

Multivariate Decomposition for Nonlinear Models

Daniel A. Powers
and
Thomas W. Pullum

Department of Sociology
and
Population Research Center
University of Texas at Austin

February 18, 2006

Abstract

This paper illustrates a multivariate decomposition method based on extending the well-known Oaxaca-Blinder approach to logit, probit, and rate, and other models that are intrinsically nonlinear in the parameters. This method has wide applicability in demographic research but has not received a great deal of exposure. Related methods that have been applied in recent work suffer from several shortcomings that are remedied with the proposed methodology. In particular, detailed decompositions are free from problems of path dependency, and the technique is valid for predicted probabilities in the tails of a distribution. We use this approach to investigate compositional and return to risk components of the differential in infant mortality due to respiratory distress syndrome between 1991 and 1998 and the compositional and return to risk contributions to racial/ethnic differences in mortality.

Introduction

This paper describes and illustrates a multivariate decomposition approach that is applicable to many demographic outcomes, which is especially useful for models that are nonlinear in parameters such as binary response, event count, and hazard rate models. The proposed decomposition method has wide applicability in demographic research but has not received a great deal of attention in the literature. To the extent that related methods have been applied in recent work, there are several notable shortcomings that are remedied by the proposed methodology.

Multivariate decomposition is used to partition the difference in mean responses between

groups or over time into components that reflect the difference in the mean levels of model predictors and difference in the effects of those predictors across groups or over time. The Oaxaca-Blinder multivariate decomposition approach is the most familiar and widely used method for linear models (Blinder 1973; Oaxaca 1973) (see also, Winsborough and Dickinson (1969)). This technique has been widely applied to the study of wage differentials to understand the relative roles played by group differences in levels of certain attributes, characteristics, or endowments, and group differences in the effects, or coefficients, on group differences in mean wage rates.¹

Although the Oaxaca-Blinder technique has wide applicability in demographic models whenever group differences in sample means or changes in sample means over time are the focus of inference, many socio-demographic outcomes involve differences in predicted rates or proportions, which may not be suitable for a technique based on the classical linear regression model. For example, in a multivariate analysis, rates and proportions are estimated using nonlinear response models. The usual Oaxaca-Blinder method of mean/coefficient substitution is not always applicable in this case, unless the average rates or proportions assume values near 0.5, and models are linear over a wide range of predictions.

This paper draws on recent work by Even and Macpherson (1993), which has been extended by Nielsen (1998), and developed more systematically by Yun (2004). These authors have addressed several weaknesses with previous approaches to nonlinear decomposition that are remedied by the proposed methodology. In particular, unlike the linear case, the results of the detailed nonlinear decomposition may be sensitive to the order in which variables are entered into the decomposition. The general technique described here is not sensitive to order in which covariates enter into the decomposition. Previous methods that have been proposed to deal with

¹It is argued that portions of the wage differential that cannot be accounted for by group differences in characteristics are a result of labor market discrimination, or differences in returns to human capital factors such as education or job experience and to differences in unmeasured factors.

this problem are far more cumbersome (see e.g., Fairlie 2003).

This technique evaluates a nonlinear response function over the entire range of covariate values (not simply at the covariate means). This implies that results are valid for many demographic applications, such as for differences in probabilities of rare events.

We apply these techniques using logit models to investigate compositional and coefficient (return to risk) components of the differential in infant mortality due to respiratory distress syndrome (RDS) between 1991 and 1998 within the population of non-Hispanic whites and non-Hispanic blacks. We also decompose racial/ethnic differential in RDS mortality in 1998. Section 1 describes the general methodology as applied to mortality differentials, Section 2 presents results of this technique applied to our specific problem, and Section 3 provides a discussion and extension of the technique to other nonlinear response models.

1 Multivariate Decomposition

The goal of decomposition is to partition a difference in mean values between two groups into components owing to group differences in observed characteristics and to group differences in the estimated effects of those characteristics based on a regression model. Multivariate decomposition provides more detail by assessing the relative contribution of specific covariates to these components.

Example Problem. This paper illustrates a general methodology for multivariate decomposition nonlinear models using differences in infant mortality rates over time and across groups as a motivating example. The substantive problem to be addressed concerns a decomposition of the mortality differential between a higher-risk group (birth cohorts of 1991) and a lower-risk group (1998 birth cohorts). We will use the same convention to decompose mortality differentials between disadvantaged and advantaged groups (i.e., between blacks and

whites), where it is expected that compositional differences play a more significant role than they do in comparisons over time. We illustrate the approach using a logistic regression model for infant mortality due to respiratory distress syndrome (RDS), in which the estimated logit of mortality for infant i in year j is

$$\text{logit}(\hat{r}_{ij}) = \mathbf{x}'_{ij} \mathbf{b}_j.$$

The difference in observed mortality rates between 1991 and 1998 is the focus of the decomposition, which is equal to the difference in the average predicted IMR from models estimated for the years 1991 and 1998, or

$$\bar{r}_{91} - \bar{r}_{98} = \overline{F(\mathbf{x}'_{i91} \mathbf{b}_{91})} - \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{98})}, \quad (1)$$

where \mathbf{b}_j is a $K \times 1$ vector of coefficients, \mathbf{x}_{ij} is a $K \times 1$ vector denoting the set of measured characteristics for the i th individual in the j th year ($j = 91, 98$), and $\overline{F(\mathbf{x}'_{ij} \mathbf{b}_j)}$ denotes the mean of the logistic cdf evaluated at \mathbf{b} and \mathbf{x} , that is²

$$\overline{F(\mathbf{x}'_{ij} \mathbf{b}_j)} = \frac{1}{N_j} \sum_{i=1}^{N_j} F(\mathbf{x}'_{ij} \mathbf{b}_j).$$

