



THE STATISTICAL PROPERTIES OF THE MUTUAL INFORMATION INDEX OF MULTIGROUP SEGREGATION

Ricardo Mora and Javier Ruiz-Castillo¹

Departamento de Economía, Universidad Carlos III de Madrid

Abstract

In this paper the Kullback-Leibler notion of discrepancy (Kullback and Leibler, 1951) is used to propose a measure of segregation within a general statistical framework. Under general conditions, this measure coincides with the Mutual Information index of segregation, M , first proposed by Theil and Finizza (1971), and fully characterized in terms of eight ordinal axioms by Frankel and Volij (2009). In this paper, two specific issues are addressed in relation to this index: the evaluation of statistical significance for observed differences in M measurements, and the control for the statistical association between demographic groups and schools and other socioeconomic variables. Among the main results of the paper it is established that M can be decomposed to isolate segregation conditional on any vector of socioeconomic characteristics. Furthermore, consistent estimators for M and the terms in its decomposition are proposed, and their asymptotic properties are obtained. As a result, the M index now stands as the only index of segregation which has been fully characterized in terms of axiomatic properties, is well embedded into a general statistical framework, and can be used when samples are finite and a multivariate framework is required. The usefulness of the approach is illustrated by looking at patterns of multigroup school segregation in the U.S. for the school years 1989-90 and 2005-06.

Keywords: Multigroup Segregation Measurement; Axiomatic Properties; Econometric Models.

¹ This is a completely new version of a 2007 Working Paper with the same title. The authors acknowledge financial support from the Spanish DGI, Grants SEJ2006-05710 and SEJ2007-67436.

I. INTRODUCTION

Social scientists have long been interested in the measurement of occupational segregation by gender, as well as in residential and educational segregation by ethnic group.² Mathematically, these problems are similar in the sense that both involve summarizing by means of a real number the information contained in the frequency of individuals (workers, residents, students) over a finite set of organizational units (occupations, neighbourhoods, schools) and a finite set of groups (defined in terms of gender, racial, or ethnic categories). Such a real number is referred to as an index of segregation. For concreteness, this paper will use the example of school segregation in the multiracial case.

Since the seminal contribution by Duncan and Duncan (1955), a plethora of segregation indices have been used in the empirical literature –for instance, Massey and Denton (1988) survey 20 indices. Unfortunately, two segregation indices may show different trends in segregation and may produce different rankings for a set of, say, school districts, in a given year.³ Thus, the design of measures with desirable properties is a central methodological issue, and the merits of competing indices are regularly debated.⁴ Since segregation measures are routinely computed using samples, the debate not only includes the attempt to provide characterizations but also the study of the statistical properties of segregation measures and the development of statistical frameworks that permit the testing of hypothesis regarding segregation.

This paper addresses two specific issues that have attracted the attention of applied researchers on segregation: the statistical significance of segregation indices, and the control for the statistical association between demographic groups and schools and other socioeconomic variables. To appreciate this paper's contribution, it is first necessary to briefly review what has been accomplished by other

² For a recent treatise on occupational segregation by gender, see Flückiger and Silber (1999), and for a recent treatise on residential and school segregation, see Reardon and Firebaugh (2002).

³ See, *inter alia*, Jonung (1984), James and Taeuber (1985), Karmel and MacLachlan (1988), Blackburn *et al.* (1993), Anker (1998), and Flückiger and Silber (1999).

⁴ See, *inter alia*, the methodological contributions by James and Taeuber (1985), Massey and Denton (1988), Siltanen (1990), Hutchens (1991, 2001, 2004), Watts (1992, 1994, 1997a, 1997b, 1998a, 1998b), Blackburn *et al.* (1993, 1995), Flückiger and Silber (1999), Reardon and Firebaugh (2002), and Frankel and Volij (2008, 2009).

researchers in these two areas.

First, it is natural to study whether two different samples are drawn from populations with the same level of segregation. For example, when indices are computed using annual census data from micro-areas such as school districts, the observed racial frequencies in each district can be regarded as a sample drawn from a population with a given multinomial distribution. Therefore, it will be of interest to test whether the populations to which the samples are associated entail the same level of segregation. In the segregation literature, the statistical significance of segregation indices is dealt with in several ways. The simplest approach involves reporting t -statistics computed by bootstrap techniques or linear approximations, as in Deutsch *et al.* (1994), Boisso *et al.* (1994), and Ransom (2000). A related approach consists of standardizing the segregation measure, using as mean and standard deviations estimates obtained from resampling under random assignment into schools and races, as in Cortese *et al.* (1976), Carrington and Troske (1997), and Aslund and Skans (2009). The validity of resampling methods assumes a statistical framework and certain regularity conditions (see, for example, Davidson and Hinkley, 1997, p. 38) which are rarely explicitly discussed in empirical applications. Furthermore, the absence of a formal statistical framework makes it difficult to study the statistical properties of the indices under different hypothesis or segregation scenarios. In contrast to these approaches, several authors have made use of a rather restricted statistical framework for the empirical analysis of segregation, either informally, as in Reardon *et al.* (2000), or more explicitly as in Charles (1992, 1998), Charles and Grusky (1995), and Kakwani (1994).

The second question addressed in this paper refers to the control for group and school differences in socioeconomic variables. The use of indices which are group and school-decomposable into between and within discrete categories, such as in Reardon *et al.* (2000) and Frankel and Volij (2009), only partially answers this question because it cannot properly deal with the situation where the controls are continuous. Other researchers have tried to develop notions of conditional segregation which should in

practice be implementable in a general multivariate framework. Frequently, their ultimate purpose is to assess to what extent segregation can be explained by the determinants of individual choice; to do so, they borrow the tools used in the econometrics literature on discrete choice. Spriggs and Williams (1996), for example, compute measures of occupational segregation by race and gender based on multinomial logit models.⁵ Although these indices sometimes have a clear intuitive appeal, they are not characterized in terms of axiomatic properties. Consequently, one cannot be certain of what the index actually measures.

In this paper, a general statistical framework to analyse multigroup school segregation is set up by borrowing the Kullback and Leibler (1951) notion of discrepancy from Information Theory. A measure of segregation, M_{KL} , is then proposed and shown to satisfy several important properties.

First, M_{KL} coincides with the Mutual Information index, M , first proposed by Theil and Finizza (1971) as a measure of racial school segregation at district level, and recently characterized by Frankel and Volij (2009) as an index which represents the unique segregation ordering satisfying eight ordinal desirable axioms.⁶

Second, Frankel and Volij (2009) show that, for any variable d which partitions the set of schools or the set of racial categories, M is strongly decomposable and the within term in this decomposition can be interpreted as segregation conditional on d . In this paper, this result is generalized to condition segregation on any vector of (possibly continuous) student and school characteristics \mathbf{x} . In particular, the M_{KL} index can be decomposed into a between term, M_{KL}^B , which captures the statistical dependence

⁵ Kalter (2000) proposes a segregation index by linking the Index of Dissimilarity with the Multinomial Logit model. Aslund and Skans (2009) propose estimating the propensity score for each group given the vector of characteristics to create the benchmark no-segregation counterfactual for any segregation index. They propose a non-parametric procedure when all characteristics are discrete, and develop a test of conditional segregation using an index of exposure. Hellerstein and Neumark (2008) apply a similar idea to an index of exposure and isolation.

⁶ In the two group case, Chakravarty and Silber (1992) characterize an index of absolute segregation. Philipson (1993) provides an axiomatic characterization of a large family of segregation orderings that have an additively separable representation. Chakravarty and Silber (2007) axiomatically derive a general class of numerical indexes of relative segregation which parallel the multidimensional Atkinson inequality indices. One member of that class is monotonically related to the square root index, independently characterized by Hutchens (2004). In the multigroup case, Frankel and Volij (2008) provide an ordinal characterization of two families of Atkinson indexes.

between race status (or school membership) and \mathbf{x} , and a within term, M_{KL}^W , which captures multigroup school segregation conditional on \mathbf{x} . Because M_{KL}^B and M_{KL}^W are independent, in the sense that it is possible to introduce changes in the population to eliminate conditional segregation M_{KL}^W keeping conditioning segregation M_{KL}^B constant, this decomposition allows us to answer questions such as “to what extent is racial segregation at school level associated with racial differences in socioeconomic variables?”⁷ Moreover, since M_{KL}^B and M_{KL}^W are functions of terms that can be interpreted as qualitative response models, the decomposition provides an intuitive unifying econometric framework for studies of segregation using segregation indices and econometric models.

Third, for any sample of size T , under weak regularity conditions, consistent estimators for both the M_{KL} index, \hat{M}_T , and also the between and within terms in its decomposition, \hat{M}_T^B and \hat{M}_T^W , can be obtained from the principle of analogy, regardless of whether the set of covariates \mathbf{x} includes continuous, countable, or discrete variables. \hat{M}_T is shown to be a monotonic transformation of the likelihood-ratio statistic for testing statistical independence between school membership and racial status. Furthermore, when the vector of covariates \mathbf{x} only include discrete variables, it is shown that \hat{M}_T^W can be interpreted as a monotone transformation of the likelihood-ratio statistic for testing statistical independence between school membership and racial status given \mathbf{x} .

Finally, sufficient conditions are provided to characterize under all segregation scenarios the asymptotic properties of \hat{M}_T , \hat{M}_T^B , and \hat{M}_T^W , both in the case when all variables are discrete and also when there is at least one continuous variable in \mathbf{x} .

To summarize, elsewhere it has been shown that M is well grounded on an axiomatic notion of segregation. In this paper, we show that it can be used when samples are finite and a multivariate

⁷ In the field of income inequality, between-groups income inequality can also be interpreted as the amount by which overall income inequality is reduced when the differences between subgroup income means are eliminated by making them equal to the population income mean (see, *inter alia*, Shorrocks, 1984). As shown by Mora and Ruiz-Castillo (2009), the corresponding interpretation is impossible in segregation studies.

framework is required. The usefulness of the approach is illustrated by applying it to the analysis of racial school segregation in the U.S.

The rest of the paper contains four sections. Section 2 sets up the general statistical framework, and defines M_{KL} and its decomposition in a multivariate framework. Section 3 proposes estimators \hat{M}_T , \hat{M}_T^B , and \hat{M}_T^W and presents the asymptotic results. Section 4 contains the empirical illustration, while section 5 offers some concluding comments.

II. A GENERAL STATISTICAL MODEL OF MULTIGROUP SCHOOL SEGREGATION

II. 1. Measures of Segregation

It is useful to refer to a specific segregation problem. For consistency with the empirical illustration in section IV, the case discussed throughout the paper is the multigroup school segregation problem. Assume a city \mathbf{X} consisting of N schools, indexed by $n = 1, \dots, N$. Each student belongs to any of G racial groups, indexed by $g = 1, \dots, G$. The data available can be organized into the following $G \times N$ matrix:

$$\mathbf{X} = \{t_{gn}\} = \begin{bmatrix} t_{11} & \dots & t_{1N} \\ \vdots & \ddots & \vdots \\ t_{G1} & \dots & t_{GN} \end{bmatrix}, \quad (1)$$

where t_{gn} is the number of individuals of racial group g attending school n , so that $t = \sum_{n=1}^N \sum_{g=1}^G t_{gn}$ is the total student population.

The information contained in the joint absolute frequencies of racial groups and schools, t_{gn} , is usually summarized by means of numerical indices of segregation. Let $\mathbf{X}(G, N)$ be the set of all cities with G groups and N schools. A segregation index S is a real valued function defined in $\mathbf{X}(G, N)$, where $S(\mathbf{X})$ provides the extent of school segregation for any city $\mathbf{X} \in \mathbf{X}(G, N)$. Let $p_{gn} = t_{gn}/t$, and denote by

$P_{gn} = \{p_{gn}\}_{g=1, n=1}^{G, N}$, the joint distribution of racial groups and schools in a city $\mathbf{X} \in \mathcal{X}(G, N)$. In the following section, the discussion will be restricted to indices that capture a *relative* view of segregation in which all that matters is the joint distribution, i.e. indices which admit a representation as a function of P_{gn} .⁸

II. 2. A Kullback-Leibler Measure of Segregation

Following Kullback (1959), consider the probability space $(\Omega, \mathcal{F}, \mathbf{m})$ where Ω is the set of possible samples $\{g, n, \mathbf{x}\} \in \Omega$ where g stands for the student's race status code and takes a finite number of values, $g \in \{1, \dots, G\}$, n stands for a particular school code, $n \in \{1, \dots, N\}$, and $\mathbf{x} \in \Lambda \subset \mathbb{R}^k$ is a vector of k covariates. \mathcal{F} is the σ -algebra of subsets of Ω , and \mathbf{m} is a measure of the probability of the events in \mathcal{F} . We assume that there are two absolutely continuous measures with respect to \mathbf{m} , \mathbf{m}_1 and \mathbf{m}_2 , and two generalized density functions $f_i(g, n, \mathbf{x})$, $i=1, 2$, such that

$$\mathbf{m}_i(E) = \int_E f_i(g, n, \mathbf{x}) d\mathbf{m}, \quad i=1, 2,$$

for all $E \in \mathcal{F}$. The elements in \mathbf{x} may be univariate or multivariate, discrete or continuous, qualitative or quantitative, and the generalized density functions f_i are known at most up to a parameter vector.

Consider the partition of Ω into GN sets $D_{gn} = \{(r, s, \mathbf{x}) \in \mathcal{F} : r = g, s = n, \mathbf{x} \in \Lambda\}$ and let

$$\mathbf{m}_i(g, n) = \int_{D_{gn}} f_i(r, s, \mathbf{x}) d\mathbf{m}, \quad i=1, 2,$$

so that the probability that a student is of race g and school n under the probability measure \mathbf{m}_1 is

$$p_{gn} = \mathbf{m}_1(g, n) = \int_{D_{gn}} f_1(r, s, \mathbf{x}) d\mathbf{m} \tag{1}$$

⁸ This property, satisfied by most segregation indices, is referred to as *Size Invariance* in James and Taeuber (1985), and as *Weak Scale Invariance* in Frankel and Volij (2008, 2009). For a study that focuses on translation invariant segregation indices that represent an absolute view of segregation, see Chakravarty and Silber (1992).

where $p_{gn} \geq 0$ and $\sum_{g=1}^G \sum_{n=1}^N p_{gn} = 1$. The marginal probabilities for race status and school membership are

$p_{g\bullet} = \sum_{n=1}^N p_{gn}$ and $p_{\bullet n} = \sum_{g=1}^G p_{gn}$, respectively. For all g and n such that $\mathbf{m}_i(g, n) > 0$, $i = 1, 2$, the

generalized conditional density given race and school status is $f_i(\mathbf{x} | g, n) = \frac{f_i(g, n, \mathbf{x})}{\mathbf{m}_i(g, n)}$. In this

framework, a Kullback-Leibler discrepancy measure between f_1 and f_2 is defined as:

$$I(1:2) = \int f_1(g, n, \mathbf{x}) \log \left(\frac{f_1(g, n, \mathbf{x})}{f_2(g, n, \mathbf{x})} \right) d\mathbf{m}$$

Let H_i , $i = 1, 2$, represent the hypothesis that (g, n, \mathbf{x}) is from the statistical population with probability

measure \mathbf{m} , and define the logarithm of the likelihood ratio, $\log \left(\frac{f_1(g, n, \mathbf{x})}{f_2(g, n, \mathbf{x})} \right)$, as the information in

(g, n, \mathbf{x}) for discrimination in favour of H_1 against H_2 .⁹ Then $I(1:2)$ can be interpreted as the mean discrepancy (or information for discrimination) in favour of H_1 against H_2 per observation from \mathbf{m} (see Kullback, 1959, p. 5).

Define the conditional probability of school membership n given race status g as $p_{n|g} = \frac{p_{gn}}{p_{g\bullet}}$, and

let $P_{n|g} = \{p_{n|g}\}_{n=1}^N$ represent the conditional distribution of students from group g across schools.

