

# Elements of Structural Modelling

## Bayesian Model

Cursillo de Verano 2020, Facultad de Matemáticas, UC

Ernesto San Martín

Laboratorio Interdisciplinario de Estadística Social LIES, UC  
The Economics School of Louvain, Université Catholique de Louvain, Belgium

January 2021



Laboratorio  
Interdisciplinario de  
Estadística Social

- We consider a statistical model:

$$\{p(x | a) : a \in A\}$$

- We consider that the density (distribution)  $p(x | a)$  is a function of  $a \in A$ : for each  $x$ ,

$$a \mapsto p(x | a)$$

- The statistical model describes the way on which the data are in fact, or supposed to be, generated by the family of distributions  $\{p(x | a) : a \in A\}$ .

- We consider a statistical model:

$$\{p(x | a) : a \in A\}$$

- We consider that the density (distribution)  $p(x | a)$  is a function of  $a \in A$ : for each  $x$ ,

$$a \mapsto p(x | a)$$

- The statistical model describes the way on which the data are in fact, or supposed to be, generated by the family of distributions  $\{p(x | a) : a \in A\}$ .
- Are we agree on that?

- We define a probability distribution on the parameter space  $A$ , denoted as  $\mu(a)$ .
- Therefore, we have  $\{p(x | a) : a \in A\}$  and  $\mu$  defined on  $(A, \mathcal{A})$ , where  $\mathcal{A}$  is the  $\sigma$ -algebra of subsets of the parameter space  $A$ .
- These two ingredients allow to define/construct a unique probability distribution on the product space “observations  $\times$  parameters” defined as

$$p(x, a) = p(x | a)\mu(a).$$

- $\mu$  is called *prior distribution*.

- Therefore, a Bayesian model is a unique probability distribution defined on the product space “observations  $\times$  parameters”.

- Therefore, a Bayesian model is a unique probability distribution defined on the product space “observations  $\times$  parameters”.
- The unique joint distribution  $p(x, a)$  can also be decomposed as

$$p(x, a) = p(x)\mu(a | x),$$

where  $\mu(a | x)$  is called *posterior distribution* and  $p(x)$  is called *predictive distribution*.

- Thus,

$$\begin{aligned} p(x, a) &= p(x | a)\mu(a) \\ &= \mu(a | x)p(x). \end{aligned}$$

- Typically, the inference is based on the posterior distribution  $\mu(x | a)$ .

- Let us consider the following example developed by Poirier (1998, Sectio 3.2) and discussed in San Martín et al. (2013): let  $(Y | \psi, \lambda) \sim \mathcal{N}(\psi, \sigma_1^2)$  be the likelihood or sampling distribution, and let

$$\begin{pmatrix} \psi \\ \lambda \end{pmatrix} \sim \mathcal{N}_2 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_2^2 + 1 & 1 \\ 1 & 1 \end{pmatrix} \right)$$

be the prior specification. For simplicity, it is assumed that  $\sigma_1^2$  and  $\sigma_2^2$  are known constants.

- It is clear that  $\psi$  is identified because the mapping  $\psi \mapsto \mathcal{N}(\psi, \sigma_1^2)$  is injective.
- Similarly,  $\lambda$  is unidentified and therefore it can not be said something about it.

# Example 1

- It is nevertheless possible to compute the posterior distribution of  $\lambda$ , which is given by

$$(\lambda | Y) \sim \mathcal{N}\left(\frac{Y}{\sigma_1^2 + \sigma_2^2 + 1}, \frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \sigma_2^2 + 1}\right).$$

This distribution provides information about  $\lambda$ .



# Example 1

- It is nevertheless possible to compute the posterior distribution of  $\lambda$ , which is given by

$$(\lambda | Y) \sim \mathcal{N}\left(\frac{Y}{\sigma_1^2 + \sigma_2^2 + 1}, \frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \sigma_2^2 + 1}\right).$$

This distribution provides information about  $\lambda$ .

- Has the situation actually been saved?

- It is nevertheless possible to compute the posterior distribution of  $\lambda$ , which is given by

$$(\lambda | Y) \sim \mathcal{N}\left(\frac{Y}{\sigma_1^2 + \sigma_2^2 + 1}, \frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \sigma_2^2 + 1}\right).$$

This distribution provides information about  $\lambda$ .