It is important to note that the linear model is a much simpler case because the mean of the predictions equals the mean response, or $\overline{\mathbf{x}'_i \mathbf{b}} = \bar{\mathbf{x}}' \mathbf{b} = \bar{y}$. For the logit model $F(\bar{\mathbf{x}}' \mathbf{b})$ is not equal to $\overline{F(\mathbf{x}'_i \mathbf{b})}$ except when the probability of the outcome is 0.5. $F(\bar{\mathbf{x}}' \mathbf{b})$ can be far from the sample proportion for certain distributions as shown in Figure 1. This implies that computations could be simplified considerably for outcome probabilities in the neighborhood of 0.5.

²If a constant is included in the logit model, the predicted differential will equal the difference in the observed proportions. The logit model offers some advantage over a probit model insofar as the sample proportion equals the average predicted probability and difference in the actual proportions is decomposed rather than the difference in predicted probabilities. The general approach outlined here will also work for linear models. This will result in familiar expressions involving only covariate means and model estimates (see e.g., Blinder 1973, Oaxaca 1973, and Winsborough and Dickinson 1969).

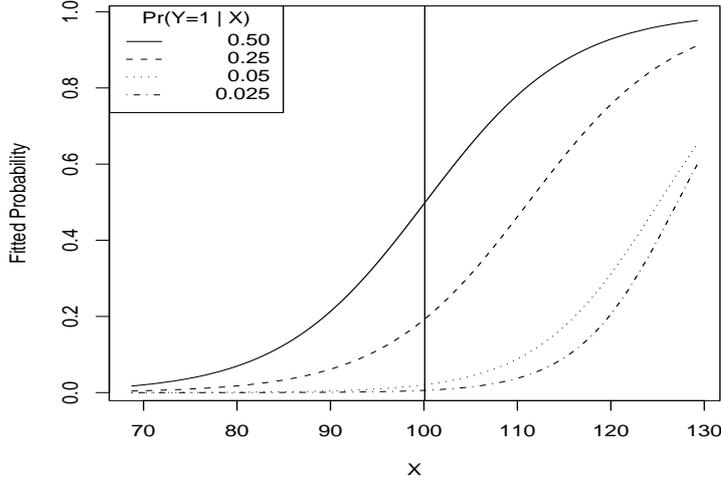


Figure 1: Means of Fitted Probabilities ($\overline{F(xb)}$) under Several Hypothetical Distributions. The vertical line indicates \bar{x} . Each curve intersects this line at $F(\bar{x}b)$. The values of $F(\bar{x}b)$ corresponding to sample the proportions ($\overline{F(xb)}$) are 0.502, 0.186, 0.015, and 0.006, respectively.

We would like to decompose the overall difference into components that reflect compositional differences between groups (differences in endowments) and differences in the effects of characteristics (differences in coefficients or returns to risk) between groups. We can rewrite Eq. 1 as

$$\bar{r}_{91} - \bar{r}_{98} = \underbrace{\{F(\mathbf{x}'_{i91} \mathbf{b}_{91}) - F(\mathbf{x}'_{i98} \mathbf{b}_{91})\}}_E + \underbrace{\{F(\mathbf{x}'_{i98} \mathbf{b}_{91}) - F(\mathbf{x}'_{i98} \mathbf{b}_{98})\}}_C \quad (2)$$

The first term appearing in the sum in Eq. 2 is the portion of the differential attributed to compositional differences, or “endowments,” E , which is the predicted mortality probability for the 1991 birth cohort minus the predicted mortality probability if the 1998 birth cohort faced the same returns to risk as the 1991 birth cohort. This component reflects the contribution to mortality differences that would have occurred if the two cohorts differed only with respect to endowments. The second term in Eq. 2 is the portion of the differential due to changes in the coefficient component, C , which assesses the contribution to mortality differences that would have occurred if the 1998 birth cohort’s returns to risk equaled those of the 1991 cohort and if group

characteristics were held fixed at 1998 levels.³ In this example, the 1991 birth cohort is the comparison group and the 1998 cohort is the reference group.

The same differential can be obtained from an alternative decomposition that switches the roles of the reference and comparison groups. This is referred to as the “indexing” problem. In the expressions above, the 1991 coefficients are used as weights in the composition component and the 1998 covariate values are used as weights in the coefficient component. By fixing the coefficients in the composition component to 1991 levels, we are assessing contribution to the mortality gap that would have occurred if the 1991 contexts had not changed, or if the returns to risk associated with the covariates in the model had remained at 1991 levels. By fixing the characteristics to 1998 levels in the coefficient component, we are assessing the contribution to the differential that is due to changing contexts between 1991 and 1998. An equivalent decomposition would reverse this procedure. That is, we could perform another decomposition that would weight the composition component by 1998 coefficient values and use observed characteristics in 1991 as weights in the coefficient component.⁴

1.1 Detailed Decomposition

The essence of a multivariate decomposition is to learn the contribution of each covariate to the components of the gap. To obtain a detailed multivariate decomposition, we break E and C into parts that are due to each of the K independent variables in the model.⁵ Following Even and

³We can motivate this decomposition in other ways. In a logit model, the total differential

$$d\hat{r} = \partial\hat{r}/\partial\mathbf{x}\mathbf{b} = \mathbf{b}f(\mathbf{x}\mathbf{b})d\mathbf{x} + \mathbf{x}f(\mathbf{x}\mathbf{b})d\mathbf{b}.$$

In terms of finite differences,

$$\Delta\hat{r} \approx \mathbf{b}f(\mathbf{x}\mathbf{b})\Delta\mathbf{x} + \mathbf{x}f(\mathbf{x}\mathbf{b})\Delta\mathbf{b}.$$

This approach is useful for obtaining the weights used in the detailed decomposition.

