

Sets, Maps & Numbers

Kollerup &amp; Hure:

Theoretical foundation for data scientists  
(HT19)SetsA set is a collection of distinct elements.e.g.  $\{0, 0\}$ . We give names to sets, e.g.  $A = \{0, 0\}$ . $A' = \{0, 0, 0\}$  is not a set (multiset) $\emptyset = \{\}$  is the empty set.An element belongs to (or doesn't belong to) a set and we write ' $\in$ ' or ' $\notin$ ' respectively. e.g.  $0 \in \{0, 0\}$  and  $0 \notin \{0, 1\}$ . etc. etc. det här heter du...Set operationsWe can add elements to an existing set by union operation.

e.g.  $\{0\} \cup \{0, 1\} = \{0, 1\}$

$\{0, 0\} \cup \{0, 1\} = \{0, 0, 1\}$

$\{0, 0\} \cup \{0, 1\} = \{0, 0, 1\}$

Def  $A \cup B := \{x \mid x \in A \text{ or } x \in B\}$

Intersection:  $A \cap B := \{x \mid x \in A \text{ and } x \in B\}$

Set difference:  $A \setminus B := \{x \mid x \in A \text{ and } x \notin B\}$

complement

Given a universal set  $U$ ,  $A^c := \{x \mid x \in U \text{ and } x \notin A\}$

## Maps

A map or a function associates each element in the set called domain with exactly one element in the set range (codomain)

formally

A function is a specific kind of relation between elements in the domain and range

## Inverse Image / pre-image

The inverse image of a function  $f: X \rightarrow Y$  is  $f^{-1}(y) = \{x \in X \mid f(x) = y\}$  or more generally, for any  $B \subseteq Y$

$$f^{-1}(B) = \{x \in X \mid f(x) \in B\} \text{ and } f^{-1} = \bigcup_{A \subseteq Y} f^{-1}(A)$$

↑  
collection of subsets of Y

## Note

In this case  $N = \{1, 2, 3, \dots\}$  and  $Z_0 = \{0, 1, 2, 3, \dots\}$

## Probability

### Language

An experiment is an activity that produces distinct observable outcomes. The set of such outcomes is called the sample space of the experiment, denoted by  $\Omega$ .

An event is a subset of the sample space.

Probability is a function

$P: \{\text{events}\} \rightarrow [0, 1]$  where  $P$  satisfies

1)  $\forall \text{ event } A, 0 \leq P(A) \leq 1$

2)  $P(\Omega) = 1$

3)  $A \cap B = \emptyset \Rightarrow P(A \cup B) = P(A) + P(B)$

4)  $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$  where  $A_i$  are pairwise disjoint

### Motivation for axioms:

Idea of long-term relative frequency of independent experiments (experimental trials). If we repeat an experiment a large number of times, the fraction of times the event  $A$  occurs will be close to  $P(A)$ .

Formally, let  $N(A, n)$  be the number of times  $A$  occurs in the first  $n$  trials.

$$P(A) = \lim_{n \rightarrow \infty} \frac{N(A, n)}{n}$$

axiom i)  $0 \leq \frac{N(A, n)}{n} \leq 1$

ii)  $\frac{N(\Omega, n)}{n} = \frac{n}{n} = 1$  "something" happens every time

iii)  $A \cap B = \emptyset \Rightarrow N(A \cup B, n) = N(A, n) + N(B, n)$

iv) If this is ignored the mathematics you need is much harder

1.1 Tossing a fair coin. Can construct reads by dyadic partitions.

$$\Omega = \{HT\}$$

$$\begin{pmatrix} H \\ T \end{pmatrix}$$

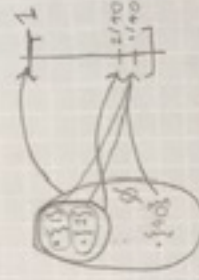
Bernoulli random variable,  $\text{Bernoulli}(\theta)$

$$2^2 \cdot F_2$$

1.2 NZ Lotto (40 balls)

Label of the 1st ball that comes out:

$$\Omega = \{1, 2, \dots, 40\}$$



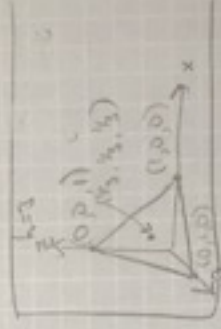
De Moivre random variable,  $2^2 = F_2$

$$\text{De Moivre}(t, \frac{1}{40}, \dots, \frac{1}{40}) = \text{De Moivre}(\theta_1, \theta_2, \dots, \theta_L)$$

Here we have  $\theta_1 = \theta_2 = \dots = \theta_L = \frac{1}{40}$ , an Equi-probable De Moivre variable

$$\begin{aligned} \text{So, what is } P(\text{"even number"}) &= P(\{2, 4, \dots, 40\}) = \\ &= P(\{2\} \cup \{4\} \cup \dots \cup \{40\}) = P(\{2\}) + P(\{4\}) + \dots + P(\{40\}) = \\ &= 20 \cdot \frac{1}{40} = \frac{1}{2} \end{aligned}$$

repeatedly



Properties:

$$1. P(A) = 1 - P(A^c)$$

$$2. A, B \text{ events. } P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

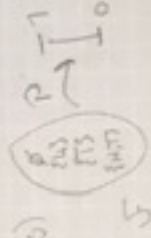
The domain of  $P$  is called a sigma field or sigma algebra, It's denoted by  $\mathcal{F}(\Omega)$  or  $F_\Omega$  or just  $\mathcal{F}$  if  $\Omega$  is clear from context.

We see that

$$i) \Omega \in \mathcal{F}, \quad ii) A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}, \quad iii) A_1, A_2, \dots \in \mathcal{F} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F} \text{ by Kolmogorov.}$$



Bernoulli experiment (coin flipping)



The triple  $(\Omega, \mathcal{F}, P)$  is called the probability space.

