**Centre for Efficiency and Productivity Analysis**

**Working Paper Series**
**No. WP03/2007**

**Date: June, 2007**

**School of Economics**
**University of Queensland**
**St. Lucia, Qld. 4072**
**Australia**

# Non-Hierarchical Bivariate Decomposition of Theil Indexes

Kam Ki TANG *, #

School of Economic, University of Queensland

Dennis PETRIE

School of Population Health, University of Queensland

**June, 2007**

**Abstract**

This paper develops a method to conduct non-hierarchical bivariate decomposition of

Theil indexes. The method has the merits that, first, it treats all variates symmetrically

and therefore facilitates the comparison of inequalities associated with different

variates; and, second, it highlights the interaction between variates in the creation of

inequality. The method is applied to measure gender and ethnic income inequality in

Australia.

**Keyword:** Theil index, hierarchical decomposition

**JEL Classification:** D30, D63

* Corresponding author. Address: School of Economics, University of Queensland, QLD 4072,
Australia. Tel: +617 3365 9796. Fax: +617 33657209. Email: kk.tang@uq.edu.au.

## 1. Introduction

The Theil index is the most commonly used entropy measure in economic studies; e.g., see Mishra and Parikh (1992), and Conceicao and Galbraith (2000). A property of the entropy family, of which the Theil index belongs to, is that its members can be decomposed into exhaustive and exclusive components. The decomposability of inequality measures has been discussed extensively in a number of studies, including Bourguignon (1979), Shorrocks (1980), Cowell (1980; 1985), Cowell and Kuga (1981), and Adelman & Levy (1984). Amongst them Adelman & Levy (1984) and Cowell (1985) are concerned with multilevel decomposition of the Theil index, an issue closely related to the theme of this paper.

The additive decomposability of the Theil index allows the examination of how overall inequality is related to subgroup characters. For instance, we can decompose the Theil measures of population-wide income inequality into between-gender and within-gender inequalities. Likewise, we can decompose the Theil measure based on other stratifications, such as ethnicity. This decomposition method allows us to slice the pie of total inequality according to either gender or ethnicity, but only one dimension at a time. However, since a population can be stratified by gender and ethnicity simultaneously, can we also decompose the Theil measure according to both variates *simultaneously*?

When the decomposition is hierarchical, the answer is simply yes. Hierarchical decomposition means that the Theil measure is decomposed first in one dimension and then in another. For instance, in Panel A of Figure 1, the Theil index is decomposed first by ethnicity and then by gender. Based on the traditional decomposition method, within-ethnicity inequality is decomposed into within-ethnicity-between-gender inequality and

within-ethnicity-within-gender inequality, but the latter is indeed the same as within-ethnicity-within-gender. In Panel B, the order of decomposition is reversed.

For hierarchical decompositions the order of decomposition matters: Panels A and B have only one common term – the within-gender-ethnicity inequality (or, in equivalent, within-ethnicity-gender inequality); all other terms are different. This is not an issue if there is a natural hierarchical order between the variates, such as the province-city stratification in Akita (2003): as city must be hierarchically under province, the decomposition is naturally done first by province and then by city. However, in many other cases, there is no natural hierarchical order, e.g. gender and ethnicity, occupation and education, and industry and region.

Considering these limitations of the hierarchical decomposition, this paper aims to develop a simple method to obtain a non-hierarchical bivariate decomposition of the Theil measure. The method has two merits as compared to hierarchical decomposition. First, it treats all variates symmetrically and therefore facilitates the comparison of inequalities associated with different variates. Second, the method highlights the interaction between variates in the creation of inequality.

The next section explains both hierarchical and non-hierarchical decomposition methods. As an illustration, Section 3 applies the method to decompose labour income inequality in Australia by gender and ethnicity.

## 2. Hierarchical and Non-Hierarchical Decomposition of Theil

### 2.1 Hierarchical Decomposition

Consider the income inequality of a population of people with both genders and mixed ethnic backgrounds. The Theil-L index for a population is expressed as[1]

$$L = \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \log \left( \frac{N_{egi}/N}{Y_{egi}/Y} \right) \tag{1}$$

where $e$ = ethnicity index, $g$ = gender index, $i$ = income division index,[2] $N_{egi}$ = the size of group $egi$, $N = \sum_e \sum_g \sum_i N_{egi}$ = the size of the whole population, $Y_{egi}$ = income of group $egi$, and $Y = \sum_e \sum_g \sum_i Y_{egi}$ = total income of the population.

