

# Likelihood in the symbolic context. Examples

Richard Emilion

<sup>1</sup>Denis Poisson Institute  
University of Orléans

Data Science, January 11<sup>th</sup>, 2019

Warmest thanks To Edwin Diday and to the organizers for this invitation.

- **I - Probabilistic setting for the symbolic paradigm**
- **II - Likelihood, examples**

## I - A Probabilistic setting for the symbolic paradigm

. E. Diday (1990).

Symbolic data expresses the variability of the data within a class of data.

SDA : Considering classes of data as new statistical units.

## I - A Probabilistic setting for the symbolic paradigm

. E. Diday (1990).

Symbolic data expresses the variability of the data within a class of data.

SDA : Considering classes of data as new statistical units.

. R.E. (ISI Conference, Rio 2015) Introducing a formalism:  $X$  (data r.v.),  $C$  (class variable),  $S$  symbolic variable

## I - A Probabilistic setting for the symbolic paradigm

. E. Diday (1990).

Symbolic data expresses the variability of the data within a class of data.

SDA : Considering classes of data as new statistical units.

. R.E. (ISI Conference, Rio 2015) Introducing a formalism:  $X$  (data r.v.),  $C$  (class variable),  $S$  symbolic variable

.Some modifications of Rio talk: Emilion-Diday 2018, Book chapter, Eds. G. Saporta

## I - A Probabilistic setting for the symbolic paradigm

. E. Diday (1990).

Symbolic data expresses the variability of the data within a class of data.

SDA : Considering classes of data as new statistical units.

. R.E. (ISI Conference, Rio 2015) Introducing a formalism:  $X$  (data r.v.),  $C$  (class variable),  $S$  symbolic variable

.Some modifications of Rio talk: Emilion-Diday 2018, Book chapter, Eds. G. Saporta

. Present talk : Examples of symbolic likelihood. Applications.

# I.1 Description Variable

4 / 17

- Population of individuals:  $(\Omega, \mathcal{F}, \mathbb{P})$  a probability space
- $(\mathbb{V}, \mathcal{V})$  a measurable space of descriptions.  
 $X : \Omega \rightarrow \mathbb{V}$ , measurable w.r.t.  $\mathcal{F}$  and  $\mathcal{V}$ , a random variable (r.v) which describes the individuals.
- Generally,  $\mathbb{V}$  is a measurable subset of  $\mathbb{R}^p$  and  
 $X = (X_1, \dots, X_p)$ ,  $p = 1, 2, \dots$
- Standard Data Analysis:  $n \times p$  numerical table of a  $n$ -sample of  $(X_1, \dots, X_p)$

## I.2 Class Variable

5 / 17

- Class variable

$$C : \Omega \longrightarrow \mathbb{C}, \text{ measurable w.r.t. } \mathcal{F} \text{ and } \mathcal{C}, \quad (1)$$

r.v. which assigns a class label to each individual.  
 $(\mathbb{C}, \mathcal{C})$  is a measurable space of class labels.



## I.2 Class Variable

- Class variable

$$C : \Omega \longrightarrow \mathbb{C}, \text{ measurable w.r.t. } \mathcal{F} \text{ and } \mathcal{C}, \quad (1)$$

r.v. which assigns a class label to each individual.

$(\mathbb{C}, \mathcal{C})$  is a measurable space of class labels.

- $X$  and  $C$  are correlated
- Class with label  $c \in \mathbb{C}$ , shortly Class  $c$  :

$$(C = c) = \{\omega \in \Omega : C(\omega) = c\} \quad (2)$$

- It is assumed that singletons  $\{c\}$  belong to  $\mathcal{C}$ ,  $\forall c \in \mathbb{C}$ , so that classes for  $c \in \text{range}(C)$  form a measurable partition of  $\Omega$

## 1.5 Symbolic variable, Symbolic data

6 / 17

**Definition** : A symbolic variable  $S$  of the context  $(X, C)$  is defined as a mapping

$$\begin{aligned} S : \mathbb{C} &\longrightarrow \mathbb{S} \\ S(c) &= f(\mathbb{P}_{X|C=c}). \end{aligned} \quad (3)$$

where

$$f : \mathcal{M}_1(\mathcal{V}) \longrightarrow \mathbb{S} \quad (4)$$

is a measurable function taking value in some measurable space of symbols  $(\mathbb{S}, \mathcal{S})$ .