If  $P$  is a set of probabilities, then  $(\Omega, \mathcal{F}, P)$  is called a statistical experiment.

Recall

events  $A, B$  are independent  $\Leftrightarrow P(A \cap B) = P(A)P(B)$

The product experiment

$X_1, X_2, \dots, X_n \stackrel{i.i.d.}{\sim} P_0$  (indep., ident. dist.) then  $P_0(\{X_1, X_2, \dots, X_n\}) = \prod_{i=1}^n P_0(X_i)$

Conditional probability

$P(A|B) := \frac{P(A \cap B)}{P(B)}$  provides  $P(B) > 0$

is conditional probability a prob.? Yes, basically  $B$  is new  $\Omega$

Constructing random graph from tossing unfair coins.

Toss coin with  $P(H) = \theta$  iid  $n$  times.

Graph:

Let  $V = \{v_1, \dots, v_n\}$  be a set of vertices and let  $E \subseteq V^2$  and  $|E| = k$

( $n$  vertices,  $k$  edges). Let

adjacency matrix  $A = \begin{bmatrix} v_1 & v_2 & \dots & v_n \\ v_1 & & & \\ v_2 & & & \\ \vdots & & & \\ v_n & & & \end{bmatrix}$  where  $A_{i,j} = \begin{cases} 1, & \text{if } H \\ 0, & \text{if } T \end{cases}$  in coin toss

so,  $E(|E|) = |V|^2 \theta$

### Exercises

6.11

$$f(t) = \begin{cases} \lambda e^{-\lambda t}, & t \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

a)  $\lambda = 0.01$

$$P(\text{"burns at least } \tau \text{ hours"}) = 1 - F(\tau) = 1 - \int_0^{\tau} \lambda e^{-\lambda t} dt = 1 - [-e^{-\lambda t}]_0^{\tau} = e^{-\lambda \tau}$$

$$= 1 + e^{-\lambda \tau} - 1 = e^{-0.01 \tau}$$

b)  $e^{-0.01 \tau} = \frac{1}{2} \Leftrightarrow -0.01 \tau = \ln\left(\frac{1}{2}\right) \Leftrightarrow \tau = \frac{\ln(2)}{-0.01} \approx 69.3 \text{ h}$

6.12

576 squares, 537 hits

$$P(0) = \frac{129}{576}, P(1) = \frac{211}{576}, P(2) = \frac{93}{576}, P(3) = \frac{35}{576}, P(4) = \frac{7}{576}, P(5) = \frac{1}{576}$$

$$\frac{129}{576} P(0) \approx \lambda e^{-\lambda} = \lambda \Rightarrow \lambda \approx \frac{129}{576}$$

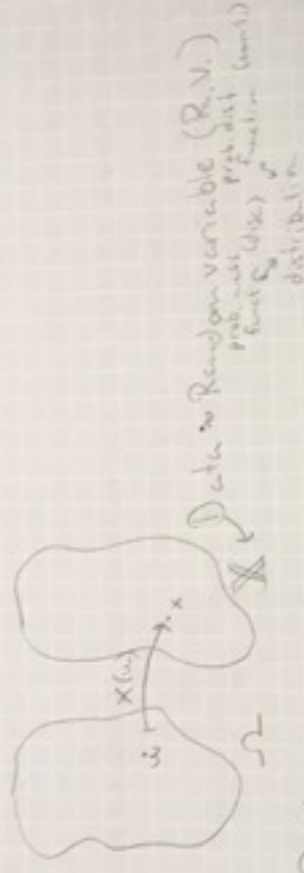
$$\lambda = \frac{537}{576}, \text{ No. of hits} \sim P_0(\lambda)$$

$$\text{approx} = [576 \cdot P(P_0(\lambda) = i) \text{ for } i \in (0, 1, 2, 3, 4, 5)] = [226.7, 211.4, 98.5, 30.6, 7.1, 1.3]$$

Which is very close to actual.

## Prelude to Decision theory

### Estimation in parametric models



### Problem

Let  $X$  be a RV with values in  $\mathbb{X}$  and  $L(\omega, X)$  is known as  $x$ .

Assume  $L(X)$  is known up to a finite dimensional parameter  $\theta$  from the parameter space  $\Theta$  (boldface):

$L(X) \in \{P_\theta \mid \theta \in \Theta\}$ , here we assume  $\Theta \subseteq \mathbb{R}^d$  for  $d \geq 1$ .

The decision problem is to estimate a function  $g(\theta)$  based on a realisation of  $X$ .

Typically (WLOG)

$X = (X_1, \dots, X_n)$ ,  $X_i \stackrel{iid}{\sim} X_1$

$n$  is called sample size, (initially)  $\mathbb{X}$  is countable or  $\subseteq \mathbb{R}^d$ .

### Def

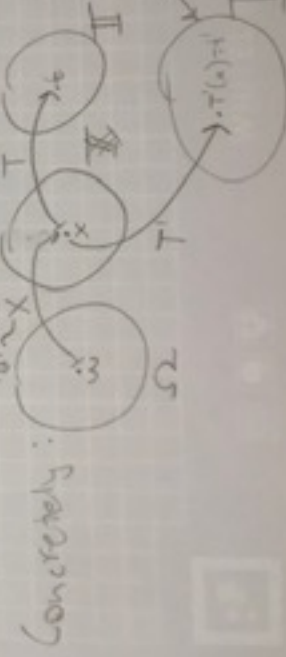
A statistic  $T$  is an arbitrary function of the observed RV  $X$  (data)

$$\{g(\theta) \mid \theta \in \Theta\}$$

### Def

As an estimator  $T$  of  $g(\theta)$  we admit any  $T: \mathbb{X} \rightarrow g(\Theta)$ .