⁴We could perform both of these decompositions to deal with the indexing problem and average the results from each approach. In the example considered here, both sets of results are very similar.

⁵We include the constant term in the coefficient component. However, we separate out this contribution when discussing the results.

Macpherson (1993), Nielsen (1998), and Yun (2004), the detailed decomposition is⁶

$$\bar{r}_{91} - \bar{r}_{98} = \sum_{k=1}^K W_{\Delta x_k} \{ \overline{F(\mathbf{x}'_{i91} \mathbf{b}_{91})} - \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{91})} \} + \sum_{k=1}^K W_{\Delta b_k} \{ \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{91})} - \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{98})} \}, \quad (3)$$

where the weights $W_{\Delta x_k}$ and $W_{\Delta b_k}$ are given by

$$W_{\Delta x_k} = \frac{(\bar{x}_{91k} - \bar{x}_{98k})b_{91k}}{\sum_{k=1}^K (\bar{x}_{91k} - \bar{x}_{98k})b_{91k}}$$

and

$$W_{\Delta b_k} = \frac{\bar{x}_{98k}(b_{91k} - b_{98k})}{\sum_{k=1}^K \bar{x}_{98k}(b_{91k} - b_{98k})},$$

where $\sum_k W_{\Delta x_k} = \sum_k W_{\Delta b_k} = 1$.

Weights for Composition and Coefficient Components. Yun (2004) derives these weights using a first-order Taylor expansion of E around $\bar{\mathbf{x}}'_{91} \mathbf{b}_{91}$ and C around $\bar{\mathbf{x}}'_{98} \mathbf{b}_{98}$, which gives

$$\bar{r}_{91} - \bar{r}_{98} = [(\bar{\mathbf{x}}_{91} - \bar{\mathbf{x}}_{98})' \mathbf{b}_{91}] f(\bar{\mathbf{x}}'_{91} \mathbf{b}_{91}) + [\bar{\mathbf{x}}'_{98} (\mathbf{b}_{91} - \mathbf{b}_{98})] f(\bar{\mathbf{x}}'_{98} \mathbf{b}_{98}) + R_M + R_T$$

where $f(\cdot, \cdot)$ is the first-order derivative of $F(\cdot, \cdot)$ (i.e., the logistic pdf) and R_M and R_T are approximation residuals resulting from substituting mean values into E and C and from the Taylor expansion, respectively. These weights are the same as those used by Even and Macpherson (1993). The weights used by Nielsen (1998) for the detailed decomposition of the coefficients component are derived using the total differential of the coefficient component, however the details on their construction are not shown. The composition component is

⁶Even and Macpherson (1993) provide a detailed decomposition of the endowment component only. Nielsen (1998) provides a decomposition for the coefficient component, but provides no details on the construction of the weights used. Yun (2004) provides a details on both components as well as a discussion of the construction of the weight factors used below.

approximated by the sum of specific terms pertaining to the k th covariate,

$$E \approx \sum_{k=1}^K f(\bar{\mathbf{x}}_j' \mathbf{b}_j) b_{jk} (\bar{x}_{ijk} - \bar{x}_{ij'k}).$$

Because $f(\bar{\mathbf{x}}_j' \mathbf{b}_j)$ is a scalar quantity, this results in a simple expression in which the k th component of the weight involves only the sample means and model estimates

$$W_{\Delta x_k} = \frac{f(\bar{\mathbf{x}}_j' \mathbf{b}_j) b_{jk} (\bar{x}_{ijk} - \bar{x}_{ij'k})}{f(\bar{\mathbf{x}}_j' \mathbf{b}_j) \sum_{k=1}^K b_{jk} (\bar{x}_{ijk} - \bar{x}_{ij'k})} = \frac{b_{jk} (\bar{x}_{ijk} - \bar{x}_{ij'k})}{\sum_{k=1}^K b_{jk} (\bar{x}_{ijk} - \bar{x}_{ij'k})}.$$

The same logic can be applied to obtain the weights for the coefficient component.

$$C \approx \sum_{k=1}^K f(\bar{\mathbf{x}}_j' \mathbf{b}_j) \bar{x}_{jk} (b_{jk} - b_{j'k}),$$

and k th component of the weight is

$$W_{\Delta b_k} = \frac{f(\bar{\mathbf{x}}_j' \mathbf{b}_j) \bar{x}_{jk} (b_{jk} - b_{j'k})}{f(\bar{\mathbf{x}}_j' \mathbf{b}_j) \sum_{k=1}^K \bar{x}_{jk} (b_{jk} - b_{j'k})} = \frac{\bar{x}_{jk} (b_{jk} - b_{j'k})}{\sum_{k=1}^K \bar{x}_{jk} (b_{jk} - b_{j'k})}.$$