Similarly, define the conditional probability of racial status g given school membership n as $p_{g|n} = \frac{p_{gn}}{p_{\bullet n}}$,

and denote by $P_{g|n} = \{p_{g|n}\}_{g=1}^G$ the racial mix within school n . Indices in the segregation literature

associate the absence of segregation with two situations. First, racial groups are not segregated if the

⁹ The base of the logarithm is immaterial, providing essentially a unit of measure. The natural logarithm is used throughout the paper.

relative frequency with which a student attends school n is constant, regardless of her racial group, i.e. $p_{n|g} = p_{\bullet n}$.¹⁰ Second, the racial composition at all schools is fully representative of the population if the relative frequency with which a student belongs to racial group g is constant regardless of the school which she attends, i.e. $p_{g|n} = p_{g\bullet}$.¹¹ These two notions of absence of segregation are equivalent and coincide with the concept of statistical independence between race status and school membership:

$$p_{g|n} = p_{g\bullet} \Leftrightarrow p_{n|g} = p_{\bullet n} \Leftrightarrow p_{ng} = p_{g\bullet} p_{\bullet n}.$$

Under the following three assumptions the Kullback-Leibler notion of discrepancy between dependence and independence of race and school membership becomes a measure of segregation. For all $g = 1, \dots, G$, $n = 1, \dots, N$, and $\mathbf{x} \in \Lambda \subset \mathbb{R}^k$:

$$\underline{A1}: p_{gn} > 0.$$

$$\underline{A2}: f_i(\mathbf{x} | g, n) = f(\mathbf{x} | g, n) > 0 \text{ as, } i = 1, 2.$$

$$\underline{A3}: \mathbf{m}_2(g, n) = p_{g\bullet} p_{\bullet n} = \left(\sum_{n=1}^N p_{gn} \right) \left(\sum_{g=1}^G p_{gn} \right).$$

A1 avoids trivialities (i.e. combinations of races and schools that are *a priori* impossible to observe), and ensures the existence of the generalized conditional density of \mathbf{x} given race and school status. A transformation of the data that considers only a marginal distribution in a multivariate situation generally results in a loss of information. However, A2 ensures that the marginal probabilities p_{gn} are sufficient statistics with respect to multigroup school segregation, so that no information is lost by disregarding \mathbf{x} . By identifying H_2 in the Kullback-Leibler general discrepancy measure with the notion of statistical independence, A3 implies that the measure of discrepancy captures the mean discrepancy (or

¹⁰ Absence of segregation in this sense is consistent with the notion of segregation as “evenness”, advocated by James and Taeuber (1985), according to which segregation is seen as the tendency of racial groups to have different distributions across schools.

¹¹ Absence of segregation in this sense follows the idea of “representativeness”, emphasized by Frankel and Volij (2009), which asks to what extent schools have different racial compositions from the population as a whole, and it is closely related to the idea of “isolation” distinguished by Massey and Denton (1988) in the two-group case

information for discrimination) in favour of the observed joint distribution against independent assignment of races and schools. Under assumptions A1, A2, and A3, the discrepancy measure $I(1:2)$, denoted by M_{KL} , takes the form:

$$\begin{aligned} M_{KL} &= \int f_1(r, s, \mathbf{x}) \log \left(\frac{p_{gn}}{p_{g\bullet} p_{\bullet n}} \right) d\mathbf{m} \\ &= \sum_{g=1}^G \sum_{n=1}^N \log \left(\frac{p_{gn}}{p_{g\bullet} p_{\bullet n}} \right) \int f_1(r, s, \mathbf{x}) d\mathbf{m}. \end{aligned}$$

Given equation (1), a relative measure of segregation –that is, a real value function which only depends on the joint distribution of race and school membership– is obtained:

$$M_{KL} = \sum_{g=1}^G \sum_{n=1}^N p_{gn} \log \left(\frac{p_{gn}}{p_{g\bullet} p_{\bullet n}} \right) \quad (2)$$

Information Theory has previously been used to motivate segregation indices. In the context of racial school segregation by school district, Theil and Finizza (1971) defined the *local index of segregation* for school n in a given district as:

$$M^n(P_g, P_{g|n}) = \sum_{g=1}^G p_{g|n} \log \left(\frac{p_{g|n}}{p_{g\bullet}} \right)$$

where $P_g = \{p_{g\bullet}\}_{g=1}^G$. Since $M^n(P_g, P_{g|n})$ measures the extent to which the racial composition in school n differs from the one for the city as a whole, it can be interpreted as a measure of deviations from representativeness. The demographically weighted average of all local school measures can be taken as an overall measure of segregation for the city:

$$M = \sum_{n=1}^N p_{\bullet n} M^n(P_g, P_{g|n}) = \sum_{n=1}^N p_{\bullet n} \sum_{g=1}^G p_{g|n} \log \left(\frac{p_{g|n}}{p_{g\bullet}} \right)$$

The M index can be interpreted as the information gained about race when the school is known. Mora and Ruiz-Castillo (2005) for the two-group case and Frankel and Volij (2009) for the multigroup case

note that, for the same data set, the M index can also be expressed as the demographically weighted averages of local indices of segregation by race:

$$M = \sum_{g=1}^G p_{g\bullet} \sum_{n=1}^N p_{n|g} \log \left(\frac{p_{n|g}}{p_{\bullet n}} \right)$$

Since it represents both the information gained about race when the school is known, as well as the information gained about the school when race is known, the M index is referred to as the *Mutual Information* index. Frankel and Volij (2009) provide a full characterization for M in terms of eight ordinal axioms. It can immediately be seen that the M_{KL} measure of segregation defined in equation (2) coincides with the M index:¹²

Proposition 1: Under assumptions A1-A3,

$$I(1:2) \equiv M_{KL} = M. \quad \blacksquare$$

Theil (1972) shows that M (and, by Proposition 1, M_{KL}) is bounded. The lower bound 0 is achieved whenever $p_{ng} = p_{g\bullet} p_{\bullet n}$ for all g and n , while the upper bound is $\min\{\log(G), \log(N)\}$.

II. 2. Multigroup School Conditional Segregation

Assumptions A1-A3 do not require independence between race status (or school membership) and any of the covariates in \mathbf{x} . A matter of interest concerns the extent to which the measure of discrepancy M_{KL} can be attributed to the discrepancy in the distributions of covariates across racial groups (or schools). As an illustration, consider the case whereby schools in a city are organized into a set of school districts. From the point of view of representativeness, school segregation by racial group is the result of differences in racial composition across schools. These differences may arise because districts differ in their racial composition, or because schools within the same district have different compositions. In terms of policy implications, this distinction is essential because, as argued by Reardon *et al.* (2000), “[between-district segregation] is generally the result of forces affecting racial groups’

¹² Proof of all propositions, lemmata, and theorems can be found in the Appendix.

differential access to housing markets (...)", whereas "[within-district segregation] may be the result of both residential patterns and school assignment practices within individual school districts". In order to address this and similar questions, organizational units and demographic groups can be partitioned into a finite set of major categories. Therefore, some indices of segregation can be decomposed into two terms measuring between and within-categories segregation.

More generally, it may be advisable to decompose overall segregation to evaluate the impact of possibly continuous covariates on it. Assume, for instance, that students can be characterized by the income level of the household to which they belong, and that there is a statistical association between student race and household income levels. Given that household income is a potential determinant of residential and school choice, multigroup school segregation may be partially due to income inequality. Therefore, it would be interesting to identify the extent to which the value of segregation arises from income and other possibly continuous socioeconomic characteristics. In the absence of a better strategy, one can obviously discretize the vector of socioeconomic controls. This option has, however, a practical and a conceptual drawback. The practical drawback stems from the curse of dimensionality. To avoid serious aggregation bias, one should consider as many discrete categories as possible for each control. As the number of necessary categories increases exponentially, in order to account for all possible combinations in the vector of controls, available datasets make this strategy viable only for vectors of very reduced dimensions. The conceptual drawback is due to the absence of a clear interpretation of the between term: the discrete categories are not real, and the between term cannot be interpreted as a relative index of segregation.

Without loss of generality, let us consider the statistical association between racial groups and covariates \mathbf{x} . Since it is always possible to factorize the generalized density $f_i(g, n, \mathbf{x})$ as

$$f_i(g, n, \mathbf{x}) = f_i(n | g, \mathbf{x}) f_i(g, \mathbf{x}),$$

where $f_i(g, \mathbf{x}) = \sum_{n=1}^N f_i(g, n, \mathbf{x})$, $i = 1, 2$, $I(1:2)$ can always be decomposed into two terms:

$$\begin{aligned} I(1:2) &= \int f_1(g, n, \mathbf{x}) \log \left(\frac{f_1(g, \mathbf{x})}{f_2(g, \mathbf{x})} \right) d\mathbf{m} \\ &\quad + \int f_1(g, n, \mathbf{x}) \log \left(\frac{f_1(n | g, \mathbf{x})}{f_2(n | g, \mathbf{x})} \right) d\mathbf{m} \end{aligned} \quad (3)$$

The first term captures the discrepancy between $f_1(g, \mathbf{x})$ and $f_2(g, \mathbf{x})$, while the second captures the discrepancy in conditional school assignment rules $f_1(n | g, \mathbf{x})$ and $f_2(n | g, \mathbf{x})$. The following assumptions are sufficient to obtain a parametric decomposition of M_{KL} and a characterization of the two terms in the decomposition, as racial discrepancy across covariates and as an expectation of conditional school segregation terms:

A4: $f_1(g, \mathbf{x}) = f(g | \mathbf{x}; \mathbf{a}) f(\mathbf{x})$ where $f(\bullet | \mathbf{x}; \mathbf{a})$ is known up to parameter vector $\mathbf{a} \in \mathbb{R}^{k_a}$.

A5: $f_2(g, \mathbf{x}) = p_{g\bullet} f(\mathbf{x})$ with $p_{g\bullet} = \int_{\mathbf{x} \in \Lambda} f(g | \mathbf{x}; \mathbf{a}) f(\mathbf{x}) d\mathbf{x}$ not uniquely identified by \mathbf{a} .

Assumptions A4 and A5 together imply that the first term in decomposition (3) captures the discrepancy between the statistical association of race status and vector \mathbf{x} under \mathbf{m}_1 and statistical independence. The parametric specification in A4 is not necessary to achieve an interpretable decomposition, but it proves useful in the derivation of the asymptotic properties for proposed estimators of the decomposition in the next section.

A6: $f_1(g, n | \mathbf{x}) = f(g, n | \mathbf{x}; \mathbf{b})$ where $f(\bullet, \bullet | \mathbf{x}; \mathbf{b})$ is known up to parameter vector $\mathbf{b} \in \mathbb{R}^{k_b}$

which is not a function of (g, \mathbf{a}) where $\mathbf{g} = (p_{1\bullet}, \dots, p_{G\bullet})'$.

A7: $f_2(g, n | \mathbf{x}) = f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})$ where $f(n | \mathbf{x}; \mathbf{b}) = \left\{ \sum_{g=1}^G f(g, n | \mathbf{x}; \mathbf{b}) \right\}$.

Taken together, assumptions A6 and A7 provide a segregation interpretation for the second term in equation (3). As with A4, the parametric assumption in A6 facilitates the derivation of the asymptotic

properties in the next section. We can now state the following result:

Proposition 2: Under assumptions A1 to A7,

$$M = M^B(\mathbf{g}, \mathbf{a}) + M^W(\mathbf{a}, \mathbf{b}) \quad (4)$$

where

$$M^B(\mathbf{g}, \mathbf{a}) = \int_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \left\{ \sum_{g=1}^G f(g | \mathbf{x}; \mathbf{a}) \log \left(\frac{f(g | \mathbf{x}; \mathbf{a})}{\hat{p}_g} \right) \right\} d\mathbf{x}$$

and

$$M^W(\mathbf{a}, \mathbf{b}) = \int_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \left\{ \sum_{g=1}^G \sum_{n=1}^N f(g, n | \mathbf{x}; \mathbf{b}) \log \left(\frac{f(g, n | \mathbf{x}; \mathbf{b})}{f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})} \right) \right\} d\mathbf{x}.$$

■

The term $M^B(\mathbf{g}, \mathbf{a})$ identifies the level of segregation which would remain if there were no segregation after controlling for the statistical dependence between the vector of covariates \mathbf{x} and racial status. Since $M^W(\mathbf{a}, \mathbf{b})$ is the level of segregation which is not related to racial discrepancy by covariates \mathbf{x} , it can be referred to as school segregation by race conditional on \mathbf{x} .

Decomposition (4) is appealing for several reasons. First, the between and within terms are independent, in the sense that it is possible to introduce changes in the densities involved to make $M^W(\mathbf{a}, \mathbf{b})$ equal to zero keeping $M^B(\mathbf{g}, \mathbf{a})$ constant. Second, decomposition (4) states that the *KL* discrepancy/segregation measure M_{KL} decomposes into $M^B(\mathbf{g}, \mathbf{a})$ a *KL* measure of discrepancy itself, and $M^W(\mathbf{a}, \mathbf{b})$ a weighted average of conditional *KL* discrepancy/segregation measures with weights equal to the marginal densities for the vector of covariates \mathbf{x} . Finally, the conditional densities $f(g | \mathbf{x}; \mathbf{a})$, $f(n | \mathbf{x}; \mathbf{b})$, and $f(g, n | \mathbf{x}; \mathbf{b})$ can be interpreted as qualitative response models which stem from economic agents' utility maximizing choices under constraints. Thus, decomposition (4) provides an intuitive, unifying, econometric framework for studies of segregation using segregation

indices and qualitative response econometric models.

Kullback (1959) points out that *KL* discrepancies can be recursively decomposed into more than two terms. This is trivially seen with decomposition (3) as the first term is itself a *KL* discrepancy measure and, thus, it can itself be decomposed. A direct application of this property to the problem of multigroup school segregation permits the decomposition of M into three terms capturing between-cities, within-cities, and within-districts school segregation.¹³ For reasons of brevity, we leave to the reader the details of decompositions of more than two terms in the model.

One final point needs to be clarified. Suppose that all covariates \mathbf{x} are discrete and that they partition the set of schools into disjoint subsets. Without loss of generality consider the district versus school segregation problem as representative of this situation. More specifically, assume that each school belongs to one of K different school districts and let p_{gnd} denote the proportion of students of racial group g at school n within district d , $p_{gnd} = p_{g|n \in d}$. Define $p_{g \bullet d}$ as the joint probability of race and district membership and let $p_{\bullet \bullet d}$ and $p_{g \bullet d}$ denote the marginal distribution of districts and the joint distribution of race and school membership conditional on district d , respectively. Finally, let $p_{g \bullet d}$ and $p_{\bullet |d}$ be the marginal distributions within district d of race and school membership. It has previously been shown that the M index is decomposable for any partition of the schools into K school districts into a between and a within term:¹⁴

$$M = M^B + M^W, \quad (5)$$

where

$$M^B = \sum_{d=1}^K \sum_{g=1}^G p_{g \bullet d} \log \left(\frac{p_{g \bullet d}}{p_{g \bullet} p_{\bullet \bullet d}} \right)$$

¹³ See also Herranz *et al.* (2005) for an application of this principle in the context of occupational segregation by gender and Frankel and Volij (2009) for sequential clustering of racial categories in multiracial school segregation.

¹⁴ See Frankel and Volij (2009) and Mora and Ruiz-Castillo (2009). For the two-group case, see Mora and Ruiz-Castillo (2003, 2004), and Herranz *et al.* (2005).

and

$$M^W = \sum_{k=1}^K \hat{p}_{\bullet\bullet d} \sum_{n \in d} \sum_{g=1}^G \hat{p}_{gnd} \log \left(\frac{\hat{p}_{gnd}}{\hat{p}_{g\bullet d} \hat{p}_{\bullet nd}} \right).$$

How does decomposition (5) relate to decomposition (4)? In order to apply Proposition 2 note that since $\mathbf{x} = d$, then: (a) the generalized density $f(\mathbf{x})$ coincides with the unconditional distribution by districts $\hat{p}_{\bullet\bullet d}$; (b) the conditional density $f(g | \mathbf{x}, \mathbf{a})$ turns out to be the conditional marginal distribution by race status within district d , $\hat{p}_{g\bullet d}$; (c) the joint density $f(g, n | \mathbf{x}, \mathbf{b})$ corresponds to the joint race and school probability within district d , \hat{p}_{gnd} ; and (d) the conditional density $f(n | \mathbf{x}, \mathbf{b})$ coincides with the marginal distribution by school membership within district d , $\hat{p}_{\bullet nd}$. Under these conditions, it can readily be shown that $M^B(\mathbf{g}, \mathbf{a}) = M^B$ and $M^W(\mathbf{a}, \mathbf{b}) = M^W$. Thus, when the vector of covariates includes only discrete variables, the general decomposition in equation (4) exactly matches the decomposition of the M index previously proposed in the literature for any partition of schools (or groups) as in equation (5).

