# A new approach to the measurement of polarization for grouped data

Eckart BOMSDORF
*University of Cologne*

Clemens A. OTTO
*Singapore Management University*, clemensotto@smu.edu.sg

# A new approach to the measurement of polarization for grouped data

**Eckart Bomsdorf · Clemens Otto**

**Summary**  In this paper we develop a measure of polarization for discrete distributions of non-negative grouped data. The measure takes into account the relative sizes and homogeneities of individual groups as well as the heterogeneities between all pairs of groups. It is based on the assumption that the total polarization within the distribution can be understood as a function of the polarizations between all pairs of groups.

The measure allows information on existing groups within a population to be used directly to determine the degree of polarization. Thus the impact of various classifications on the degree of polarization can be analysed.

The treatment of the distribution's total polarization as a function of pairwise polarizations allows statements concerning the effect of an individual pair or an individual group on the total polarization.

## 1 Introduction

Polarization describes a certain characteristic of a distribution. A distribution is polarized if it is composed of few but large groups that are highly homogeneous themselves, but highly heterogeneous with respect to each other. Thus, the extent of polarization is determined by three factors: the size and number of groups, homogeneity within these groups, and heterogeneity between the groups.

E. Bomsdorf (✉) · C. Otto (✉)

Department of Economic and Social Statistics, University of Cologne, Albertus-Magnus-Platz, 50923 Cologne, Germany

e-mail: bomsdorf@wiso.uni-koeln.de
e-mail: cotto.phd2007@london.edu

In its beginnings, the measuring of polarization was highly influenced by the analysis of income distributions. This analysis was mainly driven by the interest in the phenomenon of a 'disappearing middle class'. A population's receding middle class was associated with a rising potential of social tension or even public disturbances caused by the gap between a poor lower class and a rich upper class.

Esteban and Ray (1994) and Wolfson (1994) pointed out the unsuitability of measures of inequality for the analysis of polarized distributions[1] and presented first approaches to the measurement of polarization. Their work was of fundamental impact and credit for pioneering insights into the field of polarization measurement clearly belongs to them. In the course of the following ten years, new approaches to the measurement of polarization or derivatives of already existing measures were published.[2] These concepts can easily be applied to all distributions of non-negative data. Nevertheless, the analysis of income distributions remains the main field of application for the measurement of polarization.

The existing measures of polarization can be classified according to two criteria: the application to grouped or ungrouped data and the application to continuous or discrete distributions. Measures that are based on the assumption of continuous distributions appeal through their elegance, but as empirical income distributions tend to be discrete, measures that can be applied to discrete functions are needed for practical purposes. Comparing the measurement of polarization for grouped and ungrouped data, we note the following: In practice, criteria to group income data in a meaningful way are often lacking. In those cases, measures that can be applied to ungrouped data are needed. On the other hand, if criteria that allow a grouping exist, this additional information can be applied to the measurement of polarization with measures for grouped data. This allows the analysis of polarization with respect to explicit grouping according to a certain criterion as well as the comparison of various polarizations due to different groupings.

Unfortunately, many of the existing measures of polarization are flawed, as is demonstrated, for example, by Schmidt (2004)[3].

For a population that is made up of only two groups, the measure proposed by Esteban and Ray (1994) is defined as

$$P^{ER} = \left[\pi^{1+\alpha} \cdot (1-\pi) + \pi \cdot (1-\pi)^{1+\alpha}\right] \cdot (y-x) \ ,$$

with relative group sizes $\pi$ and $(1-\pi)$, incomes $y$ and $x$, and $a \in [1; 1,6]$. If $y$ and $x$ denote the natural logarithm of the respective income, only positive income is per-

---

[1] For example, the Pigou–Dalton axiom known in inequality measurement does not hold true for the measurement of polarization.

[2] These include, among others, D'Ambrosio (2001), Duclos et al. (2004), Esteban et al. (1999), Gradín (2000), Schmidt (2004), Wang and Tsui (2000), and Zhang and Kanbur (2001).

[3] In the subsequent discussion of the various measures of polarization we will follow the approach taken by Schmidt (2004) and define maximum polarization as the state in which the total population consists of two groups of equal size, where all individuals in one group have no income and all others equally share the total income.

Maximum inequality shall be defined as the state in which the total population consists of two groups with one individual representing the first group and possessing the total income while all other individuals representing the second group have no income, i.e. the state in which the Gini coefficient is maximized.

mitted and the measure is not defined for the cases of maximum polarization and maximum inequality. If $y$ and $x$ denote the incomes after they have been normalized by the mean income, we calculate for $\pi = 0, 5$ and $y = 2$ and $x = 0$ (i.e. for the case of maximum polarization)

$$P^{ER} = \left(\frac{1}{2}\right)^{\alpha}.$$

In the case of maximum inequality $\left(\pi = \frac{n-1}{n}; \ y = 0; \ x = n\right)$, we calculate

$$P^{ER} = \left(\frac{n-1}{n}\right)\left[\left(\frac{n-1}{n}\right)^{a} + \left(\frac{1}{n}\right)^{\alpha}\right]$$

and

$$\lim_{n\to\infty} P^{ER} = 1.$$

The measure proposed by Esteban et al. (1999) is defined as

$$P^{EGR} = P^{ER} - \beta\left[G(f) - G(\rho^*)\right],$$

with $\beta \geq 0$. $P^{ER}$ denotes the measure proposed by Esteban and Ray (1994), $G(f)$ is the Gini coefficient of the ungrouped distribution, and $G(\rho^*)$ is the Gini coefficient of the grouped distribution assuming that all individuals in one group earn the same income. Thus, in the case of maximum polarization, we calculate $G(f) = G(\rho^*) = \frac{1}{2}$ and