Statistic



### Indicator function

$X(w) \mathbb{I}_{(A)}^{(w)} := \begin{cases} 1 & \text{if } w \in A \\ 0 & \text{if } w \notin A \end{cases}$  is Bernoulli

$\text{Law}(X_i) \sim \text{Bernoulli}(\theta) = \text{Bin}(1, \theta)$

Suppose the data vector  $x = (x_1, \dots, x_n)$  is a realisation of  $X \sim \bigotimes_{i=1}^n \text{Bern}(\theta)$ , i.e.  $x \in \mathbb{R}^n = \{0, 1\}^n$ , for unknown but fixed  $\theta \in \Theta = [0, 1]$

Note that  $P_\theta(X=x) = \prod_{i=1}^n P_\theta(X_i=x_i) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}$

$$\text{Let } T(x) = \sum_{i=1}^n x_i$$

Thus, the prob. of data (i.e.  $P_\theta(X=x)$ ) only depends on the statistic  $T$

Now, consider another statistic: (sample mean)

$$T(X) = \frac{1}{n} \sum_{i=1}^n X_i =: \bar{X}_n$$

Then  $\theta$  becomes  $\theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} = \theta^{n \bar{X}_n} (1-\theta)^{n(1-\bar{X}_n)} = \theta^{n \bar{X}_n} (1-\theta)^{n(1-\bar{X}_n)}$

An estimator of  $\theta$ ,  $g(\theta)$ , (say  $g(\theta) = \theta$ ), should converge towards the "true" but unknown  $\theta$  to be estimated, as the sample size  $n \rightarrow \infty$

Def

A sequence  $T_n := T_n(\underbrace{X_1, X_2, \dots, X_n}_{X^{(n)}})$  of estimators (each based on a sample of size  $n$ ) for a parameter  $\theta$  is called (asymptotically) consistent if  $\forall \epsilon > 0, \theta \in \Theta: P_{\theta,n}(|T_n(X^{(n)}) - \theta| > \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$  or, in shorter notation:

$$T_n = T_n(X^{(n)}) \xrightarrow[n \rightarrow \infty]{P_\theta} \theta \text{ if } \text{Law}(X^{(n)}) = P_\theta$$



ex (Bon. product ex.)

The estimator  $T_n(X_1, \dots, X_n) = \bar{X}_n$  is consistent for  $\theta$

Proof

Since  $X_1, \dots, X_n$  are i.i.d  $\text{Be}(\theta)$  R.V.s, we have  $E(X_i) = \theta$  and the result follows from the law of large numbers.

So, consistency can be seen as a minimal requirement on estimators. But this still leaves a lot of consistent estimators to choose from. A quantitative comparison of estimators is made possible by the approach of statistical decision theory.

We choose a loss function  $\text{Loss}(t, \theta)$  which measures the loss (inaccuracy) with the unknown parameter  $\theta$  is estimated by  $t$ . <sup>this estimate is a realization of the estimator  $T(X^{(n)})$</sup>

Natural choices for loss, when  $\theta \in \mathbb{R}$   $\otimes \mathbb{R}$

Absolute error:  $\text{Loss}(t, \theta) = |t - \theta|$

Quadratic error:  $\text{Loss}(t, \theta) = (t - \theta)^2$

or  $\text{Loss}(t, \theta) = 1_{\{t \neq \theta\}}$  For some  $\delta > 0$  to emphasise the distance being less than  $\delta$

Note

Loss is a R.V.,  $\therefore \text{Loss}(T(X), \theta)$  needs to account for randomness.

Def

The Risk of an estimator  $T$  at parameter  $\theta$  is

$$R(t, \theta) := E_{\theta}(\text{Loss}(T(X), \theta))$$

<sup>↑</sup>  
Risk function of  $T$

Note Risk might exist

since the expectation might not exist

Idea

class of estimators

$$R(T^*, \theta) = \min_{T \in \mathcal{T}} R(T, \theta) \text{ for any fixed } \theta \in \Theta$$

We find an estimator  $T^*$  that minimizes the whole risk function simultaneously. If such a  $T^*$  can be found, then it is called a uniformly best estimator (in  $\mathcal{T}$ ).

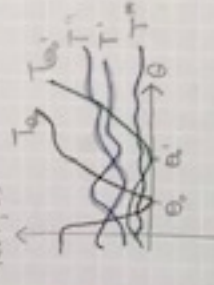
In general, such a UBE won't exist

Argument

for each  $\theta_0 \in \Theta$ , consider  $T_{\theta_0}(x) = \theta_0$

so  $R(T_{\theta_0}, \theta_0) = 0 \Rightarrow$  best if  $\theta_0$  is true

If  $T^*$  was UBE, then it would have to compete with  $T_{\theta_0}$



Unbiased estimators

Def

Consider an estimator  $T$  s.t.  $E_{\theta} T$  exists.

Then  $E_{\theta}(T) - g(\theta)$  is called bias of the estimator.

If  $E_{\theta}(T) = g(\theta)$  for any  $\theta \in \Theta$ , then the estimator  $T$  is called unbiased for  $g(\theta)$

Def

for all

A statistic  $S$  is called sufficient for  $\theta$  if  $P_{\theta}(X \in B | S(X) = s)$  is independent of  $\theta$ , for all values of  $s$  and all events  $B$ .



In words, conditional distribution of the data  $X$  given that  $S(X)$  takes any value does not depend on the parameter  $\theta$

For discrete experiments

$$P_{\theta}(X \in B | S(X) = s) = \begin{cases} \frac{P_{\theta}(X \in B \cap \{S(X) = s\})}{P_{\theta}(S(X) = s)} & , \text{ if } P_{\theta}(S(X) = s) > 0 \\ 0 & , \text{ if } P_{\theta}(S(X) = s) = 0 \end{cases}$$

Let's clarify again:

Originally, the law of  $X$  depends on  $\theta$  ( $L(x) = P_\theta(\cdot)$ ).

After the value of the sufficient statistic  $S(X)$  is known, then the conditional law  $P_\theta(\cdot | S(x))$  is no longer dependent on  $\theta$ .