III. ESTIMATION AND ASYMPTOTIC PROPERTIES

III.1 Estimation and Asymptotics of the M Index

Assume that a sample of T observations from students with information on their race status, school membership, and covariates, $\{g_i, n_i, \mathbf{x}_i\}$, $i=1, \dots, T$, is available. Let T_{gn} be the number of students of racial group g in school n so that $T = \sum_{n=1}^N \sum_{g=1}^G T_{gn}$. Let $T_{g\bullet} = \sum_{n=1}^N T_{gn} > 0$ and $T_{\bullet n} = \sum_{g=1}^G T_{gn} > 0$. Note that under assumptions A1 to A7 the probabilities \hat{p}_{gn} are fully flexible so that race and school status are jointly distributed as a nonparametric multinomial model. Without loss of

generality, denote by $GN^c = \{(g, n) : g = 1, \dots, G, n = 1, \dots, N, (g, n) \neq (G, N)\}$ the set of all race and school combinations except combination (G, N) . Then the marginal probabilities of race and school membership are fully identified by the vector of dimension $GN - 1$:

$$\mathbf{q} = (p_{11}, p_{21}, \dots, p_{G, N-2}, p_{G, N-1})' \in \Theta \equiv \left\{ \sum_{GN^c} p_{gn} < 1, p_{gn} > 0 \right\} \subset \mathbb{R}^{GN-1}$$

and $p_{GN} = 1 - \sum_{(g,n) \in GN^c} p_{gn} > 0$. The M index is bounded and continuous. Moreover, it can always be

estimated by its sample analogue:

$$\hat{M}_T = \sum_{g=1}^G \sum_{n=1}^N \hat{p}_{gn} \log \left(\frac{\hat{p}_{gn}}{\hat{p}_{g\bullet} \hat{p}_{\bullet n}} \right)$$

where $\hat{p}_{gn} = T_{gn} / T$, $\hat{p}_{g\bullet} = T_{g\bullet} / T$, $\hat{p}_{\bullet n} = T_{\bullet n} / T$, and $0 \log(0) = 0$. We first prove consistency for \hat{M}_T :

Proposition 3: Under assumptions A1 to A3,

$$\text{plim } \hat{M}_T = M. \quad \blacksquare$$

An implication of Proposition 3 is that \hat{M}_T converges in probability to 0 if and only if $p_{gn} = p_{g\bullet} p_{\bullet n}$ for all g and n . Moreover, whenever two cities are to be ranked according to M , the implicit ordering from \hat{M}_T converges in probability to the ordering induced by M as the sample size in each city becomes infinitely large. Proposition 3 also highlights that although \hat{M}_T is size invariant, in the sense that it only depends on relative frequencies (see note 8), its variance will tend to zero as sample size increases.

In the previous Section it was argued that absence of segregation under the notions of evenness and representativeness is equivalent, and that it coincides with the concept of statistical independence between race status and school membership. In statistics, testing $p_{g|n} = p_{g\bullet}$ (or testing $p_{n|g} = p_{\bullet n}$) against the two-sided alternative for any two categorical variables is sometimes presented as testing that

$N(\text{or } G)$ independent trials, each having $G(\text{or } N)$ mutually exclusive results, have been drawn from the same population, while testing $p_{gn} = p_{g\bullet} p_{\bullet n}$ is motivated as testing for the independence of categorical variables in a contingency table. The likelihood ratio test is exactly the same test for the three nulls and, therefore, it is an attractive option when, *a priori*, neither evenness nor representativeness is preferred.¹⁵ A likelihood ratio test for the hypothesis of independence against the two-sided alternative seems *a priori* an appealing statistical measure of deviations from random assignment under both the notion of evenness and the notion of representativeness. Clearly, the larger the value of the statistic, the less likely is the sample under the null hypothesis of absence of segregation.

Given the statistical framework presented in the previous Section, researchers using small samples can either compute \hat{M}_T or carry out a test for the independence of race and school membership, such as the log-likelihood test. We now study the relation between the estimator \hat{M}_T and testing for the independence of racial status and school membership using the log-likelihood ratio test in a small sample. To implement the likelihood ratio statistic, an additional assumption on the conditional density for covariates \mathbf{x} is required. It suffices to assume that this conditional density belongs to a parametric family with k_j -dimensional parameter vector \mathbf{j} which is not a function of \mathbf{q} :

$$\Delta 8 : f(\mathbf{x} | g, n) = f(\mathbf{x} | g, n; \mathbf{j}) \text{ such that } f(\mathbf{x}) = \sum_{g=1}^G \sum_{n=1}^N f(\mathbf{x} | g, n; \mathbf{j}) p_{gn} \text{ and } \mathbf{j} \in \mathbb{R}^{k_j} \text{ does not}$$

depend on \mathbf{q} .

Consider testing for the independence of race and school membership, i.e. $H_0 : p_{gn} = p_{g\bullet} p_{\bullet n}$ for all $(g, n) \in GN^c$, versus $H_1 : p_{gn} \neq p_{g\bullet} p_{\bullet n}$. Suppose that the family of densities $f(\mathbf{x} | g, n; \mathbf{j})$ is known up to parameter vector \mathbf{j} . Let $l(\hat{\mathbf{q}}, \hat{\mathbf{j}}) \equiv \log(L(\hat{\mathbf{q}}, \hat{\mathbf{j}}))$ be the log-likelihood evaluated at the maximum

¹⁵ The likelihood ratio test is the uniformly most powerful test, UMP, when testing a simple null against a simple alternative (the Neyman-Pearson Lemma). For a contingency table with $G = 2$ and $N = 2$, it is uniformly most powerful unbiased, UMPU, although for $G > 2$ or $N > 2$ there is neither a UMP nor a UMPU test (see, *inter alia*, Shao, 1998, p. 391).

likelihood estimator, and let $l(\hat{\mathbf{q}}_0, \hat{\mathbf{j}}_0) \equiv \log(L(\hat{\mathbf{q}}_0, \hat{\mathbf{j}}_0))$ be the log-likelihood for the model under H_0 evaluated at the restricted maximum likelihood estimator, so that

$$-2\log(\mathbf{I}) = -2\left(l(\hat{\mathbf{q}}_0, \hat{\mathbf{j}}_0) - l(\hat{\mathbf{q}}, \hat{\mathbf{j}})\right)$$

is the log-likelihood ratio statistic and \mathbf{I} is the likelihood ratio. The following result establishes the relationship between \mathbf{I} and \hat{M}_T .

Proposition 4: Suppose that assumptions A1, A2, A3 and A8 hold. Then:

$$\hat{M}_T = \frac{-\log(\mathbf{I})}{T}. \quad \blacksquare$$

The interpretation of \hat{M}_T as the monotonic transformation of the likelihood ratio only requires assumption A8, in addition to assumptions A1 to A3. There is no need in A8 to assume knowledge of the density $f(\mathbf{x} | g, n; \mathbf{j})$ up to parameter vector \mathbf{j} because, by assumption A3, the joint sample distribution of race and school membership is a sufficient statistic for the test. The general principle behind Proposition 4 is not new: the relation between Kullback-Leibler discrepancies measures and likelihood-ratio statistics in contingency tables is probably first stated in Kullback (1959, page 158). In the multigroup segregation literature, it has previously been used by Reardon *et al.* (2000) in relation to a normalized version of the M index known as the Entropy or H index.

Proposition 4 states that \hat{M}_T is a monotonic transformation of the likelihood-ratio statistic for testing statistical independence between school membership and racial status. This implies that the ordering across cities provided by comparisons of city-specific log-likelihood ratios, suitably divided by city-specific scalars $2T_c$, is uniquely defined by the eight ordinal properties that characterize the M index according to Frankel and Volij (2009). Why is $-\log(\mathbf{I})$ less appealing than \hat{M}_T as a measure of segregation? The ordering induced by \hat{M}_T is size invariant, while the ordering induced by $-\log(\mathbf{I})$ is

sensitive to sample size for any given set of relative frequencies.

The value $-\log(\mathbf{I})$ can be seen to be a particular case of a general *KL* divergence test for the null hypothesis that r independent samples are drawn from an identical distribution, whose functional form is known up to a vector of parameters of dimension k . Kupperman (1957) showed that, under certain regularity conditions, this general *KL* divergence test is asymptotically distributed as chi-square with $(r-1)k$ degrees of freedom.¹⁶ For the statistical model set up in the previous section, it is possible to invoke earlier well known results on the properties of the \mathbf{I} statistic under the null (see Neyman, 1949) to prove that:

Theorem 1: Suppose that A1, A2, A3, and A8 hold. If $p_{gn} = p_{g\bullet} p_{\bullet n}$, for all $(g,n) \in GN^c$, then:

$$2T\hat{M}_T \xrightarrow{d} \mathbf{c}_{(G-1)(N-1)}^2. \quad \blacksquare$$

Under the null, the finite sample distribution of $-2\log(\mathbf{I})$ can be approximated by the chi-square distribution with the appropriate degrees of freedom. Theorem 1 also implies that although \hat{M}_T is size invariant, its asymptotic confidence intervals are not. The use of $-2\log(\mathbf{I})$ or \hat{M}_T will give “similar” results under the null: the proportion of rejections of the null using the likelihood ratio tends towards the size of the test as sample size increases, whilst \hat{M}_T converges in probability to zero.

In many practical situations, the hypothesis of absence of segregation will be false, so that the relevant statistical properties for the index of segregation will be those under the true alternative. In his 1957 study of the *KL* measure of divergence, Kupperman showed that when the null hypothesis is false *KL* converges in probability to an indefinitely large number, and that the large-sample distribution may be approximated by the non-central chi-square and the same degrees of freedom as under the null. Salicrú *et al.* (1994) studied the asymptotic distribution of a family of estimators for which *KL* divided by

¹⁶ Morales *et al.* (1995) consider \hat{M}_T as a particular case of a more general family of divergence measures between two consistent estimates of a discrete distribution. They find that, under the null, the chi-square distribution is an asymptotic approximation for all members of the family.

sample size is a limiting case. Using the delta method, they find square-root convergence to a normal distribution under the alternative. The following theorem states a similar result for \hat{M}_T :

Theorem 2: Suppose that assumptions A1, A2, A3, and A8 hold. If $p_{gn} \neq p_{g\bullet} p_{\bullet n}$ for at least one $(g, n) \in GN^c$, then:

$$T^{1/2}(\hat{M}_T - M) \xrightarrow{d} N(0, \Delta m' \Sigma \Delta m)$$

where

$$\Sigma = \{\mathbf{s}_{(i,j)(g,n)}\} = \begin{cases} p_{ij}(1-p_{ij}) & \text{if } (i,j) = (g,n) \\ -p_{ij}p_{gn} & \text{if } (i,j) \neq (g,n), \end{cases}$$

$$\Delta m = \{\Delta m_{gn}\} = \log\left(\frac{p_{gn}}{p_{g\bullet} p_{\bullet n}}\right) - \log\left(\frac{p_{GN}}{p_{G\bullet} p_{\bullet N}}\right), \quad \forall (g,n) \in GN^c.$$

■

The asymptotic power of \hat{M}_T can be estimated for fixed alternatives using Theorem 2. Note, however, that the normal approximation will likely be poor if the sample is not large, and bootstrap inference may provide better approximations to the small sample distribution of \hat{M}_T .

III.2 Estimation and Asymptotics of Conditional Segregation

When a sample of *iid* observations of size T is available, estimation of decomposition (5) can be carried out using the principle of analogy. The following four estimators will be considered:

$$\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) = T^{-1} \sum_{i=1}^T \left\{ \sum_{g=1}^G f(g | \mathbf{x}_i; \hat{\mathbf{a}}) \log \left(\frac{f(g | \mathbf{x}_i; \hat{\mathbf{a}})}{\hat{p}_{g\bullet}} \right) \right\},$$

$$\hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = T^{-1} \sum_{i=1}^T \left\{ \sum_{g=1}^G \sum_{n=1}^N f(g, n | \mathbf{x}_i; \hat{\mathbf{b}}) \log \left(\frac{f(g, n | \mathbf{x}_i; \hat{\mathbf{b}})}{f(g | \mathbf{x}_i; \hat{\mathbf{a}}) f(n | \mathbf{x}_i; \hat{\mathbf{b}})} \right) \right\},$$

$$\hat{M}_T^W(\hat{\mathbf{g}}, \hat{\mathbf{a}}) = \hat{M}_T - \hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}),$$

and

$$\hat{M}_T(\hat{\mathbf{g}}, \hat{\mathbf{a}}, \hat{\mathbf{b}}) = \hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) + \hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}}),$$

where $(\hat{\mathbf{g}}, \hat{\mathbf{a}}, \hat{\mathbf{b}})$ are estimates for $(\mathbf{g}, \mathbf{a}, \mathbf{b})$. For the case in which all covariates \mathbf{x} are discrete and partition the set of schools, the sample analogues of decomposition (5) require no functional form assumptions for the densities of the variables. In this case:

$$\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) = \hat{M}_T^B = \sum_{d=1}^K \sum_{g=1}^G \hat{p}_{g \bullet d} \log \left(\frac{\hat{p}_{g \bullet d}}{\hat{p}_{g \bullet} \hat{p}_{\bullet d}} \right)$$

$$\hat{M}_T^W(\hat{\mathbf{a}}; \hat{\mathbf{b}}) = \hat{M}_T^W(\hat{\mathbf{g}}; \hat{\mathbf{b}}) = \hat{M}_T^W = \sum_{k=1}^K \hat{p}_{\bullet \bullet d} \sum_{i \in d} \sum_{g=1}^G \hat{p}_{g i d} \log \left(\frac{\hat{p}_{g i d}}{\hat{p}_{g \bullet i} \hat{p}_{\bullet i d}} \right)$$

and $\hat{M}_T = \hat{M}_T^B + \hat{M}_T^W$ always.

The remainder of this section is devoted to the study of the asymptotic properties of these estimators. Results for the case in which all covariates are discrete and for the case in which at least one covariate is not discrete are presented in III.2.1 and III.2.2, respectively.

III.2.1 Estimation and Asymptotics of District versus School Segregation

Decomposition (5) aims to answer to what extent race segregation at district level can explain a significant amount of school segregation by race. In this section, the statistical properties of the estimators for decomposition (5), \hat{M}_T^B and \hat{M}_T^W , are studied. By interchanging the notation for groups and organizational units, the results presented here can be applied to decompositions when the set of racial groups is partitioned into supergroups.

The term \hat{M}_T^B is itself a M index so that, by Proposition 3, it converges in probability to the KL measure M^B and, by Proposition 4, can be motivated as the likelihood-ratio test for the independence between race and district membership. Theorems 1 and 2 also apply to \hat{M}_T^B directly after a trivial change in notation. Therefore, under the conditions of Theorem 2, if $p_{gn} \neq p_{g \bullet} p_{\bullet n}$ for at least one

$(g,d) \in GD^c \equiv \{(g,d) : (g,d) \neq (G,D)\}$, then :

$$T^{1/2} \left(\hat{M}_T^B - M^B \right) \xrightarrow{d} N \left(0, \Delta m^B, \Sigma^B \Delta m^B \right)$$

where

$$\Sigma^B = \left\{ \mathbf{s}_{(i,j) \in (g,d)} \right\} = \begin{cases} \hat{p}_{ij} (1 - \hat{p}_{ij}) & \text{if } (i,j) = (g,d) \\ -\hat{p}_{ij} \hat{p}_{g^*} & \text{if } (i,j) \neq (g,d), \end{cases}$$

$$\Delta m^B = \left\{ \Delta m_{gd}^B \right\} = \log \left(\frac{\hat{p}_{gd}}{\hat{p}_{g^*} \hat{p}_{\bullet d}} \right) - \log \left(\frac{\hat{p}_{GD}}{\hat{p}_{G^*} \hat{p}_{\bullet D}} \right), \quad \forall (g,d) \in GD^c.$$

Similarly, under the conditions of Theorem 1, if $\hat{p}_{gd} = \hat{p}_{g^*} \hat{p}_{\bullet d}$ for all $(g,d) \in GD^c$, then

$$2T \hat{M}_T^B \xrightarrow{d} \mathbf{C}_{(G-1)(D-1)}^2.$$

The within-term \hat{M}_T^W can also be motivated as a likelihood-ratio test. Consider testing for the independence of race and school membership within any district d , i.e.