The logarithmic function in equation (1) is a measure of the deviation of the income share of the group $egi$ (i.e $Y_{egi}/Y$) from its population share (i.e. $N_{egi}/N$). If the group's income share is equal to its population share, it has its "fair share" of income and does not contribute to the inequality index. However, if the group's income share of is smaller (bigger) than its population share, it contributes positively (negatively) to the index, with its contribution weighted by its population share. In other words, the Theil-L index is a weighted sum of the deviation of income share from population share for every group in a population.[3] An important point to emphasize here is that a negative contribution, just like a positive one, indicates the existence of inequality, as with a

---

[1] The expressions for the Theil-T index can be obtained by swapping $Y$ and $N$. The discussion for Theil-T will be similar to that of Theil-L and therefore skipped.
[2] E.g. $i = 1$ for the lowest percentile of income distribution and $i = 10$ for the highest percentile.
[3] Alternatively, one can consider the logarithmic function as a measure of the deviation of the average income of the group $egi$ from the average income of the population.

negative contribution there must exist a larger positive contribution. Given this, the total weighted sum of all contribution will never be negative (see Appendix for the proof).

In Panel A of Figure 1, the Theil index is decomposed into within-ethnicity-gender $(w_{EG})$, within-ethnicity-between-gender $(w_E b_G)$, and between-ethnicity $(b_E)$ inequalities, respectively:

$$L = \sum_e \left( \frac{N_e}{N} \right) \log \left( \frac{N_e / N}{Y_e / Y} \right) \qquad (b_E)$$

$$+ \sum_e \left( \frac{N_e}{N} \right) \left[ \sum_g \left( \frac{N_{eg}}{N_e} \right) \log \left( \frac{N_{eg} / N_e}{Y_{eg} / Y_e} \right) \right] \qquad (w_E b_G) \qquad (2)$$

$$+ \sum_e \sum_g \left( \frac{N_{eg}}{N} \right) \left[ \sum_i \left( \frac{N_{egi}}{N_{eg}} \right) \log \left( \frac{N_{egi} / N_{eg}}{Y_{egi} / Y_{eg}} \right) \right] \qquad (w_{EG})$$

where $N_{eg} = \sum_i N_{egi}$, $N_e = \sum_g N_{eg}$, $Y_{eg} = \sum_i Y_{egi}$, and $Y_e = \sum_g Y_{eg}$. (See Appendix for the proof.)

$b_E$ measures the inequality between different ethnic groups, $w_E b_G$ measures the inequality between males and females across all ethnic groups, and $w_{EG}$ measures the inequality within each of the ethnic-gender groups.

In Panel B, the index is decomposed into within-gender-ethnicity inequality $(w_{GE})$, within-gender-between-ethnicity $(w_G b_E)$, and between-gender $(b_G)$ inequalities, respectively:

$$L = \sum_g \left( \frac{N_g}{N} \right) \log \left( \frac{N_g / N}{Y_g / Y} \right) \qquad (b_G)$$

$$+ \sum_g \left( \frac{N_g}{N} \right) \left[ \sum_e \left( \frac{N_{eg}}{N_g} \right) \log \left( \frac{N_{eg} / N_g}{Y_{eg} / Y_g} \right) \right] \qquad (w_G b_E) \qquad (3)$$

$$+ \sum_g \sum_e \left( \frac{N_{eg}}{N} \right) \left[ \sum_i \left( \frac{N_{egi}}{N_{eg}} \right) \log \left( \frac{N_{egi} / N_{eg}}{Y_{egi} / Y_{eg}} \right) \right]. \qquad (w_{GE})$$

$b_G$ measures the inequality between males and females, $w_G b_E$ measures the inequality

between ethnic groups across both gender groups, and $w_{GE}$ is identical to $w_{EG}$.

## 2.2 Non-Hierarchical Decomposition

Since (2) and (3) must equate each other and $w_{GE} \equiv w_{EG}$, we can state

$$w_G b_E - b_E \equiv w_E b_G - b_G \equiv residue . \qquad (4)$$

We label this residue the "gender-ethnicity interaction inequality," $i_{GE} \equiv i_{EG}$. The reason

for this will become clear later.

Using this definition of the residue, we can write

$$b_E \equiv w_G b_E - i_{GE}, \qquad (5)$$

$$b_G \equiv w_E b_G - i_{GE} . \qquad (6)$$

Substituting (5) into (3) yields a non-hierarchical decomposition of the Theil index into

four components:

$$L = w_{GE} + b_G + b_E + i_{GE} . \qquad (7)$$

The decomposition is illustrated in Figure 2.

Here $b_G$ and $b_E$ measure respectively the parts of inequality that are associated with gender and ethnicity, $w_{GE}$ measures the part of inequality that is associated with neither of them, and, as shown next, $i_{GE}$ measures the part of inequality that is associated with both gender and ethnicity.

## 2.3 Gender-Ethnicity Interaction Inequality

As in standard decomposition, total inequality is equal to the sum of within- and between-group inequalities:

$$L = \sum_e \sum_g \left( \frac{N_{eg}}{N} \right) \log \left( \frac{N_{eg}/N}{Y_{eg}/Y} \right) \qquad (b_{GE})$$
$$+ \sum_e \sum_g \left( \frac{N_{eg}}{N} \right) \left[ \sum_i \left( \frac{N_{egi}}{N_{eg}} \right) \log \left( \frac{N_{egi}/N_{eg}}{Y_{egi}/Y_{eg}} \right) \right] \qquad (w_{GE}). \tag{8}$$

where $b_{GE} \equiv b_{EG}$ measures the inequality between ethnic-gender groups.

