# A Model for Disaggregated Data Using Gini Index

Adriana Laura López Lobato[1], Martha Lorena Avendaño Garrido[2],
Johan Van Horebeek[2]

[1] Universidad Veracruzana, Facultad de Matemáticas,
Mexico

[2] Centro de Investigación en Matemáticas A.C.,
Mexico

adrilau17@gmail.com

**Abstract.** Corrado Gini developed in 1914 a methodology to measure the difference between two probability distributions, the Gini Index. In this paper, we propose the Bimodal Gini Index. We based this model on the definition of the Gini Coefficient, a model of independence between two distributions, so we set a model that approximates the Gini Index with the supposition that the searched distribution is a linear combination of independent distributions, without adding a lot of computational cost. We show some applications in political sciences concerning voting problems to illustrate the performance of the Bimodal Gini Index.

**Keywords:** Gini index, Gini coefficient, probability estimation.

## 1 Introduction

The Gini Index is a measure of the level of inequality between two probability distributions. It is applied in several fields of study like engineering, ecology, transport and economics, see [8].

The Gini Index problem is a particular case of Monge's mass transfer problem, as we will see in the following section. This problem always has a solution that is a distance between the involved probability distributions, but it can be very expensive to find it, computationally speaking, see [9] and [12]. To handle these expensive calculations, the Gini Coefficient was introduced as a natural upper bound of the Gini Index. The Gini Coefficient has several applications, many of them in economics and sociology, [2] and [8]. However, it differs a lot from the value of the Gini Index.

In this work we present the Bimodal Gini Index, a model that is a better approximation to the Gini Index than the Gini Coefficient with a low computational cost, by taking the Gini Index problem and doing the supposition that the searched probability is a linear combination of independent probability

distributions. With this model, we reduced the number of variables and the restrictions of the Gini Index problem and can be solved by using numerical optimization.

Also, it has several interesting properties, among these we highlight that it can be split in two linear programming problems, both easily solved by the simplex method.

## 2 Gini Index and Gini Coefficient

Let $X$ be a discrete random variable with $n$ elements and two probability distributions $p$ and $q$ on $X$. The Gini Index problem ($GI$) can be stated as:

$$\text{Minimize: } \sum_{i=1}^{n}\sum_{j=1}^{n} d_{ij}\pi_{ij}, \tag{1}$$

$$\text{subject to: } \pi_{ij} \geq 0, \qquad\qquad \text{for all } i,j \tag{2}$$

$$\sum_{j=1}^{n}\pi_{ij} = p_i, \qquad\qquad i = 1,2,...,n \tag{3}$$

$$\sum_{i=1}^{n}\pi_{ij} = q_j, \qquad\qquad j = 1,2,...,n \tag{4}$$

$$\sum_{i=1}^{n}\sum_{j=1}^{n}\pi_{ij} = 1, \tag{5}$$

where $p_i = p(x_i)$ y $q_i = q(x_i)$ for $i = 1,...,n$, the cost function is a distance function $d_{ij} = d(x_i, x_j)$ on $X \times X$, for all $i$ and $j$, and $\pi_{ij} = \pi(x_i, x_j)$ denotes the variables. The solution is a probability distribution $\pi^* = \{\pi_{ij}^* : i = 1,...,n, j = 1,...,n\}$. We define the Gini Index for the distributions $p$ and $q$, denoted by $GI(p,q)$, as the optimal value of the $GI$ problem. Note that there are $n^2$ no negative variables, then the solution of the problem can be expensive to find for large $n$, even with the use of computational tools. For more information about the Gini index and its problem in both forms, continuos and discrete, see [7,12,13].

On the other hand, we have a "measure of uncertainty" of a random variable, the Gini Coefficient for a discrete random variable $X$ with $n$ elements and two probability distributions $p$ and $q$ on $X$, see [1]:

$$GC(p,q) = \sum_{i=1}^{n}\sum_{j=1}^{n} d_{ij}p_i q_j.$$

With these definitions we can establish the following inequality

$$GI(p,q) \leq GC(p,q).$$

The Gini Index and the Gini Coefficient are used as indicators of social and economic inequality, as we can see in the articles [3,10,11].

