

IZA DP No. 7743

A Detailed Decomposition of Synthetic Cohort Analysis

Tavis Barr
Carl Lin

November 2013

A Detailed Decomposition of Synthetic Cohort Analysis

Tavis Barr

Beijing Normal University

Carl Lin

*Beijing Normal University
and IZA*

Discussion Paper No. 7743
November 2013

IZA

P.O. Box 7240
53072 Bonn
Germany

Phone: +49-228-3894-0
Fax: +49-228-3894-180
E-mail: iza@iza.org

Any opinions expressed here are those of the author(s) and not those of IZA. Research published in this series may include views on policy, but the institute itself takes no institutional policy positions. The IZA research network is committed to the IZA Guiding Principles of Research Integrity.

The Institute for the Study of Labor (IZA) in Bonn is a local and virtual international research center and a place of communication between science, politics and business. IZA is an independent nonprofit organization supported by Deutsche Post Foundation. The center is associated with the University of Bonn and offers a stimulating research environment through its international network, workshops and conferences, data service, project support, research visits and doctoral program. IZA engages in (i) original and internationally competitive research in all fields of labor economics, (ii) development of policy concepts, and (iii) dissemination of research results and concepts to the interested public.

IZA Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

ABSTRACT

A Detailed Decomposition of Synthetic Cohort Analysis^{*}

Social scientists are often interested in assessing relative changes between two groups over time, for example, the convergence of black-white wages from 1940 to 1990. In such situations, we need a control group for both treatment groups to remove biases resulting from time trends and unobserved factors. Smith and Welch (1989) propose a decomposition technique that handles this situation using a difference-in-differences like method, attributing relative changes to four different effects. However, a method for specifying the contribution of every variable in the equation, referred to as the detailed decomposition, has not been developed. We present a detailed decomposition and provide a Stata estimator for practitioners to implement our method. Using the convergence of black-white wage between 1960 and 1970 as an example, our detailed decomposition result shows that education accounts for 73 percent of the value of change in characteristics while experience explains only 3 percent.

JEL Classification: C20, J70

Keywords: Oaxaca-Blinder decomposition, detailed decomposition, synthetic cohort, wage gap, repeated cross sections

Corresponding author:

Carl Lin
Beijing Normal University
No.19 Xijiekou Outer St.
Haidian Dist.
Beijing
P.R. China
E-mail: csmlin@bnu.edu.cn

^{*} We are grateful to Ira Gang and Myeong-Su Yun for extended discussion and helpful comments.

1. Introduction

Decomposition techniques for linear regression models have been widely used in social research for many decades. One common such technique is the Oaxaca-Blinder decomposition (Blinder 1973; Oaxaca 1973). The technique uses the output from regression models, attributing the difference in outcomes between two groups (in a mean or proportion) to either differences in endowments or to differences in the effect of coefficients (Powers et al. 2011).

Often we are interested in assessing relative changes between two groups over time. In such situations, we need a control group for both of the treatment groups to remove biases resulting from time trends and unobserved time-constant factors. This difference-in-differences like method has been used to analyze, for example in Smith and Welch (1989, SW) and Heckman et al. (2000, HLT), the convergence of black and white wages between 1940 and 1990.¹ SW and HLT attribute the change in the black-white wage gap in two Census years to the changes due to changing characteristics and to those due to changing returns to characteristics. However, a method for specifying the contribution of every variable in the equation, referred to as the detailed decomposition, has not been developed for the four-way decomposition. In this paper, we present such a detailed decomposition and provide a Stata estimator for practitioners to implement our method.²

We begin by setting up the theoretical framework and show how it can be applied to the existing studies. Using the convergence of the black-white wage gap during the decade of 1960-1970 as an example, we find that the wages of black men grew 12.3 log points faster than the wages of white men—of which 5.15 log points (42 percent of the growth) can be attributed to changing characteristics. Our detailed decomposition further shows that out of the 5.15 log

¹ The four groups are, for example, 1940 black wage, 1940 white wage, 1990 black wage and 1990 white wage.

² The Stata estimator “quaddec” can be downloaded from www.tavisbarr.com and <https://sites.google.com/site/carlshuminglin>.

points of changing characteristics, 3.77 log points (73 percent) is due to convergence in schooling levels, 1.19 log points (23 percent) is due to changing places of residence (primarily migration of black men out of the South), and very little is due to differential growth in work experience.

The rest of the paper is organized as follows. We set up the theoretical framework and provide its application in Section 2. In Section 3, we present an empirical example. Finally, Section 4 concludes the paper.

