

Sociological Methods & Research

<http://smr.sagepub.com>

A Reply to Zax's (2002) Critique of Grofman and Migalski (1988): Double-Equation Approaches to Ecological Inference When the Independent Variable Is Misspecified

Bernard Grofman and Matt A. Barreto
Sociological Methods Research 2009; 37; 599
DOI: 10.1177/0049124109334794

The online version of this article can be found at:
<http://smr.sagepub.com/cgi/content/abstract/37/4/599>

Published by:



<http://www.sagepublications.com>

Additional services and information for *Sociological Methods & Research* can be found at:

Email Alerts: <http://smr.sagepub.com/cgi/alerts>

Subscriptions: <http://smr.sagepub.com/subscriptions>

Reprints: <http://www.sagepub.com/journalsReprints.nav>

Permissions: <http://www.sagepub.com/journalsPermissions.nav>

Citations <http://smr.sagepub.com/cgi/content/refs/37/4/599>

A Reply to Zax's (2002) Critique of Grofman and Migalski (1988)

Double-Equation Approaches to Ecological Inference When the Independent Variable Is Misspecified

Bernard Grofman

University of California, Irvine

Matt A. Barreto

University of Washington

The authors reply to Zax's critique of the double-equation method for ecological regression and of the specific extension to it proposed by Grofman and Migalski. Although Zax does correct two minor errors in Grofman and Migalski's statement of the double-equation approach, neither of those errors affected the final calculations reported in their article. Furthermore, nothing Zax reports affects their fundamental conclusion that double-equation methods can be superior to single-equation techniques if there is substantial error in the measurement of the independent variable. In particular, by analyzing an election for which, from exit polls, the "true" parameters of Hispanic and non-Hispanic levels of political cohesion are known, the authors show that double-equation ecological regression estimates derived from registration data are highly accurate in reproducing the true individual-level behavioral parameters (group means).

Keywords: *Elections; Voter Racial Turnout; Bloc Voting; Ecological Regression; Ecological Inference*

In looking at voting behavior in local elections, it is only very rarely the case that survey data on electoral behavior broken down by race are available. As a consequence, expert witnesses in litigation involving voting rights, for which information about the levels of racial bloc voting (RBV) in the electorate is almost always a key factor in determining trial outcomes (Grofman, Handley, and Niemi 1992), must rely on ecological methods to estimate

racial voting behavior. But a further complication is that data on the racial composition of the electorate at the level of voting tabulation units (precincts) are also rarely available. Thus, expert witnesses will normally have to use minority and nonminority shares of voting-age population (VAP) or of registration as proxies for the minority and nonminority composition of the actual voting electorate. But when they do so, they are obviously inputting error in the nature of the independent variable. Furthermore, it can be shown that the error in estimates so generated is of a nonlinear form.

It has been a topic of concern for several decades how to adjust ecological techniques to compensate for such systematic measurement errors in the independent variable. In an extension of standard Goodman-type ecological regression methods (Goodman 1953, 1959), various scholars (Grofman, Migalski, and Noviello 1985; Kousser 1973; Loewen and Grofman 1989) have proposed a double-equation method (commonly known as “double regression”) intended to cope with the problem. In a further extension, Grofman and Migalski (1988) showed how this two-equation regression method can be made applicable to situations in which voters may cast more than a single vote (e.g., situations in which the voting method is plurality bloc voting).

The availability of a computer program to implement King’s (1997) ecological inference methods, which represent a major improvement on Goodman’s approach, has sparked a renaissance of ecological studies, especially in political science. King’s methods of ecological inference are increasingly used by expert witnesses in voting rights cases in estimating RBV patterns.¹ An ancillary consequence of this renewed interest in ecological methods for generating inferences about individual-level patterns has been new work considering how to adjust for measurement errors in the independent variable (Cho, Judge, and Cain 2002; Zax 2002, 2005; see also Cho and Gaines 2004; King 1997:71-72). Our focus here is primarily on one of these articles, that by Zax (2002). Zax reviewed the double-equation ecological regression approach, and more

Authors’ Note: Participation in this research project was aided by funding from the Center for the Study of Democracy at the University of California, Irvine, and the University of Washington Center for Statistics in the Social Sciences. We are indebted to helpful conversations over the years with Michael Migalski, Kenneth Small, Samuel Merrill, Nicholas Noviello, Alan Lichtman, James Loewen, Gary King, Chris Adolph, Gary Segura, Anthony Salvanto, John DiNardo, and Nathan Woods about the topics discussed in this article. Errors remaining are solely the responsibility of the authors. Correspondence concerning this article should be addressed to Bernard Grofman, University of California, Irvine, Institute for Mathematical Behavioral Sciences, 2291 Social Sciences Plaza B, Irvine, CA 92697; e-mail: bgrofman@uci.edu.

specifically, he critiqued Grofman and Migalski's (1988) proposed extension of the double-equation ecological regression approach to the case in which voters may cast more than a single vote.