- Has the situation actually been saved?
- This question is meaningful because the sampling process is fully characterized by  $\psi$  in the sense that

$$Y \perp\!\!\!\perp \lambda | \psi.$$

So, what are we updating about the sampling process?

- It is nevertheless possible to compute the posterior distribution of  $\lambda$ , which is given by

$$(\lambda | Y) \sim \mathcal{N}\left(\frac{Y}{\sigma_1^2 + \sigma_2^2 + 1}, \frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \sigma_2^2 + 1}\right).$$

This distribution provides information about  $\lambda$ .

- Has the situation actually been saved?
- This question is meaningful because the sampling process is fully characterized by  $\psi$  in the sense that

$$Y \perp\!\!\!\perp \lambda | \psi.$$

So, what are we updating about the sampling process?

- Where is the trap?

# Example 1

- The posterior distribution of  $\lambda$  is fully characterized by the two first moment because it is a normal distribution.

- The posterior distribution of  $\lambda$  is fully characterized by the two first moment because it is a normal distribution.
- Well, the two first posterior moments are a function of the identified parameter  $\psi$ :

$$E[\lambda | Y] = \frac{1}{\sigma_2^2 + 1} E[\psi | Y], \quad V(\lambda | Y) = \frac{1}{\sigma_2^2 + 1} \left[ \frac{1}{\sigma_2^2 + 1} V(\psi | Y) + \sigma_2^2 \right].$$

# Example 1

- But “we are Bayesian”, so  $\psi$  and  $\lambda$  are endowed of a probability distribution and then we can marginalize the sampling distribution w.r.t  $\psi$  (by the way, this is the solution to the nuisance parameter problem from a Bayesian perspective).

- But “we are Bayesian”, so  $\psi$  and  $\lambda$  are endowed of a probability distribution and then we can marginalize the sampling distribution w.r.t  $\psi$  (by the way, this is the solution to the nuisance parameter problem from a Bayesian perspective).
- The induced sampling process is now described by the sampling distribution

$$(Y | \lambda) \sim \mathcal{N}(\lambda, \sigma_1^2 + \sigma_2^2),$$

where the prior distribution is given by  $\lambda \sim \mathcal{N}(0, 1)$ .

- In this case,  $\lambda$  is identified . . . situation saved!

# Example 1

- Where is the trap?



# Example 1

- Where is the trap?
- Which is the statistical model,

$$(Y | \psi, \lambda) \sim \mathcal{N}(\psi, \sigma_1^2)$$

or

$$(Y | \lambda) \sim \mathcal{N}(\lambda, \sigma_1^2 + \sigma_2^2)$$

?

# Example 1

- Where is the trap?
- Which is the statistical model,

$$(Y | \psi, \lambda) \sim \mathcal{N}(\psi, \sigma_1^2)$$

or

$$(Y | \lambda) \sim \mathcal{N}(\lambda, \sigma_1^2 + \sigma_2^2)$$

?

- Is it true that we accept the definition of a statistical model as the one that describes the data we are going to analyze?

- When Bayesian statisticians change arbitrarily the model specification with involved constructions, I feel that modelling is arbitrary rather than rigorous.
- Science is not objective, is rigorous.

## Example 2

- Let us consider an example widely discussed in the Bayesian literature: Poirier (1998), Carlin and Loui (2000), Xie and Carlin (2006) and Kass et al. (1998).

- Let us consider an example widely discussed in the Bayesian literature: Poirier (1998), Carlin and Loui (2000), Xie and Carlin (2006) and Kass et al. (1998).
- The data generating process is characterized by the sampling distribution

$$(Y_i | a, b) \stackrel{\text{iid}}{\sim} \mathcal{N}(a + b, \sigma_Y^2),$$

where  $\sigma_Y^2$  is known. The parameters of interest are  $(a, b)$ , but it is not identified because the mapping  $(a, b) \mapsto \mathcal{N}(a + b, \sigma_Y^2)$  is not injective.