Since we are interested in making inference about  $\theta$ , the conditional law is uninteresting for our purpose, so we can disregard it.

After taking  $S(x)$  into account, the remaining randomness does not depend on  $\theta$  anymore.

Remark (Exercise to prove formally)

The data itself is sufficient, i.e.

if  $S(X) = X$ , then for some  $\theta$ ,  $P_\theta(X \in B | X=x) = \mathbb{1}_B(x) = \begin{cases} 1, & x \in B \\ 0, & x \notin B \end{cases}$

Prop.

In  $\otimes \text{Ber}(\theta)$  exp., the sample mean  $\bar{X}_n$  is a sufficient statistic.

Proof

$$n\bar{X}_n = \sum_{i=1}^n X_i \sim \text{Bin}(n, \theta)$$

Suppose  $\bar{X}_n = t$  where  $t$  is one of the possible values. This means  $n\bar{X}_n = k$  for some  $k \in \{0, 1, \dots, n\}$ . Then, for any  $x = (x_1, \dots, x_n)$ ,

$$P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | n\bar{X}_n = k) = \frac{P_\theta(X_1 = x_1, \dots, X_n = x_n, n\bar{X}_n = k)}{P_\theta(n\bar{X}_n = k)} \quad (*)$$

If  $\theta$  R.I.F.I.s  $\sum_{i=1}^n x_i = k$ , then  $\theta$  is  $\frac{(\theta^k (1-\theta)^{n-k})}{\binom{n}{k} \theta^k (1-\theta)^{n-k}} = \frac{1}{\binom{n}{k}}$  which clearly is indep. of  $\theta$ .

If for this  $\sum_{i=1}^n x_i$ , then the numerator is 0 which is indep. of  $\theta$ , so  $\theta$  is indep. of  $\theta$ .  $\square$

Say we want to estimate  $\theta$  in this example and we limit ourselves to estimators that are functions of the sufficient statistic  $\bar{X}_n$ :

$$T(x) = h(\bar{X}_n)$$

Additionally, suppose we limit ourselves further to unbiased estimators:

$$E(T(X)) - \theta = 0 \Rightarrow E(T(X)) = \theta$$

Prop.

Under sufficiency and unbiasedness restrictions on allowed estimators of  $\theta$  in our  $\text{Bel}(\theta)$  exp., the only possible estimator is  $\bar{X}_n$ .

Proof  $E(\bar{X}_n) = \theta$

$$0 = E_{\theta} h(\bar{X}_n) - \theta \stackrel{!}{=} E_{\theta} \left( h(\bar{X}_n) - \frac{h}{n} \right) = \sum_{k=0}^n \underbrace{\left( h\left(\frac{k}{n}\right) - \frac{h}{n} \right) \binom{n}{k} \theta^k (1-\theta)^{n-k}}_{C_k} =$$

$$= (1-\theta)^n \sum_{k=0}^n C_k \frac{1}{n^k} \quad \text{for } r = \frac{\theta}{1-\theta}$$

Now, if  $\theta \in [0, 1)$ , then  $r \in [0, \infty)$ , hence the above polynomial can only be zero if every  $C_k$  is also zero.

Thus,  $0 = \binom{n}{k} \left( h\left(\frac{k}{n}\right) - \frac{h}{n} \right) \stackrel{!}{=} h\left(\frac{k}{n}\right) - \frac{h}{n}$  for  $k \in \{1, \dots, n\}$  and  $h(\bar{X}_n) = \bar{X}_n$  for all possible values of  $\bar{X}_n$ .  $\square$



### Def

A statistic  $T$  is called complete if, for all  $h: \mathbb{R} \rightarrow \mathbb{R}$   
 $[E_{\theta}(h(T(X))) = 0 \text{ for all } \theta \in \Theta] \Rightarrow$

$$\Rightarrow P(\underbrace{h(T(X)) = 0}_{\text{almost sure event}}) = 1 \text{ for all } \theta \in \Theta$$

Intuitively, completeness means there is "no superfluous information" or "no redundancy" in the complete statistic  $T$

Consider an event  $T(X) \in B$  and suppose  $P_{\theta}(T(X) \in B) = \alpha$  for  $\alpha$  indep. of  $\theta$ .  
By taking  $h(t) = 1_B(t) - \alpha$ , completeness  $\Rightarrow P_{\theta}(1_B(T(X)) - \alpha) = 0$  for all  $\theta \in \Theta$   
which means that  $\alpha$  is either 0 or 1. Thus, for any event  $T(X) \in B$   
which has a non-trivial probability ( $\neq 0, 1$ ), this prob. must depend on  $\theta$ .

$$\text{ex: } \bigotimes_{i=1}^n B_{\theta_i}(\theta)$$

The statistic  $T(X) = (X_1, \bar{X}_n)$  is sufficient but not complete:

$$h(X_1, \bar{X}_n) = X_1 - \bar{X}_n \text{ has } E(h) = 0 \text{ but } h(X_1, \bar{X}_n) \text{ is not almost surely } 0.$$

### Prop

In  $\bigotimes_{i=1}^n B_{\theta_i}(\theta)$ ,  $T(X) = \bar{X}_n$  is sufficient and complete.

### Proof

Suppose for some function  $h$ ,

$$E_{\theta}(h(\bar{X}_n)) = 0$$

This means that  $\sum_{k=0}^n h(\frac{k}{n}) \binom{n}{k} \theta^k (1-\theta)^{n-k} = 0 \quad \forall \theta \in [0, 1]$ .

Then,  $h(\frac{k}{n}) = 0$  for  $k \in \{0, 1, \dots, n\}$  as per the earlier argument used to show that  $\bar{X}_n$  is the only unbiased and sufficient estimator.

## Limits of R.V.s : Chapter 8 in CSE Book pdf 15.1-15.6

8.1 - conv. of sequence of R.V.

