

$$\frac{d}{dx} \sin^{-1} \sqrt{x}$$

$$\text{let } \sin^{-1} \sqrt{x} = \theta \Rightarrow \sqrt{x} = \sin \theta$$

$$\Rightarrow x = \sin^2 \theta$$

$$\begin{aligned} \textcircled{1} & \\ &= \frac{d\theta}{d\sin^2 \theta} = \frac{1}{\frac{d\sin^2 \theta}{d\theta}} = \frac{1}{2\sin \theta \cos \theta} = \frac{1}{2\sqrt{x}\sqrt{1-x}} \end{aligned}$$

Recap:

Maximum likelihood estimator.

① Bayes estimator ② Minimax estimator.

Frequentist:

$x_1, \dots, x_n \stackrel{iid}{\sim} f_{\theta}(x)$

Goal: provide point estimate of θ .

Frequentist

data θ are r.v.'s
and parameter θ is
fixed but unknown.

Bayesian

data are fixed,
parameter θ is a
random variable.

Bayesian: Goal is to estimate the unknown
distribution of the parameter θ .

Algorithm: ① Start with a prior distribution
on θ . This represents our prior ~~belief~~ belief on θ .

② Use this prior distribution and the data
to find the posterior distribution of θ .

③ Algorithm: Let $x_1, \dots, x_n \stackrel{iid}{\sim} f_{\theta}(x)$

Let the prior dist. θ of θ be given by $\pi(\theta)$.

The posterior dist. of $\theta | x_1, \dots, x_n$, denoted
by $\pi(\theta | x_1, \dots, x_n)$ is defined as

$$\pi(\theta | x_1, \dots, x_n) \stackrel{\text{①}}{=} \frac{\left[\prod_{i=1}^n f_{\theta}(x_i) \right] \pi(\theta)}{\int \left[\prod_{i=1}^n f_{\theta}(x_i) \right] \pi(\theta) d\theta}$$

Example: $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Ber}(p)$.

Goal: Posterior dist. of p .

① Since $0 < p < 1$, a reasonable prior dist. on p is given by $p \sim \text{Beta}(\alpha, \beta)$.

$$\pi(p) = \frac{p^{\alpha-1} (1-p)^{\beta-1}}{\text{Beta}(\alpha, \beta)}, \quad 0 < p < 1$$

$$\begin{aligned} \pi(p | x_1, \dots, x_n) &= \frac{\left[\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \right] \pi(p)}{\int \left[\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \right] \pi(p) dp} \\ &= \frac{p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i} p^{\alpha-1} (1-p)^{\beta-1}}{\int p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i} p^{\alpha-1} (1-p)^{\beta-1} dp} \end{aligned}$$

$$= \frac{p^{\alpha + \sum_{i=1}^n x_i - 1} (1-p)^{\beta + n - \sum_{i=1}^n x_i - 1}}{\int p^{\alpha + \sum_{i=1}^n x_i - 1} (1-p)^{\beta + n - \sum_{i=1}^n x_i - 1} dp}$$

$$\Rightarrow p | x_1, \dots, x_n \sim \text{Beta}\left(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i\right)$$

Example: $x_1, \dots, x_n \overset{i.i.d.}{\sim} \text{Pois}(\lambda)$, $\lambda \sim \text{Gamma}(a, b)$

$$\pi(\lambda) = \frac{\lambda^{a-1} e^{-\lambda b} b^a}{\Gamma(a)}, \quad \lambda > 0$$

$$\pi(\lambda | x_1, \dots, x_n) = \frac{\left[\prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \right] \lambda^{a-1} e^{-\lambda b} b^a}{\int \left[\prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \right] \lambda^{a-1} e^{-\lambda b} b^a d\lambda}$$

$$= \frac{e^{-(n+b)\lambda} \lambda^{\sum_{i=1}^n x_i + a - 1}}{\int e^{-(n+b)\lambda} \lambda^{\sum_{i=1}^n x_i + a - 1} d\lambda}$$

$$\lambda | x_1, \dots, x_n \sim \text{Gamma}\left(\sum_{i=1}^n x_i + a, n+b\right)$$

Conjugate family: ~~poisson~~

Let F denote the class of pdfs or pmfs. $f(x|\theta)$.
 A class π of priors distributions in a conjugate family for F if the posterior distribution is in the class π for all $f \in F$, all priors in π and all $x \in \mathcal{X}$.