$$P^{EGR} = P^{ER} = \left(\frac{1}{2}\right)^{\alpha}.$$

For maximum inequality, we calculate with $G(f) = G(\rho^*) = 1 - \frac{1}{n}$

$$P^{EGR} = P^{ER} = \left(\frac{n-1}{n}\right)\left[\left(\frac{n-1}{n}\right)^{a} + \left(\frac{1}{n}\right)^{\alpha}\right]$$

and

$$\lim_{n\to\infty} P^{EGR} = 1.$$

The measure proposed by Gradín (2000) is defined as

$$P^{G} = P^{ER} - \beta\left[G(f) - G(\rho^c) - 1\right],$$

with $\beta \geq 0$. As in the measure proposed by Esteban et al. (1999), $P^{ER}$ denotes the measure proposed by Esteban and Ray (1994), $G(f)$ is the Gini coefficient of the ungrouped distribution, and $G(\rho^c)$ is the Gini coefficient of the grouped distribution assuming that all individuals in one group earn the same income. Therefore, in the case of maximum polarization we calculate

$$P^{G} = \left(\frac{1}{2}\right)^{\alpha} + \beta$$

and for maximum inequality we have

$$P^G = \left(\frac{n-1}{n}\right)\left[\left(\frac{n-1}{n}\right)^a + \left(\frac{1}{n}\right)^\alpha\right] + \beta$$

and

$$\lim_{n\to\infty} P^G = 1 + \beta.$$

Thus, the measures proposed by Esteban and Ray (1994), by Esteban et al. (1999), and by Gradín (2000) do not reach their maximum values in the case of maximum polarization but in the case of maximum inequality.

The measure proposed by Wolfson (1994) is defined as

$$P^W = \frac{2\left[1 - 2L\left(0, 5\right) - G\right]}{\frac{m}{\mu}},$$

where $L(0, 5)$ is the value of the Lorenz curve at the 50th percentile, $G$ the Gini coefficient, $m$ the median, and $\mu$ the mean income. Therefore, the measure is not defined for the case of maximum inequality ($m = 0$) and can equal infinitely high values in the case of high income inequality.

The measure proposed by Wang and Tsui (2000) is defined as

$$P^{WT} = \left(\frac{\theta}{N}\right) \sum_{j=1}^{N} \left|\frac{x_j - m}{m}\right|^r,$$

with $\theta > 0$ and $0 < r < 1$. $N$ is the number of individuals in the population, $x_j$ the income of individual $j$ for $j = 1, \ldots, N$, and $m$ is the median income. Similar to the measure proposed by Wolfson (1994), the measure proposed by Wang and Tsui (2000) is not defined for the case of maximum inequality ($m = 0$) and can equal infinitely high values in the case of high income inequality.

The measure proposed by Zhang and Kanbur (2001) is defined as

$$P^{ZK} = \frac{\text{"between group inequality"}}{\text{"within group inequality"}}$$

and is therefore neither defined in the cases of maximum polarization and inequality nor in that of minimum polarization (*"within group inequality"* $= 0$) and can equal infinitely high values if the inequality within the groups is very low.

The measure proposed by D'Ambrosio (2001) is defined as

$$P^D = \sum_{i=1}^{N} \sum_{j=1}^{N} \pi_i^{1+\alpha} \pi_j Kov_{ij}.$$

$\pi_i$ and $\pi_j$ are the relative sizes of groups $i$ for $i = 1, \ldots, N$ and $j$ for $j = 1, \ldots, N$, and $Kov_{ij}$ is the Kolmogorov measure of variation distance between the density functions of the income distributions in groups $i$ and $j$, which are derived using a kernel density estimator. Thus, the measure is not population-invariant for discrete distributions, as the density estimation depends on the number of observations. Furthermore,

the measure does not react to changes of the heterogeneity or of the homogeneity within the groups if the income distributions of the groups do not overlap.

The measure proposed by Duclos et al. (2004) is defined as

$$P^{DER} = \int_y f(y)^\alpha a(y) dF(y) .$$

As a kernel density estimator is applied to derive the density function $f(y)$ of the income distribution that is examined, the measure is not population-invariant if applied to discrete distributions.

While the measure proposed by Schmidt (2004) seems to perform well in the examined cases, it cannot be applied to grouped data. Thus, in order to determine the polarization of discrete distributions of grouped data, a new measure is needed.

Accordingly, we will develop a new measure of polarization for grouped data in the following paper. The grouping can be based on any criteria that appear to be of interest. Examples for possible criteria could be gender, age, ethnic origin or the income itself.

For a distribution of non-negative data, the measure assigns a value describing the extent of the distribution's polarization. Polarization itself depends on the number and the size of the groups within the distribution, the homogeneity within each group, and the heterogeneity between the groups. First, we examine the homogeneity within a given group. Subsequently, we analyse the heterogeneity between two groups and after this we measure the polarization between the two groups. Finally, we introduce a new approach to measure the polarization of a distribution composed of more than two groups.

For greater convenience, the following passages will examine polarization regarding income distributions. Only non-negative income is allowed.