$X_1, \dots, X_n$  conv. to a R.V.  $X$ :

- In distribution

$X_n \rightarrow X$

- In probability

$X_n \xrightarrow{P} X$

- Markov's ineq., Chebyshev's ineq.

③ Suppose  $T$  is a sufficient statistic with values in a set  $\mathbb{T}$  and  $S: \mathbb{T} \rightarrow \mathbb{S}$  is a one-to-one mapping with values in  $\mathbb{S}$  (i.e., there exists an inverse mapping  $S^{-1}$  such that  $S^{-1}(S(t)) = t$  for each  $t \in \mathbb{T}$ ). Show that the statistic  $S(T(X))$  is sufficient.

④ In class it was claimed that the statistic  $T(X) = (X, \bar{X}_n)$  is sufficient for the  $\bigotimes_{i=1}^n \text{Bernoulli}(\theta)$  experiment. Prove this claim.

### Problem Set Week 3

- ① Show that the estimator  $\bar{X}_n$  is consistent for  $\theta = \bigotimes_{i=1}^n \text{Bernoulli}(b)$  experiment.
- ② Suppose the data  $X$  in a statistical experiment can take values in a countable set  $\bar{X}$  (i.e.,  $\text{law}(X)$  is discrete). In class it was claimed that the data itself are a sufficient statistic (i.e.,  $T(X) = X$  is sufficient). Write down the argument that proves this claim (it can be a short paragraph).

⑤ Let  $X_1, \dots, X_n$  be independent and identically distributed with  $\text{law}(X_i) \sim \text{Poisson}(\lambda)$ ,  $\lambda \in (0, \infty)$  is unknown. Show that the sample mean  $\bar{X}_n$  is a sufficient statistic.



## 8.1 - conv. of sequence of R.V.

 $X_1, \dots, X_n$  conv. to a R.V.  $X$ :

• In distribution

 $X_n \rightsquigarrow X \Leftrightarrow \forall t \text{ where } F_X \text{ cont.} \lim_{n \rightarrow \infty} F_{X_n}(t) = F_X(t)$  (pointwise convergence of dist. function)

• In probability

 $X_n \xrightarrow{P} X \Leftrightarrow \forall \varepsilon > 0: \lim_{n \rightarrow \infty} P(|X_n - X| > \varepsilon) = 0$  (uniform conv. of dist. function)

- Markov's ineq., Chebyshev's ineq.

$$\forall \varepsilon > 0: P(|X| > \varepsilon) \leq \frac{E(X)}{\varepsilon}, \quad P(|X| > \varepsilon) \leq \frac{E(X^2)}{\varepsilon^2}$$

$$P(|X| > \varepsilon) = P(X^2 > \varepsilon^2) \leq \frac{E(X^2)}{\varepsilon^2}$$

Problem set week 3  $P(|X - E(X)| > \varepsilon) \leq \frac{V(X)}{\varepsilon^2}$ 

- ① Show that the estimator  $\bar{X}_n$  is consistent for  $\theta$  in  $\mathcal{B}(\theta)$  experiment
- ② Suppose the data  $X$  in a stat. experiment can take values in a countable set  $\mathcal{X}$  (i.e.  $\text{Law}(X)$  is discrete).

In class it was claimed that the data itself is a sufficient statistic (i.e.  $T(X) = X$  is sufficient). Write down an argument for that claim that proves it.

- ③ Suppose  $T$  is a sufficient statistic with values in a set  $\mathbb{T}$  and  $S: \mathbb{T} \rightarrow \mathbb{S}$  is one-to-one (bijective)  $\exists S^{-1}: \mathbb{S} \rightarrow \mathbb{T}$  s.t.  $S^{-1}(S(t)) = t \quad \forall t \in \mathbb{T}$ . Show that the statistic  $S(T(X))$  is sufficient.

- ④ In class it was claimed that the statistic  $T(X) = (X_1, \bar{X}_n)$  is sufficient for the  $\mathcal{B}(\theta)$  experiment. Prove this claim.

- ⑤ Let  $X_1, \dots, X_n$  be independent and identically distributed (iid) with  $\text{Law}(X_i) \sim P_\theta(X)$ ,  $\lambda \in (0, \infty)$  is unknown. Show that the sample mean  $\bar{X}_n$  is a sufficient statistic for this data.



### Weak law of large numbers

$X_1, X_2, \dots \stackrel{\text{i.i.d.}}{\sim} X_1, E(X_1)$  exists. Then,  $\bar{X}_n \xrightarrow{P} E(X_1)$

Proof (In the case where  $V(X_1) < \infty$ )

$$\text{Let } \varepsilon > 0. P(|\bar{X}_n - E(\bar{X}_n)| > \varepsilon) = \frac{V(\bar{X}_n)}{\varepsilon^2} = \frac{1}{n} \frac{V(X_1)}{\varepsilon^2}.$$

Chob. ineq.  $X_1, X_2, \dots \stackrel{\text{i.i.d.}}{\sim} X_1$

We also know that  $E(\bar{X}_n) = E(X_1)$  since  $X_1, X_2, \dots \stackrel{\text{i.i.d.}}{\sim} X_1$ , so  $E(X_1) = E(X_2) = \dots$

$$\text{so } E(\bar{X}_n) = E\left(\frac{\sum_{i=1}^n X_i}{n}\right) = \frac{\sum_{i=1}^n E(X_i)}{n} = \frac{n E(X_1)}{n} = E(X_1).$$

$$\text{So } P(|\bar{X}_n - E(X_1)| > \varepsilon) = P(|\bar{X}_n - E(\bar{X}_n)| > \varepsilon) = \frac{1}{n} \frac{V(X_1)}{\varepsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Thus  $\bar{X}_n \rightarrow E(X_1)$   $\square$

### Central limit theorem (CLT)

If  $X_1, X_2, \dots \stackrel{\text{i.i.d.}}{\sim} X_1$  and  $E(X_1), V(X_1)$  exist, then

$$\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n} \rightsquigarrow X \sim \text{Normal}\left(E(X_1), \frac{V(X_1)}{n}\right).$$