Now we will learn how to create a good estimator using this Bayesian knowledge.

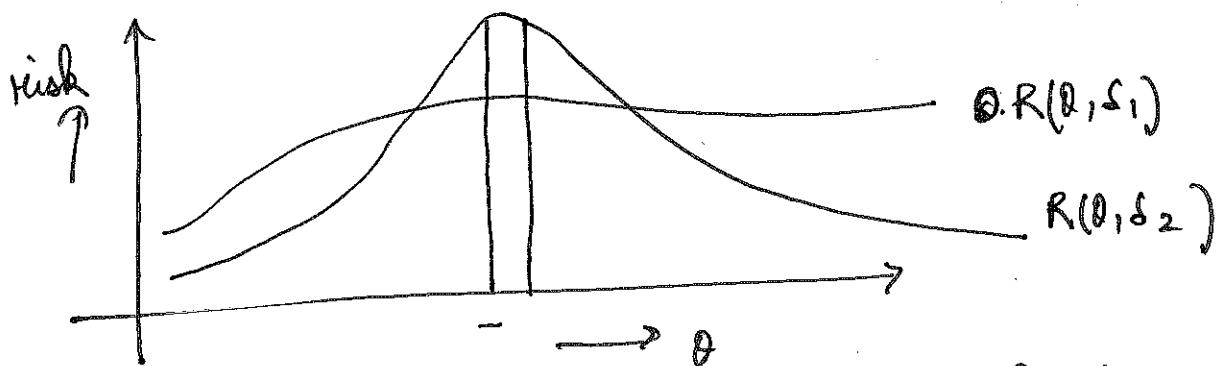
Let $\delta(x)$ be an estimator of θ . The loss in estimating θ by $\delta(x)$ is represented by a function known as the loss function. We denote the loss fn. by $L(\theta, \delta) / L(\theta, \delta(x))$

$R(\theta, \delta) = E_{x|\theta} [L(\theta, \delta(x))]$ is known as the risk function of δ .

$L(\theta, \delta(x)) = (\theta - \delta(x))^2 \rightarrow$ squared error loss function

$R(\theta, \delta) = E_{x|\theta} [L(\theta, \delta(x))] = E_{x|\theta} [(\theta - \delta(x))^2] \rightarrow$ Mean squared error.

Given two estimators δ_1 and δ_2 , we say δ_1 is better than δ_2 if $R(\theta, \delta_1) \leq R(\theta, \delta_2) \forall \theta$ and the inequality is strict at least for one θ .



We want to create a summary of the entire risk fn.

1. Average risk = $E_{\theta} [R(\theta, \delta(x))] = E_{\theta} E_{x|\theta} [L(\theta, \delta(x))]$

2. Supremum risk = $\sup_{\theta} R(\theta, \delta(x))$

If $\delta(x)$ is an estimator of θ that minimizes the average risk over all estimators, then $\delta(x)$ is called the Bayes estimator of θ .

$$\text{Average risk} = E_{\theta} E_{x|\theta} [L(\theta, \delta(x))]$$

$$= E_{x, \theta} [L(\theta, \delta(x))]$$

$$= E_x E_{\theta|x} [L(\theta, \delta(x))]$$

If \exists an estimator $\delta(x)$ that minimizes

$$E_{\theta|x} [L(\theta, \delta(x))] \quad \forall x$$

then that is going to be the minimizer of the average risk.

Finding the Bayes estimator is equivalent to finding an estimator $\delta(x)$ that minimizes

$$E_{\theta|x} [L(\theta, \delta(x))]$$

If $L(\theta, \delta(x)) = (\theta - \delta(x))^2$

argmin $E_{\theta|x} [(\theta - \delta(x))^2]$
 $\delta(x)$

$$E_{\theta|x}^* [(\theta - \delta(x))^2] = E_{\theta|x} [(\theta - E_{\theta|x}(\theta) + E_{\theta|x}(\theta) - \delta(x))^2]$$

$$= E_{\theta|x} [(\theta - E_{\theta|x}(\theta))^2] + 2 E_{\theta|x} [(\theta - E_{\theta|x}(\theta)) (E_{\theta|x}(\theta) - \delta(x))] + E_{\theta|x} [(E_{\theta|x}(\theta) - \delta(x))^2]$$

$$= E_{\theta|x} [(\theta - E_{\theta|x}(\theta))^2] + 2 (E_{\theta|x}(\theta) - \delta(x)) [E_{\theta|x}(\theta) - E_{\theta|x}(\theta)] + E_{\theta|x} [(E_{\theta|x}(\theta) - \delta(x))^2]$$

$$\geq E_{\theta|x} [(\theta - E_{\theta|x}(\theta))^2]$$

thus $\delta(x) = E_{\theta|x}(\theta)$ is the Bayes estimator in this case.