## 2 Homogeneity within a group

The more similar the incomes of different people in a group, the higher the homogeneity within the group. Let $x_1, x_2, \ldots, x_n$ be the income levels of the different individuals in the group $\left( x_k \geq 0 \text{ for } k = 1, 2, \ldots, n \text{ and } \overline{x} = \frac{1}{n} \sum_{k=1}^{n} x_k \right)$ and let $Hom = Hom(x_1, x_2, \ldots, x_n)$ be the measure of homogeneity for that group. The following requirements appear to be plausible:

1. $0 < Hom(x_1, x_2, \ldots, x_n) \leq 1$
   The measure ranges between 0 and 1.
2. $Hom(x_1, x_2, \ldots, x_n) = 1$ if $x_k = \overline{x}$ for $k = 1, 2, \ldots, n$
   The measure equals 1 if all individuals have the same income.
3. $Hom(x_1, x_2, \ldots, x_k, x_l, \ldots, x_n) > Hom(x_1, x_2, \ldots, x_k - \varepsilon, x_l + \varepsilon, \ldots, x_n)$
   for $x_k < \overline{x} < x_l$ and $0 < \varepsilon < x_k$
   $Hom(x_1, x_2, \ldots, x_k, x_l, \ldots, x_n) < Hom(x_1, x_2, \ldots, x_k + \varepsilon, x_l - \varepsilon, \ldots, x_n)$
   for $x_k < \overline{x} < x_l$ and $0 < \varepsilon < x_l - \overline{x}$

Homogeneity decreases (increases) if the differences between the individual incomes and the average income increase (decrease).

4. $Hom\,(x_1, x_2, \ldots, x_n) = Hom\,(a \cdot x_1, a \cdot x_2, \ldots, a \cdot x_n)$ for $a \in \mathbb{R}^+$

The measure is scale-invariant.

5. $Hom\,(x_1, x_2, \ldots, x_n) = Hom\,(\underbrace{x_1, \ldots, x_1}_{b\text{-times}}, \underbrace{x_2, \ldots, x_2}_{b\text{-times}}, \ldots, \underbrace{x_n, \ldots, x_n}_{b\text{-times}})$ for $b \in \mathbb{N}$

The measure is population-invariant.

In the following section we develop a measure of homogeneity for the distribution of income within a given group that satisfies the above requirements. Let us begin to determine homogeneity by looking at the differences between the incomes. Differences between the incomes mean that the individual incomes cannot all be equal to the average income. A measure for the extent of those differences is the average deviation $U$ of the incomes $x_k$ for $k = 1, 2, \ldots, n$ from the group average $\overline{x}$.

We have

$$U = \frac{1}{n} \sum_{k=1}^{n} |x_k - \overline{x}| \ .$$

By dividing $U$ by $\overline{x}$ we derive the relative measure $U^*$ with

$$U^* = \begin{cases} \frac{1}{\overline{x}} \cdot \frac{1}{n} \sum_{k=1}^{n} |x_k - \overline{x}| & \text{if} \ \ \overline{x} \neq 0 \\[2mm] 0 & \text{if} \ \ \overline{x} = 0 \end{cases}.$$

$U^*$ increases as the differences between the incomes increase. It is at its maximum if inequality among the incomes is at a maximum, i.e. if for $n \geq 2$ one individual has income $n \cdot \overline{x} > 0$ (w.l.o.g. $x_n = n \cdot \overline{x}$) and all other individuals have no income. We then have

$$U^* = \frac{1}{\overline{x}} \cdot \frac{1}{n} \sum_{k=1}^{n} |x_k - \overline{x}|$$

$$= \frac{n}{x_n} \cdot \frac{1}{n} \left[ \sum_{k=1}^{n-1} \left| 0 - \frac{x_n}{n} \right| + \left| x_n - \frac{x_n}{n} \right| \right]$$

$$= \frac{1}{x_n} \left[ (n-1)\frac{x_n}{n} + \frac{(n-1)\,x_n}{n} \right]$$

$$= 2\frac{n-1}{n} \ .$$

It follows that

$$\lim_{n \to \infty} U^* = 2 \ .$$

$U^*$ equals 0 for maximum homogeneity, and it equals 2 in the case of minimum homogeneity.

We derive a measure of homogeneity that equals 0 for minimum homogeneity and 2 for maximum homogeneity by looking at $2 - U^*$ instead of $U^*$. By dividing this by 2 we have a measure of homogeneity *Hom* for the distribution of income within a group that ranges between 0 and 1 with

$$Hom = \frac{2 - U^*}{2}$$

$$= \begin{cases} 1 - \dfrac{\frac{1}{\overline{x}} \cdot \frac{1}{n} \sum\limits_{k=1}^{n} |x_k - \overline{x}|}{2} & \text{if } \overline{x} \neq 0 \\[2mm] 1 & \text{if } \overline{x} = 0 \end{cases}.$$

*Hom* satisfies all the requirements postulated previously. It ranges between 0 and 1 and equals 1 only if all incomes are equal. The measure decreases if the differences between the individual incomes and the average income increase, i.e. if homogeneity decreases. Furthermore, the measure is invariant regarding population size and scale.[4]

For a population of $L$ groups, let $Hom_i = Hom\left(x_{i1}, x_{i2}, \ldots, x_{in_i}\right)$ be the measure of homogeneity for group $i$ ($i = 1, 2, \ldots, L$) with incomes $x_{i1}, x_{i2}, \ldots, x_{in_i}$.