$H_0 : p_{gn_d|d} = p_{g^*|d} p_{\bullet n_d|d}$, $(g, n_d) \in GN_d^c$, $d = 1, \dots, D$, versus the alternative $H_1 : p_{gn_d|d} \neq p_{g^*|d} p_{\bullet n_d|d}$ for at

least one combination (g, n_d, d) . Let $l \left(\left\{ \hat{p}_{gn_d} \right\}, \left\{ \hat{p}_{\bullet \bullet d} \right\} \right) \equiv \log \left(L \left(\left\{ \hat{p}_{gn_d} \right\}, \left\{ \hat{p}_{\bullet \bullet d} \right\} \right) \right)$ be the log-likelihood

evaluated at the maximum likelihood estimator, and let $l \left(\left\{ \hat{p}_{gn_d}^0 \right\}, \left\{ \hat{p}_{\bullet \bullet d}^0 \right\} \right) \equiv \log \left(L \left(\left\{ \hat{p}_{gn_d}^0 \right\}, \left\{ \hat{p}_{\bullet \bullet d}^0 \right\} \right) \right)$ be

the log-likelihood for the model under H_0 evaluated at the restricted maximum likelihood estimator, so

that

$$-2 \log(\mathbf{I}^W) = -2 \left(l \left(\left\{ \hat{p}_{gn_d}^0 \right\}, \left\{ \hat{p}_{\bullet \bullet d}^0 \right\} \right) - l \left(\left\{ \hat{p}_{gn_d} \right\}, \left\{ \hat{p}_{\bullet \bullet d} \right\} \right) \right)$$

is the log-likelihood ratio statistic and \mathbf{I}^W is the likelihood ratio. The following result establishes the

relation between \mathbf{I}^W and \hat{M}_T^W .

Proposition 5: Suppose that assumptions A1 to A8 hold, and that the vector \mathbf{x} includes only district code d . Then:

$$\hat{M}_T^W = \frac{-\log(\mathbf{I}^W)}{T}. \quad \blacksquare$$

Proposition 5 provides an intuitive statistical interpretation for \hat{M}_T^W . We are not aware of any other within-groups term in a decomposition of an index of segregation so closely related to a classical statistical test. The generality of this result should not pass unnoticed: Proposition 5 can be applied to any cluster of school districts, such as cities or regions, so that the within terms in the resulting decompositions can be interpreted as monotone transformations of likelihood-ratio tests for the independence between race and school membership within the districts of the corresponding cluster.

The discussion of the discrete case ends with a theorem which characterizes the asymptotic properties of \hat{M}_T^W . Let us first introduce some notation. Denote by

$$GN_d^c = \{(g, n_d) : g = 1, \dots, G, n_d \in \mathcal{S}_d, (g, n_d) \neq (G, N_d)\}$$

the set of all race and school combinations in district d except combination (G, N_d) . Let

$\mathbf{q}_d = (p_{1|d}, \dots, p_{G, N_d - 1|d})^t$ for all $d = 1, \dots, D$, and define the function $m_W(\mathbf{q}_W)$ of parameter vector $\mathbf{q}_W = (p_{\bullet\bullet 1}, \dots, p_{\bullet\bullet D-1}, \mathbf{q}_1^t, \dots, \mathbf{q}_D^t)^t$ as

$$m_W(\mathbf{q}_W) = \sum_{d=1}^{D-1} p_{\bullet\bullet d} m_d(\mathbf{q}_d) + \left(1 - \sum_{d=1}^{D-1} p_{\bullet\bullet d}\right) m_D(\mathbf{q}_D)$$

where

$$m_d(\mathbf{q}_d) = \sum_{GN_d^c} p_{gnd} \log \left(\frac{p_{gnd}}{p_{g\bullet d}(\mathbf{q}_d) p_{\bullet nd}(\mathbf{q}_d)} \right) + \left(1 - \sum_{GN_d^c} p_{gnd}\right) \log \left(\frac{1 - \sum_{GN_d^c} p_{gnd}}{p_{G\bullet d}(\mathbf{q}_d) p_{\bullet N_d}(\mathbf{q}_d)} \right)$$

for all $d = 1, \dots, D$, and $p_{g\bullet d}(\mathbf{q}_d)$ and $p_{\bullet nd}(\mathbf{q}_d)$ are defined in a similar way to $p_{g\bullet}(\mathbf{q})$ and $p_{\bullet n}(\mathbf{q})$. The following result can now be stated:

Theorem 3: Suppose that assumptions A1 to A8 hold, and that the covariates vector \mathbf{x} includes only district code d .

(a) Assume further that there is at least one district d , such that for at least two race and school combinations, (g, n) and (r, s) , we have that

$$\frac{1}{\hat{p}_{gnd}} + \frac{1}{\hat{p}_{r\bar{s}d}} \neq \frac{1}{\hat{p}_{g\bullet d}} + \frac{1}{\hat{p}_{\bullet\bar{s}d}} + \frac{1}{\hat{p}_{r\bullet d}} + \frac{1}{\hat{p}_{\bullet\bar{s}d}}. \quad (6)$$

If $\hat{p}_{gnd} = \hat{p}_{g\bullet d} \hat{p}_{\bullet n d}$ for all (g, n, d) , then

$$T\hat{M}_T^W \xrightarrow{d} ZAZ'$$

with

$$Z = N(0, \Sigma_W),$$

$$\Sigma_W = \left\{ \mathbf{s}_{(i,j,k) \neq (g,n,d)} \right\} = \begin{cases} \hat{p}_{\bullet\bullet d} (1 - \hat{p}_{\bullet\bullet d}) & \text{if } (0,0,d) = (0,0,d) \\ -\hat{p}_{\bullet\bullet d} \hat{p}_{\bullet\bullet k} & \text{if } (0,0,d) \neq (0,0,k) \\ \hat{p}_{ij\bar{d}} (1 - \hat{p}_{ij\bar{d}}) & \text{if } (i,j,d) = (g,n,d) \\ -\hat{p}_{ij\bar{d}} \hat{p}_{gnd} & \text{if } (i,j,d) \neq (g,n,d) \end{cases},$$

and

$$A = \left(\frac{1}{2} \left(\frac{\partial^2 m_W(\mathbf{q}_W)}{\partial \mathbf{q}_{W_i} \partial \mathbf{q}_{W_j}} \Big|_{\mathbf{q}_W = \mathbf{q}_W^0} \right)_{(GN-1) \times (GN-1)} \right).$$

(b) If $\hat{p}_{gnd} \neq \hat{p}_{g\bullet d} \hat{p}_{\bullet n d}$ for at least one (g, n, d) , then

$$T^{1/2} \left(\hat{M}_T^W - M^W \right) \xrightarrow{d} N \left(0, \Delta m^W \Sigma_W \Delta m^W \right)$$

where

$$\Delta m_W = \left\{ \frac{\partial m_W}{\partial \mathbf{q}_{W_i}} \right\} = \begin{cases} m_d(\mathbf{q}_d) - m_D(\mathbf{q}_D) & \text{if } \mathbf{q}_{W_i} = \hat{p}_{\bullet\bullet d} \\ \log \left(\frac{\hat{p}_{gnd}}{\hat{p}_{g\bullet d} \hat{p}_{\bullet n d}} \right) - \log \left(\frac{\hat{p}_{GN_d|d}}{\hat{p}_{G\bullet d} \hat{p}_{\bullet N|d}} \right) & \text{if } \mathbf{q}_{W_i} = \hat{p}_{gnd} \end{cases}$$

■

Regularity condition (6) is a sufficient condition to ensure that the second-order term in the Taylor approximation does not vanish at the true parameter value. Note that under absence of segregation, the

within term will not generally be approximated by a chi-square distribution. Intuitively, this term is a weighted average of terms which, by Theorem 1, each must converge under the null to a chi-square distribution.

III.2.2 Conditional Segregation with Non-Discrete Covariates

Although decomposition (4) includes decomposition (5) as a special case, the former cannot be considered a true generalization of the latter. The reason is that, while in the finite-covariates situation no restrictive functional-form assumptions for the conditional densities are required to implement decomposition (4), in the presence of countable or continuous covariates functional-form restrictions are implicit in parametric assumptions A4 to A7, and are therefore identifying. This has two important empirical implications. First, the sum $\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) + \hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ need not be equal to \hat{M}_T for small samples. Second, there is generally no monotonous relation between the likelihood-ratio test and $\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}})$ and $\hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}})$.

Clearly, the fact that $\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}})$ and $\hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ are unrelated to likelihood-ratio tests does not imply that they cannot be interpreted as statistical tests. The asymptotic properties of both estimators are next studied under different hypotheses, therefore providing their asymptotic motivation as statistical tests. Moreover, since both estimators will be shown to be consistent under general regularity conditions, the sum $\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) + \hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ will converge in probability to M . Sufficient conditions for asymptotic normality for $\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}})$ are obtained in the following result:

Theorem 4: Suppose that assumptions A1 to A5 hold and that the vector of covariates \mathbf{x} includes at least one countable or continuous variable. Let $(\mathbf{g}_0, \mathbf{a}_0)$ be the true parameter vectors of the data

generating process and define $b_B(\mathbf{x}; \mathbf{g}, \mathbf{a}) = \sum_{g=1}^G f(g | \mathbf{x}; \mathbf{a}) \log \left(\frac{f(g | \mathbf{x}; \mathbf{a})}{p_g} \right)$. Assume that:

(R1) $b_B(\mathbf{x}; \mathbf{g}, \mathbf{a})$ is differentiable with respect to (\mathbf{g}, \mathbf{a}) , with continuous partial derivatives which

are nonvanishing at $(\mathbf{g}_0, \mathbf{a}_0)$.

$$(R2) \text{Var}_{\mathbf{x}} [h_B(\mathbf{x}; \mathbf{g}_0, \mathbf{a}_0)] = \mathbf{s}_B^2 < \infty;$$

$$E_{\mathbf{x}} \left[\frac{\partial h_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \mathbf{m}_g \in \mathbb{R}^{k_g}, E_{\mathbf{x}} \left[\frac{\partial h_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \mathbf{m}_a \in \mathbb{R}^{k_a}, \text{Var}_{\mathbf{x}} \left[\frac{\partial h_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \Sigma_g, \text{Var}_{\mathbf{x}} \left[\frac{\partial h_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \Sigma_a,$$

where Σ_g and Σ_a are positive definite matrices.

$$(R3) \text{plim } \hat{\mathbf{g}} = \mathbf{g}_0 \text{ and } \text{plim } \hat{\mathbf{a}} = \mathbf{a}_0.$$

Then,

$$T^{1/2} \left(\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) - M^B(\mathbf{g}_0, \mathbf{a}_0) \right) \xrightarrow{d} N(0, \mathbf{s}_B^2). \quad \blacksquare$$

The normal approximation given in Theorem 4 covers both the case of existence of racial discrepancy across covariates (when $M^B(\mathbf{g}, \mathbf{a}) > 0$), as well as the case of no racial discrepancy across covariates \mathbf{x} ,

i.e. when $M^B(\mathbf{g}, \mathbf{a}) = 0$. The last result refers to the asymptotic properties for $\hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}})$:

Theorem 5: Suppose that assumptions A1 to A7 hold, and that the vector of covariates \mathbf{x} includes at least one countable or continuous variable. Let $(\mathbf{a}_0, \mathbf{b}_0)$ be the true parameter vectors of the

data generating process, and define $h_W(\mathbf{x}; \mathbf{a}, \mathbf{b}) = \sum_{g=1}^G \sum_{n=1}^N f(g, n | \mathbf{x}; \mathbf{b}) \log \left(\frac{f(g, n | \mathbf{x}; \mathbf{b})}{f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})} \right)$

Assume that:

(R4) $h_W(\mathbf{x}; \mathbf{a}, \mathbf{b})$ is differentiable with respect to (\mathbf{a}, \mathbf{b}) , with continuous partial derivatives which are non-vanishing at $(\mathbf{a}_0, \mathbf{b}_0)$.

$$(R5) \text{Var}_{\mathbf{x}} [h_W(\mathbf{x}; \mathbf{a}_0, \mathbf{b}_0)] = \mathbf{s}_W^2 < \infty,$$

$$E_{\mathbf{x}} \left[\frac{\partial h_W}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \mathbf{m}_a \in \mathbb{R}^{k_a}, E_{\mathbf{x}} \left[\frac{\partial h_W}{\partial \mathbf{b}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \mathbf{m}_b \in \mathbb{R}^{k_b}, \text{Var}_{\mathbf{x}} \left[\frac{\partial h_W}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \Sigma_a, \text{and } \text{Var}_{\mathbf{x}} \left[\frac{\partial h_W}{\partial \mathbf{b}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \Sigma_b,$$

where Σ_a and Σ_b are positive definite matrices.

$$(R6) \text{plim } \hat{\mathbf{a}} = \mathbf{a}_0 \text{ and } \text{plim } \hat{\mathbf{b}} = \mathbf{b}_0.$$

Then,

$$T^{1/2} \left(\hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}}) - M^W(\mathbf{a}_0, \mathbf{b}_0) \right) \xrightarrow{d} N(0, \mathbf{S}_W^2). \quad \blacksquare$$

Asymptotic normality for $\hat{M}_T^W(\hat{\mathbf{g}}, \hat{\mathbf{a}}) = \hat{M}_T - \hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}})$ and $\hat{M}_T(\hat{\mathbf{g}}, \hat{\mathbf{a}}, \hat{\mathbf{b}}) = \hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) + \hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ follows directly from Theorems 1, 2, 4, and 5. Note that these results apply to any set of consistent estimators for $(\mathbf{g}, \mathbf{a}, \mathbf{b})$. For ML estimators, consistency conditions R3 and R6 will be satisfied under general regularity conditions.

IV. U.S. SCHOOL SEGREGATION

During the past decades, the U.S. has become increasingly racially and ethnically diverse, due to higher fertility and/or immigration rates among minorities, which have led to a faster population growth than the white population. The demographic advances of Hispanics and Asians are concentrated in certain parts of the country, while scarcely apparent in many others.¹⁷ The main objective of this section is to illustrate the usefulness of the M index, the decompositions proposed in this paper, and the asymptotic results obtained to analyze racial school segregation patterns under the described changing environment.

IV.1. Data

We use the Common-Core of Data (CCD) compiled by the National Center for Educational Statistics (NCES). This dataset contains school enrolment records according to racial/ethnic group from all public schools in the United States. Results are reported for the school years 1989-90 (the first year for which complete enrolment data is available) and 2005-06. Schools are retrospectively assigned Core

¹⁷ The terms *white*, *black*, *Asian* and *Native American* are used throughout this section to refer to non-Hispanic members of these racial groups. Asians include Native Hawaiian and Pacific Islanders; Native Americans include American Indians and Alaska Natives (Innuit or Aleut). The term *Hispanic* is an ethnic rather than a racial category since Hispanic persons may belong to any race. The term racial group is used throughout to refer to each of these five racial/ethnic categories.

Based Statistical Area (CBSA) codes based on 2005 ZIP codes so that comparisons over time can be made, without changes in city boundary definitions affecting the results.¹⁸ The sample is restricted to open regular schools¹⁹ located in 960 CBSA codes –referred to as “cities” – in the 50 states and the District of Columbia. This covers approximately 74% of the student population attending U.S. public schools in 2005.

There is full information for all targeted schools in 2005. For 1989, however, a number of schools in a few states missed reporting the data. As a consequence, data for 1989 only includes 839 cities.²⁰ Unless otherwise specified, results pertain to those schools for which racial and ethnic information is available both in 1989 and in 2005. Focusing on the schools which provide information in both years probably gives a fairer comparison between the distributions observed in 1989 and in 2005, since it does not include those schools which reported in 2005 but failed to do so in 1989. However, interpretability of the results is also potentially compromised by the fact that some schools have been created whilst others have disappeared between 1989 and 2005. Nevertheless, use of all observations does not significantly change the results (available upon request), suggesting that the selection mechanisms at work are not essential to our analysis.