The last step is to separate out the contribution of each covariate to composition,

$$E_k = W_{\Delta x_k} E,$$

and coefficients

$$C_k = W_{\Delta b_k} C.$$

Thus, the composition weights reflect the relative contribution of each covariate based on the scale of the covariate, the magnitude of the difference, weighted by the effect of the covariate in the reference group. The coefficient weights reflect the relative contribution of each covariate based

on the size of the effect and the magnitude of the difference, weighted by the mean value of the covariate in the comparison group. Detailed decompositions obtained using weights constructed in this manner are free from problems of path dependency, or the order in which the independent variables are entered into the detailed decomposition. Unlike the linear model case, the independent contribution of x_k depends on the values of x_{k+1}, \dots, x_K . Thus, regardless of the ordering of variables, the relative weight associated with a particular variable remains the same.

1.2 Approximate Standard Errors

Approximate Standard Errors for the Contribution to Composition. The delta method can be used to derive standard errors of the detailed contribution to the composition component, E . Writing the first component as a weighted sum of individual contributions

$$E = \sum_{k=1}^K E_k = \sum_{k=1}^K W_{\Delta x_k} \{ \overline{F(\mathbf{x}'_{i91} \mathbf{b}_{91})} - \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{91})} \},$$

and denoting each covariate's contribution to the overall composition component as E_k ,

$$\frac{\partial E_k}{\partial b_{91_k}} = W_{\Delta x_k} \{ \overline{f(\mathbf{x}'_{i91} \mathbf{b}_{91}) x_{i91_k}} - \overline{f(\mathbf{x}'_{i98} \mathbf{b}_{91}) x_{i98_k}} \} + w_{x_k} \{ \overline{F(\mathbf{x}'_{i91} \mathbf{b}_{91})} - \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{91})} \}, \quad (4)$$

where $\overline{f(\mathbf{x}'_{ij} \mathbf{b}_j) x_{ijk}}$ is the mean of $\frac{\partial F(\mathbf{x}'_{ij} \mathbf{b}_j)}{\partial (\mathbf{x}'_{ij} \mathbf{b}_j)}$ with respect to \mathbf{b}_j , and

$$w_{x_k} = \frac{\partial W_{\Delta x_k}}{\partial b_{91_k}} = \frac{\bar{x}_{91_k} - \bar{x}_{98_k}}{\sum_k b_{91_k} (\bar{x}_{91_k} - \bar{x}_{98_k})} - \frac{b_{91_k} (\bar{x}_{91_k} - \bar{x}_{98_k})^2}{\{ \sum_k b_{91_k} (\bar{x}_{91_k} - \bar{x}_{98_k}) \}^2}.$$

Letting $\text{var}(\mathbf{b}_{91})$ denote the variance covariance matrix of \mathbf{b}_{91} and $\mathbf{E} = \{E_1, \dots, E_K\}$, the asymptotic (co)variance of the detailed composition components is then,

$$\text{var}(\mathbf{E}) = \left(\frac{\partial \mathbf{E}}{\partial \mathbf{b}_{91}} \right) \left(\frac{\partial \mathbf{E}}{\partial \mathbf{b}_{91}} \right)' \text{var}(\mathbf{b}_{91}).$$

Following the same logic, the coefficient component can be written as a sum of individual contributions as,

$$C = \sum_{k=1}^K C_k = \sum_{k=1}^K W_{\Delta b_k} \{ \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{91})} - \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{98})} \}$$

Each covariate's contribution to the overall coefficient component is more involved as this depends on two sets of parameter vectors, \mathbf{b}_{91} and \mathbf{b}_{98} .

$$\frac{\partial C_k}{\partial b_{91_k}} = W_{\Delta b_k} \overline{f(\mathbf{x}'_{i98} \mathbf{b}_{91}) x_{91_{ik}}} + w_{b_k} \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{91})} \quad (5)$$

and

$$\frac{\partial C_k}{\partial b_{98_k}} = w_{b_k} \overline{F(\mathbf{x}'_{i98} \mathbf{b}_{98})} - W_{\Delta b_k} \overline{f(\mathbf{x}'_{i98} \mathbf{b}_{98}) x_{98_{ik}}}, \quad (6)$$

where

$$w_{b_k} = \frac{\partial W_{\Delta b_k}}{\partial b_{jk}} = \frac{\bar{x}_{98k}}{\sum_k \bar{x}_{98k} (b_{91k} - b_{98k})} - \frac{\bar{x}_{98k}^2 (b_{91k} - b_{98k})}{\{\sum_k \bar{x}_{98k} (b_{91k} - b_{98k})\}^2},$$

when $j = 98$, this quantity has the opposite sign. Letting $\text{var}(\mathbf{b}_{91})$ and $\text{var}(\mathbf{b}_{98})$ denote the covariance matrix of the estimates from 1991 and 1998, respectively, and $\mathbf{C} = \{C_1, \dots, C_K\}$, the approximate (co)variance of the detailed coefficient components is

$$\text{var}(\mathbf{C}) = \left(\frac{\partial \mathbf{C}}{\partial \mathbf{b}_{91}} \right) \left(\frac{\partial \mathbf{C}}{\partial \mathbf{b}_{91}} \right)' \text{var}(\mathbf{b}_{91}) + \left(\frac{\partial \mathbf{C}}{\partial \mathbf{b}_{98}} \right) \left(\frac{\partial \mathbf{C}}{\partial \mathbf{b}_{98}} \right)' \text{var}(\mathbf{b}_{98})$$