$$\text{argmin}_{\delta(x)} E_{\theta|x} \left((\theta - \delta(x))^r \right)$$

$$\frac{d}{d\delta} E_{\theta|x} \left((\theta - \delta)^r \right) = E_{\theta|x} \left(\frac{d}{d\delta} (\theta - \delta)^r \right) = -r E_{\theta|x} \left((\theta - \delta)^{r-1} \right) = 0$$

$$\Rightarrow \delta = E_{\theta|x}(\theta)$$

$$L(\theta, \delta(x)) = w(\theta) (\theta - \delta(x))^r \rightarrow \text{weighted loss fn.}$$

$$\text{argmin}_{\delta(x)} E_{\theta|x} \left(w(\theta) (\theta - \delta(x))^r \right)$$

$$\frac{d}{d\delta} E_{\theta|x} \left(w(\theta) (\delta - \theta)^r \right) = 0$$

$$\Rightarrow 2 E_{\theta|x} \left(w(\theta) (\delta - \theta) \right) = 0 \Rightarrow \delta = E_{\theta|x} \left(\frac{\theta w(\theta)}{w(\theta)} \right) = E_{\theta|x}(\theta w(\theta))$$

$$\Rightarrow \delta(x) = \frac{E_{\theta|x}(\theta w(\theta))}{E_{\theta|x}(w(\theta))}$$

Example: x_1, \dots, x_n i.i.d. $\text{Ber}(p)$, $p \sim \text{Beta}(\alpha, \beta)$.

$$\pi(p) \pi(p|x_1, \dots, x_n) = \frac{p^{\sum_{i=1}^n x_i + \alpha - 1} (1-p)^{n - \sum_{i=1}^n x_i + \beta - 1}}{\text{Beta}(\alpha + \sum_{i=1}^n x_i, n - \sum_{i=1}^n x_i + \beta)}$$

$$E[p|x_1, \dots, x_n] = \int p \pi(p|x_1, \dots, x_n) dp$$

$$= \int p^{\sum_{i=1}^n x_i + \alpha} (1-p)^{n - \sum_{i=1}^n x_i + \beta - 1} dp$$

$$= \frac{\text{Beta}(\alpha + \sum_{i=1}^n x_i + 1, n - \sum_{i=1}^n x_i + \beta)}{\text{Beta}(\alpha + \sum_{i=1}^n x_i, n - \sum_{i=1}^n x_i + \beta)}$$

$$\begin{aligned}
 &= \frac{\Gamma(\alpha + \sum_{i=1}^n x_i + 1) \Gamma(n - \sum_{i=1}^n x_i + \beta)}{\Gamma(\alpha + 1 + n + \beta)} \\
 &= \frac{\Gamma(\alpha + \sum_{i=1}^n x_i) \Gamma(n - \sum_{i=1}^n x_i + \beta)}{\Gamma(\alpha + n + \beta)} \\
 &= \frac{\sum_{i=1}^n x_i + \alpha}{\alpha + n + \beta} \quad \left(\text{Using the fact. } \Gamma(a+1) = a \Gamma(a) \right) \\
 &= \underbrace{\bar{x}}_{\text{data mean}} \cdot \frac{n}{\alpha + \beta + n} + \frac{(\alpha + \beta)}{(\alpha + \beta + n)} \cdot \underbrace{\frac{\alpha}{(\alpha + \beta)}}_{\text{prior mean}}
 \end{aligned}$$

posterior mean is the convex combination of data and prior.

Since $E[\bar{x}] = p \Rightarrow$ the above estimator can't be an unbiased estimator of p .

