

# Investigation on income inequality of the Islands

## STAT3003 Midterm A - Survey Plan

ZHENG Weijia (William, 1155124322)

April 6, 2021

### 1 Introduction and Notations

Income inequality continues to receive world-wide attention. The most famous method to measure it to use the Gini Index proposed by Gini. Gini index gives a number inside the interval  $[0, 1]$ , the more the index close to 1, the worse the inequality problem is.

Assume we will investigate the situation of income inequality for the scale of a town (e.g., Hofn, Vardo and Takazaki are all towns).

Then denote the population of the town as  $n$ , note that this is given.

Denote the total amount of wealth of the town as  $M$ , which is also given by consulting the Hall of the town.

Denote the average amount of wealth as  $\mu = \frac{M}{n}$ .

For the  $i$ -th individual ( $i$  ranges from 1 to  $n$ ), denote the his individual amount of wealth as  $m_i$ , note that  $\sum_{i=1}^n m_i = M$ .

Then if we follow the definition of Gini index, we need to calculate

$$G = \frac{1}{2n^2\mu} \sum_{j=1}^n \sum_{i=1}^n |m_i - m_j|.$$

But this is troublesome because it is hard to estimate the  $G$  of a town properly from a sample of smaller size. Hence we switch to another method of measuring the inequality.

### 2 Mathematical Model

#### 2.1 Derived from the Concept of Entropy

As said above, hindered by the complex way to calculate the Gini index out and noted that the Islands provide the information of the total amount of a town and every individual's. We can then understand the concept of "income euqality" from the perspective of the difference between "every person's ability (as a possibility) to gain wealth".

For the  $i$ -th person, we regard the ratio between his own wealth and the total wealth, i.e.,  $p_i = \frac{m_i}{M}$  of the town as his "possibility of gaining a unit of wealth".

Note that  $\sum_{i=1}^n p_i = 1$ , hence the  $p_i$  forms a proper probability distribution. Note that the entropy

$$H(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \log_n p_i$$

actually measures the inequality of the town economy. Hence define it as the " $H$  index" of the economy.

When the town economy is the most fair, i.e., everyone has the same probability of earning money:  $p_i = \frac{1}{n}, \forall i = 1, 2, \dots, n$ , then

$$H(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \log_n p_i = - \sum_{i=1}^n \frac{1}{n} \log_n \frac{1}{n} = 1.$$

When the economy is extremely unequal, i.e., all the money go to one person, then the entropy will be 0.

The following is a simple example of the application of this Mathematical model. Consider a economy with only 3 people, each income is 10, 20, 15. And another economy with 3 people with each income 3, 7, 40. The first economy has  $H_1 = 0.9656$  and the second economy with  $H_2 = 0.5666$ , which conforms our intuition.

Note that the  $H$  index is a  $\tau$ -type metric, so all the knowledge for estimating  $\tau$  can be applied properly.

## 2.2 Some Related Properties

// TO IMPLEMENT

## 3 Sampling Method

### 3.1 Not Using Auxiliary Data

Recall that we need to estimate a  $\tau$ -type metric, and the metric  $t_i$  which we need to get from each individual (denote  $t_i = -p_i \log_n p_i$ ) is a deterministic function of  $m_i$  because  $p_i = \frac{m_i}{M}$ . Since  $m_i$  is extremely easy to obtain, we will not apply the "using auxiliary data" trick.

### 3.2 Not Using SRS

If we apply SRS, then we need a whole list which contains every one inside the town and randomly choose some from it as our sample.

Although we know the population number, but all people are grouped in houses, and we cannot know how many people are inside each house, since we cannot have such a whole list to randomly choosing elements from. Therefore we do not choose SRS.

### 3.3 Not Forming Clusters Geographically

This is an empirical decision from my investigating on some towns. I observed that the neighbourhoods of a person is likely to have similar economic status with him. In other words, people sharing similar economic status are likely to gather together geographically as a cluster. Hence I will not use Cluster Sampling under such way of forming clusters.

### 3.4 Our Sampling Method

We will adopt Cluster Sampling by forming clusters in a special way, more explicitly, in a systematic sampling manner.

First, we select a proper number  $s$  as the stepsize, and  $r$ .

We will classify all the people whose house number modulus  $s$  equals  $r$ .

Note that there are  $N = \frac{T}{s}$  clusters in totals, where  $T$  is the number of NON-EMPTY houses in the town.

We choose  $N_c$  clusters from the  $N$  clusters by SRS.

Then, for every chosen cluster  $i, i = 1, 2, \dots, N_c$  we go through all the houses and the people inside those houses and collect data from every element(person).

Using formulae

$$\hat{\tau} = \hat{H} = N \frac{1}{N_c} \sum_{i=1}^{N_c} Y_i,$$

with  $Y_i = \sum_{j=1}^{M_i} t_{ij}$ , where  $M_i, i = 1, 2, \dots, N_c$  is the number of elements (people) inside the  $i$ -th chosen cluster. We have the point estimate.

Recall that  $\hat{\tau}$  is unbiased.

Using the formulae

$$Var(\hat{\tau}) = N^2 \left(1 - \frac{N_c}{N}\right) \frac{\hat{\sigma}_c^2}{N_c},$$

we can have an estimated variance, where the  $\hat{\sigma}_c^2$  is the sample variance of the cluster totals, i.e., the variance of those  $Y_i$ 's, with  $i = 1, 2, \dots, N_c$ .

Thus an appropriate  $100(1 - \alpha)\%$  C.I.s can be given by

$$\hat{\tau} \pm t_{N_c-1, 1-\frac{\alpha}{2}} \sqrt{Var(\hat{\tau})}.$$

## 4 Preliminary Survey Result

We select the town Vardo to be under investigation.

## 5 Where to Improve