2 Results

2.0.1 Non-Hispanic Whites

We apply these methods to decompose the change in the infant mortality rate due to respiratory distress syndrome (RDS). We use the National Center for Health Statistics (NCHS) linked birth/infant death cohort files for 1991 and 1998, which include all infants born alive in the U.S.

during each of those years and mortality information and cause of death for the subset of infants in the cohort who died during the first year of life. As recommended by NCHS, mother's race/ethnicity is used to distinguish non-Hispanic whites and non-Hispanic blacks. In 1991 there were 2.5 million live births to non-Hispanic whites, with 1,154 deaths due to RDS (4.5 per 10,000). There were 2.4 million live births in 1998 with 596 deaths due to RDS (2.5 per 10,000). This reduction is thought to be largely due to surfactant therapy interventions, which were widely applied to high-risk infants beginning in 1994.

Logit models estimated separately for the years 1991 and 1998 using low birth weight (birth weight $< 2,500$ gm) and low maternal education (0-11 years of schooling) yield a change of 0.0002 in the RDS-specific IMR. This is the quantity that is subject to decomposition in components attributable to change in the distribution of low birth weight (an increase from 5.73% to 6.54% of Non-Hispanic white births over the period), low maternal education (a decrease from 15% to 12.8% over this period), and change in the coefficients (or returns to risk) associated with these variables, which included a reduction in the effect of low birth weight and a doubling of the effect of low maternal education on RDS mortality. Table 1 shows logistic regression coefficients and covariate means for 1991 and 1998.

Given the small magnitude of change, we multiply the raw components $\times 100$ and also report percentage contributions. We also decompose the coefficients effect into C , a component due solely to measured covariates, and due to U , a shift coefficient representing the change in the constant term (i.e., change that cannot be attributed to changes in model coefficients and model covariates). Table 2 shows results for non-Hispanic whites.

Decomposition of Change in RDS Mortality for Whites. Changes in low birth weight composition contribute a 31% decrease in the RDS mortality differential between 1991 and 1998 [Panel A]. This means that holding the effects of low birth weight fixed at 1991 levels, change in

the birth weight distribution (i.e., more low weight births) would have increased RDS mortality between 1991 and 1998 by about 31%. At the same time, compositional change in the distribution of low education would result in lower RDS mortality by about 1%. Taken together, compositional change accounts for a 30% decrease in the RDS mortality gap, or a 30% increase in RDS mortality between 1991 and 1998.

Change in covariate effects account for an overwhelming amount of the predicted change in RDS mortality from 1991-1998 [Panel B]. Holding composition fixed at the 1998 levels, the change in the low birth weight effect would yield an 82% decrease in RDS mortality. (or an 82% increase RDS IMR differential). Change (increase) in the maternal education effect covers zero. Change in sources of the differential that are not attributable to variables in the model (i.e., the change in the constant) are also not statistically significant. Thus, for the simple model considered here, change in the return to risk associated with low birth weight (i.e., a smaller effect in 1998) is, by far, the most significant contributor to the observed decreased in RDS mortality for Non-Hispanic Whites. This component outweighs the effect of change in composition due to increases in the low birth weight distribution. These results do not change significantly if we weight by the 1998 slopes.

2.0.2 Non-Hispanic Blacks

In 1991 there were 656,268 live births to non-Hispanic blacks, with 792 deaths due to RDS (12.1 per 10,000). There were 584,859 live births in 1998 with 379 deaths due to RDS (6.5 per 10,000). The relative reduction in RDS related mortality declined by 14% for blacks and 20% for whites. However, black/white relative risk declined by less than 1% over the period. Table 2 reports the results of the multivariate decomposition for blacks. Compositional change in low birth weight was small (13.5% in 1991 to 13.1% in 1998). Compositional change in low maternal education was greater for blacks compared to whites—a 3.6% change from 33.1% in 1991 to 26.7% in 1998 for

blacks compared to a 2.2% change for whites over the same period. The effect of low maternal education was not significantly different from zero in either year, and we observed a slight reduction in the effect of low birth weight.

Decomposition of Change in RDS Mortality for Blacks

Holding effects constant at 1998 levels, the slight decline in the percentage of low birth weight births among blacks contributed to a 6% decrease in RDS mortality between 1991 and 1998. The change in educational composition had a negligible effect on the 1991-1998 differential. The change in the low birth weight effect accounted for one third of the overall gap. Change in the low maternal education effect was not statistically significant. As in the white population, change in return to risk (i.e., a smaller effect of low birth weight in 1998) is largely responsible for lower RDS mortality among Non-Hispanic Blacks. However, this effect is nearly 2 and 1/2 times larger in the white population.

Comparing the Black–White Differential in RDS Mortality: 1998

We also apply of this technique to assess the black–white differential in RDS mortality. Table 4 shows this decomposition. The large compositional differences in the low birth weight distribution account for about 82% of the black-white RDS mortality differential. The positive sign reflects the compositional advantage of whites in terms of low birth weight (6.5%[whites] vs 13.1%[blacks] in 1998). Compositional differences in the racial/ethnic distribution of low maternal education contributed to a 5.3% reduction in black-white RDS mortality gap. This must be interpreted in light of the black coefficient for education, which is negative but not statistically different from 0. The 95% confidence interval for this component covers zero.