Thm: Under squared error loss, No unbiased estimator $\delta(x)$ of θ can be a Bayes estimator unless $E_{\theta} E_{x|\theta}((\theta - \delta(x))^2) = 0$

pf: Let $\delta(x)$ be an unbiased estimator of θ which is also a Bayes estimator. Thus,

$$E_{x|\theta}[\delta(x)] = \theta \longrightarrow \text{as } \delta(x) \text{ is unbiased}$$

and

$$E_{\theta|x}(\theta) = \delta(x) \longrightarrow \text{as } \delta(x) \text{ is Bayes estimator.}$$

$$E_{\theta} E_{x|\theta} [\delta(x)\theta] = E_{\theta} [\theta E_{x|\theta} [\delta(x)]] = E_{\theta} [\theta^2] \quad \dots \textcircled{1}$$

$$\text{Also, } E_x E_{\theta|x} [\delta(x)\theta] = E_x [\delta(x) E_{\theta|x} (\theta)] = E_x [\delta(x)^2] \quad \dots \textcircled{2}$$

$$\text{As, } E_{\theta} E_{x|\theta} [\delta(x)\theta] = E_x E_{\theta|x} [\delta(x)\theta]$$

\Rightarrow by $\textcircled{1}$ & $\textcircled{2}$

$$E_{\theta} E_{x|\theta} [\delta(x)^2] = E_{\theta} E_{x|\theta} [\delta(x)\theta] = E_{\theta} [\theta^2] \quad \dots \textcircled{3}$$

Hence,

$$\begin{aligned} E_{\theta, x} [(\delta(x) - \theta)^2] &= E_{\theta, x} [\theta^2 - 2\delta(x)\theta + \delta(x)^2] \\ &= E_{\theta} E_{x|\theta} [\delta(x)^2] + E_{\theta} [\theta^2] - 2 E_{\theta} E_{x|\theta} [\delta(x)\theta] \\ &= 0 \quad (\text{by } \textcircled{3}) \end{aligned}$$

~~both~~ \Rightarrow the above equation needs to be satisfied for $\delta(x)$ as both Bayes unbiased estimator, $\delta(x)$ to be a Bayes and unbiased estimator, the above equation needs to be satisfied.

Example: X_1, \dots, X_n i.i.d $N(\mu, \sigma^2)$, σ^2 is known,

\bar{X} is the UMVUE and

$$E_{\mu} E_{x|\mu} ((\bar{X} - \mu)^2) = E_{\mu} \left(\frac{\sigma^2}{n} \right) = \frac{\sigma^2}{n} \neq 0$$

$\Rightarrow \bar{X}$ can't be a Bayes estimator.

Minimax estimator minimizes $\sup_{\theta} E_{X|\theta} [L(\theta, S(X))]$.

Minimax estimator is difficult to find in general. But there are certain cases in which Bayes estimator becomes a minimax estimator.

Least favourable distribution

A prior dist. $\pi(\theta)$ on θ is known to be a least favourable prior if

$$E_{\theta} E_{X|\theta} [L(\theta, S_{\pi}(X))] \geq E_{\theta} E_{X|\theta} [L(\theta, S_{\pi'}(X))]$$

for all prior dist. π' on θ . Here S_{π} and $S_{\pi'}$ are Bayes estimators w.r.t. priors π and π' respectively.

Recap:

Least favourable distribution.

A prior dist. $\pi(\theta)$ of θ is a least favourable prior

if $E_{\theta} E_{X|\theta} [L(\theta, \delta_{\pi}(x))] \geq E_{\theta} E_{X|\theta} [L(\theta, \delta_{\pi'}(x))]$ for all prior dist. π' on θ . δ_{π} and $\delta_{\pi'}$ are the Bayes estimator corresponding to priors π and π' .

Result: If $\pi(\theta)$ is a prior dist. for which

$\int E_{X|\theta} [L(\theta, \delta_{\pi})] \pi(\theta) d\theta = \sup_{\delta} E_{X|\theta} [L(\theta, \delta)]$, where δ_{π} is the Bayes estimator under the prior $\pi(\theta)$, then

a) δ_{π} is minimax

b) π is least favourable.