- The prior specification is given by

$$(a | \mu_a, \sigma_a^2) \sim \mathcal{N}(\mu_a, \sigma_a^2), \quad (b | \mu_b, \sigma_b^2) \sim \mathcal{N}(\mu_b, \sigma_b^2), \quad a \perp\!\!\!\perp b | \mu_a, \mu_b, \sigma_a^2, \sigma_b^2.$$

- The Bayesian model is, consequently, characterized by a unique probability distribution given by

$$\left( \begin{array}{c} Y \\ a \\ b \end{array} \middle| \omega \right) \sim \mathcal{N}_3 \left( \left( \begin{array}{c} \mu_a + \mu_b \\ \mu_a \\ \mu_b \end{array} \right), \left( \begin{array}{ccc} \sigma_a^2 + \sigma_b^2 + \sigma_Y^2 & \sigma_a^2 & \sigma_b^2 \\ \sigma_a^2 & \sigma_a^2 & 0 \\ \sigma_b^2 & 0 & \sigma_b^2 \end{array} \right) \right),$$

where  $\omega = (\mu_a, \mu_b, \sigma_a^2, \sigma_b^2)$  are the hyper-parameters.

- The Bayesian model is, consequently, characterized by a unique probability distribution given by

$$\left( \begin{array}{c} Y \\ a \\ b \end{array} \middle| \omega \right) \sim \mathcal{N}_3 \left( \left( \begin{array}{c} \mu_a + \mu_b \\ \mu_a \\ \mu_b \end{array} \right), \left( \begin{array}{ccc} \sigma_a^2 + \sigma_b^2 + \sigma_Y^2 & \sigma_a^2 & \sigma_b^2 \\ \sigma_a^2 & \sigma_a^2 & 0 \\ \sigma_b^2 & 0 & \sigma_b^2 \end{array} \right) \right),$$

where  $\omega = (\mu_a, \mu_b, \sigma_a^2, \sigma_b^2)$  are the hyper-parameters.

- The parameters of interest are updated through their posterior distributions

$$(a | Y, \omega) \sim \mathcal{N} \left( \frac{(\sigma_b^2 + \sigma_Y^2)\mu_a + \sigma_a^2(Y - \mu_b)}{\sigma_a^2 + \sigma_b^2 + \sigma_Y^2}, \sigma_a^2 \left[ 1 - \frac{\sigma_a^2}{\sigma_a^2 + \sigma_b^2 + \sigma_Y^2} \right] \right),$$

$$(b | Y, \omega) \sim \mathcal{N} \left( \frac{(\sigma_a^2 + \sigma_Y^2)\mu_b + \sigma_b^2(Y - \mu_a)}{\sigma_a^2 + \sigma_b^2 + \sigma_Y^2}, \sigma_b^2 \left[ 1 - \frac{\sigma_b^2}{\sigma_a^2 + \sigma_b^2 + \sigma_Y^2} \right] \right).$$

- However, these posterior distributions are functions of the posterior distribution of the identified parameter since

$$E[a | Y] = \eta_{a,b} + \frac{\sigma_a^2}{\sigma_b^2 + \sigma_a^2} E[a + b | Y],$$

where  $\eta_{a,b} \doteq \frac{\sigma_b^2}{\sigma_b^2 + \sigma_a^2} \mu_a - \frac{\sigma_a^2}{\sigma_b^2 + \sigma_a^2} \mu_b$ ; and similarly for  $E[b | Y]$



- Wechsler, Izbicki and Esteves (2013), A Bayesian look at nonidentifiability: A simple example.
- They consider a problem proposed by Ross (2009):

*“A ball is in any one of  $n$  boxes and is in the  $i$ th box with probability  $p_i$ . If the ball is in box  $i$ , a search of that box will uncover it with probability  $\alpha_i$ . Show that the conditional probability that the ball is in box  $j$ , given that a search of box  $i$  did not uncover it, is*

$$\frac{p_j}{1 - \alpha_i p_i} \text{ if } j \neq i, \quad \text{and} \quad \frac{(1 - \alpha_i)p_i}{1 - \alpha_i p_i} \text{ if } j = i.”$$

- The problem rephrased:

*“A missing plane has crashed in one of  $n$  possible regions. A search for the plane is then performed in region  $i$ ;  $p_i$  is the probability that the plane has gone down in region  $i$  and  $\alpha_i$  the probability that the plane will be found in a search of the  $i$ -th region when it is, in fact, in that region. We consider  $0 < p_i < 1$  and  $0 < \alpha_i < 1$  to make the problem reasonable and avoid trivialities.”*

- The problem rephrased:

*“A missing plane has crashed in one of  $n$  possible regions. A search for the plane is then performed in region  $i$ ;  $p_i$  is the probability that the plane has gone down in region  $i$  and  $\alpha_i$  the probability that the plane will be found in a search of the  $i$ -th region when it is, in fact, in that region. We consider  $0 < p_i < 1$  and  $0 < \alpha_i < 1$  to make the problem reasonable and avoid trivialities.”*

- We note that the parameter of interest is the region in which the plane has gone down; let us denote it as  $\theta$ .
- The parameter space can be represented by a set of integers because each region is labelled by a number. Thus, the parameter space  $\Theta$  corresponds to the set  $\{1, 2, \dots, n\}$  and the prior specification is given by

$$p_j = P[\theta = j], \quad j = 1, 2, \dots, n.$$

## Example 3

- The observed data is represented by the random variable  $X$ :  $X = 1$  if the plain is found in region  $i$ ; 0 otherwise.
- Thus,  $X \in \{0, 1\}$  and, according to the conditions of the problem,

$$P[X = 0 \mid \theta = j] = 1 \text{ if } j \neq i, \quad P[X = 1 \mid \theta = i] = \alpha_i.$$

## Example 3

- The observed data is represented by the random variable  $X$ :  $X = 1$  if the plain is found in region  $i$ ; 0 otherwise.
- Thus,  $X \in \{0, 1\}$  and, according to the conditions of the problem,

$$P[X = 0 \mid \theta = j] = 1 \text{ if } j \neq i, \quad P[X = 1 \mid \theta = i] = \alpha_i.$$

- An application of Bayes theorem gives us (for  $j \neq i$ )

$$P[\theta = i \mid X = 1] = 1, \quad P[\theta = i \mid X = 0] = \frac{(1-\alpha_i)p_i}{1-\alpha_i p_i},$$

$$P[\theta = j \mid X = 1] = 0, \quad P[\theta = j \mid X = 0] = \frac{p_j}{1-\alpha_i p_i}.$$

- Sample space:  $(\{0, 1\}, \{\emptyset, \{0\}, \{1\}, \{0, 1\}\})$ .
- Sampling probabilities: for simplicity, let me suppose that the label  $i$  of the original problem corresponds to 1; then the sampling probabilities are the following:

$$P[X = 1 \mid \theta = 1] = \alpha_1, \quad P[X = 1 \mid \theta = j] = 0 \text{ for all } j \geq 2.$$

It is clear how these probabilities are indexed by the parameter  $\theta$ .

- Parameter space:  $\Theta = \{1, 2, \dots, n\}$ .
- Thus, compactly written, the statistical model corresponding to the problem at hand is given by

$$\mathcal{E} = \{ (\{0, 1\}, \{\emptyset, \{0\}, \{1\}, \{0, 1\}\}), P[\bullet \mid \theta] : \theta \in \Theta \},$$

where

$$P[X = 0 \mid \theta] = (1 - \alpha_1)\mathbb{1}_{\{\theta=1\}} + \sum_{j=2}^n \mathbb{1}_{\{\theta=j\}}.$$

- Prior specification: the corresponding measurable space is given by  $(\Theta, \mathcal{P}(\Theta))$ .
- In order to define a probability on it, it is enough to take  $n$  real numbers  $p_1, p_2, \dots, p_n$  such that

$$0 \leq p_j \leq 1, \quad \sum_{j=1}^n p_j = 1.$$

- We interpret  $p_j$  as  $P[\theta = j]$ . The prior distribution corresponds, therefore, to the vector  $(p_1, p_2, \dots, p_n)'$ : there are as many prior distributions as numbers  $p_j$ 's we choose. Each person –as Bayesian statisticians say– chooses as its initial opinion.

- These elements allow us to construct a Bayesian model, which is characterized by a *unique* probability defined on the product space “parameters  $\times$  observations”.
- In our case, the Bayesian model is given by

$$(\Theta \times \{0, 1\}, \mathcal{P}(\Theta) \otimes \mathcal{P}(\{0, 1\}), \Pi),$$

where  $\Pi$  can be fully described by the following table:

	$\{X = 0\}$	$\{X = 1\}$
$\{\theta = 1\}$	$(1 - \alpha_1)p_1$	$\alpha_1 p_1$
$\{\theta = 2\}$	$p_2$	0
$\{\theta = 3\}$	$p_3$	0
$\vdots$	$\vdots$	$\vdots$
$\{\theta = n\}$	$p_n$	0



## Example 3

- Let us consider the sampling distributions. It is easy to verify that the parameter  $\theta$  is unidentified.