$S(c), c \in \mathbb{C}$  is a 'symbolic data' representing the variability of the data  $X(\omega)$  for  $\omega \in (C = c)$ .

# I.6 Symbols in term of samples

7 / 17

- $\mathbb{P}_{X|C=c}$  is a probability distribution on  $(\mathbb{V}, \mathcal{V})$ , a complex object. It is generally estimated from an observed sample  $(x^{(1)}, c^{(1)}), \dots, (x^{(n)}, c^{(n)})$  of the pair  $(X, C)$  such that  $c^{(j)} = c$ .

# I.6 Symbols in term of samples

7 / 17

- $\mathbb{P}_{X|C=c}$  is a probability distribution on  $(\mathbb{V}, \mathcal{V})$ , a complex object. It is generally estimated from an observed sample  $(x^{(1)}, c^{(1)}), \dots, (x^{(n)}, c^{(n)})$  of the pair  $(X, C)$  such that  $c^{(j)} = c$ .
- $S(c)$  can be estimated by an aggregating function of a sample  $(x^{(1)}, c^{(1)}), \dots, (x^{(n)}, c^{(n)})$  such that  $c^{(j)} = c$ .
- BLS (Beranger - Lin - Sisson) Definition

## II - Likelihood in the symbolic context

II.1. Probability measures on  $(\mathbb{C}, \mathcal{C})$ , density 9 / 17

- $S : \mathbb{C} \rightarrow \mathbb{S}$  a symbolic variable with  $\mathbb{S} = \mathbb{N}^m$  or  $\mathbb{S} = \mathbb{R}^m$
- Problem: density of  $S$  w.r.t. to the counting (resp. the Lebesgue) measure? In the continuous case  $\mathbb{C}$  should be uncountable and even nonatomic.

$$\begin{aligned} d_S : \mathbb{S} &\longrightarrow \mathbb{R}_+ \\ \underline{s} &\longrightarrow d_S(\underline{s}) \end{aligned} \tag{5}$$

Estimating  $d_S$  given a  $n$ -sample  $\underline{s}_1, \dots, \underline{s}_n$  :

$$\underline{s}_i = (s_{i,1}, \dots, s_{i,m}) = S^{(i)}(c) \in \mathbb{R}^m, \quad S^{(i)} \stackrel{i.i.d.}{\sim} \mathbb{Q}_S, \quad i = 1, \dots, n$$

for some  $c \in \mathbb{C}$  : randomness of the sample of symbols.

## II.2. LDA model (Blei-Ng-Jordan) Specif. a 10 / 17

Model in Text Mining context. Can be applied in other domains.

. *categorical* r.v.  $X : \Omega \rightarrow \mathbb{V} = \{1, \dots, k\}$ , a finite set of  $k$  topics,

. r.v.  $N : \Omega \rightarrow \mathbb{N} = \{0, 1, 2, \dots\}$ ,

random probability vector  $\theta = (\theta_1, \dots, \theta_k) : \Omega \rightarrow \mathbb{T}_k$

$$\left\{ \begin{array}{l} (N, \theta) \sim \text{Poisson}(\xi) \otimes \text{Dirichlet}(\underline{\alpha}) \\ \mathbb{P}(X = i | \theta) = \theta_i, \quad i = 1, \dots, k. \end{array} \right. \quad (6)$$