Cor.

$$\text{i) } \bar{X}_n - E(X_1) \rightsquigarrow X - E(X_1) \sim \text{Normal}\left(0, \frac{V(X_1)}{n}\right)$$

$$\text{ii) } \sqrt{n}(\bar{X}_n - E(X_1)) \rightsquigarrow \sqrt{n}(X - E(X_1)) \sim \text{Normal}(0, V(X_1))$$

$$\text{iii) } Z_n = \frac{\bar{X}_n - E(\bar{X}_n)}{\frac{V(\bar{X}_n)}{\sqrt{n}}} = \frac{\sqrt{n}(\bar{X}_n - E(X_1))}{\sqrt{V(X_1)}} \rightsquigarrow Z \sim \text{Normal}(0, 1)$$

$$\text{iv) } \lim_{n \rightarrow \infty} P\left(\frac{\bar{X}_n - E(\bar{X}_n)}{\frac{V(\bar{X}_n)}{\sqrt{n}}} \leq z\right) = \lim_{n \rightarrow \infty} P(z_n \leq z) = P(Z \leq z) = \phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{x^2}{2}} dx$$

### Approximation

For "large"  $n$  ( $\approx n > 30$ ), we can use

$$P\left(\frac{\bar{X}_n - E(\bar{X}_n)}{V(\bar{X}_n)} \leq z\right) \approx P(Z \leq z) = \phi(z)$$



## Exercises

1. Prove that the statistic  $T(X_1, \dots, X_n) = \frac{\sum_{i=1}^n X_i}{n} = \bar{X}_n$  is consistent for  $\theta$  in  $\mathcal{B}$  Bernoulli( $\theta$ ) exp.

Sol.

We have to show that  $\bar{X}_n$  is consistent, i.e.  $\bar{X}_n \xrightarrow{P} \theta$ .

Let  $\varepsilon > 0$ .  $S_n = \sum_{i=1}^n X_i$  so  $\bar{X}_n = \frac{S_n}{n}$ .  $P(|\bar{X}_n - \theta| > \varepsilon) = P(|S_n - n\theta| > n\varepsilon)$

$= P(|S_n - n\theta| > n\varepsilon)$ . Since  $V(S_n) = n\theta(1-\theta)$ , by Chebyshev:

$$P(|\bar{X}_n - \theta| > \varepsilon) = P(|S_n - n\theta| > n\varepsilon) \leq \frac{V(S_n)}{\varepsilon^2 n} = \frac{\theta(1-\theta)}{\varepsilon^2 n} \leq \frac{\frac{1}{4}}{\varepsilon^2 n} = \frac{1}{4\varepsilon^2 n} \xrightarrow{n \rightarrow \infty} 0 \quad \square$$

Proof of  
Chebyshev

$$P(|X - E(X)| > a) = P((X - E(X))^2 > a^2) \leq \frac{1}{a^2} E((X - E(X))^2) = \frac{V(X)}{a^2}$$

Markov

2. Suppose  $X$  takes values in  $\mathcal{X}$  where  $|\mathcal{X}| \leq |\mathcal{N}|$ .

Let  $S(X) = s$  if  $X = s$ . We want to show that

$P_\theta(X \in B | S(X) = s)$  is independent of  $\theta$ .  $B \subseteq \mathcal{X}$

$$P_\theta(X \in B | S(X) = s) = P_\theta(X \in B | X = s) = \begin{cases} \frac{P_\theta(\{X \in B\} \cap \{X = s\})}{P_\theta(X = s)} & , P_\theta(X = s) > 0 \\ 0 & , P_\theta(X = s) = 0 \end{cases}$$

$$P_\theta(\{X \in B\} \cap \{X = s\}) = \begin{cases} P_\theta(X = s) & , s \in B \\ 0 & , s \notin B \end{cases}, \text{ so}$$

$$P_\theta(X \in B | S(X) = s) = \begin{cases} \frac{P_\theta(X = s)}{P_\theta(X = s)} & , P_\theta(X = s) > 0 \\ 0 & , \text{otherwise} \end{cases} = \begin{cases} 1 & \text{if } s \in B \\ 0 & \text{if } s \notin B \end{cases}$$

Alt.

Assume  $T(X) = X$ , let  $t = x$ . For an event  $B \in \mathcal{F}$

$$P_\theta(X \in B | T(X) = t) = P_\theta(X \in B | X = x) = \mathbb{1}_B(x) = \begin{cases} 1 & , x \in B \\ 0 & , x \notin B \end{cases}$$



3. Suppose  $T$  sufficient  $T: \mathcal{X} \rightarrow \mathcal{T}$  and  $S: \mathcal{T} \rightarrow \mathcal{S}$  inj.  
Show that  $S$  is sufficient.

Sol.

$\forall s \in \mathcal{S} := \{s: S(t)=s\}$  and  $\forall B \in \mathcal{F}_{\mathcal{X}}$  the cond. prob.

$$\frac{P_{\theta}(\{X \in B\} \cap \{S(T(X))=s\})}{P_{\theta}(S(T(X))=s)} = \frac{P_{\theta}(\{X \in B\} \cap \{T(X)=t\})}{P_{\theta}(T(X)=t)} \quad \text{where } t \in S^{-1}(s)$$

$= P_{\theta}(X \in B | T(X)=t)$  which is indep. of  $\theta$  because of  $T$  suff.  $\square$

Prop.

In the  $\bigotimes_{i=1}^n \text{Be}(\theta)$  experiment, the estimator  $\bar{X}_n$  is uniformly best among unbiased estimators for quadratic loss (i.e., for any unbiased estimator  $T: \mathcal{R}(\bar{X}_n, \theta) = E_{\theta}((\bar{X}_n - \theta)^2) \leq E_{\theta}((T(X) - \theta)^2) = R(T, \theta)$  for all  $\theta \in [0, 1]$ ).