- In fact,

$$\begin{aligned}P[X = 1 \mid \theta = 2] &= P[X = 1 \mid \theta = 3] = 0, \\P[X = 0 \mid \theta = 2] &= P[X = 0 \mid \theta = 3] = 1.\end{aligned}$$

- That is, we found at least two different parameters such that the respective sampling probabilities are equal.

## Example 3

- In spite of this identification problem, it is still possible to make inferences on all the parameters  $\theta$ , in particular on those that are unidentified.
- Thus, for instance,

$$P[\theta = 2 \mid X = 0] = \frac{p_2}{1 - \alpha_1 p_1}, \quad P[\theta = 2 \mid X = 1] = 0.$$

- In spite of this identification problem, it is still possible to make inferences on all the parameters  $\theta$ , in particular on those that are unidentified.
- Thus, for instance,

$$P[\theta = 2 \mid X = 0] = \frac{p_2}{1 - \alpha_1 p_1}, \quad P[\theta = 2 \mid X = 1] = 0.$$

- Those computations show the advantage of the Bayesian approach when compared to the sampling theory framework approach:

*"We reinforce that even a nonidentifiable model may bring information about the parameters of interest. Bayesian inference is easily performed in cases of nonidentifiability [...] Nonidentifiability may bring undesirable nonuniqueness to the ML estimation method" (Weshler et al., 2013, p.93).*

- Let us still focus our attention on the sampling distributions, which can be written in the following terms:

$$P[X = 0 \mid \theta] = (1 - \alpha_1)\mathbb{1}_{\{\theta=1\}} + \mathbb{1}_{\{\theta \neq 1\}},$$

where  $\{\theta \neq 1\} = \{\theta \in \{2, 3, \dots, n\}\}$ .

- Who can describe the corresponding minimal sufficient parameter?
- The sampling process is fully characterized by the measurable partition

$$\{\{\theta = 1\}, \{\theta \neq 1\}\}.$$

- This is the identified parameter.

## Example 3

- Let go to the updating process. I can update  $\{\theta = 1\}$ . This is easy.
- But I can also update  $\{\theta = j\}$  for  $j \geq 2$ , in spite that

$$\{\theta = j\} \notin \{\{\theta = 1\}, \{\theta \neq 1\}\}.$$

## Example 3

- This statement need to be carefully considered.
- In fact, when we compute the posterior probability of  $\{\theta = j\}$  with  $j = 1, 2, \dots, n$ , we consider the measurability structure given by  $\mathcal{P}(\{1, 2, \dots, n\})$ .
- In this context, we can compute the posterior distribution not only of the identified parameters  $\{\theta = 1\}$  and  $\{\theta \neq 1\}$ , but also of  $\{\theta = j\}$  for  $j = 2, 3, \dots, n$ :

$$P[\{\theta = 1\} | X] = \frac{(1 - \alpha_1)p_1}{1 - \alpha_1 p_1} \mathbb{1}_{\{X=0\}} + \mathbb{1}_{\{X=1\}};$$

$$P[\{\theta \neq 1\} | X] = \frac{p_2 + \dots + p_n}{1 - \alpha_1 p_1} \mathbb{1}_{\{X=0\}} + 0 \cdot \mathbb{1}_{\{X=1\}};$$

$$P[\{\theta = j\} | X] = \frac{p_j}{1 - \alpha_1 p_1} \mathbb{1}_{\{X=0\}} + 0 \cdot \mathbb{1}_{\{X=1\}} \text{ for } j \geq 2.$$

- Now: when we update unidentified parameters, do we learn something different from what we learn from the updating of the identified parameter?
- That is, that what we learn about unidentified parameters in *not* more than what we learn about the identified one: is it correct?
- Note the following: for  $j = 2, 3, \dots, n$ ,

$$\begin{aligned}P[\theta = j \mid X = 0] &= E[\mathbb{1}_{\{\theta=j\}} \mid \{X = 0\}] \\&= E\{E[\mathbb{1}_{\{\theta=j\}} \mid \{X = 0\}, \{\theta = 1\}, \{\theta \neq 1\}]\{X = 0\}\} \\&= E\{E[\mathbb{1}_{\{\theta=j\}} \mid \{\theta = 1\}, \{\theta \neq 1\}] \mid \{X = 0\}\}.\end{aligned}$$