## 3 Heterogeneity between two groups

Greater differences in the income levels of individuals in two groups $i$ and $j$ ($i, j = 1, 2, \ldots, L$ and $i \neq j$) result in greater heterogeneity between these two groups. Let $x_{i1}, x_{i2}, \ldots, x_{in_i}$ and $x_{j1}, x_{j2}, \ldots, x_{jn_j}$ be the incomes of the $n_i$ and $n_j$ individuals in the groups and let $Het_{i,j} = Het\left(x_{i1}, x_{i2}, \ldots, x_{in_i}; x_{j1}, x_{j2}, \ldots, x_{jn_j}\right)$ be the measure of heterogeneity for these two groups. The following requirements appear to be plausible:

1. $0 \leq Het_{i,j} \leq 1$
   The measure ranges between 0 and 1.
2. $Het\left(x_{i1}, x_{i2}, \ldots, x_{in_i}; x_{j1}, x_{j2}, \ldots, x_{jn_j}\right) = 0$ if $x_{i1} = x_{i2} = \ldots = x_{in_i} = x_{j1} = x_{j2} = \ldots = x_{jn_j}$
   Heterogeneity is 0 if all individuals have the same income.

---

[4] Potential alternatives to the measure *Hom* are the measure *Hom** that resembles Gini's mean difference with

$$Hom^* = 1 - \frac{\frac{1}{n^2} \sum\limits_{k=1}^{n} \sum\limits_{l=1}^{n} |x_k - x_l|}{\max_k (x_k) - \min_k (x_k)}$$

and the measure $Hom^G$ with

$$Hom^G = 1 - D_G ,$$

where $D_G$ is the Gini coefficient.

However, as the measure *Hom* focuses on the deviation of the incomes from the group average, i.e. on the deviation from focal point of the income distribution within that group, it appears to be more in line with the concept of identification as introduced by Esteban and Ray (1994).

3. Heterogeneity increases if the differences between the different groups' incomes increase.

4. $Het\left(x_{i1}, x_{i2}, \ldots, x_{in_i}; x_{j1}, x_{j2}, \ldots, x_{jn_j}\right) =$
$Het\left(a \cdot x_{i1}, a \cdot x_{i2}, \ldots, a \cdot x_{in_i}; a \cdot x_{j1}, a \cdot x_{j2}, \ldots, a \cdot x_{jn_j}\right)$ for $a \in \mathbb{R}^+$
The measure is scale-invariant.

5. $Het\left(x_{i1}, x_{i2}, \ldots, x_{in_i}; x_{j1}, x_{j2}, \ldots, x_{jn_j}\right) =$

$$Het\left( \underbrace{x_{i1}, \ldots, x_{i1}}_{b\text{-times}}, \underbrace{x_{i2}, \ldots, x_{i2}}_{b\text{-times}}, \ldots, \underbrace{x_{in_i}, \ldots, x_{in_i}}_{b\text{-times}}; \right.$$

$$\left. \underbrace{x_{j1}, \ldots, x_{j1}}_{b\text{-times}}, \underbrace{x_{j2}, \ldots, x_{j2}}_{b\text{-times}}, \ldots, \underbrace{x_{jn_j}, \ldots, x_{jn_j}}_{b\text{-times}} \right) \text{ for } b \in \mathbb{N}$$

The measure is population-invariant.

In order to determine the differences between the two groups' incomes, we start by comparing the average incomes of the two groups. For groups $i$ and $j$, the difference between the two averages $\overline{x}_i$ and $\overline{x}_j$ is

$$\left|\overline{x}_i - \overline{x}_j\right| .$$

By dividing this difference by the sum of the averages, we derive a measure $A_{i,j}^{\overline{x}}$ that ranges between 0 and 1 with[5]

$$A_{i,j}^{\overline{x}} = \frac{\left|\overline{x}_i - \overline{x}_j\right|}{\overline{x}_i + \overline{x}_j} .$$

$A_{i,j}^{\overline{x}}$ can be used directly to measure the extent of heterogeneity between the two groups. It satisfies all the requirements postulated above. It ranges between 0 and 1 and equals 0 if all incomes are equal. The measure increases as the difference between the two averages increases. Furthermore, the measure is invariant regarding scale and population size.[6]

This measure of heterogeneity takes into account only the average incomes of the two groups; further characteristics of the distribution of income have no effect on the measure. A possible extension of the measure is to include not only the group averages but also the minimum and maximum incomes in the groups. In this case, heterogeneity between the groups $i$ and $j$ increases if the gap between the two minima or the between the two maxima increases ceteris paribus.

$$\left| \max_{k=1,\ldots,n_i} (x_{ik}) - \max_{k=1,\ldots,n_j} \left(x_{jk}\right) \right|$$

and

$$\left| \min_{k=1,\ldots,n_i} (x_{ik}) - \min_{k=1,\ldots,n_j} \left(x_{jk}\right) \right|$$

---

[5] Let $A_{i,j}^{\overline{x}} = 0$ for $\overline{x}_i = \overline{x}_j = 0$.

[6] Alternatively, we could use $\max\left(\overline{x}_i, \overline{x}_j\right)$ in the denominator.

are the differences between the corresponding extrema.

Measures ranging between 0 and 1 for these differences are

$$A_{i,j}^{\max} = \frac{\left|\max_{k=1,\dots,n_i}(x_{ik}) - \max_{k=1,\dots,n_j}(x_{jk})\right|}{\max_{k=1,\dots,n_i}(x_{ik}) + \max_{k=1,\dots,n_j}(x_{jk})}$$

and

$$A_{i,j}^{\min} = \frac{\left|\min_{k=1,\dots,n_i}(x_{ik}) - \min_{k=1,\dots,n_j}(x_{jk})\right|}{\min_{k=1,\dots,n_i}(x_{ik}) + \min_{k=1,\dots,n_j}(x_{jk})}.$$

