

# Descriptive Statistics for Symbolic Data

Paula Brito

Fac. Economia & LIAAD-INESC TEC, Universidade do Porto

ECI 2015 - Buenos Aires  
T3: Symbolic Data Analysis:  
Taking Variability in Data into Account

# Outline

- 1 A new framework
- 2 Interval-valued variables
- 3 Histogram-valued variables

# Outline

- 1 A new framework
- Interval-valued variables
- Histogram-valued variables

# Descriptive Statistics for Symbolic Variables

No unique and straightforward definitions !

What is the variance of a set of interval observations ?

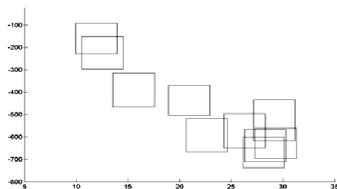
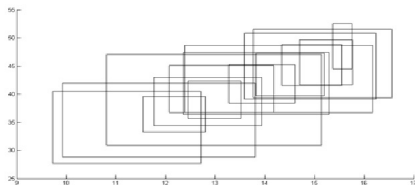
How do we measure correlation ?

- Measures based on interval parameters
- Measures based on distributional assumptions
- Measures based on distances

# Outline

- 1 A new framework
- 2 Interval-valued variables
- 3 Histogram-valued variables

# Biplots for Interval variables



# Descriptive Statistics for Interval Variables

## First option :

Using the dispersion of the interval centers

The mean value and the dispersion of all interval midpoints are given by

$$\bar{Y}_j = \frac{1}{n} \sum_{i=1}^n \frac{l_{ij} + u_{ij}}{2}$$

$$S_{Y_k}^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{l_{ij} + u_{ij}}{2} - \bar{Y}_j \right)^2$$

# Descriptive Statistics for Interval Variables

## Second option :

Using the dispersion of the interval boundaries.

The mean value and the dispersion of all interval midpoints are given by

$$\bar{Y}_j = \frac{1}{n} \sum_{i=1}^n \frac{l_{ij} + u_{ij}}{2}$$

$$S_{Y_k}^2 = \frac{1}{n} \sum_{i=1}^n \frac{(l_{ij} - \bar{Y}_j)^2 + (u_{ij} - \bar{Y}_j)^2}{2}$$



## Descriptive Statistics for Interval Variables

Under the assumption that the observed  $Y_j(s_i)$  and  $Y_{j'}(s_i)$  values,  $i = 1, \dots, n$ , are **uniformly distributed** across each interval  $I_{ik} = [l_{ik}, u_{ik}]$ ,  $k = j, j'$ , we have

$$E(Y_{ik}) = (l_{ik} + u_{ik})/2 = c_{ik} \text{ and } \text{Var}(Y_{ik}) = (u_{ik} - l_{ik})^2/12$$

- **symbolic sample mean** :

$$\bar{Y}_k = \frac{1}{2n} \sum_{i=1}^n (l_{ik} + u_{ik}) = \frac{1}{n} \sum_{i=1}^n c_{ik}$$

- **symbolic sample variance** :

$$\begin{aligned} S_{Y_k}^2 &= \frac{1}{3n} \sum_{i=1}^n [(l_{ik} - \bar{Y}_k)^2 + (l_{ik} - \bar{Y}_k)(u_{ik} - \bar{Y}_k) + (u_{ik} - \bar{Y}_k)^2] \\ &= \frac{1}{3n} \sum_{i=1}^n (l_{ik}^2 + l_{ik}u_{ik} + u_{ik}^2) - \bar{Y}_k^2 \end{aligned}$$

Bertrand and Goupil's (2000)

obtained from the empirical density function for an interval variable ▶

# Descriptive Statistics for Interval Variables

For the symbolic covariance three definitions were proposed :

- $$\text{Cov}_1(Y_j, Y_{j'}) = \frac{1}{4n} \sum_{i=1}^n (l_{ij} + u_{ij})(l_{ij'} + u_{ij'}) - \overline{Y}_j \cdot \overline{Y}_{j'}$$

Billard & Diday (2003)

obtained from the empirical joint density function

- $$\text{Cov}_2(Y_j, Y_{j'}) = \frac{1}{3n} \sum_{i=1}^n G_j G_{j'} [Q_j, Q_{j'}]^{1/2}$$

with  $Q_k = (l_{ik} - \overline{Y}_k)^2 + (l_{ik} - \overline{Y}_k)(u_{ik} - \overline{Y}_k) + (u_{ik} - \overline{Y}_k)^2$ ,

$$G_k = \begin{cases} -1 & \text{if } c_{ik} \leq \overline{Y}_k \\ 1 & \text{if } c_{ik} > \overline{Y}_k \end{cases}$$

Billard & Diday (2006)

incorporating more accurately both between and within interval variations into the overall covariance