## II.2. LDA model. Specifications b

11 / 17

.  $\{1, \dots, V\}$  a finite set of  $V$  words  
 $\beta = (\beta_{i,j}), i = 1, \dots, k, j = 1, \dots, V$  be a  $k \times V$  Markov matrix, each of its  $k$  rows being a probability vector in dimension  $V$ . A *document*, is considered as an outcome  $c$  of our class random variable  $C$  defined as a sequence of random words:

$$\left\{ \begin{array}{l} C(\omega) = (W^{(1)}(\omega), \dots, W^{(N(\omega))}(\omega)), \omega \in \Omega \\ \text{where, given } N \text{ and } \theta \\ X^{(r)} \stackrel{i.i.d}{\sim} \mathbb{P}_{(X|\theta)}, \text{ for each } r = 1, \dots, N \\ W^{(r)} : \Omega \longrightarrow \{1, \dots, V\}, r = 1, \dots, N, \text{ are independent} \\ \mathbb{P}(W^{(r)} = v | X^{(r)}) = \beta_{(X^{(r)}, v)}, \text{ for each } r = 1, \dots, N, v = 1, \dots, V. \end{array} \right. \quad (7)$$



## II.2. LDA model. Symbolic likelihood

12 / 17

Given a class label  $c = (w_1, \dots, w_N)$ , class  $c$  is defined as

$$(C = c) = \{\omega \in \Omega : W^{(1)}(\omega) = w_1, \dots, W^{(N)}(\omega) = w_N\}$$

Topic variability of a document  $c$  with  $N$  words  $(w_1, \dots, w_N)$  and unobserved topics  $(x_1, \dots, x_N)$  expressed by the *latent symbol*

$$s(c) = \left( \sum_{r=1}^N 1_{(x_r=1)}, \dots, \sum_{r=1}^N 1_{(x_r=k)} \right)$$

Random symbol  $S = s \circ C = \left( \sum_{r=1}^N 1_{(X^{(r)}=1)}, \dots, \sum_{r=1}^N 1_{(X^{(r)}=k)} \right)$  distribution, given  $(N, \theta)$ , is multinomial

$$\mathbb{P}(S = (n_1, \dots, n_k) | N = n, \theta) = \frac{n!}{n_1! \dots n_k!} \theta_1^{n_1} \dots \theta_k^{n_k}$$

if  $n_1 + \dots + n_k = n$

## II.2. LDA. Document / Corpus Probability 13 / 17

We have  $p(x_r|\theta) = \prod_{i=1}^k \theta_i^{1_{x_r=i}}$  and  $p(w_r|x_r, \beta) = \prod_{j=1}^V \beta_{x_r,j}^{1_{w_r=j}}$

$$\begin{cases} p(x_r, w_r|\theta, \beta) = \prod_{i=1}^k \theta_i^{1_{x_r=i}} \prod_{j=1}^V \beta_{x_r,j}^{1_{w_r=j}} \\ p(w_r|\theta, \beta) = \sum_{x_r} \prod_{i=1}^k \theta_i^{1_{x_r=i}} \prod_{j=1}^V \beta_{x_r,j}^{1_{w_r=j}} \end{cases} \quad (8)$$

The probability of a document is,

$$p(w_1, \dots, w_N|\theta, \beta, N) = \prod_{r=1}^N \sum_{x_r} \prod_{i=1}^k \theta_i^{1_{x_r=i}} \prod_{j=1}^V \beta_{x_r,j}^{1_{w_r=j}} \quad (9)$$

$$p(w_1, \dots, w_N|\underline{\alpha}, \beta, N) = \int Dd(\theta|\underline{\alpha}) \prod_{r=1}^N \sum_{x_{d_r}} \prod_{i=1}^k \theta_i^{1_{x_r=i}} \prod_{j=1}^V \beta_{x_r,j}^{1_{w_r=j}} d\theta$$

## II.3. BLS (Beranger-Lin-Sisson) method

14 / 17

$X : \Omega \rightarrow \mathbb{R}^p$  r.v. with density  $d_X(\cdot|\theta)$ .