Just like $A_{i,j}^{\bar{x}}$, these measures can be used directly as measures of heterogeneity. We can also combine all three measures, which leads to the measure of heterogeneity $Het_{i,j}$:

$$Het_{i,j} = \gamma_1 \cdot A_{i,j}^{\bar{x}} + \gamma_2 \cdot A_{i,j}^{\max} + \gamma_3 \cdot A_{i,j}^{\min}$$

with

$$0 \leq \gamma_r \leq 1 \quad \text{and} \quad \sum_{r=1}^{3} \gamma_r = 1 \quad \text{for} \quad r = 1, 2, 3.$$

This general measure of heterogeneity[7] also satisfies all the requirements postulated in the beginning.[8]

## 4 Polarization between two groups

Polarization between two groups $i$ and $j$ depends on the size of the groups ($n_i$ and $n_j$), the homogeneities within the groups ($Hom_i$ and $Hom_j$), and the heterogeneity between the groups ($Het_{i,j}$). Higher homogeneity, higher heterogeneity and more similar group sizes all cause polarization to increase. Let $P_{i,j} = P(n_i, n_j, Hom_i, Hom_j, Het_{i,j})$ be the measure of polarization between the groups $i$ and $j$. The following requirements appear to be plausible:

1. $0 \leq P_{i,j} \leq 1$
   The measure ranges between 0 and 1.
2. $P_{i,j} = 0$ if $Het_{i,j} = 0$
   The measure equals 0 if heterogeneity between the two groups is 0.
3. $P_{i,j} = 1$ if $n_i = n_j$ and $Hom_i = Hom_j = Het_{i,j} = 1$

---

[7] Alternatively to the measure proposed, we could also use the measure $Het_{i,j}^*$ that resembles Gini's mean difference:

$$Het_{i,j}^* = 1 - \frac{\frac{1}{n_i \cdot n_j} \sum_{k=1}^{n_i} \sum_{l=1}^{n_j} |x_k - x_l|}{\max\left(\max_k(x_k) - \min_l(x_l) ; \max_l(x_l) - \min_k(x_k)\right)}.$$

[8] For practical purposes it is plausible to choose $\gamma_1 \geq \gamma_2 + \gamma_3$, so that the impact of the differences between the extrema does not exceed the impact of the distance between the groups' averages.

The measure equals 1 if both groups are of equal size and homogeneity within both groups and heterogeneity between the groups equal 1.

4. $\frac{\partial P\left(n_i,n_j,Hom_i,Hom_j,Het_{i,j}\right)}{\partial Hom_i} > 0$ and $\frac{\partial P\left(n_i,n_j,Hom_i,Hom_j,Het_{i,j}\right)}{\partial Hom_j} > 0$

if $Het_{i,j} \neq 0$

The measure increases if ceteris paribus homogeneity in one group rises, unless heterogeneity between the groups is 0.

5. $\frac{\partial P\left(n_i,n_j,Hom_i,Hom_j,Het_{i,j}\right)}{\partial Het_{i,j}} > 0$

The measure increases if ceteris paribus heterogeneity between the groups increases.

6. $P\left(n_i, n_j, Hom_i, Hom_j, Het_{i,j}\right) < P\left(n_i - \Delta, n_j + \Delta, Hom_i, Hom_j, Het_{i,j}\right)$

$$\text{for } n_i > n_j,\ 0 < \Delta \leq \tfrac{n_i-n_j}{2},\ \Delta \in \mathbb{N}$$

$P\left(n_i, n_j, Hom_i, Hom_j, Het_{i,j}\right) < P\left(n_i + \Delta, n_j - \Delta, Hom_i, Hom_j, Het_{i,j}\right)$

$$\text{for } n_i < n_j,\ 0 < \Delta \leq \tfrac{n_j-n_i}{2},\ \Delta \in \mathbb{N}$$

The measure increases if ceteris paribus the groups' sizes become more equal.

7. $P\left(n_i, n_j, Hom_i, Hom_j, Het_{i,j}\right) = P\left(n_i, n_j, Hom_i^*, Hom_j^*, Het_{i,j}^*\right)$

with $Hom_i^* = Hom\left(a \cdot x_{i1}, a \cdot x_{i2}, \ldots, a \cdot x_{in_i}\right)$,

$Hom_j^* = Hom\left(a \cdot x_{j1}, a \cdot x_{j2}, \ldots, a \cdot x_{jn_j}\right)$ and

$Het_{i,j}^* = Het\left(a \cdot x_{i1}, a \cdot x_{i2}, \ldots, a \cdot x_{in_i}; a \cdot x_{j1}, a \cdot x_{j2}, \ldots, a \cdot x_{jn_j}\right)$

for $a \in \mathbb{R}^+$

The measure is scale-invariant.

8. $P\left(n_i, n_j, Hom_i, Hom_j, Het_{i,j}\right) = P\left(b \cdot n_i, b \cdot n_j, Hom_i^{\sim}, Hom_j^{\sim}, Het_{i,j}^{\sim}\right)$

with $Hom_i^{\sim} = Hom(\underbrace{x_{i1}, \ldots, x_{i1}}_{b\text{-times}}, \underbrace{x_{i2}, \ldots, x_{i2}}_{b\text{-times}}, \ldots, \underbrace{x_{in_i}, \ldots, x_{in_i}}_{b\text{-times}})$,

$Hom_j^{\sim} = Hom(\underbrace{x_{j1}, \ldots, x_{j1}}_{b\text{-times}}, \underbrace{x_{j2}, \ldots, x_{j2}}_{b\text{-times}}, \ldots, \underbrace{x_{jn_j}, \ldots, x_{jn_j}}_{b\text{-times}})$ and

$$Het_{i,j}^{\sim} = Het \left( \underbrace{x_{i1}, \ldots, x_{i1}}_{b\text{-times}}, \underbrace{x_{i2}, \ldots, x_{i2}}_{b\text{-times}}, \ldots, \underbrace{x_{in_i}, \ldots, x_{in_i}}_{b\text{-times}}; \right.$$

$$\left. \underbrace{x_{j1}, \ldots, x_{j1}}_{b\text{-times}}, \underbrace{x_{j2}, \ldots, x_{j2}}_{b\text{-times}}, \ldots, \underbrace{x_{jn_j}, \ldots, x_{jn_j}}_{b\text{-times}} \right) \text{ for } b \in \mathbb{N}$$

The measure is population-invariant.