# 3 Proposed Model: Bimodal Gini Index

To set the Bimodal Gini Index we consider the Gini Index problem and we add the additional assumption that the searched probability distribution $\pi$ is a linear combination of independent distributions, that means, it has the form:

$$\pi_{ij} = \alpha f_i^{(1)} f_j^{(2)} + (1 - \alpha) g_i^{(1)} g_j^{(2)}, \tag{6}$$

where $\alpha \in (0, 1)$ and $f^{(1)}, f^{(2)}, g^{(1)}$ and $g^{(2)}$ are independient probability distributions pairwise on $X$, this is that $f^{(1)}$ and $f^{(2)}$ are independent and $g^{(1)}$ and $g^{(2)}$ are independent.

As $f^{(1)}, f^{(2)}, g^{(1)}$ and $g^{(2)}$ are probability distributions, by replacing (6) in the expressions (3) and (4) we obtain:

$$p_i = \alpha f_i^{(1)} + (1 - \alpha) g_i^{(1)}, \text{ for all } i \quad \text{and} \quad q_j = \alpha f_j^{(2)} + (1 - \alpha) g_j^{(2)}, \text{ for all } j,$$

and expressing the variables $g_i^{(1)}$ and $g_j^{(2)}$ in terms of $f_i^{(1)}$ and $f_j^{(2)}$, respectively, as:

$$g_i^{(1)} = \frac{p_i - \alpha f_i^{(1)}}{1 - \alpha} \quad \text{and} \quad g_j^{(2)} = \frac{q_j - \alpha f_j^{(2)}}{1 - \alpha}, \text{ for all } i, j \text{ and } \alpha \in (0, 1),$$

we can express the variables $\pi_{ij}$ only in terms of $f_i^{(1)}$ and $f_j^{(2)}$ as:

$$\pi_{ij} = \frac{\alpha}{1 - \alpha} \left( f_i^{(1)} f_j^{(2)} - p_i f_j^{(2)} - q_j f_i^{(1)} + \frac{1}{\alpha} p_i q_j \right). \tag{7}$$

Also, we can express the values of $f_n^{(k)}$, with $k = 1, 2$, by $f_n^{(k)} = 1 - \sum_{i=1}^{n-1} f_i^{(k)}$, when we use this expressions in (7) we can define the following functions:

$$h_1(f_i^{(1)}, f_j^{(2)}) = f_i^{(1)} f_j^{(2)} - p_i f_j^{(2)} - q_j f_i^{(1)} + \frac{1}{\alpha} p_i q_j, \text{ for } i = 1, ..., n - 1, j = 1, ..., n - 1,$$

$$h_2(f^{(1)}, f_j^{(2)}) = f_j^{(2)} \left( 1 - \sum_{i=1}^{n-1} f_i^{(1)} - p_n \right) + q_j \left( \sum_{i=1}^{n-1} f_i^{(1)} - 1 + \frac{1}{\alpha} p_n \right), \text{ for } j = 1, ..., n-1,$$

$$h_3(f_i^{(1)}, f^{(2)}) = f_i^{(1)} \left( 1 - \sum_{j=i}^{n-1} f_j^{(2)} - q_n \right) + p_i \left( \sum_{j=1}^{n-1} f_j^{(2)} - 1 + \frac{1}{\alpha} q_n \right), \text{ for } i = 1, ..., n-1,$$

where $f^{(1)} = (f_1^{(1)}, ..., f_{n-1}^{(1)})$ and $f^{(2)} = (f_1^{(2)}, ..., f_{n-1}^{(2)})$, then we define

$$H(f^{(1)}, f^{(2)}) = \frac{\alpha}{1 - \alpha} \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} d_{ij} h_1(f_i^{(1)}, f_j^{(2)}) + \sum_{j=1}^{n-1} d_{ij} h_2(f^{(1)}, f_j^{(2)}) + \sum_{i=1}^{n-1} d_{ij} h_3(f_i^{(1)}, f^{(2)}).$$

This function only depends on the first $n - 1$ variables of the distributions $f^{(1)}$ and $f^{(2)}$. Also we have:

$$\frac{p_i - (1 - \alpha)}{\alpha} \le f_i^{(1)} \le \frac{p_i}{\alpha}, \quad \frac{q_j - (1 - \alpha)}{\alpha} \le f_j^{(2)} \le \frac{q_j}{\alpha}, \quad \text{for all } i, j \text{ and } \alpha \in (0, 1).$$

*Adriana Laura López Lobato, Martha Lorena Avendaño Garrido, Johan Van Horebeek*

Thus, we define the Bimodal Gini Index ($BGI$) as:

Minimize: $H(f^{(1)}, f^{(2)})$,

subject to: $\max\left\{0, \dfrac{p_i - (1-\alpha)}{\alpha}\right\} \leq f_i^{(1)} \leq \min\left\{1, \dfrac{p_i}{\alpha}\right\}$, $i = 1, ..., n-1$,

