

Chapter 1

MEASURING DISPERSION IN THE CONTEXT OF ORDERED QUALITATIVE SCALES

José Luis García-Lapresta¹ and Luis Borge²

¹*PRESAD Research Group, BORDA Research Unit, IMUVA, Departamento de Economía Aplicada, Universidad de Valladolid, Spain*

lapresta@eco.uva.es

²*Departamento de Economía Aplicada, Universidad de Valladolid, Spain*

borge@eco.uva.es

Abstract In this contribution, we introduce a family of dispersion measures in the context of ordered qualitative scales.

Keywords: Qualitative scales; Dispersion; Ordinal proximity measures.

1. Introduction

The main purpose of this contribution is to measure the dispersion associated with the different results of a variable when they form a set of linguistic terms obtained after examining different qualities of people, services, etc, ordered, but with a non-homogeneous and non-quantifiable distance between the linguistic terms.

There is in the literature a great variety of dispersion measures associated with a data set, being the most used the standard deviation. For a very long time, the normality assumption was present in most of the statistical studies and its estimator in combination with the sample mean provide excellent efficiency from the practical side due to its precision to adjust the data and from the capacity of giving rise to theoretical results in the sample distributions. However, most of the estimators of dispersion presents in literature do not hold

for our purposes, because they are based on numerical values obtained in the sample.

When the data that we want to study are obtained from a statistical behaviour different from the normal one, there are atypical observations and losses of efficiency on the classical estimators. The growth of statistical methods and their increasing use in different areas of research poses the study of alternative estimators which are less sensible to deviations to the normal. In the middle of the last century there was an important revival of the robust statistics. The concept of robustness is related to the fact that an estimator is less sensible to small changes in the statistical hypothesis about the models or about the presence of atypical observations, being the breakdown point one of the ways to measure the robustness of an estimator which is related to the percentage of polluted observations of a sample.

The most well known measure of localization of the robust statistical is the median which is estimated ordering the values obtained in a sample, while the sample median is obtained as the central value of the ordered data, if the number of data is odd, and as the average of the central data, if it is even. When the data of a sample is like the ones we are interested in, we may use this localization measurement when the number of data n is odd. When we are dealing with the robust dispersion the most used measurement is the interquartile range. This method is not the best one because an estimator is more robust if it is obtained as a median of the absolute deviation of the median of the sample, also known as MAD. If we have an ordered sample of a variable $(x_{(1)}, \dots, x_{(n)})$, its median, $\text{med}_i(x_i)$ is defined

$$\text{MAD} = \text{med}_j |x_j - \text{med}_i(x_i)|.$$

To obtain this estimator, it is necessary to evaluate the median two times. The first one between all the linguistic terms of the data. The second median, after ordering the different proximities between the linguistic terms, we choose the median of these proximities.

Rousseeuw and Croux [8] presented two alternative robust estimators to the MAD and, as the interquartile range, they do not require to center the data to respect to the sample median. The first of them is

$$S_n = \text{med}_i \{ \text{med}_j |x_i - x_j| \},$$

which for any x_i value we obtain the median of the differences of the form $\{|x_i - x_j|, j = 1, \dots, n\}$, which gives us n medians, being the estimator S_n the median of this set of medians.

The second estimator proposed by Rousseeuw and Croux, which they call Q_n , is based only on the differences between data values. If $\{x_{(1)}, \dots, x_{(n)}\}$ is the ordered sample, let $\mathcal{D} = \{x_{(i)} - x_{(j)}, i > j\}$ be the set of interpoint distances in which the number of elements is $m = \binom{n}{2}$. Then, the estimate is

defined as the k -th order statistic of \mathcal{D} , where $k = \binom{[n/2]+1}{2}$ and that is roughly half the number of observations. These estimator has attractive properties (see Maronna, et al.**)

**Some absolute dispersion measures used in Statistics are the range, the variance, the mean deviation, the standard deviation, the absolute Gini index, etc. (see Martínez-Panero et al. [7]).