2. Theoretical Framework

Suppose that a dependent variable Y is a function $F(\bullet)$ of a linear combination of independent variables X , that is $Y = F(X\beta)$, where β is a vector of coefficients. Suppose $Y_1 = Y_A - Y_B$ and $Y_2 = Y_C - Y_D$ are two *relative* outcome variables where Y_B and Y_D are the reference group variables for Y_A and Y_C , respectively. For example, in studying racial wage differentials Y_A and Y_C can be regarded as current-year and base-year black male wages, while Y_B and Y_D are the corresponding current-year and base-year wages for whites (SW; HLT). Or in the synthetic cohort analysis of immigration, Y_A and Y_C can be referred to current-census earnings of an old immigrant cohort A and a new immigrant cohort C , while Y_B and Y_D are the corresponding current-census earnings for natives (Borjas 1985). The mean *relative* outcome difference between the two groups, e.g., $\bar{Y}_1 - \bar{Y}_2$, can be written as

$$\bar{Y}_1 - \bar{Y}_2 = (\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D) = \left[\overline{F(X_A\beta_A)} - \overline{F(X_B\beta_B)} \right] - \left[\overline{F(X_C\beta_C)} - \overline{F(X_D\beta_D)} \right]. \quad (1)$$

Eq. (1) can be further decomposed as

$$\begin{aligned}
(\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D) &= \left[\overline{F(X_A \beta_D)} - \overline{F(X_B \beta_D)} \right] - \left[\overline{F(X_C \beta_D)} - \overline{F(X_D \beta_D)} \right] & \text{(i)} \\
&+ \left[\overline{F(X_A \beta_C)} - \overline{F(X_C \beta_C)} \right] - \left[\overline{F(X_A \beta_D)} - \overline{F(X_C \beta_D)} \right] & \text{(ii)} \\
&+ \left[\overline{F(X_A \beta_B)} - \overline{F(X_A \beta_D)} \right] - \left[\overline{F(X_B \beta_B)} - \overline{F(X_B \beta_D)} \right] & \text{(iii)} \\
&+ \left[\overline{F(X_A \beta_A)} - \overline{F(X_A \beta_B)} \right] - \left[\overline{F(X_A \beta_C)} - \overline{F(X_A \beta_D)} \right]. & \text{(iv)}
\end{aligned} \tag{2}$$

In a linear regression framework, $Y = F(X\beta) = X\beta + \varepsilon$. Under the least squares assumptions, $E(\beta) = \beta$ and $E(\varepsilon) = 0$. Using the sample means \bar{Y} and \bar{X} as estimates for $E(Y)$ and $E(X)$ and the least squares estimates $\hat{\beta}$ for β , Eq. (2) can be shown as

$$\begin{aligned}
(\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D) &= \hat{\beta}_D \left[(\bar{X}_A - \bar{X}_B) - (\bar{X}_C - \bar{X}_D) \right] & \text{(i)} \\
&+ (\hat{\beta}_C - \hat{\beta}_D) \left[(\bar{X}_A - \bar{X}_C) \right] & \text{(ii)} \\
&+ (\bar{X}_A - \bar{X}_B) \left[(\hat{\beta}_B - \hat{\beta}_D) \right] & \text{(iii)} \\
&+ \bar{X}_A \left[(\hat{\beta}_A - \hat{\beta}_B) - (\hat{\beta}_C - \hat{\beta}_D) \right], & \text{(iv)}
\end{aligned} \tag{3}$$

where the first term (i) in Eq. (3) captures the predicted change in $(\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D)$, the mean *relative* outcome difference, due to changes in mean characteristics valued at group D 's parameter values. The second term (ii) captures the additional change in the mean *relative* outcome difference due to an increase/decrease in mean characteristics levels, fixing the difference of returns between group C and its reference group. The third term (iii) measures the effect of a change in $(\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D)$ due to an increase/decrease in the returns to reference groups' characteristics, valued at the difference between group A and its reference group's mean characteristic levels. The last term (iv) captures the change in $(\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D)$ that occurs because the relative returns to characteristics (e.g., wage structure) of two groups are changing, valued at the mean characteristics of reference group A .

The above decomposition studied at the aggregate level in a linear regression framework is set forth in SW. It is of interest to find the contribution of each variable to the four effects in the mean *relative* outcome difference. The next step is to develop the detailed decomposition.

2.1. Detailed Decomposition

In developing the detailed decomposition, the key question is how to properly weight the contribution of each variable to the four effects in Eq. (2). To obtain a proper weight, we employ the method of Yun (2004) by evaluating the value of the functions in Eq. (2) using mean characteristics and then using a first order Taylor expansion to linearize the four effects around $\bar{X}_D\beta_D$ for (i), $\bar{X}_C\beta_C$ for (ii), $\bar{X}_B\beta_B$ for (iii), and $\bar{X}_A\beta_A$ for (iv), respectively. Eq. (2) can be re-written as

$$\begin{aligned}
(\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D) &= \beta_D [(\bar{X}_A - \bar{X}_B) - (\bar{X}_C - \bar{X}_D)] f(\bar{X}_D\beta_D) \\
&+ [(\beta_C - \beta_D)(\bar{X}_A - \bar{X}_C)] f(\bar{X}_C\beta_C) \\
&+ [(\bar{X}_A - \bar{X}_B)(\beta_B - \beta_D)] f(\bar{X}_B\beta_B) \\
&+ \bar{X}_A [(\beta_A - \beta_B) - (\beta_C - \beta_D)] f(\bar{X}_A\beta_A) \\
&+ R_{(i)} + R_{(ii)} + R_{(iii)} + R_{(iv)},
\end{aligned} \tag{4}$$

where $f(\bar{X}_i\beta_i) = \frac{dF(\bar{X}_i\beta_i)}{d(\bar{X}_i\beta_i)}$, $i = A, B, C, D$. $R_{(i)}$ to $R_{(iv)}$ are the approximation remainder