$\qquad\quad \max\left\{0, \dfrac{q_j - (1-\alpha)}{\alpha}\right\} \leq f_j^{(2)} \leq \min\left\{1, \dfrac{q_j}{\alpha}\right\}$, $j = 1, ..., n-1$,

$\qquad\quad \max\left\{0, \dfrac{p_n - (1-\alpha)}{\alpha}\right\} \leq 1 - \sum_{i=1}^{n-1} f_i^{(1)} \leq \min\left\{1, \dfrac{p_n}{\alpha}\right\}$,

$\qquad\quad \max\left\{0, \dfrac{q_n - (1-\alpha)}{\alpha}\right\} \leq 1 - \sum_{j=1}^{n-1} f_j^{(2)} \leq \min\left\{1, \dfrac{q_n}{\alpha}\right\}$.

If $f^* = (f^{(1)*}, f^{(2)*})$ is the optimal solution of the previous problem, then we define the Bimodal Gini Index as:

$$BGI(p,q) = H(f^{(1)*}, f^{(2)*}).$$

Note that the Gini Index problem has $n^2$ no negative variables and $2n + 1$ equality restrictions. With the proposed model we can reduce this amount to $2(n-1)$ variables, $2(n-1)$ box restrictions and 2 linear box restrictions.

Moreover, the Bimodal Gini Index is a better bound for the Gini Index than the Gini Coefficient, that is, the following inequality is fulfilled:

$$GI(p,q) \leq BGI(p,q) \leq GC(p,q).$$

So, we add the additional assumption that the searched probability distribution $\pi$ is a linear combination of independent distributions, as in 6, based on the idea of independence given by the Gini Coefficient, to make it more complex without adding a lot of computational cost:

- If $\alpha$ takes the value 0 or 1 in (6), then the optimal value of the BGI problem and the value of the Gini Coefficient will be the same.

- The objective function $H$ of the $BGI$ problem is a convex and symmetric function of $\alpha$ and reaches its minimum value in $1/2$ (or in $\alpha$ close to $1/2$).

- We can separate the BGI problem in two linear programming problems, both solved by the simplex method, as we will see in the following section.

### 3.1 Approximation to the Bimodal Gini Index

We can express the function $H(f^{(1)}, f^{(2)})$ as:

$$H(f^{(1)}, f^{(2)}) = H_L(f^{(1)}) + H_L(f^{(2)}) + H_C(f^{(1)}, f^{(2)}) + C,$$

where:

$$H_L(f^{(1)}) = \frac{\alpha}{1-\alpha} \left[ \sum_{i=1}^{n-1}\sum_{j=1}^{n} d_{ij}(-q_j f_i^{(1)}) + \sum_{i=1}^{n-1}\sum_{j=1}^{n-1} d_{nj}q_j f_i^{(1)} + \sum_{i=1}^{n-1} d_{in}f_i^{(1)} \right],$$

$$H_L(f^{(2)}) = \frac{\alpha}{1-\alpha} \left[ \sum_{i=1}^{n}\sum_{j=1}^{n-1} d_{ij}(-p_i f_j^{(2)}) + \sum_{i=1}^{n-1}\sum_{j=1}^{n-1} d_{in}p_i f_j^{(2)} + \sum_{j=1}^{n-1} d_{nj}f_j^{(2)} \right],$$

$$H_C(f^{(1)},f^{(2)}) = \frac{\alpha}{1-\alpha} \left[ \sum_{i=1}^{n-1}\sum_{j=1}^{n-1} d_{ij}f_i^{(1)}f_j^{(2)} - \sum_{i=1}^{n-1}\sum_{j=1}^{n-1} d_{nj}f_i^{(1)}f_j^{(2)} - \sum_{i=1}^{n-1}\sum_{j=1}^{n-1} d_{in}f_i^{(1)}f_j^{(2)} \right],$$

$$C = \frac{1}{1-\alpha} \left[ \sum_{i=1}^{n}\sum_{j=1}^{n} d_{ij}p_i q_j \right] + \frac{\alpha}{1-\alpha} \left[ \sum_{j=1}^{n-1} d_{nj}q_j + \sum_{i=1}^{n-1} d_{in}p_i \right].$$

The linear functions $H_L(f^{(1)})$ and $H_L(f^{(2)})$ depends on $f^{(1)}$ and $f^{(2)}$, respectively. The value of $C$ is known. The quadratic function $H_C(f^{(1)},f^{(2)})$ only have negative values bounded by $-2d$, where $d$ is the maximum distance between the elements of $X$. We can move the elements of $X$ to a specific range, so $d$ is as small as we desired. Then, we only consider the linear functions, leaving the following separate problems.