- But

$$\begin{aligned}P[\theta = j \mid \{\theta = 1\}, \{\theta \neq 1\}] &= P[\theta = j \mid \theta = 1]\mathbb{1}_{\{\theta=1\}} + P[\theta = j \mid \theta \neq 1]\mathbb{1}_{\{\theta \neq 1\}} \\&= 0 \cdot \mathbb{1}_{\{\theta=1\}} + \frac{p_j}{1 - p_1} \mathbb{1}_{\{\theta \neq 1\}}.\end{aligned}$$

## Example 3

- Therefore,

$$P[\theta = j \mid X = 0] = E \left[ \frac{p_j}{1 - p_1} \mathbb{1}_{\{\theta \neq 1\}} \mid X = 0 \right] = \frac{p_j}{1 - p_1} P[\theta \neq 1 \mid X = 0].$$



- Therefore,

$$P[\theta = j \mid X = 0] = E \left[ \frac{p_j}{1 - p_1} \mathbb{1}_{\{\theta \neq 1\}} \mid X = 0 \right] = \frac{p_j}{1 - p_1} P[\theta \neq 1 \mid X = 0].$$

- This type of statements should be complemented by a remark previously done, namely that the original Bayesian model can be replaced by a *reduced* model without losing relevant information: we can reduce our original Bayesian model to one which only includes the identified parameters.

## Example 3

- Therefore,

$$P[\theta = j \mid X = 0] = E \left[ \frac{p_j}{1 - p_1} \mathbb{1}_{\{\theta \neq 1\}} \mid X = 0 \right] = \frac{p_j}{1 - p_1} P[\theta \neq 1 \mid X = 0].$$

- This type of statements should be complemented by a remark previously done, namely that the original Bayesian model can be replaced by a *reduced* model without losing relevant information: we can reduce our original Bayesian model to one which only includes the identified parameters.
- By doing so, we avoid redundant information! The original Bayesian model can, therefore, be reduced to the following one:

	$\{X = 0\}$	$\{X = 1\}$
$\{\theta = 1\}$	$(1 - \alpha_1)p_1$	$\alpha_1 p_1$
$\{\theta \neq 1\}$	$p_2 + \dots + p_n$	0

(1)

- This reduction can be done *without losing relevant information*.
- Consequently, the likelihood is *exhaustively* described by the identified parameter, and the learning-by-observing process is fully concentrated on the minimal sufficient parameter in the sense that, conditionally on it, *there is nothing more to learn on  $\{\theta = j\}$  for  $j = 2, \dots, n$* .
- Moreover, the posterior probability of the unidentified parameters are less than the posterior probability of the identified ones: for each  $j = 2, \dots, n$ ,

$$P[\theta = j \mid X = 0] < P[\theta \neq 1 \mid X = 0], \quad P[\theta = j \mid X = 1] < P[\theta \neq 1 \mid X = 1].$$

- In other words, it is more informative the posterior distribution of an identified parameter than that of the unidentified parameter.

- The reduced Bayesian model is now defined on the product space

$$(\{1, \dots, n\} \times \{0, 1\}, \mathcal{A} \otimes \mathcal{S}, \Pi),$$

where

- $\mathcal{A} = \{\emptyset, \{\theta = 1\}, \{\theta \neq 1\}, \Theta\}$ ;
  - $\mathcal{S} = \{\emptyset, \{X = 0\}, \{X = 1\}, \{0, 1\}\}$ ;
  - and  $\Pi$  is the matrix that Alpha computed.
- In this model, the events

$$\{\theta = j\}, \quad j = 2, 3, \dots, n$$

are *not* measurable and, consequently, they don't provide information on the observations are not informative: all the relevant information regarding the data is concentrated on the identified parameter.