We find relatively small contributions due to race/ethnic differences in coefficients. For example, the difference in the effect of low maternal education accounts for 2.3% of the

black-white RDS mortality differential. However, this also must be interpreted in light of the black coefficient, which is *negative* and not significant. We can conclude that the black-white difference in the low birth weight distribution, or compositional differences, rather than differences in return to risk, are most important in explaining differences in RDS mortality rates.

3 Discussion

This paper describes a multivariate decomposition procedure for logit models and an example of RDS mortality differentials over time and across racial groups. With some minor modifications, the same technique can be used for other nonlinear models. Previous work has considered a decomposition of the probit model (Even and Macpherson 1993). The decomposition of multivariate models for rates is relevant for many areas of demographic research and is a straightforward extension of the methods described above. For example, consider a proportional hazards model of the form

$$F(\mathbf{x}'_i \mathbf{b}) = h_0(t) \exp(\mathbf{x}'_i \boldsymbol{\beta}),$$

where $h_0(t)$ is a piecewise constant baseline hazard. Then, absorbing the baseline hazard component into $\exp(\mathbf{x}'_i \boldsymbol{\beta})$, then other relevant quantities for the decomposition are,

$$\frac{\partial F(\mathbf{x}'_i \mathbf{b})}{\partial \mathbf{b}} = \mathbf{x}'_i \exp(\mathbf{x}'_i \mathbf{b}) \quad \text{and} \quad \frac{\partial F(\mathbf{x}'_i \mathbf{b})}{\partial \mathbf{x}} = \mathbf{b} \exp(\mathbf{x}'_i \mathbf{b})$$

The same expressions would result from log probability and count models. The expressions for the weights used in the detailed decomposition are not affected by the choice of the model.

This general technique can be extended beyond the decomposition of empirical differences. For example, at the core of this technique is a type of simulation. We ask what is the expected contribution to a group difference that would occur if we could equalize group characteristics and,

what is the expected contribution to a group difference that would occur if each group's effects were equal? We could then simulate the differentials that would accrue from equalizing group characteristics through policy interventions or could examine changes in differentials that might result from increases or decreases in the effects of certain characteristics.

References

- Blinder A. S. 1973. "Wage Discrimination: Reduced Form and Structural Variables." *Journal of Human Resources*, 8, 436-455.
- Even, W. E., and D. A. Macpherson. 1993, "The Decline of Private-Sector Unionization and the Gender Wage Gap, *Journal of Human Resources*, 279-296
- Fairlie, R. 2003. "An Extension of the Blinder-Oaxaca Decomposition Technique to Logit and Probit Models," Discussion Paper #873, Economic Growth Center, Yale University.
- Nielsen, H. S. 1998 "Discrimination and Detailed Decomposition in a Logit Model," *Economics Letters*, 61: 115-120.
- Oaxaca, R. L. 1973. "Male-Female Wage Differentials in Urban Labor Markets." *International Economic Review*, 14, 693-709.
- R Development Core Team (2005). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org>.
- Winsborough, H. H. and P. Dickinson. 1969. "Components of Negro-White Income Difference." Center for Demography and Ecology, University of Wisconsin, Madison, WI.
- Yun, M-S. 2004. "Decomposing Differences in the First Moment," *Economics Letters*, 82, 275-280.

Table 1: Coefficients, Std. Errors, and Means 1991-1998

<u>Non-Hispanic Whites</u>						
Variable	1991			1998		
	b	$se(b)$	\bar{x}	b	$se(b)$	\bar{x}
Low Birth Weight	6.130	0.161	0.057	5.399	0.172	0.065
Low Maternal Education	0.070	0.070	0.150	0.148	0.103	0.128
Constant	-11.011	0.158		-11.049	0.167	

<u>Non-Hispanic Blacks</u>						
Variable	1991			1998		
	b	$se(b)$	\bar{x}	b	$se(b)$	\bar{x}
Low Birth Weight	5.812	0.075	0.135	5.323	0.293	0.131
Low Maternal Education	-0.017	0.261	0.303	-0.163	0.116	0.267
Constant	-10.535	0.259		-10.614	0.290	

Table 2: Decomposition of Change in IMR Due to RDS 1991–1998: Non-Hispanic Whites

Scaled Raw Contributions	Composition	Coefficients	Attributed
	E	C	$E + C$
Low Birth Weight (1)	-0.0063	0.0167	0.0104
Low Maternal Education (2)	0.0002	-0.0034	-0.0033
Total	-0.0061	0.0132	0.0071
Shift Coefficient (U)	0.0132		
Differential ($R = E + C + U$)	0.0203		
Percentage Contributions			
(A) Composition as % of Differential (E/R)	-29.91±11.56		
Due to Low Birth Weight (E_1/R)	-30.86±9.64		
Due to Low Maternal Education (E_2/R)	0.96±1.92		
(B) Coefficients as % of Differential (C/R)	129.91±689.96		
Due to Low Birth Weight (C_1/R)	82.10±68.18		
Due to Low Maternal Education (C_2/R)	-17.02±66.01		
Due to Constant (U/R)	64.82±555.77		
Total (A + B)	100.00		

Table 3: Decomposition of Change in IMR Due to RDS 1991–1998: Non-Hispanic Blacks