**Linear problem with respect to $f^{(1)}$:**

Minimize: $H_L(f^{(1)})$

subject to: $\max\left\{0, \dfrac{p_i - (1-\alpha)}{\alpha}\right\} \leq f_i^{(1)} \leq \min\left\{1, \dfrac{p_i}{\alpha}\right\}, \, i = 1,...,n-1,$

$\max\left\{0, \dfrac{p_n - (1-\alpha)}{\alpha}\right\} \leq 1 - \sum_{i=1}^{n-1} f_i^{(1)} \leq \min\left\{1, \dfrac{p_n}{\alpha}\right\}.$

**Linear problem with respect to $f^{(2)}$**

Minimize: $H_L(f^{(2)})$

subject to: $\max\left\{0, \dfrac{q_j - (1-\alpha)}{\alpha}\right\} \leq f_j^{(2)} \leq \min\left\{1, \dfrac{q_j}{\alpha}\right\}, \, j = 1,...,n-1$

$\max\left\{0, \dfrac{q_n - (1-\alpha)}{\alpha}\right\} \leq 1 - \sum_{j=1}^{n-1} f_j^{(2)} \leq \min\left\{1, \dfrac{q_n}{\alpha}\right\}.$

Then we define the Separated Bimodal Gini Index as:

$$BGIs(p,q) = H(f^{(1)*}, f^{(2)*}),$$

where $f^{(1)*} = (f_1^{(1)*},...,f_{n-1}^{(1)*})$ y $f^{(2)*} = (f_1^{(2)*},...,f_{n-1}^{(2)*})$ are the points where the optimal results are reached in the linear problems with respect to distributions $f^{(1)}$ y $f^{(2)}$, respectively, and $H$ is the objective function previously expressed.

We can obtain the Separated Bimodal Gini Index solving two linear problems by the simplex method, each of one with $n-1$ variables, $n-1$ box restrictions and a linear box restriction. Solving these two problems is much less expensive, computationally speaking, than solving the original one.

*Adriana Laura López Lobato, Martha Lorena Avendaño Garrido, Johan Van Horebeek*

- The Bimodal Gini Index and the Separated Bimodal Gini Index take the same value in $\alpha = 1/2$ in computational experiments.

- If we obtain the values of the distributions $f^{(1)}$ and $f^{(2)}$ we can obtain the values of the distribution $\pi$ of the form (6). The distribution $\pi$ is of great importance for the application in the following section.

## 4 Aplication in Political Science

### 4.1 Voting Data Ohio, 1990

We can see in the Table (1) the data of race of voting-age person and the voting decision for the 1990 election in the Ohio State House, District 42, see [4]. The unobservable values in the empty cells must be found from the observed values in the marginals.

**Table 1.** Aggregate data for the 1990 election in the Ohio State House, District 42.

| Race | Voting decision | | |
| --- | --- | --- | --- |
| | Democrat | Republican | No vote |
| African american | | | 221 (0.313) |
| White | | | 484 (0.687) |
| | 130 (0.184) | 92 (0.131) | 483 (0.685) 705 (1.000) |

We want to fill this table using the problems raised in the previous section by taking the value of $\alpha = 1/2$, the random variable $X =\{$African american, White, Democrat, Republican, No vote$\}$ and the probability distributions $p = \{0.313, 0.687, 0, 0, 0\}$ and $q = \{0, 0, 0.184, 0.131, 0.685\}$. Since the values of the random variable $X$ are categorical, we will use the discrete metric. So, the problems are:

Minimize: $0.315f_1^{(1)} + 0.315f_2^{(1)}$

subject to: $0 \leq f_1^{(1)} \leq 0.626,$

$0.374 \leq f_2^{(1)} \leq 1,$

$f_1^{(1)} + f_2^{(1)} = 1.$

Minimize: $f_3^{(2)} + f_4^{(2)}$

subject to: $0 \leq f_3^{(2)} \leq 0.368,$

$0 \leq f_4^{(2)} \leq 0.262,$

$0 \leq f_3^{(2)} + f_4^{(2)} \leq 0.63.$

We found the searching value $BGIs$ in points of the form

$$(f_1^{(1)*}, f_2^{(1)*}, f_3^{(1)*}, f_4^{(1)*}, f_1^{(2)*}, f_2^{(2)*}, f_3^{(2)*}, f_4^{(2)*}) = (f_1^{(1)}, 1 - f_1^{(1)}, 0, 0, 0, 0, 0, 0),$$

with $f_1^{(1)} \in [0, 0.626]$. We analize the solution in the extreme point with $f_1^{(1)} = 0$. So the Table 2 shows the data of interest, the probability distribution $\pi$.