Equating (7) and (8) give

$$b_{GE} \equiv w_E b_G + w_G b_E - i_{GE}. \tag{9}$$

Substituting (5) and (6) into this yield

$$i_{GE} \equiv b_{GE} - b_G - b_E. \tag{10}$$

Here we can express the gender-ethnicity interaction term as (see Appendix for the proof)

$$i_{GE} = \sum_e \sum_g \left( \frac{N_{eg}}{N} \right) \log \left( \frac{\sigma_{Neg}}{\sigma_{Yeg}} \right)$$
$$\sigma_{Neg} = \frac{N_{eg}/N}{(N_e/N)(N_g/N)}, \sigma_{Yeg} = \frac{Y_{eg}/Y}{(Y_e/Y)(Y_g/Y)} \tag{11}$$

$N_j / N$ is equal to the probability that a person randomly selected from the population belongs to group $j$, $j = e, g, eg$. If the event that a person belongs to ethnic group $e$ is independent of the event that a person belongs to gender group $g$, $\log(\sigma_{Neg})$ will be equal to zero; otherwise, it will be non-zero. Therefore, $\log(\sigma_{Neg})$ is a measure of the dependency of the two events, or more explicitly, the interrelationship (or interaction) between ethnicity $e$ and gender $g$ in the allocation of the population into the ethnicity-gender group $eg$. Similarly, $\log(\sigma_{Yeg})$ is a measure of the interaction between ethnicity $e$ and gender $g$ in the allocation of the income into the ethnicity-gender group $eg$. Hence $i_{GE}$ is a weighted sum of the derivation of the interaction of $e$ and $g$ in the allocation of income into group $eg$ from that of population.

To foster a better understanding of this interaction inequality, we consider a numerical example of two ethnic groups: native and non-native. Table 1 shows the value of the total income and the total number of individuals in each of the four ethnic-gender groups. Ethnicity and gender are already independent of each other in the allocation of the population, as the number of males is 2.5 times that of females for both ethnic groups, and the number of non-natives is twice that of natives for both genders. The total income of male native is presented by $Y_{na,m}$, which is a variable in the following simulation.

In the simulation, we change the value of $Y_{na,m}$ from 5 to 200 while keeping all other figures constant. Everyone with an ethnic-gender group is assumed to earn the same income and, thus, $w_{GE} = 0$. The impacts on the Theil-L index and its various components are shown on Figure 3. The total inequality falls first and then rises again as

the total income of male native increases. The four conventional components of Theil-L index, namely $b_E$, $b_G$, $w_E b_G$ and $w_G b_E$ show a similar skewed U-shape trajectory. On the other hand, $i_{GE}$, while of relatively much smaller values, is highly non-linear. The schedule of $i_{GE}$ crosses the x-axis three times at $Y_{na,m}$ equal to 20, around 8.4 and 153.

When $Y_{na,m}$ is equal to 20, $\sigma_{Neg} = \sigma_{Yeg} = 1$ for all $e$ and $g$. That is, ethnicity and gender are completely independent of each other in the allocation of both population and income. As a result, $i_{GE} = 0$. On the other hand, when $Y_{na,m}$ is close to 8.4 and 153, $\sigma_{Neg} \neq \sigma_{Yeg}$ for individual $e$ and $g$; however, the weighted values of $\log(\sigma_{Neg} / \sigma_{Yeg})$ for various pairs of $\{e, g\}$ cancel each other out, leaving no net effect on the total inequality. This demonstrates that the independence of ethnicity and gender in the allocation of income and population is a sufficient but not necessary condition for the interaction inequality to be equal to zero.

A unique feature of $i_{GE}$, as against the conventional inequality components, is that it can be negative, due to its structural difference. When $i_{GE}$ is negative, it represents the overlapping part of $b_E$ and $b_G$; when it is positive, it represents the 'gap' between the two.

## 3. Labour Income Inequality in Australia

This section applies the proposed decomposition method to estimate gender and ethnic labour income inequality in Australia. The data are sourced from the 1998-99 Household Expenditure Survey (HES) (Australia Bureau of Statistics 2000). The data set has been used to examine trends in household income and consumption inequality in

Australia (e.g. see Harding & Greenwell 2002), but not gender and ethnic inequality in individual labour income. The magnitudes used are weekly gross wages and salaries.[4] Due to data limitation, the country of birth is used as a proxy of ethnicity. The HES categorizes countries of birth into 10 regions. Table 2 provides the summary statistics of the income data by sex and country of birth. There are totally 218,187 observations in the sample. Those who were born in Australia represent over 75 percent of the sample, well ahead of the 10 percent share of the next group – North-West Europe.