It is immediately evident that

$$n_i \cdot n_j \cdot Hom_i \cdot Hom_j \cdot Het_{i,j}$$

satisfies all but the first and third requirements, as it can equal values greater than 1. Using relative group sizes instead of absolute sizes we have a measure $P_{i,j}^*$ with

$$P_{i,j}^* = \frac{n_i}{n_i + n_j} \cdot \frac{n_j}{n_i + n_j} \cdot Hom_i \cdot Hom_j \cdot Het_{i,j}.$$

This measure additionally satisfies the first requirement, but does not equal 1 for maximum polarization. In that case, we have

$$P_{i,j}^* = \frac{1}{2} \cdot \frac{1}{2} \cdot 1 \cdot 1 \cdot 1$$
$$= \frac{1}{4}.$$

We derive a measure $P_{i,j}$ ranging between 0 and 1 with

$$P_{i,j} = 4 \cdot P_{i,j}^*$$
$$= 4 \cdot \frac{n_i}{n_i + n_j} \cdot \frac{n_j}{n_i + n_j} \cdot Hom_i \cdot Hom_j \cdot Het_{i,j}.$$

$P_{i,j}$ satisfies all the requirements postulated above. Thus, $P_{i,j}$ is a measure of polarization between the groups $i$ and $j$.[9] What we lack is a measure of polarization for a population of more than two groups.

## 5  Polarization in the entire population

If a population is made up of $N$ individuals and $L$ groups of size $n_i$ for $i = 1, \ldots, L$, polarization in the entire population can be understood as a function of the polarizations between all pairs of groups. Let $P = P\left(P_{1,2}, P_{1,3}, \ldots, P_{L-1,L}\right)$ be the measure of polarization for the entire population. The following requirements appear to be plausible:

1. $0 \le P \le 1$
   The measure ranges between 0 and 1.
2. $P = 0$ if $P_{i,j} = 0$ for $i = 1, 2, \ldots, L$ and $j = 1, 2, \ldots, L$
   The measure equals 0 if the polarizations between all pairs of groups are at their minimums.
3. $P = 1$ if $L = 2$ and $P_{1,2} = 1$
   The measure equals 1 if the population is composed of only two groups and polarization between the groups is at its maximum.
4. $P\left(P_{1,2}, P_{1,3}, \ldots, P_{L-1,L}\right) = P\left(P_{1,2}^a, P_{1,3}^a, \ldots, P_{L-1,L}^a\right)$
   with $P_{i,j}^a = P\left(n_i, n_j, Hom_i^a, Hom_j^a, Het_{i,j}^a\right)$
   with $Hom_i^a = Hom\left(a \cdot x_{i1}, a \cdot x_{i2}, \ldots, a \cdot x_{in_i}\right)$,
   $Hom_j^a = Hom\left(a \cdot x_{i1}, a \cdot x_{i2}, \ldots, a \cdot x_{in_i}\right)$ and

---

[9] Alternative ways to take homogeneity into account could be

$$P_{i,j} = 4 \cdot \frac{n_i}{n_i + n_j} \cdot \frac{n_j}{n_i + n_j} \cdot \frac{n_i \cdot Hom_i + n_j \cdot Hom_j}{n_i + n_j} \cdot Het_{i,j}$$

and

$$P_{i,j} = 4 \cdot \frac{n_i}{n_i + n_j} \cdot \frac{n_j}{n_i + n_j} \cdot \left(Hom_i^{n_i} \cdot Hom_j^{n_j}\right)^{\frac{1}{n_i + n_j}} \cdot Het_{i,j}.$$

$$Het^a_{i,j} = Het\left(a \cdot x_{i1}, a \cdot x_{i2}, \ldots, a \cdot x_{in_i}; a \cdot x_{j1}, a \cdot x_{j2}, \ldots, a \cdot x_{jn_j}\right)$$

for $a \in \mathbb{R}^+$ and $i = 1, 2, \ldots, L$ and $j = 1, 2, \ldots, L$

The measure is scale-invariant.

5. $P\left(P_{1,2}, P_{1,3}, \ldots, P_{L-1,L}\right) = P\left(P^b_{1,2}, P^b_{1,3}, \ldots, P^b_{L-1,L}\right)$

with $P^b_{i,j} = P\left(b \cdot n_i, b \cdot n_j, Hom^b_i, Hom^b_j, Het^b_{i,j}\right)$

with $Hom^b_i = Hom\,(\underbrace{x_{i1}, \ldots, x_{i1}}_{b\text{-times}}, \underbrace{x_{i2}, \ldots, x_{i2}}_{b\text{-times}}, \ldots, \underbrace{x_{in_i}, \ldots, x_{in_i}}_{b\text{-times}})$,

$Hom^b_j = Hom\,(\underbrace{x_{j1}, \ldots, x_{j1}}_{b\text{-times}}, \underbrace{x_{j2}, \ldots, x_{j2}}_{b\text{-times}}, \ldots, \underbrace{x_{jn_j}, \ldots, x_{jn_j}}_{b\text{-times}})$ and

$$Het^b_{i,j} = Het\left(\underbrace{x_{i1}, \ldots, x_{i1}}_{b\text{-times}}, \underbrace{x_{i2}, \ldots, x_{i2}}_{b\text{-times}}, \ldots, \underbrace{x_{in_i}, \ldots, x_{in_i}}_{b\text{-times}};\right.$$

$$\left. \underbrace{x_{j1}, \ldots, x_{j1}}_{b\text{-times}}, \underbrace{x_{j2}, \ldots, x_{j2}}_{b\text{-times}}, \ldots, \underbrace{x_{jn_j}, \ldots, x_{jn_j}}_{b\text{-times}}\right)$$

for $b \in \mathbb{N}$ and all $i = 1, 2, \ldots, L$ and $j = 1, 2, \ldots, L$

The measure is population-invariant.