- What we prove is that all the information generated by the sampling process is represented by the partition

$$\{\{\theta = 1\}, \{\theta \neq 1\}\}.$$

- Let us define the functions  $Y_0$  and  $Y_1$  as follows:

$$P[X = 0 \mid \theta] = (1 - \alpha_1)\mathbb{1}_{\{\theta=1\}} + \mathbb{1}_{\{\theta \neq 1\}} \doteq Y_0$$

$$P[X = 1 \mid \theta] = \alpha_1\mathbb{1}_{\{\theta=1\}} + 0 \cdot \mathbb{1}_{\{\theta \neq 1\}} \doteq Y_1.$$

Therefore,

$$\{\theta = 1\} = Y_0^{-1}[\{1 - \alpha_1\}] \cap Y_1^{-1}[\{\alpha_1\}], \quad \{\theta \neq 1\} = Y_0^{-1}[\{1\}] \cap Y_1^{-1}[\{0\}].$$

## Example 3

- It can be seen that the parameters  $\{\theta = j\}$  for  $j = 2, \dots, n$  can not be expressed as functions of the sampling distributions.
- The parameter  $\{\theta = 1\}$  that the plain is in region 1, provides the same information as the probability of finding it and of not finding it there. The parameter  $\{\theta \neq 1\}$  that the plain is not in region 1 provides the same information as the probability of finding and of not finding the plain in a region different from 1.

- Consider a Bayesian model  $p(X, \theta)$ , which can be decomposed as follows

$$p(X, \theta) = m(\theta)q(X | \theta) = q(X) m(\theta | X).$$

- Let us decompose both the observations and the parameters into two components:  $X = (Y, Z)$  and  $\theta = (\psi, \lambda)$ , where  $\psi$  is a nuisance parameter and  $\lambda$  is a parameter of interest. Therefore, we only have interest in the posterior distribution of  $\lambda$  which can be decomposed as

$$m(\lambda | X) \propto m(\lambda) q(Z | \lambda) q(Y | Z, \lambda).$$

- We are interesting in the following question: how far can we forget the term  $p(z | \lambda)$  for the inference on  $\lambda$ ? Let us introduce three sufficient conditions in a decreasing order of generality.

- If  $Z$  and  $\lambda$  are *mutually ancillary*, that is

$$Z \perp\!\!\!\perp \lambda,$$

then  $q(Z | \lambda) = q(Z)$  and  $m(\lambda | Z) = m(\lambda)$  and, therefore,

$$m(\lambda | X) \propto m(\lambda) q(Y | Z, \lambda).$$



- However, the use of

$$m(\lambda | X) \propto m(\lambda) q(Y | Z, \lambda).$$

implies integrating the conditional sampling distribution of  $(Y | Z)$  with respect to the conditional posterior distribution of  $\psi$  given  $Z$  and  $\lambda$ , namely

$$q(Y | Z, \lambda) = \int_{\Psi} q(Y, \psi | Z, \lambda) d\psi = \int_{\Psi} q(Y | Z, \theta) m(\psi | Z, \lambda) d\psi.$$

- But if  $\lambda$  is a sufficient parameter of the conditional sampling process  $(Y | Z)$ , i.e.,

$$Y \perp\!\!\!\perp \psi | Z, \lambda,$$

then  $q(Y | Z, \lambda) = q(Y | Z, \psi, \lambda)$  and, therefore,

$$\begin{aligned} q(Y | Z, \lambda) &= \int_{\Psi} q(Y | Z, \theta) m(\psi | Z, \lambda) d\psi \\ &= \int_{\Psi} q(Y | Z, \lambda) m(\psi | Z, \lambda) d\psi \\ &= q(Y | Z, \lambda) \int_{\Psi} m(\psi | Z, \lambda) d\psi = q(Y | Z, \lambda). \end{aligned}$$

- When both

$$Z \perp\!\!\!\perp \lambda, \quad Y \perp\!\!\!\perp \psi \mid Z, \lambda$$

are verified,  $\lambda$  and  $Z$  are said to be *mutually exogenous*.

- Thus, if only  $\lambda$  is of interest, and if  $\lambda$  and  $Z$  are mutually exogenous, the process generating  $Z$  becomes irrelevant, and the data-generating process may be directly conditioned on  $Z$ .
- This is precisely the meaning of the so-called “exogenous variables” used in econometric models and introduced by Koopmans (1950).