Table 3 summarizes the percentage shares of various decomposed items of Theil-L and -T measures of the labour income inequality. It can be seen that for all ages combined (Theil-L), $w_{EG}$ accounts for nearly 90 percent of the total inequality, distantly followed by $b_G$ at around 10 percent. The values of $b_E$ is less than one percent and $i_{GE}$ is negligible. The figures for Theil-T are very similar so we concentrate our discussion on Theil-L. These results indicate that while gender inequality is substantial, ethnic inequality is not as an important issue. Moreover, the bivariate decomposition shows that the interaction between ethnicity and gender has contributed little to income inequality. In other words, without losing much, one can comfortably approximate the value of the Theil index as $L \approx w_{EG} + b_E + b_G$.

Since labour income increases with experience (age), if a large amount of $w_{EG}$ is due to the income gap between workers of different ages within each ethnic-gender group, it could disguise the inequality effects of gender and ethnicity. To control for the age effect, we break down the sample into five age groups; the results are shown in columns 3 to 7. A noticeable result is that the share of $b_E$ increases substantially for the last two

---

[4] There is no information on the taxes on wages and salaries.

age groups at about 1.8 percent and 4.78 percent respectively, indicating that ethnic

inequality is more prominent amongst more experienced workers. The share of $i_{GE}$,

while remaining small in absolute term for all age groups, has increased substantially in

proportional terms, confirming the hypothesis about the masking effect of age on gender

and ethnicity inequalities.

Furthermore, gender inequality measured by $b_G$ is below one percent for the youngest

age group of 15-24 but quickly rises through child bearing and family caring ages

before starting to fall for those aged 55-64. Also, for the age group of 15-24 although

the gender-ethnicity interaction inequality is very small, it is more than half the size of

gender or ethnic inequality. This suggests that compared with gender and ethnic

inequalities a large amount of inequality is due to interaction between gender and

ethnicity for the 15-24 year olds.

Since Australia accounts for over 75 percent of the sample, we have experimented with

first grouping all other nine regions together as a single group, and second excluding

Australia from the sample. The results for these two cases are reported in the last two

columns of Table 3. The results are largely intact, indicating that the findings of the

base line case are robust to region grouping and to the migrant sub-sample. The only

noticeable difference is that in the case of Australia against all other regions together,

$i_{GE}$ is negative, indicating that $b_G$ and $b_E$ overlap and the overlapping inequality cannot

be attributed solely to either gender or ethnicity.

## 4.  Concluding Remarks

In the above empirical example of labour income inequality in Australia, the gender-ethnicity interaction inequality is found to be very small, compared with other inequality components. One may then question the practical value of conducting such decomposition. We would like to point out that, although the interaction inequality could be very small in practice, knowing its actual value allows us to approximate the total inequality by the remaining non-hierarchical components, which makes the decomposition results even easier to interpret. Moreover, for some other variates, such as occupation and education, ethnicity and region, the interaction is likely to be much stronger.

Lastly, although we focus on bi-variate decomposition here, the method can be generalized to handle decompositions of higher dimensions. The number of interaction terms increases with the number of variates. For example, in the three variate case, there will be totally four interaction terms, three corresponding to the interaction of every two variates and one to the interaction of all three variates. Despite the increasing number of interaction terms, the merit of non-hierarchical decomposition as compared with hierarchical decomposition is also greater. If the number of variates is equal to m, the total number of non-hierarchical, asymmetric decompositions is equal to m factorial (i.e. m!). In comparison, using the hierarchical decomposition, we only need to focus on a single decomposition in which all variates are treated symmetrically.

## References

Adelman, I & Levy, A 1984, 'Decomposing Theil's Index of Inequality into Between and Within Components: a note', *Review of Income and Wealth*, vol. 30, no. 1, pp. 119-21.

Akita, T 2003, 'Decomposing regional income inequality in China and Indonesia using two-stage nested Theil decomposition method', *The Annuals of Regional Science*, vol. 37, pp. 55-77.

Australia Bureau of Statistics 2000, *Household Expenditure Survey: User Guide 1998-99*, Canberra.

Bourguignon, F 1979, 'Decomposable Income Inequality Measures', *Econometrica*, vol. 47, no. 4, pp. 901-20.

Conceicao, P & Galbraith, JK 2000, 'Constructing Long and Dense Time-Series of Inequality Using the Theil Index', *Eastern Economic Journal*, vol. 26, no. 1, pp. 61-74.

Cowell, FA 1980, 'On the Structure of Additive Inequality Measure', *The Review of Economic Studies*, vol. 47, no. 3, pp. 521-31.

---- 1985, 'Multilevel Decomposition of Theil's Index of Inequality: a note', *Review of Income and Wealth*, vol. 31, no. 2, pp. 201-5.