The rest of the contribution is organized as follows. Section 2 is devoted to introduce the notation and some basic notions that are necessary for defining the two families of dispersion measures we propose in the setting of ordered qualitative scales. Section 3 present these two families of dispersion measures and includes some illustrative examples. Section 4 contains some properties. Finally, Section 5 concludes with some remarks.

2. Notation and Basic Notions

First we present the notion of ordinal proximity measure.

2.1 Ordinal Proximity Measures

Let $\mathcal{L} = \{l_1, \dots, l_g\}$ be an ordered qualitative scale arranged from the lowest to the highest linguistic terms, $l_1 < l_2 < \dots < l_g$, with $g \geq 3$.

In order to recall the notion of ordinal proximity measure on \mathcal{L} , introduced by García-Lapresta and Pérez-Román [4], we shall use a linear order $\Delta = \{\delta_1, \dots, \delta_h\}$, with $\delta_1 \succ \dots \succ \delta_h$, for representing different degrees of proximity among the terms of \mathcal{L} , being δ_1 and δ_h the maximum and minimum degrees, respectively.

It is important noticing that the elements of Δ have no meaning and they only represent different degrees of proximity.

As usual in the setting of linear orders, $\delta_r \succeq \delta_s$ means $\delta_r \succ \delta_s$ or $\delta_r = \delta_s$; $\delta_r \prec \delta_s$ means $\delta_s \succ \delta_r$; and $\delta_r \preceq \delta_s$ means $\delta_r \prec \delta_s$ or $\delta_r = \delta_s$.

Given a weak order \trianglelefteq on \mathcal{L}^n , with \triangleleft we denote the asymmetric part of \trianglelefteq , i.e., $\mathbf{x} \triangleleft \mathbf{y} \Leftrightarrow \text{not } \mathbf{y} \trianglelefteq \mathbf{x}$.

DEFINITION 1.1 ([4]) *An ordinal proximity measure on \mathcal{L} with values in Δ is a mapping $\pi: \mathcal{L}^2 \rightarrow \Delta$, where $\pi(l_r, l_s) = \pi_{rs}$ means the degree of proximity between l_r and l_s , satisfying the following conditions:*

- 1 Exhaustiveness: For every $\delta \in \Delta$, there exist $l_r, l_s \in \mathcal{L}$ such that $\delta = \pi_{rs}$.
- 2 Symmetry: $\pi_{sr} = \pi_{rs}$, for all $r, s \in \{1, \dots, g\}$.
- 3 Maximum proximity: $\pi_{rs} = \delta_1 \Leftrightarrow r = s$, for all $r, s \in \{1, \dots, g\}$.
- 4 Monotonicity: $\pi_{rs} \succ \pi_{rt}$ and $\pi_{st} \succ \pi_{rt}$, for all $r, s, t \in \{1, \dots, g\}$ such that $r < s < t$.

The previous conditions are independent (see García-Lapresta and Pérez-Román [4, Prop. 1]).

Every ordinal proximity measure $\pi : \mathcal{L}^2 \rightarrow \Delta$ can be represented by a $g \times g$ symmetric matrix with coefficients in Δ , where the elements in the main diagonal are $\pi_{rr} = \delta_1$, for every $r = 1, \dots, g$:

$$\begin{pmatrix} \pi_{11} & \cdots & \pi_{1s} & \cdots & \pi_{1g} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \pi_{r1} & \cdots & \pi_{rs} & \cdots & \pi_{rg} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \pi_{g1} & \cdots & \pi_{gs} & \cdots & \pi_{gg} \end{pmatrix} = (\pi_{rs}).$$

This matrix is called the *proximity matrix associated with π* .

Taking into account the conditions appearing in Definition 1.1, it is only necessary to show the upper half proximity matrix

$$\begin{pmatrix} \delta_1 & \pi_{12} & \pi_{13} & \cdots & \pi_{1(g-1)} & \pi_{1g} \\ & \delta_1 & \pi_{23} & \cdots & \pi_{2(g-1)} & \pi_{2g} \\ & & & \cdots & \cdots & \cdots \\ & & & & \delta_1 & \pi_{(g-1)g} \\ & & & & & \delta_1 \end{pmatrix}.$$

As shown in García-Lapresta and Pérez-Román [4, Prop. 2], the minimum proximity between linguistic terms is only reached when comparing the extreme linguistic terms: $\pi_{rs} = \delta_h \Leftrightarrow (r, s) \in \{(1, g), (g, 1)\}$.