**Table 2.** Results given by the Separated Binomial Gini Index for the 1990 election in the Ohio State House, District 42.

|  | Democrat | Republican | No vote |
|---|---|---|---|
| African american | 0.115184 | 0.082006 | 0.11581 |
| White | 0.068816 | 0.048994 | 0.56919 |

Note that with the problems of the Separate Gini Index we can obtain the wanted probabilities and a scenario of how the votes in Ohio could have been distributed with respect to the race of the voters.

In [5] three types of results obtained for this same problem given by King are shown, with the particularity that this solutions are interest intervals. Thus, when making a comparison of the puntual results obtained by the IGAs problems, we can notice that these are within the corresponding intervals.

### 4.2 Elections in the republic of Weimar, 1932

One of the most studied questions in the history is "who voted by Hitler?". In [6] identify some factors that could explain why certain groups of voters joined the Nazi party, concluding that a determining factor was the economic great depression, so the occupations of voters are studied. In the Table (3) we observe the marginals obtained for this problem, the left column of the table denotes each occupational group while the upper row indicates the different political parties.

**Table 3.** Aggregate data for elections in 1932 in the republic of Weimar.

|  | Far Left | Left/Center | Far Right | Nazi | Liberal | No vote/ Other |  |
|---|---|---|---|---|---|---|---|
| Self-employed |  |  |  |  |  |  | 0.164 |
| Blue collar |  |  |  |  |  |  | 0.314 |
| White collar |  |  |  |  |  |  | 0.144 |
| Domestic |  |  |  |  |  |  | 0.197 |
| Unemployed |  |  |  |  |  |  | 0.181 |
|  | 0.120 | 0.311 | 0.049 | 0.311 | 0.018 | 0.191 |  |

The objetive of this problem is filling the Table to answer questions like "what fraction of independent people voted for the Nazi party?". Analyzing historically this type of questions, it is expected that the results related to the working class (blue collar) will be those that favor the Nazi party, since they feared losing their jobs if the centralist party remained in power, see [6]. There are no statistical references for the solution to this problem, our results would be a way to confirm the hypothesis made by researchers in Social Sciences.

We solved the Separated Bimodal Gini Index problems with $\alpha = 1/2$ and the random variable $X =$ {Self-employed, Blue collar, White collar, Domestic, Unemployed, Far Left, Left/Center, Far Right, Nazi, Liberal, No vote/Other} and the

probability distributions $p = \{0.164, 0.314, 0.144, 0.197, 0.181, 0, 0, 0, 0, 0, 0\}$ and $q = \{0, 0, 0, 0, 0, 0.120, 0.311, 0.049, 0.311, 0.018, 0.191\}$. So, we have the following problems:

Minimize: $0.809 f_1^{(1)} + 0.809 f_2^{(1)} + 0.809 f_3^{(1)} + 0.809 f_4^{(1)} + 0.809 f_5^{(1)}$

subject to: $0 \leq f_1^{(1)} \leq 0.328, \qquad 0 \leq f_2^{(1)} \leq 0.628,$

$\qquad\qquad 0 \leq f_3^{(1)} \leq 0.288, \qquad 0 \leq f_4^{(1)} \leq 0.394,$

$\qquad\qquad 0 \leq f_5^{(1)} \leq 0.362, \qquad f_1^{(1)} + f_2^{(1)} + f_3^{(1)} + f_4^{(1)} + f_5^{(1)} = 1.$

Minimize: $f_6^{(2)} + f_7^{(2)} + f_8^{(2)} + f_9^{(2)} + f_{10}^{(2)}$

subject to: $0 \leq f_6^{(2)} \leq 0.240, \qquad 0 \leq f_7^{(2)} \leq 0.622,$

$\qquad\qquad 0 \leq f_8^{(2)} \leq 0.098, \qquad 0 \leq f_9^{(2)} \leq 0.622,$

$\qquad\qquad 0 \leq f_{10}^{(2)} \leq 0.036, \qquad 0.618 \leq f_6^{(2)} + f_7^{(2)} + f_8^{(2)} + f_9^{(2)} + f_{10}^{(2)} \leq 1.$

The minimum value is reached in the points of the form

$$(f_1^{(1)}, f_2^{(1)}, f_3^{(1)}, f_4^{(1)}, f_5^{(1)}, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, f_6^{(2)}, f_7^{(2)}, f_8^{(2)}, f_9^{(2)}, f_{10}^{(2)})$$

where the values of this variables meet the constraints of the previous problems.