resulting from evaluating the functions $F(\bullet)$ for the four effects in Eq. (2) at the mean values and by using the first order Taylor expansion. After linearization, the j th weight component for effect (i), (ii), (iii) and (iv) are

$$W_{\Delta X_j}^{(i)} = \frac{\beta_{Dj} [(\bar{X}_{Aj} - \bar{X}_{Bj}) - (\bar{X}_{Cj} - \bar{X}_{Dj})] f(\bar{X}_D\beta_D)}{\sum_{j=1}^J \beta_{Dj} [(\bar{X}_{Aj} - \bar{X}_{Bj}) - (\bar{X}_{Cj} - \bar{X}_{Dj})] f(\bar{X}_D\beta_D)} = \frac{\beta_{Dj} [(\bar{X}_{Aj} - \bar{X}_{Bj}) - (\bar{X}_{Cj} - \bar{X}_{Dj})]}{\sum_{j=1}^J \beta_{Dj} [(\bar{X}_{Aj} - \bar{X}_{Bj}) - (\bar{X}_{Cj} - \bar{X}_{Dj})]},$$

$$W_{\Delta X_j}^{(ii)} = \frac{[(\beta_{Cj} - \beta_{Dj})(\bar{X}_{Aj} - \bar{X}_{Cj})] f(\bar{X}_C \beta_C)}{\sum_{j=1}^J [(\beta_{Cj} - \beta_{Dj})(\bar{X}_{Aj} - \bar{X}_{Cj})] f(\bar{X}_C \beta_C)} = \frac{[(\beta_{Cj} - \beta_{Dj})(\bar{X}_{Aj} - \bar{X}_{Cj})]}{\sum_{j=1}^J [(\beta_{Cj} - \beta_{Dj})(\bar{X}_{Aj} - \bar{X}_{Cj})]},$$

$$W_{\Delta \beta_j}^{(iii)} = \frac{[(\bar{X}_{Aj} - \bar{X}_{Bj})(\beta_{Bj} - \beta_{Dj})] f(\bar{X}_B \beta_B)}{\sum_{j=1}^J [(\bar{X}_{Aj} - \bar{X}_{Bj})(\beta_{Bj} - \beta_{Dj})] f(\bar{X}_B \beta_B)} = \frac{[(\bar{X}_{Aj} - \bar{X}_{Bj})(\beta_{Bj} - \beta_{Dj})]}{\sum_{j=1}^J [(\bar{X}_{Aj} - \bar{X}_{Bj})(\beta_{Bj} - \beta_{Dj})]},$$

and

$$W_{\Delta \beta_j}^{(iv)} = \frac{\bar{X}_{Aj} [(\beta_{Aj} - \beta_{Bj}) - (\beta_{Cj} - \beta_{Dj})] f(\bar{X}_A \beta_A)}{\sum_{j=1}^J \bar{X}_{Aj} [(\beta_{Aj} - \beta_{Bj}) - (\beta_{Cj} - \beta_{Dj})] f(\bar{X}_A \beta_A)} = \frac{\bar{X}_{Aj} [(\beta_{Aj} - \beta_{Bj}) - (\beta_{Cj} - \beta_{Dj})]}{\sum_{j=1}^J \bar{X}_{Aj} [(\beta_{Aj} - \beta_{Bj}) - (\beta_{Cj} - \beta_{Dj})]},$$

where $\sum_{j=1}^J W_{\Delta X_j}^{(i)} = \sum_{j=1}^J W_{\Delta X_j}^{(ii)} = \sum_{j=1}^J W_{\Delta \beta_j}^{(iii)} = \sum_{j=1}^J W_{\Delta \beta_j}^{(iv)} = 1$.

The mean *relative* outcome difference can now be expressed in terms of the overall components as a sum of weighted sums of the unique contributions

$$\begin{aligned} (\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D) &= (i) + (ii) + (iii) + (iv) \\ &= \sum_{j=1}^J W_{\Delta X_j}^{(i)} (i) + \sum_{j=1}^J W_{\Delta X_j}^{(ii)} (ii) + \sum_{j=1}^J W_{\Delta \beta_j}^{(iii)} (iii) + \sum_{j=1}^J W_{\Delta \beta_j}^{(iv)} (iv) \\ &= \sum_{j=1}^J (i)_j + \sum_{j=1}^J (ii)_j + \sum_{j=1}^J (iii)_j + \sum_{j=1}^J (iv)_j. \end{aligned} \quad (5)$$

2.2. Identification Problem: Normalization of Dummy Variables

It is well-known that detailed decomposition suffers from an identification problem (Oaxaca and Ransom 1999) as long as there are dummy variables in the regression equation. That is, the detailed coefficients effect attributed to dummy variables is not invariant to the choice of reference groups. It turns out that detailed decomposition of Eq. (5) cannot escape from this identification problem. Gardeazabal and Ugidos (2004) and Yun (2005) have proposed an

intuitive solution to the problem. We follow their methods by constructing a normalized equation and use it to perform decomposition analysis.