The cardinality of Δ is located between the cardinality of \mathcal{L} and a polynomial of degree 2 of that cardinality (see García-Lapresta and Pérez-Román [4, Prop. 4]):

$$g \leq h \leq \frac{g \cdot (g-1)}{2} + 1.$$

2.2 Medians

Following García-Lapresta and Pérez-Román [5], we now introduce the median operator in the setting of ordinal degrees of proximity.

Given a vector of ordinal degrees of proximity $\delta = (\delta_1, \dots, \delta_p) \in \Delta^p$, we arrange its components in a decreasing fashion, from the highest to the lowest degrees. If p is odd, then the median of δ is unique, say $\delta_r \in \Delta$. However, if p is even, then δ has two medians, say $\delta_s, \delta_t \in \Delta$ such that $s \leq t$, i.e., $\delta_s \succeq \delta_t$. In order to unify the assignment of medians, we consider the pair of medians (δ_r, δ_r) and (δ_s, δ_t) whenever p is odd and even, respectively.

More formally, given the *set of feasible medians* $\Delta_2 = \{(\delta_r, \delta_s) \in \Delta^2 \mid r \leq s\}$, the *median operator* is the mapping

$$M : \bigcup_{p=1}^{\infty} \Delta^p \longrightarrow \Delta_2$$

that assigns the corresponding pair of medians to each vector of ordinal degrees.

For ordering the pairs of medians of ordinal proximities, consider the linear order \succeq_2 on Δ_2 defined as

$$(\delta_r, \delta_s) \succeq_2 (\delta_t, \delta_u) \Leftrightarrow \begin{cases} r + s < t + u \\ \text{or} \\ r + s = t + u \text{ and } s - r \leq u - t, \end{cases} \quad (1)$$

for all $(\delta_r, \delta_s), (\delta_t, \delta_u) \in \Delta_2$.

It is easy to see that if $r + s = t + u$, then $s - r \leq u - t \Leftrightarrow r \geq t \Leftrightarrow s \leq u$.

3. Dispersion Measures

In this section we present two families of dispersion measures in the setting of ordered qualitative scales equipped with ordinal proximity measures. The first family generalizes the most basic dispersion measure, the range.

3.1 Range-based dispersion measures

Given $n \geq 2$, let $D_R : \mathcal{L}^n \longrightarrow \Delta$ be the mapping defined as

$$D_R(\mathbf{x}) = \pi(\min \mathbf{x}, \max \mathbf{x}), \quad (2)$$

for every $\mathbf{x} \in \mathcal{L}^n$.

Based on the linear order \succeq on Δ , we introduce the weak order \trianglelefteq_R on \mathcal{L}^n defined as

$$\mathbf{x} \trianglelefteq_R \mathbf{y} \Leftrightarrow D_R(\mathbf{x}) \succeq D_R(\mathbf{y}),$$

with the meaning of the dispersion in \mathbf{x} is lower than or equal to in \mathbf{y} (with respect to D_R).

EXAMPLE 1.2 Consider the ordered qualitative scale $\mathcal{L} = \{l_1, l_2, l_3, l_4\}$ and the vectors $\mathbf{x} = (l_1, l_2, l_2, l_3)$, $\mathbf{y} = (l_3, l_3, l_4, l_4) \in \mathcal{L}^4$. We want to compare the dispersion in these vectors with respect to three different ordinal proximity measures.