Cowell, FA & Kuga, K 1981, 'Additivity and the Entropy Concept: An Axiomatic Approach to Inequality Measurement', *Journal of Economic Theory*, vol. 25, no. 1, pp. 131-43.

Harding, A & Greenwell, H 2002, 'Trends in Income and Consumption Inequality in Australia', paper presented to The 27th General Conference of The International Association for Research in Income and Wealth, Stockholm, August 18-24.

Mishra, P & Parikh, A 1992, 'Household Consumer Expenditure Inequalities in India: A Decomposition Analysis', *Review of Income and Wealth*, vol. 38, no. 2, pp. 225-36.

Shorrocks, AF 1980, 'The Class of Additively Decomposable Inequality Measures', *Econometrica*, vol. 48, no. 3, pp. 613-25.

**Table 1 Numerical Example**

| Total income, total number of individuals | Males | Females |
|---|---|---|
| Native | $Y_{na,m}$, 5 | 4, 2 |
| Non-native | 80, 10 | 16, 4 |

**Table 2 Summary Statistics of Weekly Personal Gross Labour Income in**

**Australia, 1998-99**

|  | Country of Birth | Mean | Median | Minimum | Maximum | Standard Deviation | Population share in the sample (%) |
|---|---|---|---|---|---|---|---|
| Male | Australia | 808 | 730 | 7 | 6284 | 496 | 26.37 |
|  | Other Oceania and Antarctica | 898 | 800 | 41 | 3694 | 593 | 1.16 |
|  | North-West Europe | 946 | 825 | 1 | 5709 | 547 | 3.82 |
|  | Southern and Eastern Europe | 746 | 693 | 56 | 2412 | 383 | 1.12 |
|  | North Africa and Middle East | 633 | 634 | 27 | 1575 | 410 | 0.32 |
|  | South-East Asia | 694 | 651 | 15 | 2053 | 381 | 0.90 |
|  | North-East Asia | 804 | 700 | 80 | 1942 | 366 | 0.50 |
|  | Southern and Central Asia | 1000 | 770 | 259 | 5709 | 1055 | 0.63 |
|  | Americas | 841 | 722 | 50 | 2288 | 471 | 0.42 |
|  | Sub-Saharan Africa | 937 | 752 | 114 | 2832 | 481 | 0.46 |
|  | Total | 824 | 742 | 1 | 6284 | 515 | 35.69 |
| Female | Australia | 523 | 500 | 2 | 2970 | 307 | 49.55 |
|  | Other Oceania and Antarctica | 600 | 524 | 40 | 2541 | 397 | 2.06 |
|  | North-West Europe | 589 | 560 | 50 | 2235 | 324 | 6.45 |
|  | Southern and Eastern Europe | 535 | 487 | 20 | 1123 | 245 | 1.64 |
|  | North Africa and Middle East | 494 | 480 | 115 | 1200 | 259 | 0.18 |
|  | South-East Asia | 616 | 550 | 30 | 2100 | 367 | 1.93 |
|  | North-East Asia | 558 | 549 | 35 | 1627 | 356 | 0.67 |
|  | Southern and Central Asia | 602 | 528 | 70 | 1067 | 231 | 0.47 |
|  | Americas | 475 | 458 | 12 | 1001 | 231 | 0.74 |
|  | Sub-Saharan Africa | 626 | 550 | 150 | 1500 | 312 | 0.63 |
|  | Total | 537 | 507 | 2 | 2970 | 313 | 64.31 |

**Table 3 Percentage Shares of Various Components of Theil Indexes**

| | Ten regions | | | | | | Australia vs other regions together | All regions excluding Australia |
|---|---|---|---|---|---|---|---|---|
| | All ages | 15-24 | 25-34 | 35-44 | 45-54 | 55-64 | All ages | All ages |
| $L$ | | | | | | | | |
| $w_{GE} = w_{EG}$ | 88.81 | 98.02 | 88.77 | 83.42 | 81.90 | 81.64 | 89.29 | 87.75 |
| $b_E$ | 0.88 | 0.54 | 0.85 | 0.77 | 1.80 | 4.78 | 0.52 | 1.58 |
| $b_G$ | 10.30 | 0.84 | 9.89 | 15.32 | 15.72 | 13.31 | 10.30 | 10.10 |
| $i_{GE} = i_{EG}$ | 0.01 | 0.60 | 0.49 | 0.50 | 0.59 | 0.27 | -0.12 | 0.57 |
| | | | | | | | | |
| $T$ | | | | | | | | |
| $w_{GE} = w_{EG}$ | 86.66 | 97.23 | 86.76 | 80.77 | 79.20 | 78.97 | 87.32 | 85.99 |
| $b_E$ | 1.06 | 0.81 | 1.01 | 0.88 | 2.09 | 5.83 | 0.62 | 1.75 |
| $b_G$ | 12.24 | 1.21 | 11.51 | 17.75 | 17.92 | 14.47 | 12.24 | 11.41 |
| $i_{GE} = i_{EG}$ | 0.04 | 0.76 | 0.72 | 0.60 | 0.79 | 0.74 | -0.17 | 0.86 |