Suppose there is a set of K dummy variables, indexed by d_{ik} , and a vector of J continuous variables, x_{ij} , in the function $F(X\beta)$. Under the linear regression framework, $Y_i = F(X_i\beta_i) = X_i\beta_i + \varepsilon_i$ where $i = A, B, C, D$. Without loss of generality, let the first category of dummy variables, e.g., d_{i1} , be the reference category. The equation is written as

$$Y_i = F(X_i\beta_i) = \alpha_i + \sum_{k=2}^K d_{ik}\gamma_{ik} + \sum_{j=1}^J x_{ij} + e_i. \quad (6)$$

The next step is to construct a normalized equation. Once the normalized equation is constructed, the identification problem of the decomposition analysis is automatically resolved. Following the method generalized by Jann (2008) and Yun (2008), define $\gamma_{ik}^* = \gamma_{ik} + c$ and impose restriction $\sum_{k=1}^K \gamma_{ik}^* = 0$; the normalized equation for Eq. (6) is then written as:

$$Y_i = F(X_i^*\beta_i^*) = \alpha_i^* + \sum_{k=1}^K d_{ik}\gamma_{ik}^* + \sum_{j=1}^J x_{ij}\theta_{ij}^* + e_i, \quad (7)$$

where $i = A, B, C, D$. The solution for the constraint is $c = -\bar{\gamma} = -\frac{1}{K} \sum_{k=1}^K \gamma_{ik}$, where γ_{i1} is set to zero. The transformed normalized equation then is

$$Y_i = (\alpha_i + \bar{\gamma}) + \sum_{k=1}^K d_{ik}(\gamma_{ik} - \bar{\gamma}) + \sum_{j=1}^J x_{ij}\theta_{ij} + e_i, \quad (8)$$

which is mathematically equivalent to the untransformed Eq. (6). As shown in Yun (2008), using the normalized Eq. (8) can resolve the identification problem in the detailed decomposition without altering the sizes and asymptotic variances. Hence, inferences about the effects are not affected by normalization.

2.3. Significance tests for aggregate effect and individual factors

Since \bar{X}_i and $\hat{\beta}_i$ are uncorrelated by assumption and assuming that Y_1 and Y_2 are independent, the mean and the variance for $Y_1 - Y_2$ are

$$\begin{aligned}
E(\bar{Y}_1 - \bar{Y}_2) &= E((\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D)) \\
&= E(\hat{\beta}_D [(\bar{X}_A - \bar{X}_B) - (\bar{X}_C - \bar{X}_D)]) & (i) \\
&+ E([\hat{\beta}_C - \hat{\beta}_D](\bar{X}_A - \bar{X}_C)) & (ii) \\
&+ E([\bar{X}_A - \bar{X}_B](\hat{\beta}_B - \hat{\beta}_D)) & (iii) \\
&+ E(\bar{X}_A [(\hat{\beta}_A - \hat{\beta}_B) - (\hat{\beta}_C - \hat{\beta}_D)]), & (iv)
\end{aligned} \tag{9}$$

and

$$\begin{aligned}
V(\bar{Y}_1 - \bar{Y}_2) &= V((\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D)) \\
&= V(\hat{\beta}_D [(\bar{X}_A - \bar{X}_B) - (\bar{X}_C - \bar{X}_D)]) & (i) \\
&+ V([\hat{\beta}_C - \hat{\beta}_D](\bar{X}_A - \bar{X}_C)) & (ii) \\
&+ V([\bar{X}_A - \bar{X}_B](\hat{\beta}_B - \hat{\beta}_D)) & (iii) \\
&+ V(\bar{X}_A [(\hat{\beta}_A - \hat{\beta}_B) - (\hat{\beta}_C - \hat{\beta}_D)]), & (iv)
\end{aligned} \tag{10}$$

Following the method in Jann (2005), the variance estimator for the four effects in Eq. (10) can be separately shown as

$$\begin{aligned}
\hat{\sigma}_{(i)}^2 &= \hat{\beta}_D' (\hat{V}(\bar{X}_A) + \hat{V}(\bar{X}_B)) \hat{\beta}_D + (\bar{X}_A - \bar{X}_B)' \hat{V}(\hat{\beta}_D) (\bar{X}_A - \bar{X}_B) + tr(\hat{V}(\bar{X}_A - \bar{X}_B) \hat{V}(\hat{\beta}_D)) \\
&+ \hat{\beta}_D' (\hat{V}(\bar{X}_C) + \hat{V}(\bar{X}_D)) \hat{\beta}_D + (\bar{X}_C - \bar{X}_D)' \hat{V}(\hat{\beta}_D) (\bar{X}_C - \bar{X}_D) + tr(\hat{V}(\bar{X}_C - \bar{X}_D) \hat{V}(\hat{\beta}_D)),
\end{aligned} \tag{11}$$