Before we prove this, let's revisit the notion of cond. expectation:

Let  $Y$  be a R.V. with finite support in  $\mathcal{Y}$ , i.e.  $|\mathcal{Y}| < \infty$ .

Let  $U$  be a statistic with values in  $\mathcal{U}$  and  $w \in \mathcal{U}$ ,  $h$  a real-valued fct. of  $Y$ .

Def.

The cond. expectation of  $h(Y)$  given  $U=w$ , written

$E(h(Y) | U(Y)=w)$  is defined as the expectation of  $h(Y)$  under the cond. distribution of  $Y$  given  $U(Y)=w$ :  $P(Y=y | U(Y)=w) = \frac{P(Y=y, U(Y)=w)}{P(U(Y)=w)}$  if  $P(U(Y)=w) \neq 0$  (0 otherwise).

$$E(h(Y) | U(Y)=w) = \sum_{y \in \mathcal{Y}} h(y) P(Y=y | U(Y)=w).$$

In special case  $h(y) = 1_B(y)$  for some  $B \subseteq \mathcal{Y}$ , then

$$E(1_B(y) | U(Y)=w) = \sum_{y \in B} P(Y=y | U(Y)=w) = P(y \in B | U(Y)=w).$$



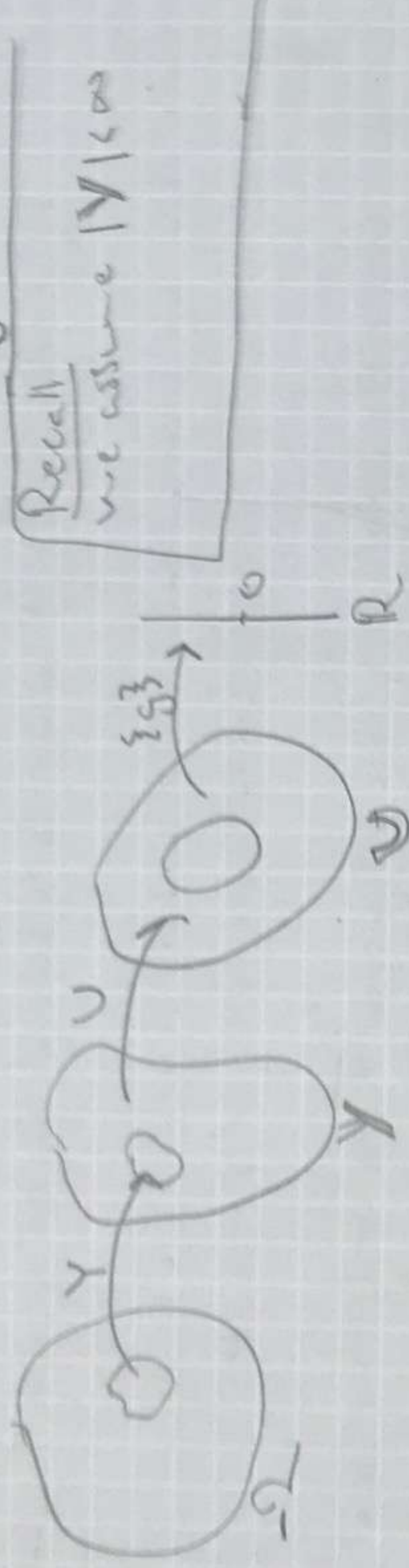
Def

The Conditional expectation can be a R.V.

$$E(h(Y)|U)$$

Note the following properties of cond. exp. as a R.V.;

Let  $M_U$  be the set of all real-valued R.V.s  $Z$  that are functions of  $U$ , i.e.  $Z = g(U)$  for some function  $g: \mathcal{U} \rightarrow \mathbb{R}$ .



$$i) E(E(h(Y)|U)) = E(h(Y))$$

proof

Exercise 0

ii) For any  $Z \in M_U$ ,  $E(Zh(Y)|U) = Z E(h(Y)|U)$  "almost surely" /

Intuition: R.V.s  $X, Y$  equal w.p.1 means  $1 = P(X=Y)$  / "with prob. 1" w.p.1

iii) For any function  $h: \mathcal{Y} \rightarrow \mathbb{R}$  s.t.  $h = h_1 + h_2$ , then  $E(h(Y)|U) = E(h_1(Y)|U) + E(h_2(Y)|U)$  w.p.1

iv) For any  $Z \in M_U$ ,  $E(Z \cdot h(Y)) = E(Z \cdot E(h(Y)|U))$

proof

Exercise 1

⊗

v) For any  $Z \in M_U$ ,  $E((h(Y) - Z)^2) = E((E(h(Y)|U) - Z)^2) + E(h(Y) - E(h(Y)|U))^2$

proof

Exercise 2 (assuming iii) and iv))

hint: add and subtract  $E(h(Y)|U)$  in ⊗ and

$$E(E(h(Y)|U)|U) = E(h(Y)|U) E(1|U) = E(h(Y)|U)$$



### Remark

The conditional expectation can be seen as an operator.

Let  $M_Y$  be all  $\mathbb{R}$ -valued R.V.s which are functions of  $Y$ .

Then, since  $U$  is  $U(Y)$  we have  $M_U \subseteq M_Y$

$\Omega \xrightarrow{Y} \mathcal{Y} \xrightarrow{U} \mathcal{U}$   
 $M_Y = \{g(Y)\} \quad \downarrow \{h(U)\} \subseteq M_U$   
 $\mathbb{R} \quad \mathbb{R}$

and both  $M_U$  and  $M_Y$  are vector spaces.

Let  $H \in M_Y$ , then the cond. expect.  $E(H|U) \in M_U$  and we define the operator  $\Pi(H) := E(H|U)$ , i.e.  $\Pi: M_Y \rightarrow M_U$

The " $\leq$ " in prop. v) tells us that  $E(H - \Pi(H))^2 = \min_{Z \in M_U} E(H - Z)^2$   
 $E[(H(Y) - E(H(Y)|U))^2] = E[(h(Y) - z)^2]$ .