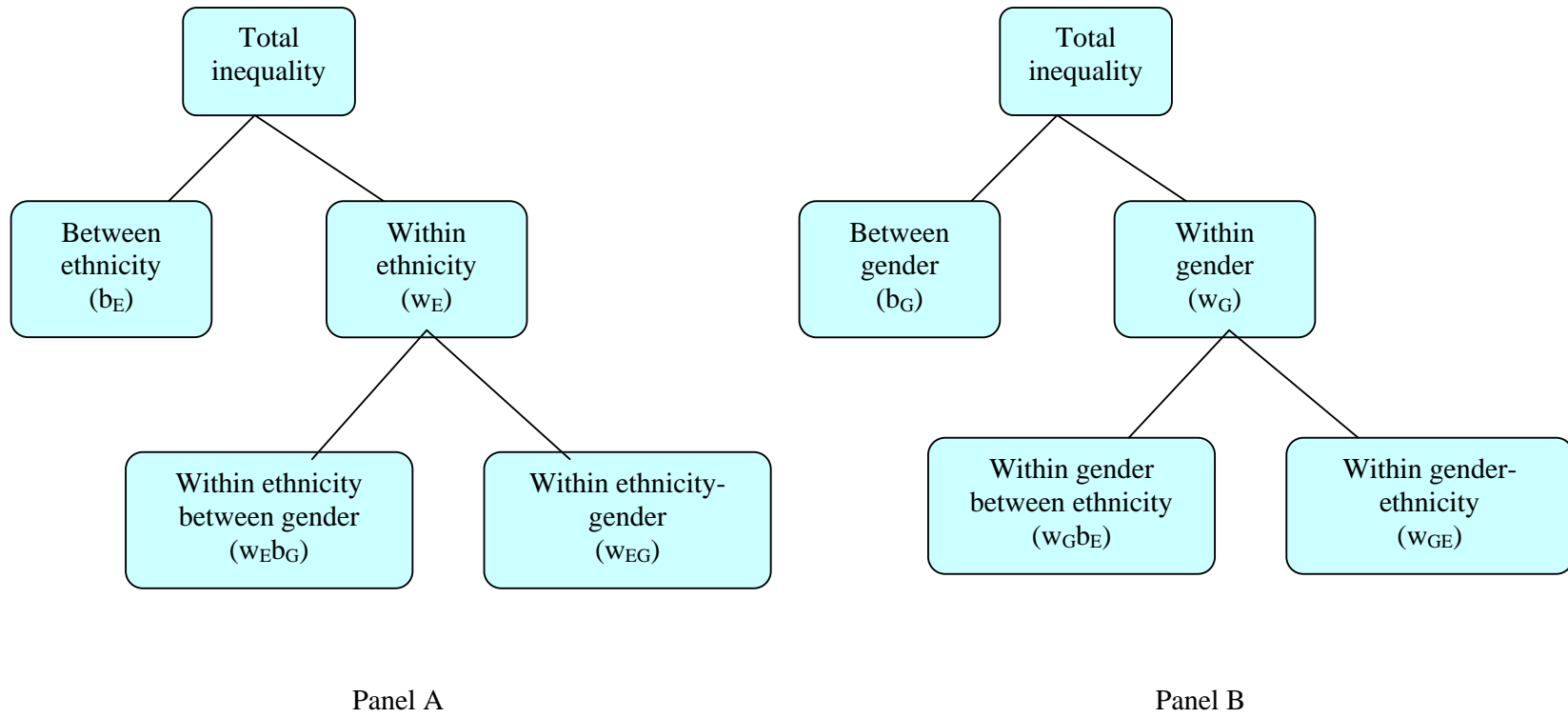**Figure 1 Hierarchical Bivariate Decomposition of Theil Indexes**