IV.2. Racial School Segregation in the U.S.: 1989 and 2005

Table 1 presents the 1989 and 2005 school enrolment by race, the overall racial compositions in the U.S. urban public schools, and the *M* index of overall racial school segregation. Native American, Asian, black, and Hispanic students already made up 34.8% of the total enrolment in 1989. Since growth rates in minority enrolment were larger than among whites, by 2005 minorities accounted for almost half, 48.05%, of total enrolment. Although all minority racial groups have increased their share at the

¹⁸ CBSAs were published by the Office of Management and Budget in 2003 and refer collectively to urban clusters of at least 10,000 people. They replace the Metropolitan Statistical Areas (MSAs) which were used during the 1950-2000 period.

¹⁹ These are all operational schools, except those focused on vocational, special, or other alternative types of education.

²⁰ Reardon *et al.* (2000) study the public school population between 1989 and 1995 in 217 out of 323 MSAs, as defined by the Census Bureau in 1993. Frankel and Volij (2009) present results only for the 2005/2006 school year, restricting the sample to districts in CBSAs with at least two schools which serve grades K-12. Thus, the three papers study a similar phenomenon, although ours covers a larger population during a longer period.

expense of white students, the largest increases are by far from Hispanics, who, in 2005, were already the largest minority group in U.S. public schools.

The last row of Table 1 presents the M index, which measures the expected information of the message that transforms the set of U.S. racial shares presented in the previous panel of Table 1 to the set of schools racial shares. In 1989, the index of segregation (multiplied by 100) is 43.92 and the index increases by 11.3% to 48.90 between 1989 and 2005.

Table 1

IV.3. District vs. School Segregation in the U.S.

Schools are organized into a set of school districts which are themselves organized into a set of cities.²¹ Thus, the overall index of segregation can be decomposed into three terms. The first term results from differences in racial shares between the cities and the national racial shares, so that it can be referred to as BC (Between Cities segregation). The second term captures differences in racial shares between the cities and the educational districts, and is referred to as WC (Within Cities segregation). Finally, the last term in the decomposition captures differences in racial shares between the districts and the schools, and is referred to as WD (Within Districts segregation). Table 2 presents the decomposition of overall racial school segregation into the three components both for 1989 and 2005. The results are in line with those reported in previous empirical studies. First, the BC term, closely linked to parental choices of residence at city level, contributes the most to overall school segregation. Second, the WD term, closely linked to the district educational authorities' decisions, is a relatively small part—around 19%—of overall segregation.

Table 2

Theorems 2 and 3 from the previous sections can be invoked to justify the use of resampling methods. Table 2 presents 5% confidence intervals based on the normal approximation. Upper and

²¹ The data originally consist of 5,834 districts in 1989 and 7,704 districts in 2005. For the common sample there are 5,429 districts.

lower limits are obtained from bootstrap estimates of the variance, using 250 bootstrap samples of each individual student racial status within schools. Given the very high level of aggregation and the large sample sizes, it is hardly surprising to confirm that all terms are significantly different from zero. Looking at differences between the 1989 and 2005 results, there is supportive evidence that the increases observed in all terms are also significant, and of the same order of magnitude: near 12% for BC, 11.5% for WC, and 10% for WD.²²

Aggregation at national level may mask large differences in segregation at city and district level. It is important to note that, by Proposition 5, both WC and WD can be interpreted as likelihood ratios for the null of absence of segregation in any of the more than 800 cities, or the more than 5,000 districts. Clearly, this does not imply that there is segregation in all cities and all districts. A direct way to find out how many cities and districts have significant levels of segregation is to look directly at each city and each district's local index of racial segregation. By Proposition 4, these indices are formal tests for the independence of racial and organizational unit status within each geographical cluster. By Theorem 1, the distribution of these local indices under absence of segregation can be approximated using the chi-square distribution with the appropriate degrees of freedom. A naïve procedure to assess how many cities and districts present significant levels of segregation applies the test to each city and district and then counts those cities and districts for which the null cannot be rejected at a given confidence level. Using this procedure with a 1% confidence interval, segregation is found to be significant in all cities and the vast majority (99%) of districts in both years. This approach has the well-known drawback that, by design, we should expect a positive number of rejections even if the null is always true, simply because of the large number of times the test is performed. Several corrections have been proposed in the literature (see, for example, Romano *et al.*, 2008). Using the Holm correction, segregation remains significant in all cities and in most districts, although the percentage of districts for which segregation is significantly different from zero slightly decreases (98%).

²² Using all schools in 1989 and 2005 results in slightly larger increases for WC: around 16%.

A related question is whether segregation levels are very different among cities and districts. Given that the M index represents a unique ordering of clusters of the organizational units satisfying a set of desirable properties, it is useful to address this question by assessing whether the ranking of cities and the ranking of districts is significant. There are several ways to define ranking significance. For brevity, here we only mean whether the position in the ranking for each of the cities and each of the districts is precisely estimated. One simple way to address this issue is by bootstrapping the rankings and reporting basic bootstrap confidence limits for the ranking for each city and district. Figure 1 presents this information graphically. The y -axis for each plot shows both the index values and 10% bootstrap confidence limits for each city and district by year. Cities and districts with the lowest levels of segregation are ranked first so that they are represented to the left on the x -axis. Thus, all graphs present a positive slope by construction. Figure 1 shows that school districts with large segregation values tend to be ranked more precisely than school districts with low segregation values. The rank of those districts with the lowest levels of segregation is, in fact, very poorly estimated and its confidence intervals often range in the hundreds of positions. Regarding cities, however, the availability of large samples allows us to obtain precise estimates of the rank in most cases. Finally, a note of caution is due regarding the interpretation of Figure 1: since the ordering is specific for each year, Figure 1 does not show ranking dynamics between 1989 and 2005.

IV.4. Multigroup Conditional School Segregation: The Role of Income, Wages, and Teachers per Pupil.

This subsection considers to what extent the measures of WC and WD presented so far are due to the statistical association between racial group membership and socioeconomic covariates using the methodological framework developed in subsection III.2.2. We focus on two sets of controls. First, it has been argued in section II.2 that, given that household income is a potential determinant of residential and school choice, it would be interesting to identify the extent to which multigroup school

segregation arises from income differences across races. In addition, residential choices may potentially be affected by the composition of earnings into wage and non-wage income in the presence of credit market restrictions. Therefore, we would hope to identify the extent to which multigroup school segregation arises from race differentials in the wage to income ratio. Second, the impact of class size, or its inverse –the number of teachers per pupil– on academic performance and other outcomes has long been subject to debate in academic studies and political circles, where the reduction of class sizes is frequently seen as an operational way for educational authorities to effectively increase resources in schools with special needs. At the same time, parents aware of the potential positive effects of small class size on their children’s educational achievements will likely make their residential and school choices dependent on how schools differ in this dimension. Thus, it would be interesting to identify the extent to which multigroup school segregation arises from class size differences across schools. Our empirical illustration tentatively addresses these two issues by merging the CCD data with aggregated measures of income and wages at county level. We specifically study the contribution to the measurements of WC and WD in 2005 of the discrepancy in the racial mix by city and by district for different values of average annual per capita income, annual wages, and teachers per pupil at county or school level.²³

Both for WC and WD, we estimate components $M^B(\mathbf{g}, \mathbf{a})$ and $M^W(\mathbf{g}, \mathbf{a}) = M - M^B(\mathbf{g}, \mathbf{a})$ using estimates of conditional densities $f(g | \mathbf{x}; \mathbf{a})$ based on logistic regressions carried out at district and school level. In particular, for each racial group at city (district) level we assume that $f(g_i | \mathbf{x}_i; \mathbf{a}) = \left(1 + e^{-\mathbf{a}_i' \mathbf{x}_i}\right)^{-1}$ where \mathbf{x}_i includes average annual per capita income, annual wages, and teachers per pupil at county (school) level in addition to dummy variables for city (district) to control for

²³ Since income includes non-wage income, income and wages are not perfectly collinear. The variable *teachers per pupil* at school level can be constructed using the information on the number of teachers and pupils reported by most schools since 2002. County codes, also available from 2002 onwards, allow us to merge the 2005 dataset with the 2004 annual per capita personal income and average wage per job by county published by the Bureau of Economic Activity of the U.S. Department of Commerce. In our county sample, the correlation between the two variables is 0.73.

between city (district) segregation. Logistic regressions for each of the five racial groups are run, using as the dependent variable in each of the regressions the logistic transformation of the observed frequency of students of a given race in a given district (school), and as controls the averages at district (school) level for wages, income, and teachers per pupil, in addition to the city (district) dummies. Table 3 presents a summary of the results.

Table 3

Estimated Marginal effects can be interpreted as the expected change in probability (in percentage terms) associated with a one-percentage increase in each of the controls. For example, a 1% increase in per capita personal income at district level is associated with a 0.83% expected increase in the probability of a student being white, and a 0.32% expected decrease in the probability of a student being black. At district level, increases in per capita income *ceteris paribus* are associated to increases in whites and Asians and decreases in blacks and Hispanics, whilst increases in wages per job are associated to decreases in whites and increases in blacks and Hispanics. These results arguably reflect both the higher probability of black and Hispanic students of having parents who have lower overall per capita income and who are more likely to be salaried workers. Increases in teachers per pupil are associated with decreases in all minority groups. With respect to whites, the point estimate of the relation is positive, although statistically not significant. At school level, increases in county per capita income are associated again with significant increases in whites and significant decreases in blacks. The signs of the estimates for the other groups are similar to those obtained for the district level regression, but the estimates are not significant. With respect to wages per job, results are again similar for whites, Hispanic and blacks, while the parameter estimates for Native American and Asians are not significant. Finally, the effect of increases in teachers per pupil is reversed for blacks at school level: a 1% increase in the teachers per pupil in a school increases by 0.05% the probability that a given student is black. Obviously, a causality interpretation should not be attached to these estimates. Nevertheless, the strong significance of these

effects suggests that a significant part of WC and WD stems from the statistical association of these covariates with race. This central issue is addressed in the last panel of Table 3.

Once estimates $\hat{\mathbf{a}}$ are obtained using the logistic regressions carried out at district and school level, the term $M^B(\mathbf{g}, \mathbf{a})$ can be estimated using $\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}})$ with $\hat{\mathbf{g}} = (\hat{p}_{1\bullet}, \dots, \hat{p}_{G\bullet})'$. The term “All controls” represents $\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}})$ at city and district level (segregation at city and district level stemming from the statistical association between race membership and per-capita personal income, wages per job, and teachers per pupil). Using the results from Theorem 4, asymptotic standard errors are shown in parenthesis. Results show that most (around 64%) of WC while over 20% of WD is accounted for by these three covariates. These effects are significant even in the within districts case. Finally, to evaluate the potentially attenuating effect on segregation of teachers per pupil, the conditional segregation terms are simulated as if this control had no effect. The average effect remains the same for WD, while it decreases very slightly for WC.

V. CONCLUSIONS

The starting point of this paper is the use of the Kullback-Leibler notion of discrepancy (Kullback and Leibler, 1951) to propose a measure of segregation within a general statistical framework. Under general conditions, this measure coincides with the Mutual Information index of segregation, M , first introduced by Theil and Finizza (1971). Elsewhere, it has been shown that M is well grounded on an axiomatic notion of segregation (Frankel and Volij, 2009). In this paper, two specific issues are addressed in relation to this index: the evaluation of statistical significance for observed differences in M measurements, and the control for the statistical association between demographic groups and schools and other socioeconomic variables. Among the main results of the paper, it is shown that M can be decomposed to isolate segregation conditional on any vector of socioeconomic characteristics.

Furthermore, consistent estimators for M and the terms in its decomposition are proposed, and their asymptotic properties are obtained. As a result, the M index now stands as the only index of segregation which has been fully characterized in terms of axiomatic properties, is well embedded into a general statistical framework, and can be used when samples are finite and a multivariate framework is required.

The usefulness of the approach is illustrated by looking at patterns of multigroup school segregation in the U.S. for the school years 1989-90 and 2005-06. Several interesting results stem from direct application of the tools developed in the paper.

Overall multigroup school segregation, which is measured as the discrepancy between the set of U.S. racial shares to the set of schools' racial shares, is significantly positive and has significantly increased during the 15-year period. In the decomposition of overall segregation into between city, within city, and within district segregation, the findings are in line with previous studies: between city segregation, closely linked to parental choices of residence at city level, contributes the most to overall school segregation. In contrast, within districts segregation, potentially linked to policies by the district educational authorities, represents around 19% of overall segregation. All terms in the decomposition of overall segregation are significantly different from zero, and evidence is found that all of them significantly increased during the period.

Aggregation at national level may mask large differences in segregation at city and district level. However, when segregation is studied recursively by city and district, it is found to be significant in all cities and the vast majority of districts in both years. A related question is whether the ranking of cities and the ranking of districts is significant. Using bootstrap techniques, it is found that the rank of those districts with the lowest levels of segregation is, in fact, very poorly estimated and confidence intervals often range in the hundreds of positions. Regarding the ranking of cities, however, the availability of large samples allows us to obtain precise rank estimates for most cities. Finally, we study to what extent the measures of within city and within district segregation are due to the statistical association between

racial group membership and three continuous variables: annual per capita county income, hourly wages at county level, and teachers per pupil at district and school level. Results show that around 64% and 20% of within city and within district segregation is accounted for by these three covariates, and that the effects are strongly significant.

APPENDIX

Proposition 1: If, for all $g = 1, \dots, G$, and $n = 1, \dots, N$,

A1 $p_{gn} > 0$.

A2 $f_i(\mathbf{x} | g, n) = f(\mathbf{x} | g, n) > 0$ as, $i = 1, 2$.

A3 $\mathbf{m}_2(g, n) = p_{g\bullet} \cdot p_{\bullet n} = \left(\sum_{n=1}^N p_{gn} \right) \left(\sum_{g=1}^G p_{gn} \right)$.

then $M_{KL} = M$.

Proof. Under assumptions A1-A3 $\frac{f_1(e, s, \mathbf{x})}{f_2(e, s, \mathbf{x})} = \frac{p_{gn}}{p_{g\bullet} \cdot p_{\bullet n}} = \frac{p_{n|g}}{p_{\bullet n}}$ so that

$$\begin{aligned} M_{KL} &= \sum_{g=1}^G \sum_{n=1}^N p_{gn} \log \left(\frac{p_{n|g}}{p_{\bullet n}} \right) \\ &= \sum_{g=1}^G \sum_{n=1}^N p_{n|g} p_{g\bullet} \log \left(\frac{p_{n|g}}{p_{\bullet n}} \right) \\ &= \sum_{g=1}^G p_{g\bullet} \sum_{n=1}^N p_{n|g} \log \left(\frac{p_{n|g}}{p_{\bullet n}} \right) = M. \end{aligned}$$

■

Proposition 2: If, in addition to assumptions A1-A3, assume that

A4 $f_1(g, \mathbf{x}) = f(g | \mathbf{x}; \mathbf{a}) f(\mathbf{x})$ where $f(\bullet | \mathbf{x}; \mathbf{a})$ is known up to parameter vector $\mathbf{a} \in \mathbb{R}^{k_a}$.

A5 $f_2(g, \mathbf{x}) = p_{g\bullet} f(\mathbf{x})$ with $p_{g\bullet} = \int_{\mathbf{x} \in \Lambda} f(g | \mathbf{x}; \mathbf{a}) f(\mathbf{x}) d\mathbf{x}$ not uniquely identified by \mathbf{a} .

A6 $f_1(g, n | \mathbf{x}) = f(g, n | \mathbf{x}; \mathbf{b})$ where $f(\bullet, \bullet | \mathbf{x}; \mathbf{b})$ is known up to parameter vector $\mathbf{b} \in \mathbb{R}^{k_b}$ which is not a function of (\mathbf{g}, \mathbf{a}) where $\mathbf{g} = (p_{1\bullet}, \dots, p_{G\bullet})'$.

A7 $f_2(g, n | \mathbf{x}) = f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})$ where $f(n | \mathbf{x}; \mathbf{b}) = \left\{ \sum_{g=1}^G f(g, n | \mathbf{x}; \mathbf{b}) \right\}$.