1 If \mathcal{L} is equipped with the ordinal proximity measure

$$\pi : \mathcal{L}^2 \longrightarrow \Delta = \{\delta_1, \dots, \delta_7\}$$

with associated proximity matrix

$$\begin{pmatrix} \delta_1 & \delta_2 & \delta_4 & \delta_7 \\ & \delta_1 & \delta_3 & \delta_6 \\ & & \delta_1 & \delta_5 \\ & & & \delta_1 \end{pmatrix},$$

we have $D_R(\mathbf{x}) = \pi_{13} = \delta_4 \succ \delta_5 = \pi_{34} = D_R(\mathbf{y})$. Thus, $\mathbf{x} \triangleleft_R \mathbf{y}$.

2 If \mathcal{L} is equipped with the ordinal proximity measure

$$\pi : \mathcal{L}^2 \longrightarrow \Delta = \{\delta_1, \dots, \delta_7\}$$

with associated proximity matrix

$$\begin{pmatrix} \delta_1 & \delta_4 & \delta_6 & \delta_7 \\ & \delta_1 & \delta_3 & \delta_5 \\ & & \delta_1 & \delta_2 \\ & & & \delta_1 \end{pmatrix},$$

we have $D_R(\mathbf{x}) = \pi_{13} = \delta_6 \prec \delta_2 = \pi_{34} = D_R(\mathbf{y})$. Thus, $\mathbf{y} \triangleleft_R \mathbf{x}$.

3 If \mathcal{L} is equipped with the ordinal proximity measure

$$\pi : \mathcal{L}^2 \longrightarrow \Delta = \{\delta_1, \dots, \delta_4\}$$

with associated proximity matrix

$$\begin{pmatrix} \delta_1 & \delta_2 & \delta_3 & \delta_4 \\ & \delta_1 & \delta_2 & \delta_3 \\ & & \delta_1 & \delta_2 \\ & & & \delta_1 \end{pmatrix},$$

we have $D_R(\mathbf{x}) = \pi_{13} = \delta_3 \prec \delta_2 = \pi_{34} = D_R(\mathbf{y})$. Thus, $\mathbf{y} \triangleleft_R \mathbf{x}$.

3.2 Gini-based dispersion measures

We now introduce a new family of dispersion measures in the mentioned framework. It is based on the Gini index ([6]) and it is closely related to the scale estimator appearing in Shamos [9, p. 260] and Bickel and Lehmann [3, p. 38] in the setting of real numbers (see Rousseeuw and Croux [8, p. 1277]).

Given $n \geq 2$, let $D_G : \mathcal{L}^n \longrightarrow \Delta_2$ be the mapping defined as

$$D_G(\mathbf{x}) = M(\pi(x_i, x_j)_{i < j}), \quad (3)$$

for every $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{L}^n$.

Based on the linear order \succeq_2 on Δ_2 defined in (1), we introduce the weak order \triangleleft_G on \mathcal{L}^n defined as

$$\mathbf{x} \triangleleft_G \mathbf{y} \Leftrightarrow D_G(\mathbf{x}) \succeq_2 D_G(\mathbf{y}),$$

with the meaning of the dispersion in \mathbf{x} is lower than or equal to in \mathbf{y} (with respect to D_G).

Since some vectors can share the same pair of medians, it is necessary to devise a tie-breaking process for ordering the vectors. We propose to use a sequential procedure based on Balinski and Laraki [1] (see Balinski and Laraki [2] for practical examples). It consists of withdrawing the pair of medians of the vectors that are in a tie, and then selecting the new pairs of medians of the remaining proximity degrees for the corresponding vectors. The process continues until the ties are broken. It is important to note that different vectors never are in a final tie.

EXAMPLE 1.3 Consider Example 1.2 and the same three ordinal proximity measures.

1 We have

$$D_G(\mathbf{x}) = M(\pi_{12}, \pi_{12}, \pi_{13}, \pi_{22}, \pi_{23}, \pi_{23}) = M(\delta_2, \delta_2, \delta_4, \delta_1, \delta_3, \delta_3) = (\delta_2, \delta_3)$$

and

$$D_G(\mathbf{y}) = M(\pi_{33}, \pi_{34}, \pi_{34}, \pi_{34}, \pi_{34}, \pi_{44}) = M(\delta_1, \delta_5, \delta_5, \delta_5, \delta_5, \delta_1) = (\delta_5, \delta_5).$$

Since $(\delta_2, \delta_3) \succ_2 (\delta_5, \delta_5)$, we have $\mathbf{x} \triangleleft_G \mathbf{y}$.