Panel A

Panel B

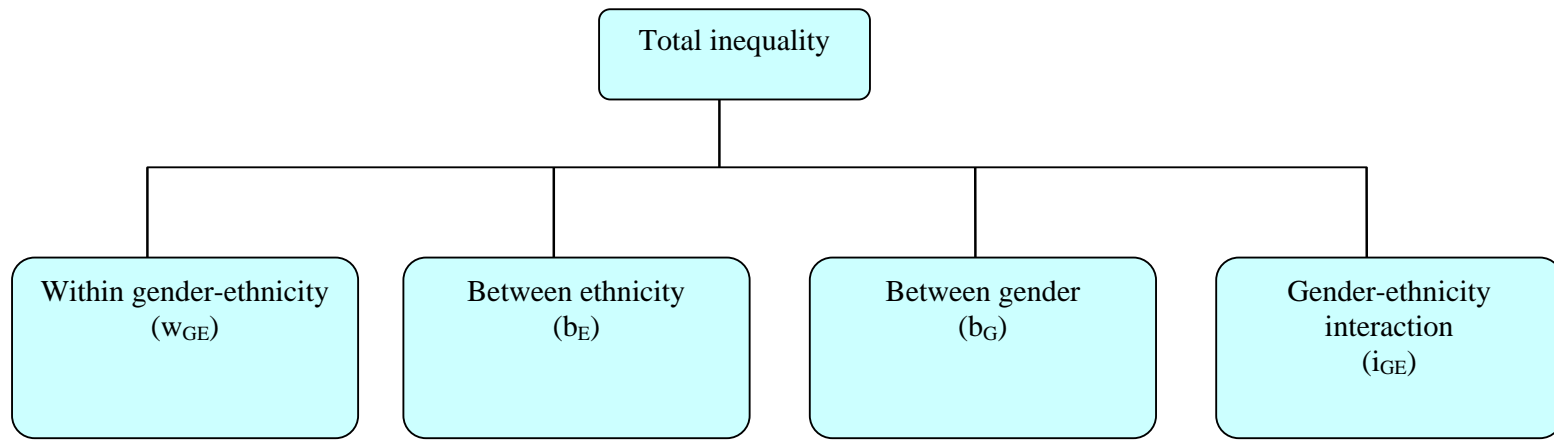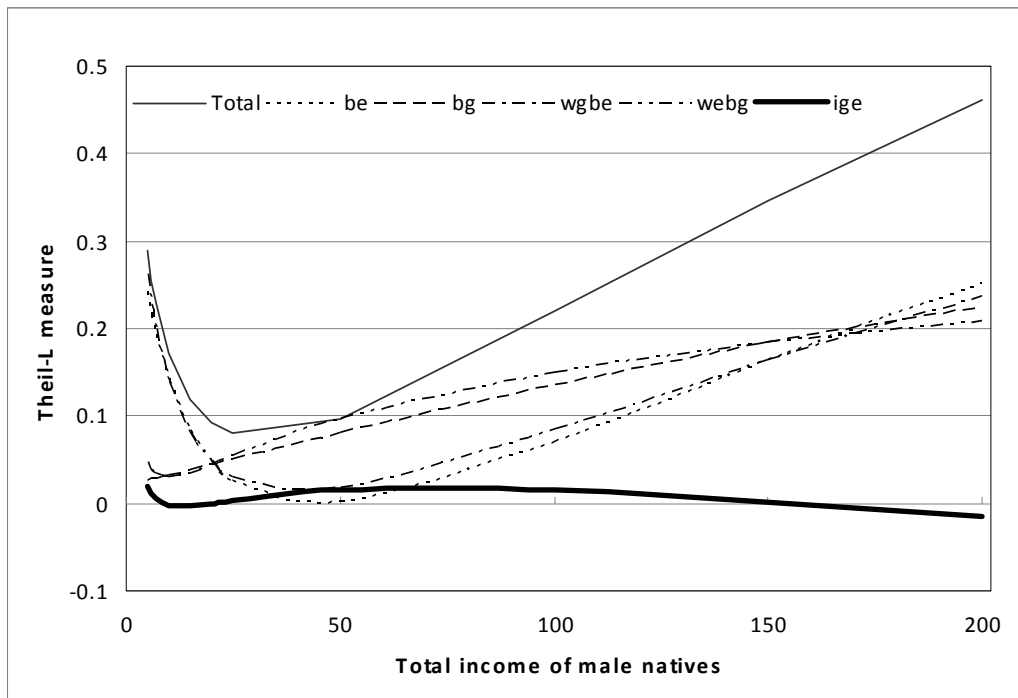**Figure 2 Non-hierarchical Bivariate Decomposition of Theil indexes**

**Figure 3 Changes in Theil-L components with the total income of male natives**

**Appendix (not for publication, but available to readers on request)**

Derivation of equation (2).

In Panel A of Figure 1, the Theil-L index is first composed into:

$$L = \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \left[ \log \left( \frac{N_{egi}/N_e}{Y_{egi}/Y_e} \right) + \log \left( \frac{N_e/N}{Y_e/Y} \right) \right]$$

$$= \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \log \left( \frac{N_{egi}/N_e}{Y_{egi}/Y_e} \right) + \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \log \left( \frac{N_e/N}{Y_e/Y} \right)$$

$$= \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \log \left( \frac{N_{egi}/N_e}{Y_{egi}/Y_e} \right) + \sum_e \left( \frac{N_e}{N} \right) \log \left( \frac{N_e/N}{Y_e/Y} \right)$$

where $N_e = \sum_g \sum_i N_{egi}$ = the size of group $e$, $N_{eg} = \sum_i N_{egi}$ = the size of group $eg$, and

so forth.