$N \geq 2$  any very large integer.  $(X^{(1)}, \dots, X^{(N)})$ ,  $X^{(r)} \stackrel{i.i.d}{\sim} \mathbb{P}_X$ .

$(x^{(1)}, \dots, x^{(N)})$  an observed large sample

$c = (B_k)_{k \in K}$  a finite partition covering the support of  $\mathbb{P}_{(X^{(1)}, \dots, X^{(N)})}$

Given partition  $c$ , the joint distribution of

$$(1_{(X^{(1)} \in B_1)}, \dots, 1_{(X^{(N)} \in B_1)}, \dots, 1_{(X^{(1)} \in B_K)}, \dots, 1_{(X^{(N)} \in B_K)})$$

is captured by the symbolic variable  $S$  defined as

$$S(c) = \left( \sum_{r=1}^N 1_{(X^{(r)} \in B_1)}, \dots, \sum_{r=1}^N 1_{(X^{(r)} \in B_k)}, \dots \right),$$

whose distribution is multinomial  $(N, p_1, \dots, p_K)$  with

$$p_k = \int_{B_k} d_X(x|\theta) dx$$

## II.3. BLS (Beranger-Lin-Sisson) result

15 / 17

The main observation in BLS ( Beranger B., Lin H., Scott A. S., New models for symbolic data analysis, ArXiv e-prints (2018)) is that the total probability formula

$$d_{S|\nu,\theta} = \int_{t \in (\mathbb{R}^p)^N} d_{S|(X^{(1)}, \dots, X^{(N)})=t, \nu} (d_X)^{\otimes N} (t|\theta) \quad (10)$$

yields an inference on parameter  $\theta$  from an inference on the symbolic likelihood with parameter  $\nu$ .

This considerably reduces inference complexity and seems to be a significant application of the symbolic approach

## II.4 Dirichlet Process Mixture (DPM)

16 / 17

- $h_i$  : symbol (e.g. histogram  $h_i =$  bins, frequencies) Bayesian parametric:  $h_i|\theta \sim F(\theta)$ ,  $\theta \sim D$ : apriori  $D$ , shape

## II.4 Dirichlet Process Mixture (DPM)

16 / 17

- $h_i$  : symbol (e.g. histogram  $h_i = \text{bins, frequencies}$ ) Bayesian parametric:  $h_i|\theta \sim F(\theta)$ ,  $\theta \sim D$ : apriori  $D$ , shape
- DPM: Bayesian nonparametric, flexible, infinite mixture model

$$\begin{cases} h_i|\theta_i \stackrel{\text{ind}}{\sim} F(\cdot, \theta_i), i = 1, \dots, n, \quad \theta_i \in \Theta \\ \theta_i|P = p \stackrel{\text{i.i.d.}}{\sim} p, i = 1, \dots, n \\ P \sim DP(c, P_0) \text{ a Dirichlet Process on } \Theta \end{cases} \quad (11)$$

Draw  $p$  from  $DP(c, P_0)$ ,  $\theta_i$  from  $p$  and  $h_i$  from  $F(\cdot, \theta_i)$ : the distribution on  $\Theta$  is the mixture  $\int_{\Theta} F(\cdot, \theta) dp(\theta)$

## II.5 Mixture of Dirichlet Processes (MDP) 17 / 17

- Antoniak (Ann. Stat. 1974): If the  $h_i$ 's and  $P$  are as in (11) then the posterior

$$P|h_1, \dots, h_n \sim \int DP(cP_0 + \sum_{i=1}^n \delta_{\theta_i}) d\mathbb{P}_{\theta_1, \dots, \theta_n | h_1, \dots, h_n} \quad (12)$$

- In other words the posterior is a Mixture of Dirichlet Processes (MDP)
- The posterior MDP provides a classification of the histogram data without any apriori number of classes, mixture component = fuzzy class.
- A mixture of DD estimated from an histogram dataset converges to a MDP as the bin width goes to  $0^+$  (R.E. *Stat. Anal. & Data Mining* 2012)
- 'DPpackage' in R