$$\begin{aligned}
\hat{\sigma}_{(ii)}^2 &= (\hat{\beta}_C - \hat{\beta}_D)' \hat{V}(\bar{X}_A) (\hat{\beta}_C - \hat{\beta}_D) + \bar{X}_A' (\hat{V}(\hat{\beta}_C) + \hat{V}(\hat{\beta}_D)) \bar{X}_A + tr(\hat{V}(\hat{\beta}_C - \hat{\beta}_D) \hat{V}(\bar{X}_A)) \\
&+ (\hat{\beta}_C - \hat{\beta}_D)' \hat{V}(\bar{X}_C) (\hat{\beta}_C - \hat{\beta}_D) + \bar{X}_C' (\hat{V}(\hat{\beta}_C) + \hat{V}(\hat{\beta}_D)) \bar{X}_C + tr(\hat{V}(\hat{\beta}_C - \hat{\beta}_D) \hat{V}(\bar{X}_C)),
\end{aligned} \tag{12}$$

$$\begin{aligned}
\hat{\sigma}_{(iii)}^2 &= (\bar{X}_A - \bar{X}_B)' \hat{V}(\hat{\beta}_B) (\bar{X}_A - \bar{X}_B) + \hat{\beta}_B' (\hat{V}(\bar{X}_A) + \hat{V}(\bar{X}_B)) \hat{\beta}_B + tr(\hat{V}(\bar{X}_A - \bar{X}_B) \hat{V}(\hat{\beta}_B)) \\
&+ (\bar{X}_A - \bar{X}_B)' \hat{V}(\hat{\beta}_D) (\bar{X}_A - \bar{X}_B) + \hat{\beta}_D' (\hat{V}(\bar{X}_A) + \hat{V}(\bar{X}_B)) \hat{\beta}_D + tr(\hat{V}(\bar{X}_A - \bar{X}_B) \hat{V}(\hat{\beta}_D)),
\end{aligned} \tag{13}$$

and

$$\begin{aligned} \hat{\sigma}_{(iv)}^2 = & \bar{X}_A' (\hat{V}(\hat{\beta}_A) + \hat{V}(\hat{\beta}_B)) \bar{X}_A + (\hat{\beta}_A - \hat{\beta}_B)' \hat{V}(\bar{X}_A) (\hat{\beta}_A - \hat{\beta}_B) + \text{tr}(\hat{V}(\bar{X}_A) \hat{V}(\hat{\beta}_A - \hat{\beta}_B)) \\ & + \bar{X}_A' (\hat{V}(\hat{\beta}_C) + \hat{V}(\hat{\beta}_D)) \bar{X}_A + (\hat{\beta}_C - \hat{\beta}_D)' \hat{V}(\bar{X}_A) (\hat{\beta}_C - \hat{\beta}_D) + \text{tr}(\hat{V}(\bar{X}_A) \hat{V}(\hat{\beta}_C - \hat{\beta}_D)). \end{aligned} \quad (14)$$

Combining (11), (12), (13) and (14), we get

$$\hat{V}(\bar{Y}_1 - \bar{Y}_2) = \hat{\sigma}_{(i)}^2 + \hat{\sigma}_{(ii)}^2 + \hat{\sigma}_{(iii)}^2 + \hat{\sigma}_{(iv)}^2, \quad (15)$$

which constitutes the basis of the significance test for the aggregate effect $\bar{Y}_1 - \bar{Y}_2$. The test hypothesis and the test statistic under the null are

$$H_0 : \bar{Y}_1 - \bar{Y}_2 = 0 \text{ vs. } H_1 : \bar{Y}_1 - \bar{Y}_2 \neq 0,$$

and

$$t_{\bar{Y}_1 - \bar{Y}_2} = \frac{\bar{Y}_1 - \bar{Y}_2}{\sqrt{\hat{V}(\bar{Y}_1 - \bar{Y}_2)}},$$

where the test statistic is asymptotically normally distributed.

Let E_l represent the four effects in Eq. (9) and $l = (i), (ii), (iii), (iv)$. Whether E_l are statistically meaningful in explaining the mean *relative* outcome difference between $Y_1 = Y_A - Y_B$ and $Y_2 = Y_C - Y_D$, four null hypotheses are tested:

$$H_0 : E_l = 0 \text{ vs. } H_1 : E_l \neq 0,$$

where $l = (i), (ii), (iii), (iv)$. The test statistics, $t_{E_l} = E_l / \hat{\sigma}_l$, are asymptotically normally distributed.

For significance tests of the four effects at the individual factor level, let

$$\begin{aligned} E_{(i)}^j = & W_{\Delta X_j}^{(i)} \hat{\beta}_{Dj} [(\bar{X}_{Aj} - \bar{X}_{Bj}) - (\bar{X}_{Cj} - \bar{X}_{Dj})] \quad , \quad E_{(ii)}^j = W_{\Delta X_j}^{(ii)} [(\hat{\beta}_{Cj} - \hat{\beta}_{Dj})(\bar{X}_{Aj} - \bar{X}_{Cj})] \quad , \\ E_{(iii)}^j = & W_{\Delta \beta_j}^{(iii)} [(\bar{X}_{Aj} - \bar{X}_{Bj})(\hat{\beta}_{Bj} - \hat{\beta}_{Dj})] \quad , \quad \text{and} \quad E_{(iv)}^j = W_{\Delta \beta_j}^{(iv)} \bar{X}_{Aj} [(\hat{\beta}_{Aj} - \hat{\beta}_{Bj}) - (\hat{\beta}_{Cj} - \hat{\beta}_{Dj})] \quad \text{be} \end{aligned}$$