Then,

$$M = M^B(\mathbf{g}, \mathbf{a}) + M^W(\mathbf{a}, \mathbf{b}) \tag{App.1}$$

where

$$M^B(\mathbf{g}, \mathbf{a}) = \int_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \left\{ \sum_{g=1}^G f(g | \mathbf{x}; \mathbf{a}) \log \left(\frac{f(g | \mathbf{x}; \mathbf{a})}{p_{g\bullet}} \right) \right\} d\mathbf{x}$$

and

$$M^W(\mathbf{a}, \mathbf{b}) = \int_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \left\{ \sum_{g=1}^G \sum_{n=1}^N f(g, n | \mathbf{x}; \mathbf{b}) \log \left(\frac{f(g, n | \mathbf{x}; \mathbf{b})}{f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})} \right) \right\} d\mathbf{x}.$$

Proof. From assumptions A4 and A5 we directly see that the first term in equation (App.1) equals $M^B(\mathbf{g}, \mathbf{a})$ since

$\frac{f_1(g, \mathbf{x})}{f_2(g, \mathbf{x})} = \frac{f(g | \mathbf{x}; \mathbf{a})}{p_{g\bullet}}$ and thus

$$\begin{aligned} \int_{\mathbf{x} \in \Lambda} \left\{ \sum_{g=1}^G f_1(g, \mathbf{x}) \log \left(\frac{f(g | \mathbf{x}; \mathbf{a})}{p_{g\bullet}} \right) \right\} d\mathbf{x} &= \int_{\mathbf{x} \in \Lambda} \left\{ \sum_{g=1}^G f(\mathbf{x}) f(g | \mathbf{x}; \mathbf{a}) \log \left(\frac{f(g | \mathbf{x}; \mathbf{a})}{p_{g\bullet}} \right) \right\} d\mathbf{x} \\ &= \int_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \left\{ \sum_{g=1}^G f(g | \mathbf{x}; \mathbf{a}) \log \left(\frac{f(g | \mathbf{x}; \mathbf{a})}{p_{g\bullet}} \right) \right\} d\mathbf{x}. \end{aligned}$$

From assumptions A5 and A6, $f_1(n | g, \mathbf{x}) = \frac{f(n, g | \mathbf{x}; \mathbf{b})}{f(g | \mathbf{x}; \mathbf{a})}$ and $f_1(g, n, \mathbf{x}) = f(g, n | \mathbf{x}; \mathbf{b}) f(\mathbf{x})$. Given assumption A7 we directly see that:

$$\begin{aligned} \int f_1(g, n, \mathbf{x}) \log \left(\frac{f_1(n | g, \mathbf{x})}{f_2(n | g, \mathbf{x})} \right) d\mathbf{m} &= \int_{\mathbf{x} \in \Lambda} \left\{ \sum_{g=1}^G \sum_{n=1}^N f(g, n | \mathbf{x}; \mathbf{b}) f(\mathbf{x}) \log \left(\frac{f(g, n | \mathbf{x}; \mathbf{b})}{f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})} \right) \right\} d\mathbf{x} \\ &= \int_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \left\{ \sum_{g=1}^G \sum_{n=1}^N f(g, n | \mathbf{x}; \mathbf{b}) \log \left(\frac{f(g, n | \mathbf{x}; \mathbf{b})}{f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})} \right) \right\} d\mathbf{x} \\ &= M^W(\mathbf{a}, \mathbf{b}). \end{aligned}$$

■

Proposition 3: Under assumptions A1-A3, $\text{plim } \hat{M}_T = M = M_{KL}$.

Proof. We first note that by direct application of Lemma 1 in Rao (1957), the sample frequencies $\hat{p}_{gn} = \frac{T_{gn}}{T}$ converge in probability to the actual probabilities, i.e. $\text{plim } \hat{p}_{gn} = p_{gn}$. Define the mutual information index as a function of the parameter vector:

$$m(\mathbf{q}) = \sum_{GN^c} p_{gn} \log \left(\frac{p_{gn}}{p_{g\bullet}(\mathbf{q}) p_{\bullet n}(\mathbf{q})} \right) + \left(1 - \sum_{GN^c} p_{gn} \right) \log \left(\frac{1 - \sum_{GN^c} p_{gn}}{p_{G\bullet}(\mathbf{q}) p_{\bullet N}(\mathbf{q})} \right)$$

where $\mathbf{q} \in \Theta$,

$$p_{g\bullet}(\mathbf{q}) = \begin{cases} \sum_{n=1}^N p_{gn} & \text{if } g \neq G \\ \sum_{n=1}^{N-1} p_{Gn} + \left(1 - \sum_{(j,n) \in GN^c} p_{jn} \right) & \text{if } g = G \end{cases}$$

and

$$p_{\bullet n}(\mathbf{q}) = \begin{cases} \sum_{g=1}^G p_{gn} & \text{if } n \neq N, \\ \sum_{g=1}^{G-1} p_{gN} + \left(1 - \sum_{(g,j) \in GN^c} p_{gj} \right) & \text{if } n = N. \end{cases}$$

Let \mathbf{q}^0 be the vector containing the true probabilities. Since $m(\mathbf{q})$ is continuous at $\mathbf{q} = \mathbf{q}^0$, by the Slutsky theorem it follows that $\text{plim } \hat{M}_T = \text{plim } m(\hat{\mathbf{q}}) = m(\mathbf{q}^0) = M$.

■

Proposition 4: Suppose that assumptions A1, A2, A3 hold. Further assume that

A8 $f(\mathbf{x} | g, n) = f(\mathbf{x} | g, n; \mathbf{j})$ such that $f(\mathbf{x}) = \sum_{g=1}^G \sum_{n=1}^N f(\mathbf{x} | g, n; \mathbf{j}) p_{gn}$ and $\mathbf{j} \in \mathbb{R}^{k_j}$ does not depend on \mathbf{q} .

Then: $\hat{M}_T = \frac{-\log(\mathbf{I})}{T}$.

Proof. From assumption A8 and the first order conditions from ML estimation we have $\hat{\mathbf{j}}_0 = \hat{\mathbf{j}}$ so that

$$L(\hat{\mathbf{q}}, \hat{\mathbf{j}}) = \left\{ \prod_{i=1}^T f(\mathbf{x}_i | g, n; \hat{\mathbf{j}}) \right\} \left\{ \prod_{g=1}^G \prod_{n=1}^N (\hat{p}_{gn})^{T \hat{p}_{gn}} \right\}$$

and

$$L(\hat{\mathbf{q}}_0, \hat{\mathbf{j}}_0) = \left\{ \prod_{i=1}^T f(\mathbf{x}_i | g, n; \hat{\mathbf{j}}) \right\} \left\{ \prod_{g=1}^G \prod_{n=1}^N (\hat{p}_{g\bullet} \hat{p}_{\bullet n})^{T \hat{p}_{gn}} \right\}.$$

The result of the proposition follows after some simple algebraic manipulations.

■

Theorem 1: Suppose that assumptions A1, A2, A3, and A8 hold. If $p_{gn} = p_{g\bullet} p_{\bullet n}$, for all $(g, n) \in GN^c$:

$$2T\hat{M}_T \xrightarrow{d} \mathbf{c}_{(G-1)(N-1)}^2.$$

Proof. We first note that the \mathbf{I} statistic is also the likelihood ratio for testing $H_0: p_{g|n} = p_{g\bullet}$, $g = 1, \dots, G-1$, $n = 1, \dots, N$, versus the two-sided alternative $H_1: p_{g|n} \neq p_{g\bullet}$ where the quantities $T_{\bullet n}$ are assumed to be constants. Direct application of Theorem 7 in Neyman (1949) implies that $-2\log(\mathbf{I}) \xrightarrow{d} \mathbf{c}_{df}^2$ where $df = (G-1)N - (G-1) = (G-1)(N-1)$. The proof then follows from Proposition 3.

■

Theorem 2: Suppose that assumptions A1, A2, A3, and A8 hold. If $p_{gn} \neq p_{g\bullet} p_{\bullet n}$ for at least one $(g, n) \in GN^c$ then:

$$T^{1/2}(\hat{M}_T - M) \xrightarrow{d} N(0, \Delta m' \Sigma \Delta m)$$

where

$$\Sigma = \left\{ \mathbf{s}_{(i,j)(g,n)} \right\} = \begin{cases} p_{ij}(1-p_{ij}) & \text{if } (i, j) = (g, n) \\ -p_{ij}p_{gn} & \text{if } (i, j) \neq (g, n), \end{cases}$$

$$\Delta m = \left\{ \Delta m_{gn} \right\} = \log \left(\frac{\hat{p}_{gn}}{\hat{p}_{g\bullet} \hat{p}_{\bullet n}} \right) - \log \left(\frac{\hat{p}_{GN}}{\hat{p}_{G\bullet} \hat{p}_{\bullet N}} \right), \quad \forall (g, n) \in GN^c.$$

Proof. We first note that $T^{1/2}(\hat{\mathbf{q}} - \mathbf{q}^0) \xrightarrow{d} N(0, \Sigma)$ (see, for example, Serfling, 1980, Theorem 2.7, p.109). To prove Theorem 2, we will use two lemmata:

Lemma 1: Suppose that $\hat{\mathbf{q}}$ is $AN(\mathbf{q}, T^{-1}\Sigma)$. If $m(\mathbf{q})$ has a non-zero partial derivative at $\mathbf{q} = \mathbf{q}^0$,

$$\Delta m \equiv \left(\left. \frac{\partial m(\mathbf{q})}{\partial \mathbf{q}_1} \right|_{\mathbf{q}=\mathbf{q}^0}, \dots, \left. \frac{\partial m(\mathbf{q})}{\partial \mathbf{q}_{NG-1}} \right|_{\mathbf{q}=\mathbf{q}^0} \right)' \neq 0_{NG-1}, \text{ then } m(\mathbf{q}) \text{ is } AN(m(\mathbf{q}^0), T^{-1}\Delta m' \Sigma \Delta m).$$

Proof. This is a direct from Theorem 3.3A in Serfling (1980).

■

Lemma 2: For $\mathbf{q}_i = \hat{p}_{gn}$:

$$\frac{\partial m(\mathbf{q})}{\partial \mathbf{q}_i} = \log\left(\frac{\hat{p}_{gn}}{\hat{p}_{g\bullet}\hat{p}_{\bullet n}}\right) - \log\left(\frac{\hat{p}_{GN}}{\hat{p}_{G\bullet}\hat{p}_{\bullet N}}\right).$$

Proof. Define $h_{gn} = \hat{p}_{gn} \log\left(\frac{\hat{p}_{gn}}{\hat{p}_{g\bullet}\hat{p}_{\bullet n}}\right)$, where $\hat{p}_{g\bullet} = \sum_{n=1}^N \hat{p}_{gn}$, $\hat{p}_{\bullet n} = \sum_{g=1}^G \hat{p}_{gn}$, and $\hat{p}_{GN} = 1 - \sum_{n=1}^{N-1} \sum_{g=1}^{G-1} \hat{p}_{gn} - \sum_{g=1}^{G-1} \hat{p}_{gN} - \sum_{n=1}^{N-1} \hat{p}_{GN}$. Then, for any race and school combination e and s the partial derivative of h_{es} with respect to \hat{p}_{gn} is $\frac{\partial h_{es}}{\partial \hat{p}_{gn}} = \left[1 + \log\left(\frac{\hat{p}_{es}}{\hat{p}_{e\bullet}\hat{p}_{\bullet s}}\right)\right] \frac{\partial \hat{p}_{es}}{\partial \hat{p}_{gn}} - \left(\frac{\hat{p}_{es}}{\hat{p}_{e\bullet}\hat{p}_{\bullet s}}\right) \frac{\partial (\hat{p}_{e\bullet}\hat{p}_{\bullet s})}{\partial \hat{p}_{gn}}$. The result follows after some algebraic manipulation after noticing that $\frac{\partial m}{\partial \hat{p}_{gn}} = \sum_{e=1}^G \sum_{s=1}^N \frac{\partial h_{es}}{\partial \hat{p}_{gn}}$. ■

In view of Lemmata 1 and 2, to prove Theorem 2 we only need to show that if $\hat{p}_{gn} \neq \hat{p}_{g\bullet}\hat{p}_{\bullet n}$ for at least one $(g, n) \in GN^c$, then $\Delta m \neq 0_{NG-1}$. Since $\hat{p}_{GN} = 1 - \sum_{g=1}^{G-1} \sum_{n=1}^{N-1} \hat{p}_{gn} - \sum_{n=1}^{N-1} \hat{p}_{gN} - \sum_{g=1}^{G-1} \hat{p}_{gN}$, if $\hat{p}_{gn} = \hat{p}_{g\bullet}\hat{p}_{\bullet n}$ for all $(g, n) \in GN^c$, then $\hat{p}_{GN} = \hat{p}_{G\bullet}\hat{p}_{\bullet N}$. Thus, if $\hat{p}_{gn} \neq \hat{p}_{g\bullet}\hat{p}_{\bullet n}$ for at least one $(g, n) \in GN^c$, then there must be at least another combination (\hat{g}, \hat{n}) such that $\hat{p}_{\hat{g}\hat{n}} \neq \hat{p}_{\hat{g}\bullet}\hat{p}_{\bullet\hat{n}}$. Assume, wlog, that $\hat{p}_{GN} \neq \hat{p}_{G\bullet}\hat{p}_{\bullet N}$. We can now prove by contradiction that $\Delta m \neq 0_{NG-1}$. Assume otherwise that

$$\log\left(\frac{\hat{p}_{gn}}{\hat{p}_{g\bullet}\hat{p}_{\bullet n}}\right) = \log\left(\frac{\hat{p}_{GN}}{\hat{p}_{G\bullet}\hat{p}_{\bullet N}}\right) \quad \forall (g, n) \in GN^c.$$

Then

$$\hat{p}_{gn} = \left(\frac{\hat{p}_{GN}}{\hat{p}_{G\bullet}\hat{p}_{\bullet N}}\right) \hat{p}_{g\bullet}\hat{p}_{\bullet n}$$

and summing over all $(g, n) \in GN^c$, $\sum_{GN^c} \hat{p}_{gn} = \left(\frac{\hat{p}_{GN}}{\hat{p}_{G\bullet}\hat{p}_{\bullet N}}\right) \sum_{GN^c} \hat{p}_{g\bullet}\hat{p}_{\bullet n}$. Given that $\sum_{GN^c} \hat{p}_{g\bullet}\hat{p}_{\bullet n} = 1 - \hat{p}_{G\bullet}\hat{p}_{\bullet N}$ always

then $1 - \hat{p}_{GN} = \left(\frac{\hat{p}_{GN}}{\hat{p}_{G\bullet}\hat{p}_{\bullet N}}\right) (1 - \hat{p}_{G\bullet}\hat{p}_{\bullet N})$, which contradicts the initial assumption $\hat{p}_{GN} \neq \hat{p}_{G\bullet}\hat{p}_{\bullet N}$. So it must be true that if $\hat{p}_{gn} \neq \hat{p}_{g\bullet}\hat{p}_{\bullet n}$ for at least one $(g, n) \in GN^c$ then $\Delta m \neq 0_{NG-1}$. ■

Proposition 5: Suppose that assumptions A1 to A8 hold and that covariates vector \mathbf{x} includes only district code d . Then:

$$\hat{M}_T^W = \frac{-\log(\mathbf{I}^W)}{T}.$$

Proof. From ML first order conditions we have that

$$L(\{\hat{p}_{gnd}\}, \{\hat{p}_{\bullet\bullet d}\}) = \prod_{d=1}^D (\hat{p}_{\bullet\bullet d})^{T\hat{p}_{\bullet\bullet d}} \left\{ \prod_{g=1}^G \prod_{n=1}^N (\hat{p}_{gnd})^{T\hat{p}_{gnd}} \right\}$$

and

$$L(\{\hat{p}_{gnd}^0\}, \{\hat{p}_{\bullet\bullet d}^0\}) = \prod_{d=1}^D (\hat{p}_{\bullet\bullet d})^{T_{\hat{p}_{\bullet\bullet d}}} \left\{ \prod_{g=1}^G \prod_{n=1}^N (\hat{p}_{g\bullet d} \hat{p}_{\bullet n d})^{T_{\bullet\bullet d} \hat{p}_{gnd}} \right\}$$

so that

$$\mathbf{I}^W = \prod_{d=1}^D \prod_{g=1}^G \prod_{n=1}^N \left(\frac{\hat{p}_{g\bullet d} \hat{p}_{\bullet n d}}{\hat{p}_{gnd}} \right)^{T_{\bullet\bullet d} \hat{p}_{gnd}}.$$

The result of the proposition follows after some simple algebraic manipulations. ■

Theorem 3: Suppose that assumptions A1 to A8 hold and that covariates vector \mathbf{x} includes only district code d .