Scaled Raw Contributions	Composition	Coefficients	Attributed
	E	C	$E + C$
Low Birth Weight (1)	0.0037	0.0185	0.0222
Low Maternal Education (2)	-0.0001	0.0113	0.0112
Total	0.0036	0.0298	0.0334
Shift Coefficient (U)	0.0225		
Differential ($R = E + C + U$)	0.0559		
Percentage Contributions			
(A) Composition as % of Differential (E/R)	6.46±5.17		
Due to Low Birth Weight (E_1/R)	6.61±0.22		
Due to Low Maternal Education (E_2/R)	-0.16±4.94		
(B) Coefficients as % of Differential (C/R)	93.54±540.25		
Due to Low Birth Weight (C_1/R)	33.15±25.42		
Due to Low Maternal Education (C_2/R)	20.21±132.18		
Due to Constant (U/R)	40.19±382.64		
Total (A + B)	100.00		

Table 4: Decomposition of Black/White Differential in IMR Due to RDS 1998

Scaled Raw Contributions	Composition	Coefficients	Attributed
	E	C	$E + C$
Low Birth Weight (1)	0.03275	-0.00012	0.03264
Low Maternal Education (2)	-0.00212	-0.00093	-0.00305
Total	0.03063	-0.00104	0.02959
Shift Coefficient (U)	0.01012		
Differential ($R = E + C + U$)	0.03971		
Percentage Contributions			
(A) Composition as % of Differential (E/R)	77.13±73.19		
Due to Low Birth Weight (E_1/R)	82.47±61.76		
Due to Low Maternal Education (E_2/R)	-5.34±11.43		
(B) Coefficients as % of Differential (C/R)	22.87±101.84		
Due to Low Birth Weight (C_1/R)	-0.29±7.24		
Due to Low Maternal Education (C_2/R)	-2.33±6.06		
Due to Constant (U/R)	25.49±88.54		
Total (A + B)	100.00		

Appendix A

Programming. Below is a program written in the R programming language (R Development Core Team [2005]) use to carry out a logit decomposition.

```
# decompLogit.R
#
#
dat <- read.table( 'NHW_rds.raw' )
  yr   <- dat[,1]
  bwtlow <- dat[,2]
  edu011 <- dat[,3]
  Yrds   <- dat[,4]
  Yother <- dat[,5]
  Yfact  <- dat[,6]
#####
Y <- Yrds
#####
x0   <- rep(1, length(yr))
IMR91 <- mean(Y[yr==1])
IMR98 <- mean(Y[yr==2])
year91 <- yr==1
year98 <- yr==2
  N1 <- sum(year91)
  N2 <- sum(year98)
m191 <- mean(bwtlow[yr==1])
m291 <- mean(edu011[yr==1])
m198 <- mean(bwtlow[yr==2])
m298 <- mean(edu011[yr==2])
m91  <- cbind(1,m191,m291)
```

```

m98      <- cbind(1,m198,m298)
# 1991
mod1 <- glm(Y ~ bwtlow + edu011 , family="binomial", sub=year91, x=TRUE)
# extract coef vector and x
x91      <- mod1$x
b91      <- mod1$coef
b091     <- b91[1]
varb91   <- summary(mod1)$cov.scaled
b91L     <- b91 - 1.96*sqrt(diag(varb91))
b91H     <- b91 + 1.96*sqrt(diag(varb91))
# extract fitted probs
p91 <- mod1$fitted
#1998
mod2 <- glm(Y ~ bwtlow + edu011, family="binomial", sub=year98, x=TRUE)
# extract coef vector and x
x98      <- mod2$x
b98      <- mod2$coef
b098     <- b98[1]
varb98   <- summary(mod2)$cov.scaled
b98L     <- b98 - 1.96*sqrt(diag(varb98))
b98H     <- b98 + 1.96*sqrt(diag(varb98))

# logistic CDF
CDF.lgt <- function(b,x) {
  xb <- x%*%b
  F <- exp(xb)/(1 + exp(xb))
return(F)
}

# logistic pdf
pdf.lgt <- function(b,x) {
  xb <- x%*%b
  f <- exp(xb)/(1 + exp(xb))^2
return(f)
}

# weight function (composition)
Wdx.F <- function(b,x1,x2){
  A <- (x1-x2)%*%b
  Wdx <- NULL
  for (i in 1:length(b)){
    Wdx[i] <- (x1[i] - x2[i])*b[i] / A
  }
return(Wdx)
}

# weight function (coefficient)
Wdb.F <- function(b1,b2,x){

A <- x%*%(b1-b2)
Wdb <- NULL
for (i in 1:ncol(x)){
  Wdb[i] <- (x[i]*(b1[i] - b2[i])) / A
}