By looking at all pairs of groups and the corresponding polarizations in a population of $L$ groups we derive the following matrix:

|  | Group 1 | Group 2 | ... | Group i | Group j | ... | Group L |
|---|---|---|---|---|---|---|---|
| Group 1 | $P_{1,1}$ | $P_{1,2}$ | ... | $P_{1,i}$ | $P_{1,j}$ | ... | $P_{1,L}$ |
| Group 2 | $P_{2,1}$ | $P_{2,2}$ |  | $P_{2,i}$ | $P_{2,j}$ | ... | $P_{2,L}$ |
| ... | ... | ... | ... | ... | ... | ... | ... |
| Group i | $P_{i,1}$ | $P_{i,2}$ | ... | $P_{i,i}$ | $P_{i,j}$ | ... | $P_{i,L}$ |
| Group j | $P_{j,1}$ | $P_{j,2}$ | ... | $P_{j,i}$ | $P_{j,j}$ | ... | $P_{j,L}$ |
| ... | ... | ... | ... | ... | ... | ... | ... |
| Group L | $P_{L,1}$ | $P_{L,2}$ | ... | $P_{L,i}$ | $P_{L,j}$ | ... | $P_{L,L}$ |

It is immediately evident that this matrix is symmetrical. For the values on the main diagonal we have

$$P_{i,i} = 0 \quad \text{for} \quad i = 1, 2, \ldots, L.$$

Calculating the polarization between a group and itself makes little sense. By definition, the value of the measure would always be 0. Thus, we consider all values on the main diagonal as irrelevant for the determination of polarization.

If we calculate a weighted average of all $P_{i,j}$ with $i < j$,[10] we derive a measure $P^*$ with

$$P^* = \frac{\displaystyle\sum_{i=1}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} \lambda_{i,j} \cdot P_{i,j}}{\displaystyle\sum_{i=1}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} \lambda_{i,j}} \, .$$

In the easiest case, all weights $\lambda_{i,j}$ could be equal. However, failing to take into account the size of each group would lead to unpleasant results. Let us consider, for example, the case where the entire population is made up of $N = 2n+1$ individuals and 3 groups of size $n_i$ for $i = 1, 2, 3$ with $n_1 = n_2 = n$ and $n_3 = 1$. If all individuals in group 1 have no income, all individuals in group 2 have the income 1, and the individual of group 3 has the income $n+1$, we have

$$Hom_1 = Hom_2 = Hom_3 = 1$$
$$Het_{1,2} = Het_{1,3} = 1$$
$$Het_{2,3} = \frac{n}{n+2}$$

and

$$P_{1,2} = 4 \cdot \frac{n}{n+n} \cdot \frac{n}{n+n} \cdot 1 \cdot 1 \cdot 1 = 1$$
$$P_{1,3} = 4 \cdot \frac{n}{n+1} \cdot \frac{1}{n+1} \cdot 1 \cdot 1 \cdot 1 = \frac{4n}{(n+1)^2}$$
$$P_{2,3} = 4 \cdot \frac{n}{n+1} \cdot \frac{1}{n+1} \cdot 1 \cdot 1 \cdot \frac{n}{n+2} = \frac{4n^2}{(n+2)(n+1)^2} \, .$$

Thus, with all weights equal, we calculate

$$P^* = \frac{1}{3}\left[ 1 + \frac{4n}{(n+1)^2} + \frac{4n^2}{(n+2)(n+1)^2} \right]$$

and see that

$$\lim_{n \to \infty} P^* = \frac{1}{3} \, .$$

This is quite disturbing as for $n \to \infty$ the examined case resembles that of maximum polarization.

However, if we require the measure of polarization to take into account the effect of the size of a group on the total polarization, we can use $\lambda_{i,j} = n_i \cdot n_j$ for

---

[10] Due to the symmetry of the matrix, it is sufficient to take into account only the $P_{i,j}$ with $i < j$.

$i, j = 1, \ldots, L$ in analogy to Gini's mean difference.[11] This leads to the measure $P$ with

$$P = \frac{\displaystyle\sum_{\substack{i=1}}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} n_i \cdot n_j \cdot P_{i,j}}{\displaystyle\sum_{\substack{i=1}}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} n_i \cdot n_j}$$

$$= \frac{2 \displaystyle\sum_{\substack{i=1}}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} n_i \cdot n_j \cdot P_{i,j}}{N^2 - \displaystyle\sum_{i=1}^{L} n_i^2},$$

and for the distribution described above we calculate

$$P = \left[ n \cdot n \cdot 1 + n \cdot 1 \cdot \frac{4n}{(n+1)^2} + n \cdot 1 \cdot \frac{4n^2}{(n+2)(n+1)^2} \right] \frac{1}{n^2 + 2n}$$

$$= \left[ n^2 + \frac{4n^2}{(n+1)^2} + \frac{4n^3}{(n+2)(n+1)^2} \right] \frac{1}{n^2 + 2n}.$$

We now see that $\lim\limits_{n \to \infty} P = 1$, which is quite desirable because as mentioned earlier, for $n \to \infty$ the examined case resembles that of maximum polarization.[12]

Due to the fact that the polarizations between all pairs of groups are taken into account, we can easily determine the impact of any single pair of groups or any single group on the total polarization.