2 We have

$$D_G(\mathbf{x}) = M(\pi_{12}, \pi_{12}, \pi_{13}, \pi_{22}, \pi_{23}, \pi_{23}) = M(\delta_4, \delta_4, \delta_6, \delta_1, \delta_3, \delta_3) = (\delta_3, \delta_4)$$

and

$$D_G(\mathbf{y}) = M(\pi_{33}, \pi_{34}, \pi_{34}, \pi_{34}, \pi_{34}, \pi_{44}) = M(\delta_1, \delta_2, \delta_2, \delta_2, \delta_2, \delta_1) = (\delta_2, \delta_2).$$

Since $(\delta_3, \delta_4) \prec_2 (\delta_2, \delta_2)$, we have $\mathbf{y} \triangleleft_G \mathbf{x}$.

3 We have

$$D_G(\mathbf{x}) = M(\pi_{12}, \pi_{12}, \pi_{13}, \pi_{22}, \pi_{23}, \pi_{23}) = M(\delta_2, \delta_2, \delta_3, \delta_1, \delta_2, \delta_2) = (\delta_2, \delta_2)$$

and

$$D_G(\mathbf{y}) = M(\pi_{33}, \pi_{34}, \pi_{34}, \pi_{34}, \pi_{34}, \pi_{44}) = M(\delta_1, \delta_2, \delta_2, \delta_2, \delta_2, \delta_1) = (\delta_2, \delta_2).$$

Consequently, in \mathbf{x} and \mathbf{y} the dispersion is the same. If we apply the tie-breaking procedure, then we have $D_G(\mathbf{x}) = M(\delta_1, \delta_2, \delta_2, \delta_3) = (\delta_2, \delta_2)$ and $D_G(\mathbf{y}) = M(\delta_1, \delta_1, \delta_2, \delta_2) = (\delta_1, \delta_2)$. Since $(\delta_2, \delta_2) \prec_2 (\delta_1, \delta_2)$, we finally have $\mathbf{y} \triangleleft_G \mathbf{x}$.

4. Properties

Let \mathcal{L} be an ordered qualitative scale equipped with an ordinal proximity measure $\pi : \mathcal{L}^2 \rightarrow \Delta$. We say that π is *totally uniform* if $\pi_{r(r+t)} = \pi_{s(s+t)}$ for all $r, s, t \in \{1, \dots, g-1\}$ such that $r+t \leq g$ and $s+t \leq g$.

Let $N : \mathcal{L} \rightarrow \mathcal{L}$ be the *negation operator* defined as $N(l_r) = l_{g+1-r}$, for every $r \in \{1, \dots, g\}$.

Given $k \in \{1-g, \dots, g-1\}$, let $T_k : \mathcal{L} \rightarrow \mathcal{L}$ be the *translation operator* defined as $T_k(l_r) = l_{r+k}$, for every $r \in \{1, \dots, g\}$ such that $r+k \leq g$.

In the following proposition we establish some properties of the mappings introduced in (2) and (3). They are related to the ones considered in Martínez-Panero et al. [7] in a quantitative context.

PROPOSITION 1.1 *Consider the mappings $D_R : \mathcal{L}^n \rightarrow \Delta$ and $D_G : \mathcal{L}^n \rightarrow \Delta_2$ defined in (2) and (3), respectively, and its extensions $\tilde{D}_R : \bigcup_{n=2}^{\infty} \mathcal{L}^n \rightarrow \Delta$ and $\tilde{D}_G : \bigcup_{n=2}^{\infty} \mathcal{L}^n \rightarrow \Delta_2$. The following properties hold:*