contributions of factor j to $(\bar{Y}_A - \bar{Y}_B) - (\bar{Y}_C - \bar{Y}_D)$ in the four effects E_l , respectively. The asymptotic variances of E_l are

$$\begin{aligned}\hat{\sigma}_{(i)j}^2 &= \hat{\beta}_D^{j'} (\hat{V}(\bar{X}_A) + \hat{V}(\bar{X}_B)) \hat{\beta}_D^j + (\bar{X}_A^j - \bar{X}_B^j)' \hat{V}(\hat{\beta}_D^j) (\bar{X}_A^j - \bar{X}_B^j) + \text{tr}(\hat{V}(\bar{X}_A - \bar{X}_B) \hat{V}(\hat{\beta}_D^j)) \\ &\quad + \hat{\beta}_D^{j'} (\hat{V}(\bar{X}_C) + \hat{V}(\bar{X}_D)) \hat{\beta}_D^j + (\bar{X}_C^j - \bar{X}_D^j)' \hat{V}(\hat{\beta}_D^j) (\bar{X}_C^j - \bar{X}_D^j) + \text{tr}(\hat{V}(\bar{X}_C - \bar{X}_D) \hat{V}(\hat{\beta}_D^j)),\end{aligned}$$

$$\begin{aligned}\hat{\sigma}_{(ii)j}^2 &= (\hat{\beta}_C^j - \hat{\beta}_D^j)' \hat{V}(\bar{X}_A) (\hat{\beta}_C^j - \hat{\beta}_D^j) + \bar{X}_A^{j'} (\hat{V}(\hat{\beta}_C^j) + \hat{V}(\hat{\beta}_D^j)) \bar{X}_A^j + \text{tr}(\hat{V}(\hat{\beta}_C^j - \hat{\beta}_D^j) \hat{V}(\bar{X}_A)) \\ &\quad + (\hat{\beta}_C^j - \hat{\beta}_D^j)' \hat{V}(\bar{X}_C) (\hat{\beta}_C^j - \hat{\beta}_D^j) + \bar{X}_C^{j'} (\hat{V}(\hat{\beta}_C^j) + \hat{V}(\hat{\beta}_D^j)) \bar{X}_C^j + \text{tr}(\hat{V}(\hat{\beta}_C^j - \hat{\beta}_D^j) \hat{V}(\bar{X}_C)),\end{aligned}$$

$$\begin{aligned}\hat{\sigma}_{(iii)j}^2 &= (\bar{X}_A^j - \bar{X}_B^j)' \hat{V}(\hat{\beta}_B^j) (\bar{X}_A^j - \bar{X}_B^j) + \hat{\beta}_B^{j'} (\hat{V}(\bar{X}_A) + \hat{V}(\bar{X}_B)) \hat{\beta}_B^j + \text{tr}(\hat{V}(\bar{X}_A - \bar{X}_B) \hat{V}(\hat{\beta}_B^j)) \\ &\quad + (\bar{X}_A^j - \bar{X}_B^j)' \hat{V}(\hat{\beta}_D^j) (\bar{X}_A^j - \bar{X}_B^j) + \hat{\beta}_D^{j'} (\hat{V}(\bar{X}_A) + \hat{V}(\bar{X}_B)) \hat{\beta}_D^j + \text{tr}(\hat{V}(\bar{X}_A - \bar{X}_B) \hat{V}(\hat{\beta}_D^j)),\end{aligned}$$

and

$$\begin{aligned}\hat{\sigma}_{(iv)j}^2 &= \bar{X}_A^{j'} (\hat{V}(\hat{\beta}_A^j) + \hat{V}(\hat{\beta}_B^j)) \bar{X}_A^j + (\hat{\beta}_A^j - \hat{\beta}_B^j)' \hat{V}(\bar{X}_A) (\hat{\beta}_A^j - \hat{\beta}_B^j) + \text{tr}(\hat{V}(\bar{X}_A) \hat{V}(\hat{\beta}_A^j - \hat{\beta}_B^j)) \\ &\quad + \bar{X}_A^{j'} (\hat{V}(\hat{\beta}_C^j) + \hat{V}(\hat{\beta}_D^j)) \bar{X}_A^j + (\hat{\beta}_C^j - \hat{\beta}_D^j)' \hat{V}(\bar{X}_A) (\hat{\beta}_C^j - \hat{\beta}_D^j) + \text{tr}(\hat{V}(\bar{X}_A) \hat{V}(\hat{\beta}_C^j - \hat{\beta}_D^j)).\end{aligned}$$

So the null hypotheses for factor j are

$$H_0 : E_l^j = 0 \text{ vs. } H_1 : E_l^j \neq 0, \quad l = (i), (ii), (iii), (iv).$$

The test statistics under the null, $t_{E_l^j} = E_l^j / \hat{\sigma}_{E_l^j}$, are asymptotically normally distributed.

2.4. Application

In this section we show how our method in the preceding section can be applied to repeated cross section data, particularly to synthetic cohort analysis. The first application is the aggregate decomposition method set forth in SW and re-examined in HLT.