# Descriptive Statistics for Interval Variables

$$\begin{aligned}
 \bullet \text{Cov}_3(Y_j, Y_{j'}) &= \underbrace{\frac{1}{n} \sum_{i=1}^n \frac{(u_{ij} - l_{ij})(u_{ij'} - l_{ij'})}{12}}_{\text{WithinSP}} + \\
 &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \left( \frac{l_{ij} + u_{ij}}{2} - \bar{Y}_j \right) \left( \frac{l_{ij'} + u_{ij'}}{2} - \bar{Y}_{j'} \right)}_{\text{BetweenSP}} \\
 &= \frac{1}{6n} \sum_{i=1}^n [2(l_{ij} - \bar{Y}_j)(l_{ij'} - \bar{Y}_{j'}) + (l_{ij} - \bar{Y}_j)(u_{ij'} - \bar{Y}_{j'}) \\
 &\quad + (u_{ij} - \bar{Y}_j)(l_{ij'} - \bar{Y}_{j'}) + 2(u_{ij} - \bar{Y}_j)(u_{ij'} - \bar{Y}_{j'})]
 \end{aligned}$$

Billard (2008)

considering a decomposition into

Within observations Sum of Products (WithinSP) and

Between observations Sum of Products (BetweenSP)

## Interval-valued variables: Distance measures

Many measures proposed in the litterature

Hausdorff distance :

$$d_H(l_i, l_j) = \max \{ \{|l_i - l_j|, |u_i - u_j|\} \}$$

Euclidean distance :

$$d_2(l_i, l_j) = \sqrt{(l_i - l_j)^2 + (u_i - u_j)^2}$$

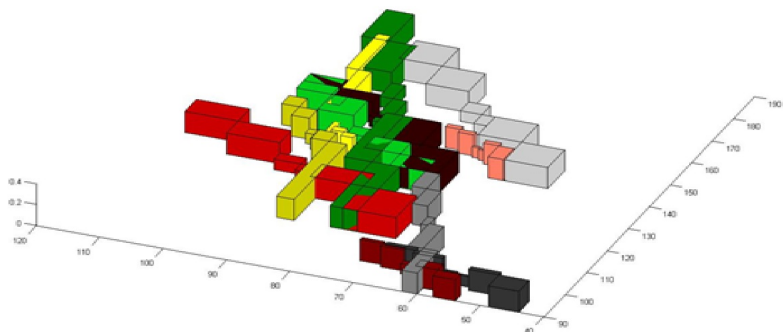
City-Block distance :

$$d_1(l_i, l_j) = |l_i - l_j| + |u_i - u_j| .$$

# Outline

- 1 A new framework
- 2 Interval-valued variables
- 3 Histogram-valued variables

# Biplots for Histogram variables



# Descriptive Statistics for Histogram Variables

Assuming an Uniform distributon within each sub-interval of  $Y_k(s_j)$ ,  $i = 1, \dots, n$ ,  $l_{ik\ell} = [l_{ik\ell}, u_{ik\ell}]$ ,  $\ell = 1, \dots, K_j$ ,  $k = j, j'$  we have

- **symbolic sample mean :**

$$\bar{Y}_k = \frac{1}{2n} \sum_{i=1}^n \sum_{\ell=1}^{K_j} ((l_{ik\ell} + u_{ik\ell}) p_{ik\ell})$$

- **symbolic sample variance :**

$$S_{Y_k}^2 = \frac{1}{3n} \sum_{i=1}^n \sum_{\ell=1}^{K_j} ((l_{ik}^2 + l_{ik} u_{ik} + u_{ik}^2) p_{ik\ell}) - \bar{Y}_k^2$$

Billard and Diday (2003)

# Descriptive Statistics for Histogram Variables

And for the symbolic covariance three definitions :

$$\text{Cov}_1(Y_j, Y_{j'}) = \frac{1}{4n} \sum_{i=1}^n \sum_{\ell=1}^{K_j} p_{ij\ell} p_{ij'\ell} (l_{ij} + u_{ij})(l_{ij'} + u_{ij'}) - \overline{Y_j} \cdot \overline{Y_{j'}}$$

Billard & Diday (2003)

obtained from the empirical joint density function



## Correlation Between Symbolic Variables

As in the classic variables: the **correlation coefficient** is defined as :

$$r_{Y_j Y_{j'}} = \frac{\text{Cov}(Y_j, Y_{j'})}{S_{Y_j} S_{Y_{j'}}$$

where

$\text{Cov}(Y_j, Y_{j'})$  is the covariance function between  $Y_j$  and  $Y_{j'}$

$S_{Y_j}, S_{Y_{j'}}$  the symbolic standard deviation of the variables  $Y_j$  and  $Y_{j'}$ , respectively.

In the particular case of interval variables the descriptive statistics depend on the assumed distribution within each interval.

Results already obtained for other distributions, e.g., the triangular distribution.