- 1 **Symmetry:** $D_R(x_{\sigma(1)}, \dots, x_{\sigma(n)}) = D_R(x_1, \dots, x_n)$ and $D_G(x_{\sigma(1)}, \dots, x_{\sigma(n)}) = D_G(x_1, \dots, x_n)$, for every permutation $\sigma : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ and every $(x_1, \dots, x_n) \in \mathcal{L}^n$.
- 2 **Invariance for replications:** $\tilde{D}_R(\overbrace{\mathbf{x}, \dots, \mathbf{x}}^m) = D_R(\mathbf{x})$ and $\tilde{D}_G(\overbrace{\mathbf{x}, \dots, \mathbf{x}}^m) = D_G(\mathbf{x})$, for every $\mathbf{x} \in \mathcal{L}^n$ and any number $m \in \mathbb{N}$ of replications of \mathbf{x} .
- 3 **Minimum dispersion:** $D_R(x_1, \dots, x_n) = \delta_1 \Leftrightarrow x_1 = \dots = x_n$; and $D_G(l_r, \dots, l_r) = (\delta_1, \delta_1)$, for every $l_r \in \mathcal{L}$.
- 4 **Anti-self-duality:** if π is totally uniform, then $D_R(N(x_1), \dots, N(x_n)) = D_R(x_1, \dots, x_n)$ and $D_G(N(x_1), \dots, N(x_n)) = D_G(x_1, \dots, x_n)$, for every $(x_1, \dots, x_n) \in \mathcal{L}^n$.
- 5 **Invariance for translations:** if π is totally uniform, then $D_R(T_k(x_1), \dots, T_k(x_n)) = D_R(x_1, \dots, x_n)$ and $D_G(T_k(x_1), \dots, T_k(x_n)) = D_G(x_1, \dots, x_n)$, for every $(x_1, \dots, x_n) \in \mathcal{L}^n$ and every $k \in \{1, \dots, g-1\}$ such that $(T_k(x_1), \dots, T_k(x_n)) \in \mathcal{L}^n$.

5. Concluding Remarks

Acknowledgments

This contribution is dedicated to the memory of Pedro Gil. The first author gratefully acknowledges the funding support of the Spanish *Ministerio de Economía y Competitividad* (project ECO2016-77900-P) and ERDF.

References

- [1] Balinski, M., Laraki, R. (2007). A theory of measuring, electing and ranking. *Proceedings of the National Academy of Sciences of the United States of America* 104, pp. 8720–8725.
- [2] Balinski, M., Laraki, R. (2013). How best to rank wines: Majority Judgment. In: *Wine Economics: Quantitative Studies and Empirical Observations*, pp. 149–172, Palgrave-MacMillan.
- [3] Bickel, P.J., Lehmann, E.L. (1979). Descriptive Statistics for nonparametric models IV: Spread. In: J. Jurečková (ed.) *Contributions to Statistics, Hájek Memorial Volume*, pp. 33–40, Academia, Prague.
- [4] García-Lapresta, J.L., Pérez-Román, D. (2015). Ordinal proximity measures in the context of unbalanced qualitative scales and some applications to consensus and clustering. *Applied Soft Computing* 35, pp. 864–872.
- [5] García-Lapresta, J.L., Pérez-Román, D. (2017). A consensus reaching process in the context of non-uniform ordered qualitative scales. *Fuzzy Optimization and Decision Making*. Forthcoming, DOI: 10.1007/s10700-016-9256-6.
- [6] Gini, C. (1912). *Variabilità e Mutabilità*, Tipografia di Paolo Cuppini, Bologna.
- [7] Martínez-Panero, M., García-Lapresta, J.L., Meneses, L.C. (2016). Multidistances and dispersion measures. In: T. Calvo Sánchez, J. Torrens Sastre (eds.) *Fuzzy Logic and Information Fusion*, Studies in Fuzziness and Soft Computing, pp. 123–134, Springer.
- [8] Rousseeuw, P.J., Croux, C. (1993). Alternatives to the median absolute deviation. *Journal of the American Statistical Association* 88, pp. 1273–1283.
- [9] Shamos, M.I. (1976). Geometry and Statistics: Problems at the interface. In: J.F. Traub (ed.) *New Directions and Recent Results in Algorithms and Complexity*, pp. 251–280, New York Academic Press.