Prf: a) For any estimator $\delta(x)$,

$$\sup_{\delta} E_{X|\theta} [L(\theta, \delta)] \geq E_{\theta} E_{X|\theta} [L(\theta, \delta)]$$

$$\geq E_{\theta} E_{X|\theta} [L(\theta, \delta_{\pi})] = \sup_{\delta} E_{X|\theta} [L(\theta, \delta_{\pi})]$$

θ thus δ_{π} minimizes supremum risk over all estimators. Hence δ_{π} is minimax estimator.

b) Note that

$$E_{\theta} E_{X|\theta} [L(\theta, \delta_{\pi})] = \sup_{\delta} E_{X|\theta} [L(\theta, \delta_{\pi})]$$

$$\geq \int E_{X|\theta} [L(\theta, \delta_{\pi})] \pi'(\theta) d\theta \quad (\text{For any other prior dist. } \pi')$$

$$\geq \int E_{X|\theta} [L(\theta, \delta_{\pi'})] \pi'(\theta) d\theta$$

$$= E_{\theta} E_{X|\theta} [L(\theta, \delta_{\pi'})]$$

①

⇒ ~~π~~ in π is the least favourable dist.

Example: $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Ber}(p)$, $p \sim \text{Beta}(\alpha, \beta)$.

Bayes estimator of p under squared error loss

$$\delta_{\pi}(x) = \frac{\alpha + \sum_{i=1}^n x_i}{n + \alpha + \beta}$$

$$E_{x|p} \left((\delta_{\pi}(x) - p)^2 \right) = \frac{1}{(\alpha + \beta + n)^2} \left[\alpha^2 + \underbrace{\{n - 2\alpha(\alpha + \beta)\}}_{\text{constant}} p + \underbrace{\{(\alpha + \beta)^2 - n\}}_{\text{constant}} p^2 \right]$$

We want to make this risk fn. constant fn. of p

thus we set,

$$n = 2\alpha(\alpha + \beta) \quad \text{and} \quad (\alpha + \beta)^2 = n$$

$$\Rightarrow \alpha = \frac{\sqrt{n}}{2}, \quad \beta = \frac{\sqrt{n}}{2}$$

Thus, under the prior dist. $p \sim \text{Beta}\left(\frac{\sqrt{n}}{2}, \frac{\sqrt{n}}{2}\right)$,

Bayes estimator produces constant risk.

Now use the result to argue that this Bayes

estimator $\frac{\frac{\sqrt{n}}{2} + \sum_{i=1}^n x_i}{n + \sqrt{n}}$ is ~~the~~ a minimax

estimator under the prior $\text{Beta}\left(\frac{\sqrt{n}}{2}, \frac{\sqrt{n}}{2}\right)$.

Of course, by the same result $\text{Beta}\left(\frac{\sqrt{n}}{2}, \frac{\sqrt{n}}{2}\right)$ is the least favourable dist.

3. Testing of Hypothesis

Statistical hypothesis testing is all about

① Beginning with a tentative idea about the unknown parameter.

② Want to test validity of this tentative idea based on sample information.

③ Existing tentative idea : H_0 (null hypothesis),
 new idea : H_1 (alternative hypothesis).

④ We begin by assuming that the null hypothesis is true. Only when there is an overwhelming evidence contradicting the null do we reject it in favour of alternative.

	H_0 is true	H_0 is false
Do not reject H_0	Correct	Type 2 error.
reject H_0	Type 1 error	Correct.

Our goal is to minimize

$$P(\text{Type 1 error}) = P(\text{reject } H_0 \mid H_0 \text{ is true})$$

$$P(\text{Type 2 error}) = P(\text{not rejecting } H_0 \mid H_0 \text{ is false}).$$

They can't be minimized together. Thus we fix $P(\text{type 1 error})$ and minimize $P(\text{Type 2 error})$

Minimizing $P(\text{Type 2 error})$

$$\Leftrightarrow \text{Maximizing } 1 - P(\text{Type 2 error}) \\ = P(\text{rejecting } H_0 \mid H_0 \text{ is false})$$

$P(\text{type 1 error})$ is called the level of the test and $1 - P(\text{type 2 error})$ is called the power of the test.

Parametric tests: Let $X_1, \dots, X_n \stackrel{iid}{\sim} f(x|\theta)$

We test $H_0: \theta \in \mathcal{H}_0$ vs. $H_1: \theta \in \mathcal{H}_1$, where

\mathcal{H}_0 and \mathcal{H}_1 are disjoint sets. If \mathcal{H}_0 is a singleton set, then the null hypothesis is called a simple null hypothesis. If \mathcal{H}_0 has multiple elements, then the null hypothesis is called a composite null hypothesis.