Following the notations in HLT, let t be the current year and τ a base year. Let $\bar{\mathbf{z}}_t^w$, $\bar{\mathbf{z}}_t^b$, $\bar{\mathbf{z}}_\tau^w$, $\bar{\mathbf{z}}_\tau^b$ denote the mean vectors of black and white characteristics included in the earnings model and $\boldsymbol{\gamma}_t^w$, $\boldsymbol{\gamma}_t^b$, $\boldsymbol{\gamma}_\tau^w$, $\boldsymbol{\gamma}_\tau^b$ denote the associated vectors of coefficients. HLT and SW show the change

in log black wages minus log white wages between time periods t and τ can be decomposed in the following way,

$$\begin{aligned}
[(\bar{\mathbf{z}}_t^b \gamma_t^b - \bar{\mathbf{z}}_t^w \gamma_t^w) - (\bar{\mathbf{z}}_\tau^b \gamma_\tau^b - \bar{\mathbf{z}}_\tau^w \gamma_\tau^w)] &= [(\bar{\mathbf{z}}_t^b - \bar{\mathbf{z}}_t^w) - (\bar{\mathbf{z}}_\tau^b - \bar{\mathbf{z}}_\tau^w)] \gamma_\tau^w && \text{(Main Effect)} \\
&+ (\bar{\mathbf{z}}_t^b - \bar{\mathbf{z}}_\tau^b) (\gamma_\tau^b - \gamma_\tau^w) && \text{(Race Interaction Effect)} \\
&+ (\bar{\mathbf{z}}_t^b - \bar{\mathbf{z}}_t^w) (\gamma_t^w - \gamma_\tau^w) && \text{(Year Interaction Effect)} \\
&+ \bar{\mathbf{z}}_t^b [(\gamma_t^b - \gamma_t^w) - (\gamma_\tau^b - \gamma_\tau^w)]. && \text{(Race-Year interaction Effect)}
\end{aligned} \tag{16}$$

The first two terms measure the contribution of changing mean characteristics, valued at base-year returns. The second two terms measure the contribution of changing returns. In addition to the four effects, social scientists are most interested in what factor can account for the changing black-white wage gap between two periods. Our detailed decomposition method provides a solution.

3. An Empirical Example: Smith and Welch (1989)

SW and HLT examine the change in the gap between the wages of black and white men between 1940 and 1980 using decennial U.S. Census data from each of the intervening census years; HLT focuses specifically on the decade of 1960-1970.³ Their overall findings are that convergence between black and white men in both the quantity of schooling attained, and in how schooling is rewarded, can explain a large part of the improvement in the racial wage gap, while migration of black men out of the South in this period is also a significant contributor. Table 14 of SW decomposes the effect of education on the change in the racial wage gap between 1940 and 1980 using the formula in Eq.(16), while Table 2 of HLT shows a similar decomposition for

³ As described in SW, the sample is restricted to black and white men, aged 15-64, who did not reside in group quarters, were not in the military, were not attending school (unless they worked at least 50 weeks), were not self-employed (unless they worked in agriculture), worked at least 27 weeks, earned between \$6.25 and \$625 per week (\$10 and \$1250 per week in 1970), and did not earn the top-coded wage value unless they worked at least 50 weeks. Weekly wages are calculated from annual wages, and top-coded wages are adjusted, as described in their Footnote 13.

the decade of 1960-1970. In our example, we also use the decade of 1960-1970, but for simplicity's sake, we use the regression specification of SW.⁴

As Table 1 shows, during the decade of the 1960s, the wages of black men grew 12.3 log points faster than the wages of white men. The Main Effect shows what would have happened to the wage gap based on changing characteristics, if all men were treated the way white men were in 1960. The effect would have been a closing in the wage gap by 5.15 log points, of which 3.77 log points (73 percent) is due to convergence in schooling levels, 1.19 log points (23 percent) is due to changing places of residence (primarily migration of black men out of the South), and very little is due to differential growth in work experience. The Race Interaction effect shows the degree to which the Main Effect would be attenuated because black men were treated worse in 1960. This overall attenuation was fairly small (roughly 1.6 log points), in which a worse return to education was balanced out by a slightly positive effect of urban residence. The Race-Year Interaction effect shows how black workers in 1970 were better off simply because the rewards for their attributes grew faster than the rewards for the attributes of white workers. Had the characteristics of black workers simply stayed constant at the 1970 level, the racial wage gap would have closed by 9.2 log points. About 24 percent of this can be explained by an improved return to education, while the constant term (reflecting an improvement in the relative pay of black workers unrelated to any of the variables in the regression) simply dwarfs the other terms, showing a 21-log-point unexplained improvement in relative pay. On the other hand, their relative wages would have deteriorated about 3 log points from the changing treatment of

⁴ The main difference between the specifications is that HLT include dummy variables for each of the nine census regions, and for state of birth, interacted with years of schooling. For our purposes, this makes the uninteracted state and region dummies more difficult to interpret as an effect of geography, since they may pick up some of the effect of education. We were able to replicate Heckman et al.'s results. (There are several PUMS microsamples available, and we are not sure which combination they used, so we cannot match their exact numbers, but the differences between our numbers using their specification and the numbers they report in Table 2 are not statistically significant.) However, in addition to the difficulty of interpreting the dummies, using this many dummy variables in a regression also makes the estimated geographic effects somewhat unstable.

location (mainly a worsening treatment of urban residents) and by about 11 log points from a worsening treatment of work experience (reflecting the fact that older black cohorts were losing out relative to their younger counterparts during this period). The Year Interaction shows the degree to which this improving treatment of black workers was attenuated because black workers had less valuable characteristics in 1970 and were therefore less able to benefit from the growing rewards to their human capital (taking the improvement in the treatment of white men as the numeraire). Again, the effect is small (about -0.46 log points), in which a negative effect of education is balanced out by a positive effect of geographic location (that is, black men were hurt less than white men by the declining status of urban location).