In English: the cond. expect. of  $H$  is the element of  $M_U$  which is closest to  $H$ . You can see this as the "projection"  $\Pi$  of  $H$  onto the space  $M_U$ .

The "=" in prop. v) is the orthogonal decomposition

$$E((H - Z)^2) = E((\Pi(H) - Z)^2) + E((H - \Pi(H))^2)$$

We are now ready to prove that  $\bar{X}_n$  is uniformly best among all unbiased estimators of  $\theta$  in  $\hat{\Theta}_{Be}(\theta)$  exp.:



## Proof

N.T.S. (need to show) that the risk of  $\bar{X}_n$  is the lowest.

i.e.  $R(\bar{X}_n, \theta) := E_{\theta}[(\bar{X}_n - \theta)^2] \leq E_{\theta}[(T(X_1, \dots, X_n) - \theta)^2] =: R(T, \theta)$   
for all  $\theta \in \Omega$ , and all unbiased estimators  $T$ .

Define a R.V.:  $g(\bar{X}_n) = E_{\theta}(T | \bar{X}_n) \stackrel{\text{K}}{=} E(T | \bar{X}_n)$  and regard it as  
an estimator of  $\theta$

by prop. (i),  $E(g(\bar{X}_n)) = E(E(T | \bar{X}_n)) = E(T) = \theta$   
so,  $g$  is unbiased.

Before, we showed that  $g(\bar{X}_n) = \bar{X}_n$  is the only unbiased  
estimator of  $\theta$  ( $\bar{X}_n$  is complete)

Hence,  $g(\bar{X}_n) = \bar{X}_n$

In prop. (v), set  $Z = \theta$ , (point-mass R.V.  $(P(Z=z) = \begin{cases} 1, & z=0 \\ 0, & z \neq 0 \end{cases})$ , so  $Z \in \mathcal{W}$ ),  
s.t.  $h(Y) = T$ ,  $U = \bar{X}_n$ ,  $E(h(Y) | U) = g(\bar{X}_n)$

to get  $E_{\theta}(T - \theta)^2 \geq E_{\theta}(g(\bar{X}_n) - \theta)^2$

so, the estimator  $g(\bar{X}_n) = \bar{X}_n$  is at least as good as any  $T$ .  $\square$

For unbiased estimators, the quadratic risk is also the estimator's

variance:

$$E_{\theta}(T - \theta)^2 = V(T) := E_{\theta}(T - E_{\theta}(T))^2$$

$\therefore \bar{X}_n$  for  $\bigotimes_{i=1}^n \text{Be}(\theta)$  exp is also called a

uniform minimum variance unbiased estimator (UMVUE)



## Bayes estimators

In the Bayesian approach to estimation of the param.  $\theta$  underpinning the law of data  $X = (X_1, \dots, X_n)$  in an exp., one assumes that prior distribution is given over the parameter space  $\Theta \ni \theta$

ex In the  $\bigotimes_{i=1}^n \text{Be}(\theta)$  exp., assume that the prior distribution is given in the form of density on  $\Theta = [0, 1]$ .

## Def (Integrated risk)

For an estimator  $T$  of  $\theta$  the prior  $g(\theta)$  can be used to reduce the risk function  $R(T, \theta)$  for each  $\theta \in [0, 1]$  to a single number,

by integration:  $B(R(T)) = B_r(T) := \int_0^1 R(T, \theta) g(\theta) d\theta = \int_0^1 E_\theta(T - \theta)^2 g(\theta) d\theta$

Suppose we have quadratic loss of  $T$  under risk  $R$ .

$B_r(T)$  is called integrated risk or mixed risk

## Def

A Bayes estimator  $T_B$  of  $\theta$  is the estimator that minimizes the

integrated risk, i.e.  $T_B := \arg \min B_r(T)$

and  $B_r(T_B)$ , the minimal integrated risk, is called Bayes Risk

$$\arg \min_T B_r(T) := T \text{ s.t. } B_r(T) = \min B_r(T)$$

The name "Bayesian" comes from Bayes formula:

$\{D_1, \dots, D_k\}$  partition  $\Omega$ , then  $P(A) = \sum_{i=1}^k P(A|D_i)P(D_i)$

Motivation: Consider the case when  $P$  is the joint distribution of

$(X, U)$  where  $X$  is data and  $U$  is a R.V. that takes  $k$

possible values in  $\Theta := \{\theta_1, \dots, \theta_k\}$ . Then  $A = "X \in A"$  and

observed data

$$D_i = "U = \theta_i"$$

$$\text{Posterior prob. of } \theta_i \text{ given data} \rightarrow P(U = \theta_i | X \in A) = \frac{P(X \in A | U = \theta_i) P(U = \theta_i)}{\sum_{j=1}^k P(X \in A | U = \theta_j) P(U = \theta_j)} \leftarrow \begin{matrix} \text{prior prob.} \\ \text{of param } \theta_i \end{matrix}$$



The Bayesian approach views  $\{P_\theta: \theta \in \Theta\}$  as a family of conditional distributions given  $\theta$ , exhibits a prior dist.

~~ex~~ In  $\tilde{\otimes}_{i=1}^n \text{Be}(\theta)$  exp., consider the family of prior densities for  $\theta \in [0, 1]$ :

for each  $(\alpha, \beta) \in (0, \infty)^2 =: \mathbb{R}_{>0}^2$  ( $\alpha > 0, \beta > 0$ )

$$g_{\alpha, \beta}(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}, \quad \theta \in [0, 1] \quad \text{where } B \text{ is beta function:}$$

$$B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha + \beta)}, \quad \text{and } \Gamma(\alpha) := \int_0^\infty x^{\alpha-1} e^{-x} dx.$$

Consider the whole family  $g_{\alpha, \beta}$  for specification of prior density.