(a) Assume further that there is at least one district d such that for at least two race and school combinations, (g, n) and (r, s) , we have that

$$\frac{1}{\hat{p}_{gnd}} + \frac{1}{\hat{p}_{rnd}} \neq \frac{1}{\hat{p}_{gnd}} + \frac{1}{\hat{p}_{rnd}} + \frac{1}{\hat{p}_{snd}} + \frac{1}{\hat{p}_{snd}}.$$

If $\hat{p}_{gnd} = \hat{p}_{g\bullet d} \hat{p}_{\bullet nd}$ for all (g, n, d) , then

$$TM_T^W \xrightarrow{d} ZAZ'$$

with $Z = N(0, \Sigma_W)$,

$$\Sigma_W = \left\{ \mathbf{s}_{(i,j,k) \in \{g,n,d\}} \right\} = \begin{cases} \hat{p}_{\bullet\bullet d} (1 - \hat{p}_{\bullet\bullet d}) & \text{if } (0,0,d) = (0,0,d) \\ -\hat{p}_{\bullet\bullet d} \hat{p}_{\bullet\bullet k} & \text{if } (0,0,d) \neq (0,0,k) \\ \hat{p}_{i\bullet d} (1 - \hat{p}_{i\bullet d}) & \text{if } (i,j,d) = (g,n,d) \\ -\hat{p}_{i\bullet d} \hat{p}_{gnd} & \text{if } (i,j,d) \neq (g,n,d), \end{cases}$$

and

$$A = \left(\frac{1}{2} \left(\frac{\partial^2 m_W(\mathbf{q}_W)}{\partial \mathbf{q}_{Wi} \partial \mathbf{q}_{Wj}} \right) \Big|_{\mathbf{q}_W = \mathbf{q}_W^0} \right)_{(GN-1) \times (GN-1)}.$$

(b) If $\hat{p}_{gnd} \neq \hat{p}_{g\bullet d} \hat{p}_{\bullet nd}$ for at least one (g, n, d) , then

$$T^{1/2} (\hat{M}_T^W - M^W) \xrightarrow{d} N(0, \Delta m^W, \Sigma_W \Delta m^W)$$

where

$$\Delta m_W = \left\{ \frac{\partial m_W}{\partial \mathbf{q}_{Wi}} \right\} = \begin{cases} m_d(\mathbf{q}_d) - m_D(\mathbf{q}_D) & \text{if } \mathbf{q}_{Wi} = \hat{p}_{\bullet\bullet d} \\ \log \left(\frac{\hat{p}_{gnd}}{\hat{p}_{g\bullet d} \hat{p}_{\bullet nd}} \right) - \log \left(\frac{\hat{p}_{GNd}}{\hat{p}_{G\bullet d} \hat{p}_{\bullet Nd}} \right) & \text{if } \mathbf{q}_{Wi} = \hat{p}_{gnd}. \end{cases}$$

Proof. We first prove part (a) and then prove part (b). In both cases, we exploit the fact that the within term M^W can be expressed as a well-behaved function of the sample relative frequencies. Since these are the ML estimators for the actual probabilities, we have (see, for example, Serfling, 1980, pp.109) that $T^{1/2} (\hat{\mathbf{q}}_W - \mathbf{q}_W^0) \xrightarrow{d} N(0, \Sigma_W)$ where \mathbf{q}_W^0 is the vector with the actual probabilities. To prove Theorem 3a, we will use Lemma 2 and the following corollary from Theorem 3.3B in Serfling (1980):

Lemma 3: Suppose that $\hat{\mathbf{m}} \in \mathbb{R}^K$ is $\mathcal{AN}(\mathbf{m}^0, T^{-1} \Sigma_m)$. If $g(\mathbf{m})$ is a real-valued function possessing continuous partial derivatives of second order in a neighbourhood of $\mathbf{m} = \mathbf{m}^0$, with the first order partial derivatives vanishing at $\mathbf{m} = \mathbf{m}^0$, but with the second order partial derivatives not all vanishing at $\mathbf{m} = \mathbf{m}^0$. Then

$$T\left(g(\hat{\mathbf{m}}_T) - g(\mathbf{m}^0)\right) \xrightarrow{d} ZAZ' \text{ with } Z = (Z_1, \dots, Z_K) \sim N(0, \Sigma_{\mathbf{m}}) \text{ and } A = \left(\frac{1}{2} \left(\frac{\partial^2 g(\mathbf{m})}{\partial \mathbf{m}_i \partial \mathbf{m}_j} \Big|_{\mathbf{m}=\mathbf{m}^0} \right)_{K \times K}\right).$$

In view of Lemma 3, to prove Theorem 3a we only need to show that: (A) $m_W(\mathbf{q}_W^0) = 0$; (B)

$$\frac{\partial m_W(\mathbf{q}_W)}{\partial \mathbf{q}_{W_i}} \Big|_{\mathbf{q}_W = \mathbf{q}_W^0} = 0 \text{ for all } i = 1, \dots, GN - 1; \text{ (C) } \frac{\partial^2 m_W(\mathbf{q}_W)}{\partial \mathbf{q}_{W_i} \partial \mathbf{q}_{W_j}}$$

(D) at least one second partial derivative $\frac{\partial^2 m_W(\mathbf{q}_W)}{\partial \mathbf{q}_{W_i} \partial \mathbf{q}_{W_j}}$ does not vanish at $\mathbf{q}_W = \mathbf{q}_W^0$. Condition (A) follows

immediately from the definition of $m_W(\mathbf{q}_W)$ and the assumption that $\dot{p}_{gnd} = \dot{p}_{g \bullet | d} \dot{p}_{\bullet | nd}$ for all (g, n, d) . To show that condition (B) holds, first note that Lemma 2 can be used to show that the partial derivative with respect to \dot{p}_{gnd} is zero when $\dot{p}_{gnd} = \dot{p}_{g \bullet | d} \dot{p}_{\bullet | nd}$ for all (g, n, d) . In addition, since

$$\frac{\partial m_W(\mathbf{q}_W)}{\partial \dot{p}_{\bullet \bullet | d}} = m_d(\mathbf{q}_d) - m_D(\mathbf{q}_D) \text{ for all } d = 1, \dots, D - 1, \text{ then these derivatives also vanish when } \dot{p}_{gnd} = \dot{p}_{g \bullet | d} \dot{p}_{\bullet | nd}$$

since then $m_d(\mathbf{q}_d) = m_D(\mathbf{q}_D) = 0$. It is then straightforward to show that the second order partial derivatives $\frac{\partial^2 m_W(\mathbf{q}_W)}{\partial \mathbf{q}_{W_i} \partial \mathbf{q}_{W_j}}$ are continuous in a neighbourhood of $\mathbf{q}_W = \mathbf{q}_W^0$. Condition (D) follows directly given that for any g and

$$n \text{ in any district } d \text{ we have that } \frac{\partial^2 m_W(\mathbf{q}_W)}{\partial \dot{p}_{gnd}^2} = \left(\frac{1}{\dot{p}_{gnd}} + \frac{1}{\dot{p}_{GN_d|d}} \right) - \left(\frac{1}{\dot{p}_{g \bullet | d}} + \frac{1}{\dot{p}_{G \bullet | d}} + \frac{1}{\dot{p}_{\bullet | nd}} + \frac{1}{\dot{p}_{\bullet | N_d|d}} \right). \text{ Therefore, a}$$

sufficient condition for (D) to hold is that we choose $G = r$ and $N_d = s$.

To prove Theorem 3b, we first note that by Theorem 3.3.A in Serfling (1980), Lemma 2, and the fact that

$$\Delta m_W = \left\{ \frac{\partial m_W}{\partial \mathbf{q}_{W_i}} \right\} = \begin{cases} m_d(\mathbf{q}_d) - m_D(\mathbf{q}_D) & \text{if } \mathbf{q}_{W_i} = \dot{p}_{\bullet \bullet | d} \\ \log \left(\frac{\dot{p}_{gnd}}{\dot{p}_{g \bullet | d} \dot{p}_{\bullet | nd}} \right) - \log \left(\frac{\dot{p}_{GN_d|d}}{\dot{p}_{G \bullet | d} \dot{p}_{\bullet | N_d|d}} \right) & \text{if } \mathbf{q}_{W_i} = \dot{p}_{gnd}, \end{cases}$$

the result follows if it is shown that $m_W(\mathbf{q}_W)$ has a non-zero partial derivative at $\mathbf{q}_W = \mathbf{q}_W^0$. Now, using a similar argument to the argument used in the proof of Theorem 2, given that

$$\dot{p}_{gnd} \neq \dot{p}_{g \bullet | d} \dot{p}_{\bullet | nd} \text{ for at least one } (g, n, d), \text{ then } \frac{\partial m_d(\mathbf{q}_d)}{\partial \mathbf{q}_d} \Big|_{\mathbf{q}_d = \mathbf{q}_d^0} \neq 0_d. \text{ Since}$$

$$\frac{\partial m_W(\mathbf{q}_W)}{\partial \dot{p}_{gnd}} = \dot{p}_{\bullet \bullet | d} \frac{\partial m_d(\mathbf{q}_d)}{\partial \dot{p}_{gnd}}$$

$$\text{and } \dot{p}_{\bullet \bullet | d} > 0 \text{ for all } d, \text{ then } \frac{\partial m_W(\mathbf{q}_W)}{\partial \mathbf{q}_W} \Big|_{\mathbf{q}_W = \mathbf{q}_W^0} \neq \mathbf{0}_W.$$

■

Theorem 4: Suppose that assumptions A1 to A5 hold and that the vector of covariates \mathbf{x} includes at least one countable or continuous variable. Let $(\mathbf{g}_0, \mathbf{a}_0)$ be the true parameter vectors of the data

generating process and define $b_B(\mathbf{x}; \mathbf{g}, \mathbf{a}) = \sum_{g=1}^G f(g | \mathbf{x}; \mathbf{a}) \log \left(\frac{f(g | \mathbf{x}; \mathbf{a})}{p_{g^*}} \right)$. Assume that:

(R1) $b_B(\mathbf{x}; \mathbf{g}, \mathbf{a})$ is differentiable with respect to (\mathbf{g}, \mathbf{a}) , with continuous partial derivatives which are nonvanishing at $(\mathbf{g}_0, \mathbf{a}_0)$.

(R2) $\text{Var}_x [b_B(\mathbf{x}; \mathbf{g}_0, \mathbf{a}_0)] = \mathbf{S}_B^2 < \infty$;

$$E_x \left[\frac{\partial b_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \mathbf{m}_g \in \mathbb{R}^{k_g}, E_x \left[\frac{\partial b_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \mathbf{m}_a \in \mathbb{R}^{k_a}, \text{Var}_x \left[\frac{\partial b_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \Sigma_g, \text{Var}_x \left[\frac{\partial b_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}_0, \mathbf{a}_0} \right] = \Sigma_a,$$

where Σ_g and Σ_a are positive definite matrices.

(R3) $\text{plim } \hat{\mathbf{g}} = \mathbf{g}_0$ and $\text{plim } \hat{\mathbf{a}} = \mathbf{a}_0$.

Then,

$$T^{1/2} \left(\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) - M^B(\mathbf{g}_0, \mathbf{a}_0) \right) \xrightarrow{d} N(0, \mathbf{S}_B^2).$$

Proof. We first note that by the Lindeberg-Levy central limit theorem,

$$T^{1/2} \left(\tilde{M}_T^B(\mathbf{g}_0, \mathbf{a}_0) - M^B(\mathbf{g}_0, \mathbf{a}_0) \right) \xrightarrow{d} N(0, \mathbf{S}_B^2)$$

where $\tilde{M}_T^B(\mathbf{g}_0, \mathbf{a}_0) = T^{-1} \sum_{i=1}^T b_B(\mathbf{x}_i; \mathbf{g}_0, \mathbf{a}_0)$. By the mean value theorem, we have that

$$\hat{M}_T^B(\hat{\mathbf{g}}, \hat{\mathbf{a}}) = \tilde{M}_T^B(\mathbf{g}_0, \mathbf{a}_0) + T^{-1} \sum_{i=1}^T \left\{ \frac{\partial b_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}^*, \mathbf{a}^*} (\hat{\mathbf{g}} - \mathbf{g}_0) + \frac{\partial b_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}^*, \mathbf{a}^*} (\hat{\mathbf{a}} - \mathbf{a}_0) \right\}$$

satisfying

$$\|\mathbf{g}^* - \mathbf{g}_0\| \leq \|\hat{\mathbf{g}} - \mathbf{g}_0\|$$

$$\|\mathbf{a}^* - \mathbf{a}_0\| \leq \|\hat{\mathbf{a}} - \mathbf{a}_0\|.$$

Therefore, to prove the theorem we only need to prove that

$$T^{-1} \sum_{i=1}^T \left\{ \frac{\partial b_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}^*, \mathbf{a}^*} (\hat{\mathbf{g}} - \mathbf{g}_0) + \frac{\partial b_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}^*, \mathbf{a}^*} (\hat{\mathbf{a}} - \mathbf{a}_0) \right\} = o_p(1).$$

Given (R3), this condition is satisfied if $T^{-1} \sum_{i=1}^T \frac{\partial b_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}^*, \mathbf{a}^*}$ and $T^{-1} \sum_{i=1}^T \frac{\partial b_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}^*, \mathbf{a}^*}$ are asymptotically normal, which

follows by the squeeze theorem from asymptotic normality for $T^{-1} \sum_{i=1}^T \frac{\partial b_B}{\partial \mathbf{g}} \Big|_{\hat{\mathbf{g}}, \hat{\mathbf{a}}}$ and $T^{-1} \sum_{i=1}^T \frac{\partial b_B}{\partial \mathbf{a}} \Big|_{\hat{\mathbf{g}}, \hat{\mathbf{a}}}$, (this is due to R3

and continuity in the partial derivatives) and asymptotic normality for $T^{-1} \sum_{i=1}^T \frac{\partial b_B}{\partial \mathbf{g}} \Big|_{\mathbf{g}_0, \mathbf{a}_0}$ and $T^{-1} \sum_{i=1}^T \frac{\partial b_B}{\partial \mathbf{a}} \Big|_{\mathbf{g}_0, \mathbf{a}_0}$ (due to the central limit theorem).

■

Theorem 5: Suppose that assumptions A1 to A7 hold, and that the vector of covariates \mathbf{x} includes at least one countable or continuous variable. Let $(\mathbf{a}_0, \mathbf{b}_0)$ be the true parameter vectors of the

data generating process, and define $h_W(\mathbf{x}; \mathbf{a}, \mathbf{b}) = \sum_{g=1}^G \sum_{n=1}^N f(g, n | \mathbf{x}; \mathbf{b}) \log \left(\frac{f(g, n | \mathbf{x}; \mathbf{b})}{f(g | \mathbf{x}; \mathbf{a}) f(n | \mathbf{x}; \mathbf{b})} \right)$

Assume that:

(R4) $h_W(\mathbf{x}; \mathbf{a}, \mathbf{b})$ is differentiable with respect to (\mathbf{a}, \mathbf{b}) , with continuous partial derivatives which are nonvanishing at $(\mathbf{a}_0, \mathbf{b}_0)$.

(R5) $\text{Var}_x [h_W(\mathbf{x}; \mathbf{a}_0, \mathbf{b}_0)] = \mathbf{s}_W^2 < \infty$,

$\text{E}_x \left[\frac{\partial h_W}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \mathbf{m}_a \in \mathbb{R}^{k_a}$, $\text{E}_x \left[\frac{\partial h_W}{\partial \mathbf{b}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \mathbf{m}_b \in \mathbb{R}^{k_b}$, $\text{Var}_x \left[\frac{\partial h_W}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \Sigma_a$, and $\text{Var}_x \left[\frac{\partial h_W}{\partial \mathbf{b}} \Big|_{\mathbf{a}_0, \mathbf{b}_0} \right] = \Sigma_b$,

where Σ_a and Σ_b are positive definite matrices.

(R6) $\text{plim } \hat{\mathbf{a}} = \mathbf{a}_0$ and $\text{plim } \hat{\mathbf{b}} = \mathbf{b}_0$.

