2nd Ed., TBp. 587-588).

- Let  $A$  be an event (e.g., image that  $A \subset S$ : sample space).  $\rightarrow$  regarded as parameter.
- For a game in which if  $A$  occurs, the Bayesian will be rewarded \$1,  $P(A)$  is the maximum amount of money the Bayesian would be willing to pay to buy into the game.  $\rightarrow$  a fair game, average reward = pay,  $1 \cdot p - R = 0 \Leftrightarrow R = p$ .  $\rightarrow$  Frequentist approach
- Example: if the Bayesian is willing to pay at most 50 cents to buy in,  $P(A) = 0.5$ .
- Note:** Bayesian probability is personal.  $P(A)$  may vary from person to person.  $\rightarrow$  cf.
- Bayesian probability is a model for quantifying the strength of personal opinions or personal beliefs. Bayesian: 每人心心中皆可有不同的平行世界, Frequentist: 每人心中的平行世界皆相同

### Evolution of Personal Opinion (2nd Ed., TBp. 588).

- Bayes Thm describes how personal opinions evolve with experience.

Suppose that the prior probability of  $\theta$  is  $P(\theta)$ .

On observation of an event  $A$ ,  $\leftarrow$  every person observe same  $A$

the opinion about  $\theta$  changes to:

can result in consistent probability with the long-run average interpretation

posterior dist.

$P(\theta | A)$

$P(A | \theta)$

$P(\theta)$

$P(A)$

Bayes' Thm (LN, CHI~6, P.23)

lot to lot in Ex10.1 (LNp.1)

### Difference btw Frequentist and Bayesian Approaches – Point Estimation (TBp. 588)

- Data  $X$  (random) has a joint pdf/pmf  $f(x | \theta)$ .  $\leftarrow$  parameter  $\leftarrow$  Freq.: dist. of  $X$  Baye.: cond. dist. of  $X$
- Frequentist:**  $\theta$  has some fixed value that is unknown. Since  $\theta$  is not random, it makes no sense to assign  $\theta$  a pdf/pmf.  $\rightarrow$  cf.
- Bayesian:** The prior opinion is  $g(\theta)$  for  $\theta$  (random), and  $f(x | \theta) \rightarrow f_{X|\theta}(x | \theta)$ . Having observed the data  $X = x$ , the new opinion about  $\theta$  is  $\leftarrow$  conditional dist.  $\rightarrow$  cf.

Bayes est'or  $\rightarrow$  update  $\rightarrow$  posterior

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

$f_x(x)$  or  $P_x(x)$ : marginal

prior can reflect personal experience. However, because it's personal, it can be biased

- Note 1:** As a function of  $\theta$ , the posterior density  $h(\theta | x)$  is proportional to  $f(x | \theta) g(\theta)$ , i.e.,  $h(\theta | x) \propto f(x | \theta) g(\theta)$ .

no particular preference on some specific  $\theta$ 's

if  $g(\theta) \approx c$ , a constant

Check Ex10.7 (LNp.15)  $\rightarrow$  usually centering  $\rightarrow$  posterior (Bayesian)  $\leftarrow$  cf.  $\rightarrow$  usually maximizing  $\rightarrow$  likelihood (Frequentist)

- Note 2:** If  $g(\theta)$  is nearly uniform in the region of  $\theta$  where  $f(x | \theta)$ , regarded as a function of  $\theta$ , has almost all its mass, then  $h(\theta | x)$  is nearly proportional to the likelihood function  $\mathcal{L}(\theta) = f(x | \theta)$ .  $\leftarrow$  In likelihood  $\leftarrow$   $x$ : fixed  $\theta$ : variable

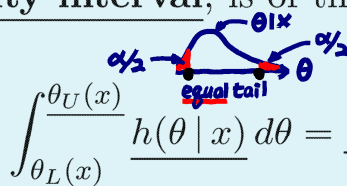
### Difference btw Frequentist and Bayesian Approaches – Interval Estimation (TBp. 588)

- Frequentist:** The confidence interval is a random interval  $(\theta_L(X), \theta_U(X))$ . A  $100(1 - \alpha)\%$  C.I. covers the true, fixed, unknown value of  $\theta$  with probability  $1 - \alpha$ . Once  $X = x$  is observed,  $P(\theta \in (\hat{\theta}_L(x), \hat{\theta}_U(x)))$  is 0 or 1.  $\rightarrow$  Bayesian approach  $\rightarrow$  cf.

- **Bayesian:** A Bayesian  $100(1 - \alpha)\%$  confidence interval, also called a **credibility interval**, is of the form  $(\theta_L(x), \theta_U(x))$  where  $\theta_L(x)$  and  $\theta_U(x)$

highest posterior density

satisfies



$$\int_{\theta_L(x)}^{\theta_U(x)} h(\theta | x) d\theta = 1 - \alpha = P(\theta \in (\theta_L(x), \theta_U(x)) | X = x)$$

Frequentist approach

Once  $X = x$  is observed, the interval is fixed and  $\theta$  is random.