We calculated the values in the Table 4 for the point

$$(0.198, 0.12, 0.258, 0.254, 0.17, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0.18, 0.074, 0.098, 0.23, 0.036).$$

This point meets the aforementioned restrictions.

**Table 4.** Results given by the Separated Binomial Gini Index for the 1932 elections in the republic of Weimar.

|  | Far Left | Left/Center | Far Right | Nazi | Liberal | No vote/ Other |
|---|---|---|---|---|---|---|
| Self-employed | 0.022 | 0.043 | 0.009 | 0.048 | 0.004 | 0.038 |
| Blue collar | 0.026 | 0.144 | 0.006 | 0.113 | 0.002 | 0.023 |
| White collar | 0.024 | 0.018 | 0.013 | 0.035 | 0.005 | 0.049 |
| Domestic | 0.027 | 0.047 | 0.013 | 0.057 | 0.004 | 0.049 |
| Unemployed | 0.021 | 0.059 | 0.008 | 0.058 | 0.003 | 0.032 |

As we can see, it is true that the working class, blue collar, is the most likely to belong to the Nazi party or the centralist party, as expected.

## 5  Conclusions and Future Work

The Bimodal Gini Index is a better bound for the Gini Index than the Gini Coefficient. The Bimodal Gini Index has many favorable properties like the Separated Bimodal Gini Index problems. This is possible because of the specific form

given to the searched distribution, that reduces the feasible set of the problem. Also, because this model is based in the Gini Coefficient, the computational cost does not increase as much. In this way we reduced the problem in terms of the number of variables and we found a simpler way to solve it by means of two linear problems with box constraints using the simplex method.

We can also observe in the given examples that the problems of the separated Bimodal Gini Index are very useful to solve problems where we have grouped information and we want to obtain data at a disaggregated level. The solved examples are current problems pertinent to political science and history, and their solutions are of great importance for these fields of science.

As future work, we want to use this model in other data bases in different areas of science and in any type of problems that involved disaggregated data or lack of information.

# References

1. Bassetti, F., Bodini, A., Regazzini, E.: On Minimum Kantorovich Distance Estimators. Statistics and probability letters 76(12), 1298–1302 (2006)
2. Chakravarty, S.: Ethical Social Index Numbers. Springer Berlin Heidelberg (2012)
3. Han, J., Zhao, Q., Zhang, M.: China's income inequality in the global context. Perspectives in Science 7, 24–29 (2016)
4. King, G.: A Solution to the Ecological Inference Problem: Reconstructing Individual Behavior from Aggregate Data. Princeton University Press (1997)
5. King, G.: A solution to the ecological inference problem: Reconstructing individual behavior from aggregate data. Princeton University Press (2013)
6. King, G., Rosen, O., Tanner, M., Wagner, A.F.: Ordinary Economic Voting Behavior in the Extraordinary Election of Adolf Hitler. The Journal of Economic History 68(4), 951–996 (2008)
7. Peyré, G., Cuturi, M., et al.: Computational optimal transport. Foundations and Trends® in Machine Learning 11(5-6), 355–607 (2019)
8. Rachev, S., Klebanov, L., Stoyanov, S., Fabozzi, F.: The Methods of Distances in the Theory of Probability and Statistics. SpringerLink : Bucher, Springer New York (2013)
9. Rubner, Y., Tomasi, C., Guibas, L.J.: The Earth Mover's Distance as a Metric for Image Retrieval. International journal of computer vision 40(2) (2000)
10. Schneider, M., et al.: Measuring Inequality: The Origins of the Lorenz Curve and the Gini Coefficient. La Trobe University, School of Business (2004)
11. Sturm, J.E., De Haan, J.: Income inequality, capitalism, and ethno-linguistic fractionalization. American Economic Review 105(5), 593–97 (2015)
12. Villani, C.: Topics in Optimal Transportation. American Mathematical Society (2003)
13. Villani, C.: Optimal transport: old and new, vol. 338. Springer Science & Business Media (2008)