The within-ethnicity inequality component can be further decomposed:

$$w_E = \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \log \left( \frac{N_{egi}/N_e}{Y_{egi}/Y_e} \right) = \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \left[ \log \left( \frac{N_{egi}/N_{eg}}{Y_{egi}/Y_{eg}} \right) + \log \left( \frac{N_{eg}/N_e}{Y_{eg}/Y_e} \right) \right]$$

$$= \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \log \left( \frac{N_{egi}/N_{eg}}{Y_{egi}/Y_{eg}} \right) + \sum_e \sum_g \sum_i \left( \frac{N_{egi}}{N} \right) \log \left( \frac{N_{eg}/N_e}{Y_{eg}/Y_e} \right)$$

$$= \sum_e \sum_g \sum_i \left( \frac{N_e}{N} \frac{N_{eg}}{N_e} \frac{N_{egi}}{N_{eg}} \right) \log \left( \frac{N_{egi}/N_{eg}}{Y_{egi}/Y_{eg}} \right) + \sum_e \sum_g \left( \frac{N_{eg}}{N} \right) \log \left( \frac{N_{eg}/N_e}{Y_{eg}/Y_e} \right)$$

$$= \sum_e \left( \frac{N_e}{N} \right) \left[ \sum_g \sum_i \left( \frac{N_{eg}}{N_e} \frac{N_{egi}}{N_{eg}} \right) \log \left( \frac{N_{egi}/N_{eg}}{Y_{egi}/Y_{eg}} \right) \right] + \sum_e \sum_g \left( \frac{N_e}{N} \frac{N_{eg}}{N_e} \right) \log \left( \frac{N_{eg}/N_e}{Y_{eg}/Y_e} \right)$$

$$= \sum_e \left( \frac{N_e}{N} \right) \left\{ \sum_g \left( \frac{N_{eg}}{N_e} \right) \left[ \sum_i \left( \frac{N_{egi}}{N_{eg}} \right) \log \left( \frac{N_{egi}/N_{eg}}{Y_{egi}/Y_{eg}} \right) \right] \right\} + \sum_e \left( \frac{N_e}{N} \right) \left[ \sum_g \left( \frac{N_{eg}}{N_e} \right) \log \left( \frac{N_{eg}/N_e}{Y_{eg}/Y_e} \right) \right]$$

Hence we obtain equation (2).

To derive (11), we make use of (10):

$$i_{GE} = b_{EG} - b_E - b_G$$

$$= \sum_g \sum_e \left(\frac{N_{eg}}{N}\right) \log\left(\frac{N_{eg}/N}{Y_{eg}/Y}\right) - \sum_e \left(\frac{N_e}{N}\right) \log\left(\frac{N_e/N}{Y_e/Y}\right) - \sum_g \left(\frac{N_g}{N}\right) \log\left(\frac{N_g/N}{Y_g/Y}\right)$$

$$= \sum_g \sum_e \left(\frac{N_{eg}}{N}\right) \log\left(\frac{N_{eg}/N}{Y_{eg}/Y}\right) - \sum_e \sum_g \left(\frac{N_{eg}}{N}\right) \log\left(\frac{N_e/N}{Y_e/Y}\right) - \sum_e \sum_g \left(\frac{N_{eg}}{N}\right) \log\left(\frac{N_g/N}{Y_g/Y}\right)$$

$$= \sum_e \sum_g \left(\frac{N_{eg}}{N}\right) \log\left\{\left[\frac{N_{eg}/N}{(N_e/N)(N_g/N)}\right]\left[\frac{(Y_e/Y)(Y_g/Y)}{Y_{eg}/Y}\right]\right\}$$

**Non-negativity of Theil indexes**

A basic Theil-L index has a structure of

$$L = \sum_i \frac{N_i}{N} \log\left(\frac{Y/N}{Y_i/N_i}\right) \tag{12}$$

where $Y_i$ is the total income of group $i$, $N_i$ is the total population of group $i$,

$Y = \sum_i Y_i$, and $N = \sum_i N_i$.

The Theil-L index can be rewritten as

$$L = \sum_i \frac{N_i}{N}\left[\log(Y/N) - \log(Y_i/N_i)\right]$$

$$= \log\left[\frac{(Y/N)}{\prod_i (Y_i/N_i)^{N_i/N}}\right] \tag{13}$$

$$= \log\left[\frac{\sum_i (Y_i/N_i)}{\prod_i (Y_i/N_i)^{N_i/N}}\right]$$

Inside the logarithmic function, the numerator and denominator are the arithmetic

mean and geometric mean of the group level average incomes. Since the arithmetic

mean of non-negative real numbers must be greater than or equal to the geometric

mean, and the logarithmic of a value greater than or equal to one must be non-

negative, hence the Theil-L index must be non-negative. The reason for the non-negativity of the Theil-T index is the same.