### Difference btw Frequentist and Bayesian Approaches - Testing (2nd Ed., TBp. 589)

- **Frequentist:** probabilities of Type I and Type II errors ( $\alpha$  and  $\beta$ )
- **Bayesian:** deciding between two hypotheses (i.e., two subsets of  $\Omega$ ) in light of data reduces to comparing their posterior probabilities.  $\Omega_0$  &  $\Omega_A$

check Lnp.32~35

$$h(\theta | x), \theta \in \Omega_0$$

$$h(\theta | x), \theta \in \Omega_A$$

❖ Reading: textbook (2nd ed.), 15.3

### Bayesian Inference for the Normal distribution

#### Theorem 10.5 (2nd Ed., TBp. 590)

dist. of data in Frequentist  
cond'nal dist. of data in Bayesian

- Suppose that a single observation  $X$  is taken, and
    - random —  $X | \mu \sim N(\mu, \sigma^2)$ , where  $\mu$  is parameter and  $\sigma^2$  known
    - fixed —  $\mu \sim N(\mu_0, \sigma_0^2)$  ← Bayesian
- assumed known for simplicity

precision  $\xi = 1/\sigma^2$  and  $\xi_0 = 1/\sigma_0^2$

The posterior distribution of  $\mu$  is normal with mean

best predictor of  $\mu$  under prior  $\Rightarrow$  best estimate of  $\mu$  before observing data

Frequentist estimator of  $\mu$ 

update

Bayes est'or of  $\mu \rightarrow \mu_1$

and variance  $\sigma_1^2 = 1/\xi_1$ , where  $\xi_1 = \xi + \xi_0$ .

$$\mu_1 = \frac{\xi_0 \mu_0 + \xi x}{\xi + \xi_0} = \frac{\xi_0}{\xi + \xi_0} \mu_0 + \frac{\xi}{\xi + \xi_0} x$$

Sum=1

if  $\xi \gg \xi_0$   
if  $\xi \ll \xi_0$

**Proof.** The posterior distribution of  $\mu$  is

$$h(\mu | x) \propto f(x | \mu) g(\mu) \propto \exp \left[ -\frac{1}{2\sigma^2} (x - \mu)^2 \right] \times \exp \left[ -\frac{1}{2\sigma_0^2} (\mu - \mu_0)^2 \right]$$

similar form:  
2nd-order polynomial of  $\mu$

$$\sigma_1^2 = \frac{\sigma^2 \sigma_0^2}{\sigma^2 + \sigma_0^2}$$

When will  $\sigma_1^2$  small?

This is a pdf of  $\mu$

That is, it is known  $\int h(\mu | x) d\mu = 1$

$$\propto \exp \left\{ -\frac{1}{2} \left[ \mu^2 \left( \frac{1}{\sigma^2} + \frac{1}{\sigma_0^2} \right) - 2\mu \left( \frac{x}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) + \frac{x^2}{\sigma^2} + \frac{\mu_0^2}{\sigma_0^2} \right] \right\}$$

$$\propto \exp \left[ -\frac{1}{2a^{-1}} \left( \mu - \frac{b}{a} \right)^2 \right]$$

$\rightarrow$  It has the form of normal pdf.

where  $a = \sigma^{-2} + \sigma_0^{-2}$  and  $b = x\sigma^{-2} + \mu_0\sigma_0^{-2}$ . Thus the posterior distribution of  $\mu$  is normal with mean

$$\mu_1 = \frac{b}{a} = \frac{x(1/\sigma^2) + \mu_0(1/\sigma_0^2)}{1/\sigma^2 + 1/\sigma_0^2} = \frac{x\xi + \mu_0\xi_0}{\xi + \xi_0}$$

and variance  $a^{-1} = (\sigma^{-2} + \sigma_0^{-2})^{-1} = (\xi + \xi_0)^{-1}$ .



## Notes (2nd Ed., TBp. 590).

- For a normal distribution  $N(\mu, \sigma^2)$ ,  $1/\sigma^2$  is called its precision.  $\uparrow$  when  $\sigma^2 \downarrow$
- $\xi$  ( $= 1/\sigma^2$ ),  $\xi_0$  ( $= 1/\sigma_0^2$ ), and  $\xi_1$  ( $= 1/\sigma_1^2$ ) are precisions of  $X|\mu$ , prior, and posterior distributions, respectively. Notice that  $\xi_1 = \xi + \xi_0$ .  $\therefore$  combine information

biased estimator under  $X|\mu$ 

$$\sigma_1^2 < \sigma_0^2$$

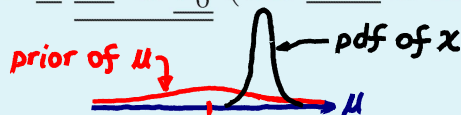
posterior

data

prior

- The posterior mean is weighted average of prior mean  $\mu_0$  and data  $x$ , with weights proportional to the respective precisions.

- If  $\sigma^2 \ll \sigma_0^2$  (the data is much more informative than the prior), then



$$\xi \gg \xi_0$$

$$\text{and } \xi_1 \approx \xi$$

$$\Rightarrow \mu_1 \approx x$$

Bayes est'or

Thus,  $h(\mu|x) \approx f(x|\mu)$ , i.e.,  $\mu|x$  is nearly distributed as  $N(x, 1/\xi)$ 

posterior

likelihood

update

**Exercise:** What if  $\sigma^2 \gg \sigma_0^2$ ?  $\rightarrow$  prior dominates

$$N(\mu_0, 1/\xi_0)$$

- If prior distribution is quite flat relative to  $f(x|\mu)$ , as a function of  $\mu$ ,

likelihood

1. the prior distribution has little influence on the posterior,2. the posterior distribution is approximately proportional to the likelihood function. $\rightarrow$  Bayesian inference  $\leftarrow$  consistent  $\rightarrow$  Frequentist inference  
(may not be identical)Such a prior is often called a vague, or noninformative, prior.