Then,

$$T^{1/2} \left(\hat{M}_T^W(\hat{\mathbf{a}}, \hat{\mathbf{b}}) - M^W(\mathbf{a}_0, \mathbf{b}_0) \right) \xrightarrow{d} N(0, \mathbf{s}_W^2).$$

Proof. It is formally similar to the proof of Theorem 4. ■

REFERENCES

- Anker, R. (1998), *Gender and jobs: Sex segregation of occupations in the world*. Geneva, ILO.
- Aslund, O., and O.N. Skans (2009), "How to measure segregation conditional on the distribution of covariates," *Journal of Population Economics* 22(4):971-981.
- Blackburn, R.M., J. Jarman, and J. Siltanen (1993), "The Analysis of Occupational Gender Segregation over Time and Place: Considerations of Measurement and Some New Evidence," *Work, Employment and Society* 7: 335-36.
- Blackburn, R.M., J. Jarman, and J. Siltanen (1995), "The Measurement of Occupational Gender Segregation: Current Problems and a New Approach," *Journal of the Royal Statistical Society A*, Part 2 **158**: 319-331.
- Boisso, D., K. Hayes, J. Hirschberg, and J. Silber (1994), "Occupational Segregation in the Multidimensional Case: Decomposition and Test of Significance," *Journal of Econometrics* **61**: 161-171.
- Carrington, W.J., and K.R. Troske (1997), "On Measuring Segregation in Samples with Small Units," *Journal of Business and Economic Statistics*, **15**: 402-409.
- Chakravarty, S.R., and J. Silber (1992), "Employment Segregation Indices: An Axiomatic Characterization" In Eichhorn, W. (ed), *Models and Measurement of Welfare and Inequality*, New York: Springer-Verlag.
- Chakravarty, S.R., and J. Silber (2007), "A generalized index of employment segregation," *Mathematical Social Sciences* **53(2)**:185-195.
- Charles, M. (1992), "Cross-National Variation in Occupational Sex Segregation," *American Sociological Review* **57**: 483-502.
- Charles, M. (1998), "Structure, Culture, and Sex Segregation in Europe," *Research in Social Stratification and Mobility* **16**: 89-116.
- Charles, M. and D. Grusky (1995), "Models for Describing the Underlying Structure of Sex Segregation," *American Journal of Sociology* **100**: 931-971.
- Cortese, C.F., R.F. Falk, and J.K. Cohen (1976) "Further Considerations on the Methodological Analysis of Segregation Indices," *American Sociological Review* **41**:630-637.
- Davison, A.C. and D.V. Hinkley (1997), *Bootstrap Methods and their Applications*, Cambridge University Press, Cambridge.
- Deutsch, J, Flückiger, Y. and Silber, J. (1994), "Measuring Occupational Segregation," *Journal of Econometrics* **61**: 133-146.
- Duncan, O. and Duncan, B. (1955), "A Methodological Analysis of Segregation Indices," *American Sociological Review* **20**: 210-217.

- Frankel, D and O. Volij (2008), "Scale Invariant Measures of Segregation," mimeo, August 2008.
- Frankel, D and O. Volij (2009), "Measuring School Segregation," mimeo, March 2009.
- Flückiger, Y. and Silber, J. (1999), *The Measurement of Segregation in the Labor Force*, Heidelberg, Physica-Verlag.
- Hellerstein, J.K., and D. Neumark (2008), "Workplace Segregation in the United States: Race, Ethnicity, and Skill," *Review of Economics and Statistics* **90(3)**: 459-477.
- Herranz, N., Mora, R., and Ruiz-Castillo, J. (2005), "An Algorithm to Reduce the Occupational Space in Gender Segregation Studies," *Journal of Applied Econometrics* **20**: 25-37
- Hutchens, R. M. (1991), "Segregation Curves, Lorenz Curves and Inequality in the Distribution of People Across Occupations," *Mathematical Social Sciences* **21**: 31-51.
- Hutchens, R. M. (2001), "Numerical Measures of Segregation: Desirable Properties and Their Implications," *Mathematical Social Sciences* **42**: 13-29.
- Hutchens, R. M. (2004), "One Measure of Segregation," *International Economic Review* **45**: 555-578.
- James, D.R. and K.E. Taeuber (1985), "Measures of Segregation," in G. Schmid and R. Weitzel (eds.), *Sex Discrimination and Equal Opportunity: The Labor Market and Employment Policy*, London, Gower Publishing Company.
- Jonung, C. (1984), "Patterns of Occupational Segregation by Sex in the Labor Market," in N.B. Tuma (ed.), *Sociological Methodology*, San Francisco, Jossey-Bass.
- Kakwani, N.C. (1994), "Segregation by Sex: Measurement and Hypothesis Testing," *Research on Economic Inequality* **5**: 1-26.
- Kalter, F. (2000), "Measuring Segregation and Controlling for Independent Variables," Working Paper 19-2000, Manheimer Zentrum für Europäische Sozialforschung.
- Karmel, T. and M. MacLachlan (1988), "Occupational Sex Segregation: Increasing or Decreasing?" *Economic Record* **64**:187-195.
- Kullback, S. (1959), *Information Theory and Statistics*, John Wiley and Sons, NY.
- Kullback, S., and R.A. Leibler (1951), "On Information and Sufficiency," *The Annals of Mathematical Statistics* **22(1)**:79-86.
- Kupperman, M. (1957), "Further applications of information theory to multivariate analysis and statistical inference," *Dissertation*, Graduate Council of George Washington Univ.
- Massey, D. and N. Denton (1988), "The Dimensions of Residential Segregation," *Social Forces* **67**: 281-315.

- Mora, R. and Ruiz-Castillo, J. (2003), "Additively Decomposable Segregation Indices. The Case of Gender Segregation By Occupations and Human Capital Levels In Spain," *Journal of Economic Inequality* **1**: 147-179.
- Mora, R. and Ruiz-Castillo, J. (2004), "Gender Segregation by Occupations in the Public and the Private Sectors. The Case of Spain In 1977 and 1992," *Investigaciones Económicas* **XXVIII**: 399-428.
- Mora, R. and Ruiz-Castillo, J. (2005), "Axiomatic Properties of an Entropy Based Index of Segregation," Working Paper 05-62, Economics Series 31, Universidad Carlos III.
- Mora, R. and Ruiz-Castillo, J. (2009), "Entropy-based Segregation Indices," Working Paper 09-32, Economics Series 18, Universidad Carlos III.
- Morales, D., L. Pardo, I. Vajda (1995), "Asymptotic divergence of estimates of discrete distributions," *Journal of Statistical Planning and Inference* **48**(3): 347-369.
- Neyman, J. (1949), "Contribution to the theory of the χ^2 test," in J. Neyman (ed.) *Proceedings of the First Berkeley Symposium on Mathematical Statistics and Probability*, Univ. of California Press, Berkeley, 547-566.
- Philipson, T. (1993), "Social Welfare and Measurement of Segregation," *Journal of Economic Theory* **60**(2):322-334.
- Rao, C.R. (1957), "Maximum Likelihood Estimation for the Multinomial Distribution," *The Indian Journal of Statistics (1933-1960)* **18**:139-148.
- Ransom, M.R. (2000), "Sampling Distributions of Segregation Indexes," *Sociological Methods & Research* **28**: 454-475.
- Reardon, S., J. Yun and T. McNulty, (2000), "The Changing Structure of School Segregation: Measurement and Evidence of Multi-racial Metropolitan Area School Segregation, 1989-1999," *Demography* **37**: 351-364.
- Reardon, S. and G. Firebaugh, (2002), "Measures of Multigroup Segregation," *Sociological Methodology* **32**: 33-67.
- Romano, J.P., A.M. Shaikh, and M. Wolf (2008), "Formalized Data Snooping Based On Generalized Error Rates," *Econometric Theory* **24**(02):404-447.
- Salicrú, M., D. Morales, M.L. Menéndez and L. Pardo (1994), "On the applications of divergence type measures in testing statistical hypotheses," *Journal of Multivariate Analysis* **51**: 372-391.
- Serfling, R.J. (1980) *Approximation Theorems of Mathematical Statistics*, New York: John Wiley & Sons.
- Shao, J. (1998), *Mathematical Statistics*, Springer-Verlag New York, Inc.
- Siltanen, J. (1990), "Social Change and the Measurement of Occupational Segregation by Sex: An Assessment of the Sex-Ratio Index," *Work, Employment and Society* **4**: 1-29.

- Spriggs, W.E., and R.M. Williams (1996), "Logit Decomposition Analysis of Occupational Segregation: Results for the 1970s and 1980s," *The Review of Economics and Statistics* **78**: 348-355.
- Theil, H. and A.J. Finizza (1971), "A Note on the Measurement of Racial Integration of Schools by Means of Information Concepts," *Journal of Mathematical Sociology* **1**: 187-194.
- Theil, H. (1972), *Statistical Decomposition Analysis*, Amsterdam: North Holland.
- Watts, M. (1992), "How Should Occupational Segregation Be Measured?" *Work, Employment and Society* **6**: 475-487.
- Watts, M. (1994), "A Critique of Marginal Matching," *Work, Employment and Society* **8**: 421-431.
- Watts, M. (1997a), "Multidimensional Indices of Occupational Segregation," *Evaluation Review* **21**: 461-482.
- Watts, M. (1997b), "The Measurement of Occupational Gender Segregation," *Journal of the Royal Statistical Society A* **160**: 141-145.
- Watts, M. (1998a), "Occupational Gender Segregation: Index Measurement and Econometric Modelling," *Demography* **35**: 489-496.
- Watts, M. (1998b), "The Analysis of Sex Segregation: When is Index Measurement Not Index Measurement?" *Demography* **35**: 505-508.

Table 1. Urban Public School Enrolment, Racial Mix, and School Segregation in the U.S., 1989:2005

	No. of students (millions)			Racial Shares (%)		
	1989	2005	Change (%)	1989	2005	Change (%)
Native American	0.18	0.31	70.36	0.68	0.89	0.20
Asian	1.12	2.03	82.14	4.15	5.49	1.34
Black	4.55	6.94	52.63	16.10	17.80	1.70
Hispanic	3.75	8.39	123.95	13.85	23.87	10.02
Non-white	9.59	17.68	84.26	34.78	48.05	13.27
White	16.98	18.47	8.77	65.22	51.95	-13.27
Total	26.57	36.14	36.03	100.00	100.00	0.00
	Mutual Information index of Segregation					
	1989		2005		Change (%)	
M	43.92		48.90		11.33%	

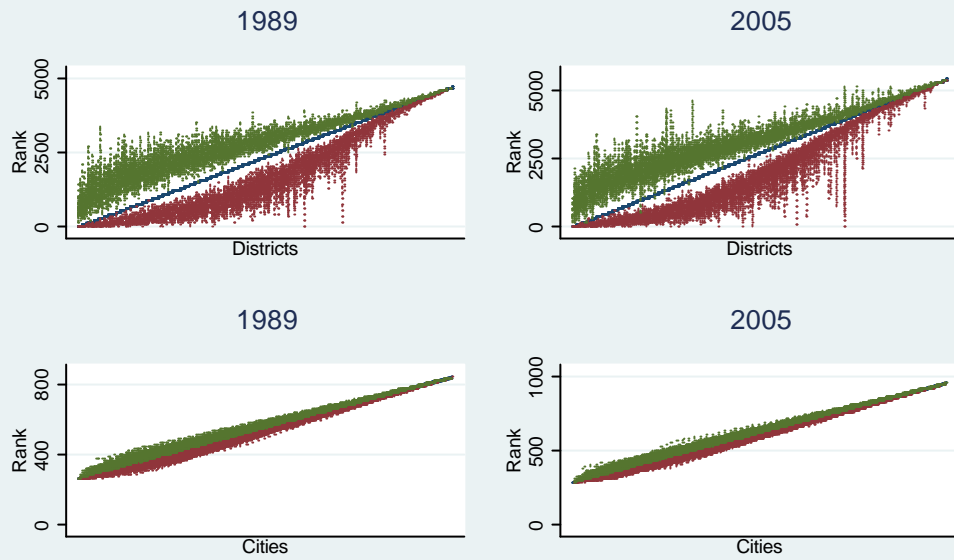
Note: Ethnic shares are the percentages of students from every race/ethnic group. The terms Native American, Asian, Black, and White refer to non-Hispanic members of these racial groups. Asian includes Native Hawaiians and Pacific Islanders; Native American includes American Indians and Alaska Natives (Innuut or Aleut). The term Hispanic is an ethnic rather than a racial category since Hispanic persons may belong to any race. Total Non-white includes all categories except White.

Table 2. Between Cities, Within Cities, and Within Districts Segregation in the U.S.

	1989			2005		
	Index	Lower Bound	Upper Bound	Index	Lower Bound	Upper Bound
Between Cities	21.04	21.02	21.06	23.50	23.48	23.52
(% over total)	47.91			48.06		
Within Cities	14.71	14.69	14.72	16.40	16.38	16.42
(% over total)	33.49			33.54		
Within Districts	8.18	8.16	8.19	9.00	8.99	9.01
(% over total)	18.62			18.40		
Total	43.92	43.90	43.95	48.90	48.87	48.92

Note: Ethnic shares are the percentages of students from every race/ethnic group. The terms Native American, Asian, Black, and White refer to non-Hispanic members of these racial groups. Asian includes Native Hawaiians and Pacific Islanders; Native American includes American Indians and Alaska Natives (InnuIt or Aleut). The term Hispanic is an ethnic rather than a racial category since Hispanic persons may belong to any race. Total Non-white includes all categories except White.

Rankings based on the Mutual Information Index Levels and Upper and Lower Bounds



Notes: 5% and 95% bounds obtained from 250 bootstrap replications

Table 3. Multigroup Conditional School Segregation: The Role of Income, Wages, and Teachers per Pupil, 2005

Fixed-Effects Logistics Regressions: Marginal Effects ^a					
District Level Data (7419 districts, 923 cities)					
	Native American	Asian	Black	Hispanic	White
Per-capita personal income	0.00	0.04***	-0.32***	-0.09***	0.83***
Wages per job	-0.01***	-0.003	0.62***	0.14***	-1.48***
Teachers per pupil	-0.02***	-0.08***	-0.15***	-0.15***	0.10
School Level Data (53174 schools, 7419 districts)					
	Native American	Asian	Black	Hispanic	White
Per-capita personal income	0.003	0.01	-0.05	-0.14***	0.18***
Wages per job	-0.004	0.03*	0.14***	0.19***	-0.35***
Teachers per pupil	-0.01***	-0.06***	0.05***	-0.01	-0.14***
The Decomposition of Within-Cities and Within-Districts Segregation ^b					
	Within Cities		Within Districts		
	Index	% over Total	Index	% over Total	
Total	16.72	100.00	9.07	100.00	
All Controls	10.67 (0.003)	63.78	1.93 (0.001)	21.25	
Income and Wages	10.69 (0.003)	63.92	1.93 (0.001)	21.32	
Conditional Segregation	6.06	36.22	7.14	78.75	

Notes:

^a District level regressions include city fixed effects. School level regressions include district fixed effects. Marginal effects are sample averages of the estimated partial derivative of each of the controls over the probability of belonging to each of the races and can be interpreted as the expected change in probability (in percentage terms) brought about by a one-percentage increase in the control. ***, **, and * denote parameter significant at 1, 5, and 10 significance level.

^b Total-Within Cities Index measures the expected information of the message that transforms the set of city ethnic shares to the set of district ethnic shares for the regressions sample. All controls-Within Cities captures segregation stemming from the statistical association between per-capita personal income, wages per job, and teachers per pupil and race membership at district level. Income and Wages-Within Cities simulates the value of segregation stemming from the statistical association between income and wages as if teachers per pupil played no role. Conditional Segregation-Within Cities reports the difference between Total-Within Cities and All Controls-Within Cities. Total-Within District Index measures the expected information of the message that transforms the set of district ethnic shares to the set of school ethnic shares for the regressions sample. All controls-Within Districts captures segregation stemming from the statistical association between per-capita personal income, hourly wages, and teachers per pupil and race membership at school level. Income and Wages-Within Districts simulates the value of segregation stemming from the statistical association between income and wages as if teachers per pupil played no role at school level. Conditional Segregation-Within Districts reports the difference between Total-Within Districts and All Controls-Within Districts. Asymptotic Standard Errors are shown in parenthesis.