For the impact of the polarization between groups $k$ and $l$ we have

$$\frac{n_k \cdot n_l \cdot P_{k,l}}{\displaystyle\sum_{\substack{i=1}}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} n_i \cdot n_j \cdot P_{i,j}}$$

and for the impact of group $k$ on the total polarization we have

$$\frac{n_k \cdot \displaystyle\sum_{\substack{j=1 \\ j \neq k}}^{L} n_j \cdot P_{j,k}}{\displaystyle\sum_{\substack{i=1}}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} n_i \cdot n_j \cdot P_{i,j}}.$$

---

[11] See Mosler and Schmid (2005), page 46.

[12] Note that the value of the measure will usually change if the grouping is changed even if the overall income distribution remains the same. This is well in line with the assumption that the polarization in the entire population not only depends on the income distribution itself, but also on the number and the size of the groups. Splitting a group into two, for example, will increase the number of groups and decrease the sizes of the two new groups in comparison with the old group and, therefore, create a situation that is different from the original one.

If polarization is at a minimum[13], i.e. if all individuals in the population have the same income, polarization $P_{i,j}$ between all pairs of groups $i$ and $j$ with $i \neq j$ does not depend on the groups' sizes and we have

$$
\begin{aligned}
P_{i,j} &= 4 \cdot \frac{n_i}{n_i + n_j} \cdot \frac{n_j}{n_i + n_j} \cdot Hom_i \cdot Hom_j \cdot Het_{i,j} \\
&= 4 \cdot \frac{n_i}{n_i + n_j} \cdot \frac{n_j}{n_i + n_j} \cdot 1 \cdot 1 \cdot 0 \\
&= 0
\end{aligned}
$$

and thus

$$
P = \frac{2 \sum_{\substack{i=1}}^{L-1} \sum_{\substack{j=2 \\ j>i}}^{L} n_i \cdot n_j \cdot P_{i,j}}{N^2 - \sum_{i=1}^{L} n_i^2}
$$

$$
= 0.
$$

For maximum polarization, i.e. if there are exactly two groups of equal size, and all individuals of one group have no income, and all individuals of the other group have the same income, we have

$$
\begin{aligned}
n_1 &= n_2, \\
Hom_1 &= Hom_2 = 1, \\
Het_{1,2} &= 1,
\end{aligned}
$$

and thus we calculate

$$
\begin{aligned}
P &= P_{1,2} \\
&= 4 \cdot \frac{n_1}{n_2 + n_1} \cdot \frac{n_1}{n_1 + n_2} \cdot Hom_1 \cdot Hom_2 \cdot Het_{1,2} \\
&= 4 \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot 1 \cdot 1 \cdot 1 \\
&= 1.
\end{aligned}
$$

For maximum inequality, assuming $n - 1$ individuals have no income and compose one group and that a second group consists of only one individual with income $n \cdot \overline{x}$, we have

$$
\begin{aligned}
P &= P_{1,2} \\
&= 4 \cdot \frac{n_1}{n_2 + n_1} \cdot \frac{n_1}{n_1 + n_2} \cdot Hom_1 \cdot Hom_2 \cdot Het_{1,2} \\
&= 4 \cdot \frac{n-1}{n} \cdot \frac{1}{n} \cdot 1 \cdot 1 \cdot 1 \\
&= \frac{4(n-1)}{n^2}.
\end{aligned}
$$

---

[13] This case corresponds to minimum inequality.

In this case, we see that $\lim_{n \to \infty} P = 0$.

Thus, the measure $P$ satisfies all the requirements for a measure of polarization as postulated in the beginning and can, therefore, be applied in the desired way.

## References

D'Ambrosio, C. (2001) Household characteristics and the distribution of income in Italy: an application of social distance measures. The Review of Income and Wealth **47**(1), 43–64

Duclos, J.-Y., Esteban, J., Ray, D. (2004) Polarization: Concepts, measurement, estimation. Econometrica **72**(6), 1737–1772

Esteban, J., Gradín, C., Ray, D. (1999) Extensions of a measure of Polarization with an Application to the income distribution of five OECD countries. Working Paper No. 218, Maxwell School of Citizenship and Public Affairs, Syracuse University, Syracuse, New York

Esteban, J., Ray, D. (1994) On the measurement of polarization. Econometrica **62**(4), 819–851

Esteban, J., Ray, D. (1999) Conflict and distribution. Journal of Economic Theory **87**, 379–415

Gradín, C. (2000) Polarization by sub-populations in Spain, 1973–91. The Review of Income and Wealth **46**(4), 457–474

Mosler, K., Schmid, F. (2005) Beschreibende Statistik und Wirtschaftsstatistik. Springer, Berlin

Schmidt, A. (2004) Statistische Messung der Einkommenspolarisation. Eul Verlag, Lohmar-Köln

Wang, Y.-Q., Tsui, K.-Y. (2000) Polarization orderings and new classes of polarization indices. Journal of Public Economic Theory **2**(3), 349–363

Wolfson, M.C. (1994) When inequalities diverge. The American Economic Review **84**(2), 353–358

Zhang, X., Kanbur, R. (2001) What difference do polarisation measures make? An application to China. Journal of Development Studies **37**, 85–98