```

```

    }
    return(Wdb)
  }

dW.F <- function(b,x1,x2) {
  dW <- NULL
  A <- (x1-x2)%*%b
  for (i in 1:length(b)) {
    dW[i] <- (x1[i] - x2[i])/A - (b[i]*(x1[i]-x2[i])^2)/A^2
    cat(dW[i], "\n")
  }
  return(dW)
}

dWA.F <- function(b1,b2,x2) {
# derivative of Wdb wrt b1 = -derivative of Wdb wrt b2
  dWA1 <- NULL
  A <- x2*%*(b1-b2)
  for (i in 1:length(b1)){
    dWA1[i] <- x2[i]/A - (x2[i]^2*(b1[i]-b2[i]))/A^2
  }
  return(dWA1)
}

wb <- dWA.F(b91,b98,m98)
Wdx <- Wdx.F(b91, m91, m98)
Wdb <- Wdb.F(b91, b98, m98)
# check-sum to 1
sum(Wdx)
sum(Wdb)
#Convention: Yhi - Ylo = 1st moment higher group - 1st moment lower group
#decomp total:
#Composition or Endowments: E = F(Bhi,Mhi)-F(Bhi,Mlo) [use 91 coefs as weights]
E <- mean(CDF.lgt(b91,x91)) - mean(CDF.lgt(b91,x98))
#Coefficients + Unexplained: C = F(Bhi,Mlo)-F(Blo,Mlo) [use 98 means as weights]
C <- mean(CDF.lgt(b91, x98)) - mean(CDF.lgt(b98, x98))
# get dEdb for variance estimator
dWx <- dW.F(b91, m91, m98)

# gradient (composition)
dEdb <- NULL
for (k in 1:length(b91)){
dEdb[k] <- Wdx[k]*(mean(pdf.lgt(b91,x91)*x91[,k]) -
                    mean(pdf.lgt(b91,x98)*x98[,k])) + dWx[k]*E
}

# gradient (coefficients)
dCdb1 <- NULL
dCdb2 <- NULL
for (k in 1:length(b98)){
  dCdb1[k] <- Wdb[k]*mean(pdf.lgt(b91,x91)*x91[,k]) + wb[k]*mean(CDF.lgt(b91,x98))
  dCdb2[k] <- wb[k]*mean(CDF.lgt(b98,x98)) - Wdb[k]*mean(pdf(x98,b98)*x98[,k])
}

varb.b1 <- varb91

```

```

varb.b2 <- varb98
varb.b12 <- varb91 + varb98

### Variances
#Composition
K <- length(b91)
Var.E.k <- matrix(0,K,K)
  for (k in 1:K){
    for (l in 1:K){
      Var.E.k[k,l] <- dEdb[k]*dEdb[l]*varb.b1[k,l]
    }
  }
###
seWdx <- sqrt(diag(Var.E.k))
# or
seWdx <- sqrt(diag((dEdb%*%t(dEdb)*varb.b1)))

### Variances
#Coefficients
K <- length(b91)
Var.C.k <- matrix(rep(0,K*K),K,K)
  for (k in 1:K){
    for (l in 1:K){
      Var.C.k[k,l] <- (dCdb1[k])*(dCdb1[l])*varb.b1[k,l]
                        +
                        (dCdb2[k])*(dCdb2[l])*varb.b2[k,l]
    }
  }
###
seWdb <- sqrt(diag(Var.C.k))

#detailed decomp (due to composition) %
E*Wdx
#detailed decomp (due to coefficients)
C*Wdb
R <- E + C
CompH <- (E*Wdx) + 1.96*seWdx
Comp <- (E*Wdx)
CompL <- (E*Wdx) - 1.96*seWdx
CoefH <- (C*Wdb) + 1.96*seWdb
Coef <- (C*Wdb)
CoefL <- (C*Wdb) - 1.96*seWdb

gap <- IMR91-IMR98
cat("OBSERVED gap in IMR91 - IMR98=", gap, "\n")
cat("PREDICTED gap in model prediction=" , R, "\n")
cat("Total Amount due to composition:", "\n",
    "Raw Amount          =", E, "\n",
    "Percentage of predicted  =", (E/R)*100, "\n")
cat("Total Amount due to coefficients:", "\n",
    "Raw Amount =", C, "\n",
    "Percentage of predicted  =", (C/R)*100, "\n")
cat("Amount due to low birth weight composition:", "\n",

```

```

" Lower95                = ", CompL[2], "\n",
" Estimate                = ", E*Wdx[2], "\n",
" Upper95                = ", CompH[2], "\n",
" Percentage of predicted = ", ((E*Wdx[2])/R)*100, "\n")
cat("Amount due to low education composition:", "\n",
" Lower95                = ", CompL[3], "\n",
" Estimate                = ", E*Wdx[3], "\n",
" Upper95                = ", CompH[3], "\n",
" Percentage of predicted = ", ((E*Wdx[3])/R)*100, "\n")
cat("Amount due to birth weight coefficient =", C*Wdb[2], "\n",
" Lower95                = ", CoefL[2], "\n",
" Estimate                = ", C*Wdb[2], "\n",
" Upper95                = ", CoefH[2], "\n",
"Percentage of predicted = ", ((C*Wdb[2])/R)*100, "\n")
cat("Amount due to education coefficient   =", C*Wdb[3], "\n",
" Lower95                = ", CoefL[3], "\n",
" Estimate                = ", C*Wdb[3], "\n",
" Upper95                = ", CoefH[3], "\n",
"Percentage of predicted = ", ((C*Wdb[3])/R)*100, "\n")
cat("Amount due to Constant =                = ", C*Wdb[1], "\n",
" Lower95                = ", CoefL[1], "\n",
" Estimate                = ", C*Wdb[1], "\n",
" Upper95                = ", CoefH[1], "\n",
"Percentage of predicted = ", ((C*Wdb[1])/R)*100, "\n")
comp <- E*Wdx*100
coef <- C*Wdb*100
cbind(comp, coef)
cbind(sum(comp), sum(coef))
CompL/gap*100 -> EL
E*Wdx/gap*100 -> EP
CompH/gap*100 -> EH
EL
EP
EH

CoefL/gap*100 -> CL
C*Wdb/gap*100 -> CP
CoefH/gap*100 -> CH

CL
CP
CH

```