- When  $\lambda$  is a sufficient parameter of the sampling distribution conditional on  $Z$ , a rather natural way of obtaining the mutual exogeneity of  $\lambda$  and  $Z$  is through the structure of a *Bayesian cut* which means that along with

$$Y \perp\!\!\!\perp \psi \mid Z, \lambda,$$

$\psi$  is furthermore a sufficient parameter of the sampling marginal process of  $Z$ , i.e.

$$Z \perp\!\!\!\perp \lambda \mid \psi,$$

and is also a priori independent of  $\lambda$ :

$$\psi \perp\!\!\!\perp \lambda.$$

- The structure of a Bayesian cut is given by the following decomposition:

$$\begin{aligned}q(Y, Z, \lambda, \psi) &= p(Y | Z, \lambda, \psi) q(Z | \lambda, \psi) m(\lambda, \psi) \\ &= [p(Y | Z, \lambda) m(\lambda)] [q(Z | \psi) m(\psi)]\end{aligned}$$

- In a sampling theory framework, the corresponding concept of cut (Barndorff-Nielsen, 1978; Engle, Hendry and Richard, 1980) replaces the prior independence between  $\psi$  and  $\lambda$  by the condition of being variation-free, that is,

$$(\psi, \lambda) \in \Lambda \times \Psi.$$

- The Bayesian cut implies different consequences:
  - 1 Conditions  $Z \perp\!\!\!\perp \lambda \mid \psi$  and  $\psi \perp\!\!\!\perp \lambda$  are jointly equivalent to  $(Z, \psi) \perp\!\!\!\perp \lambda$ , which implies condition  $Z \perp\!\!\!\perp \lambda$ . In other words, Bayesian cut implies mutual ancillarity between  $Z$  and  $\lambda$ .
  - 2  $(Z, \psi) \perp\!\!\!\perp \lambda$  also implies that  $\psi \perp\!\!\!\perp \lambda \mid Z$ , that is, the posterior independence of  $\lambda$  and  $\psi$  conditionally on  $Z$ .
  - 3 Condition  $Z \perp\!\!\!\perp \lambda \mid \psi$ , along with  $\psi \perp\!\!\!\perp \lambda \mid Z$ , are jointly equivalent to  $\psi \perp\!\!\!\perp (Y, \lambda) \mid Z$ . This latter condition implies that  $\psi \perp\!\!\!\perp \lambda \mid (Y, Z)$ . In other words, the Bayesian cut implies the posterior independence of  $\psi$  and  $\lambda$  given  $(Y, Z)$ .

- Let us consider bivariate parameters and observations:  $a = (b, c)$  and  $X = (T, U)$ . Suppose now that each coordinate may assume two values only:  $b \in \{b_1, b_2\}$ ,  $c \in \{c_1, c_2\}$ ,  $T \in \{T_1, T_2\}$  and  $U \in \{U_1, U_2\}$ . The Bayesian model characterized by the joint probability  $\pi(a, X)$  is suitably defined by the following 15 numbers (assumed to be different from zero):

$$\begin{aligned}\mu(c_1) &= \mu_0 \\ \mu(b_1 | c_i) &= \mu_i, & i &= 1, 2 \\ p(T_1 | c_i, b_j) &= p_{ij}, & i, j &= 1, 2 \\ p(U_1 | T_i, c_j, b_k) &= q_{ijk}, & i, j, k &= 1, 2.\end{aligned}$$

- Mutual ancillarity between  $c$  and  $T$  is equivalent to  $p(T_1 | c_1) = p(T_1 | c_2)$ , that is,

$$\mu_1 p_{11} + (1 - \mu_1) p_{12} = \mu_2 p_{21} + (1 - \mu_2) p_{22}. \quad (2)$$

- For mutual exogeneity the supplementary condition ( $a \perp\!\!\!\perp X | T, c$ ) is equivalent to the four equalities

$$q_{ij1} = q_{ij2} \quad i, j = 1, 2. \quad (3)$$

- For a cut, it is necessary to verify condition (3) plus the following three:

$$\mu_1 = \mu_2, \quad \text{i.e., } b \perp\!\!\!\perp c, \quad (4)$$

$$p_{1j} = p_{2j}, \quad j = 1, 2, \quad \text{i.e., } a \perp\!\!\!\perp T \mid b. \quad (5)$$

- Even under prior independence (4), it should be clear that mutual exogeneity (conditions (2) and (3)) does not imply a cut (conditions (3), (4) and (5)).