4. Conclusion

Decomposing the relative wage gains of two groups into a main effect, group interaction, year interaction, and group-year interaction is a well-known and popular means of studying the reasons for these gains. Breaking down this decomposition to separately study the effect of each variable requires careful attention to the treatment of categorical variables, and benefits from a closed analytic formula for the standard errors. Both of these problems have been recently solved for the two-way decomposition, and in this paper we show how to adapt this solution to a four-way decomposition. We illustrated the benefits of this decomposition using the convergence of black-white wage between 1960 and 1970 as an example. We were able to separately quantify the contributions of education, experience, and geography to the overall convergence of black and white wages in the 1960s, where we illustrated both the negative impact of the treatment of aging cohorts, and the overall large size of the residual (unexplained) component of wages in shaping the convergence of black and white wages over the 1960s.

We hope that this contribution makes it easier for other researchers to study relative wage gains using the detailed decomposition methodology.

References

- Blinder, Alan S.** 1973. "Wage Discrimination: Reduced Form and Structural Estimates." *The Journal of Human Resources*, 8(4), 436-55.
- Borjas, George J.** 1985. "Assimilation, Changes in Cohort Quality, and the Earnings of Immigrants." *Journal of Labor Economics*, 3(4), 463-89.
- Gardeazabal, Javier and Arantza Ugidos.** 2004. "More on Identification in Detailed Wage Decompositions." *The Review of Economics and Statistics*, 86(4), 1034-36.
- Jann, Ben.** 2005. "Standard Errors for the Blinder-Oaxaca Decomposition," Stata Users Group, _____. 2008. "A Stata Implementation of the Blinder-Oaxaca Decomposition," Chair of Sociology, ETH Zurich,
- Oaxaca, Ronald L.** 1973. "Male-Female Wage Differentials in Urban Labor Markets." *International Economic Review*, 14(3), 693-709.
- Oaxaca, Ronald L. and Michael R. Ransom.** 1999. "Identification in Detailed Wage Decompositions." *The Review of Economics and Statistics*, 81(1), 154-57.
- Powers, Daniel A., Hirotoshi Yoshioka and Myeong-Su Yun.** 2011. "Mvdcmp: Multivariate Decomposition for Nonlinear Response Models." *Stata Journal*, 11(4), 556-76.
- Smith, James P. and Finis R. Welch.** 1989. "Black Economic Progress after Myrdal." *Journal of Economic Literature*, 27(2), 519-64.
- Yun, Myeong-Su.** 2004. "Decomposition Differences in the First Moment." *Economics Letters*, 82(2), 275-80.
- _____. 2005. "A Simple Solution to the Identification Problem in Detailed Wage Decompositions." *Economic Inquiry*, 43(4), 766-72.
- _____. 2008. "Identification Problem and Detailed Oaxaca Decomposition: A General Solution and Inference." *Journal of Economic and Social Measurement*, 33(1), 27-38.

Table 1 Detailed Decomposition of Black-White Wage Gap, 1960–1970

Log Wage		Coeff.	S.E.	Observations
Current Year: 1970				
	Black (A)	9.61	0.00	33,742
	White (B)	10.06	0.00	333,693
Base Year: 1960				
	Black (C)	9.24	0.00	28,656
	White (D)	9.81	0.00	289,689
Log Wage Gap		Coeff.	S.E.	
	Current Year (A–B)	-44.64	0.36	
	Base Year (C–D)	-56.93	0.40	
	Change in Log Wage Gap (A–B)–(C–D)	12.30	0.54	
Detailed Decomposition		Coeff.	S.E.	Contribution (%)
1. Main Effect		5.15	0.28	100
	Education	3.77	0.00	73.32
	Experience	0.18	0.39	3.48
	Geography	1.19	0.11	23.20
2. Race Interaction Effect		-1.59	0.18	100
	Education	-3.09	0.00	193.86
	Experience	0.17	0.07	-10.44
	Geography	1.33	0.12	-83.41
3. Year Interaction Effect		-0.46	0.12	100
	Education	-1.07	0.00	234.74
	Experience	0.10	0.00	-22.96
	Geography	0.51	0.09	-111.77
4. Race-Year Interaction Effect		9.20	0.49	100
	Education	2.24	0.01	24.39
	Experience	-11.10	2.76	-120.68
	Geography	-3.03	0.82	-32.89
	Constant	21.09	2.27	229.18

Note: the unit is basis point in the log wage gap and detailed decomposition.