Rejection region: Let $R = \{x \in \mathcal{X} \mid H_0 \text{ is rejected for } x\}$

be known as the rejection region or critical region of a test.

Let $\phi(x) = \text{Prob. of rejecting } H_0 \text{ when } x \text{ is observed.}$
The power function of a test is given by

$\beta(\theta) = E_{x|\theta}[\phi(x)]$. Power function is a function of the parameter θ .

Consider the situation $H_0: \theta = \theta_0$ vs. $H_1: \theta = \theta_1$.

$\beta(\theta_0) = E_{x|\theta_0}[\phi(x)] = \text{level of the test.}$

$\beta(\theta_1) = E_{x|\theta_1}[\phi(x)] = \text{power of the test.}$

Under the composite null hypothesis, $H_0: \theta \in \mathcal{H}_0$.
the level of the test α implies

$$\sup_{\theta \in \mathcal{H}_0} \beta(\theta) \leq \alpha.$$

Given a certain level, we want to find the most powerful test.

X	0	1	2	3
f_0	$1/8$	$1/8$	$1/4$	$1/2$
f_1	$1/2$	$1/4$	$1/8$	$1/8$

$$X \sim f_0 \quad \theta = 0, 1$$

How to find the best (most powerful level $1/8$) test).
level $1/8$ means rejection region R must have prob. $1/8$ under H_0 .

$$R_1 = \{0\} \text{ test 1}$$

$$R_2 = \{1\} \text{ test 2}$$

$$\begin{aligned} \text{power of test 1} &= P(\text{rejecting } H_0 \mid H_0 \text{ is false}) \\ &= P(X \in R_1 \mid H_0 \text{ is false}) = 1/2 \end{aligned}$$

$$\begin{aligned} \text{power of test 2} &= P(X \in R_2 \mid H_0 \text{ is false}) = 1/4 \\ \text{test 1 is more powerful than test 2 under} \\ &\text{the same level.} \end{aligned}$$

f_1 is higher \odot at 0 than at 1.

Let us decide to reject H_0 with a rejection region that contains high values of f_1 .

$$\text{power} = P(X \in R \mid f_1 \text{ is true})$$

Theorem: Neyman-Pearson Lemma:

Consider testing $H_0: \theta = \theta_0$ vs. $H_1: \theta = \theta_1$, where pdf or pmf corresponding to θ_i is $f(x|\theta_i)$, $i=0,1$. Using a test with rejection region R that satisfies

$$\phi(x) = \begin{cases} 1 & \text{if } f(x|\theta_1) > k f(x|\theta_0) \Leftrightarrow \frac{f(x|\theta_1)}{f(x|\theta_0)} > k \\ 0 & \text{o.w.} \end{cases}$$

for some $k \geq 0$, and $\alpha = P_{\theta_0}(X \in R)$. Then

- Any test that satisfies the above is the most powerful level α -test.
- If there exists a test satisfying the above, then it is a most powerful level α -test of its level.

Note:

$$\frac{f(x|\theta_1)}{f(x|\theta_0)} \Leftrightarrow \frac{g(T(x)|\theta_1) h(x)}{g(T(x)|\theta_0) h(x)}$$

MP test can be written as

$$\phi(x) = \begin{cases} 1 & \text{if } g(T(x)|\theta_1) > k g(T(x)|\theta_0) \\ 0 & \text{o.w.} \end{cases}$$

Example: $X_1, X_2 \stackrel{iid}{\sim} \text{Ber}(\theta)$. Want to test $H_0: \theta = \frac{1}{2}$
vs. $H_1: \theta = \frac{3}{4}$.

$\sum_{i=1}^2 X_i =$ sufficient stat.

$\sum_{i=1}^2 X_i \sim \text{Bin}(2, \theta)$.

$$\frac{f(0|\theta=3/4)}{f(0|\theta=1/2)} = \frac{1}{4}, \quad \frac{f(1|\theta=3/4)}{f(1|\theta=1/2)} = \frac{3}{4}, \quad \frac{f(2|\theta=3/4)}{f(2|\theta=1/2)} = \frac{9}{4}.$$

if we choose $\frac{3}{4} < k < \frac{9}{4} \Rightarrow R = \{2\}$

MP level α test for $\alpha = P\left(\sum_{i=1}^2 X_i = 2 \mid \theta = \frac{1}{2}\right) = \frac{1}{4}.$