Zax (2002) made three major claims about the double-equation ecological regression technique as expounded in Grofman and Migalski (1988). First, he regarded the double-regression technique as unnecessary, because King's (1997) method of ecological inference is now available. Second, he argued that the attempt of Grofman and Migalski to expand the scope of applicability of the double-regression technique to the multiseat case (plurality bloc voting), in which voters have multiple votes to cast, fails to go through because of the existence of a critical error in Grofman and Migalski's equation 5, an error that is perpetuated in most of the subsequent equations in that article. Third, Zax asserted that even for the basic case in which voters have a single vote to cast, the double-regression technique is statistically misguided, because there exists no way to reliably estimate the standard errors of its parameter estimates, and he stated that the attempt in Grofman and Migalski to validate the approach through the use of the seemingly unrelated regressions (SUR) approach is statistically inappropriate.²

The last two of these critiques are reasonable ones, but even these criticisms are not 100 percent accurate and, in any case, point to only minor problems that do not in any way affect the substantive conclusions of Grofman and Migalski's (1988) article. On the other hand, the first of Zax's (2002) critiques, his more fundamental criticism of the double-equation approach, is simply flat wrong.

Accuracy of Double Regression and Other Double-Equation Techniques Versus Goodman Single-Equation Ecological Regression and Basic King Single-Equation Ecological Inference

Zax's (2002) claims about double-equation approaches not being needed now that a superior approach, King's (1997) ecological inference method, is available are quite misleading. If Zax were correct, we would not need double-equation methods to improve the accuracy of estimates about bloc voting patterns in situations in which the independent variable is misspecified but could simply use the basic King single-equation methodology for ecological inference. However, reading King would lead us to a quite

different expectation. King was well aware of the potential problems for ecological inference caused by measurement error in the independent variable:

Most methodological discussions of ecological inference avoid this problem by “assuming” that x_i (the proportion of those voting who are black) is known and used in place of X_i (the proportion of VAP that is black). However, although data on x_i are available in a few data sets, these data are quite rare in real voting and most other applications. *Thus, almost any practical use of aggregate data in race and voting studies to make inferences about individuals should include the insights from the double regression procedure* [italics added]. (P. 71)

Thus, contra Zax, we should not expect that single-equation ecological inference would necessarily be better than double-equation methods. Indeed, King himself has instantiated in his computer program EZI a double-equation method (EI2) that can be used for calculating polarization in situations in which the independent variable is misspecified.³

Using hypothetical data, it is possible to construct examples in which single-equation ecological regression outperforms single-equation ecological inference, or vice versa, in mimicking the “true” RBV parameters used to construct the simulated data. Similarly, we can easily construct hypothetical scenarios in which single-equation ecological regression or double regression will produce “out of bounds” results (i.e., estimates in which voting percentages are calculated as below 0 percent or above 100 percent). If, however, we want to make sensible practical judgments about the relative accuracy of different ecological estimation methods for particular types of applications, such as expert witness analysis of RBV patterns, what we want to know is (1) how different are the estimates produced by those methods in practice, and ideally, too, we would like (2) evidence comparing the results of the application of these methods to voting data from real-world elections involving minority versus nonminority candidates that we may reasonably regard as relatively typical of the sorts of elections to which ecological methods to estimate RBV will customarily be applied, but for which we “know” the right answer.

Here, we consider evidence from one real-world municipal election, the 2001 runoff contest between James K. Hahn and Antonio Villaraigosa for mayor of Los Angeles, which has been the subject of considerable research on RBV (Abrajano, Nagler, and Alvarez 2005; Barreto, Villarreal, and Woods 2005; Sonenshein and Pinkus 2002).⁴ For this election, we have

high-quality exit-poll data from a survey of election-day voters, conducted under the auspices of the *Los Angeles Times*, which includes information asked of respondents about their race and whether they were Hispanic or Latino.⁵ In our ecological analyses, we focus on Hispanic versus non-Hispanic voting patterns for the Hispanic candidate, Villaraigosa.⁶ We look at two different proxies for Hispanic share of the actual voting electorate: Hispanic VAP⁷ and Spanish surname percentage among registered voters.⁸ We then compare five different ecological methods for generating estimates of bloc voting patterns in this contest against one another and against the exit-poll results⁹:

Goodman: Standard single-equation estimate of Goodman (1953, 1959) ecological regression.

King E-I: Standard single-equation estimate of King's (1997) ecological inference using King's EZI software.

Goodman DE: Goodman double-equation ecological regression (Grofman et al. 1985), intended to allow for differential turnout rates for minority and nonminority populations. One Goodman regression is used to predict candidate 1's vote share as a fraction of eligible voters; then a second Goodman regression is used to predict candidate 2's vote share as a fraction of eligible voters. Simple algebra is then used to calculate the proportion of nonvoters (abstainers) and to recalculate bloc voting levels as a proportion of the actual (rather than eligible) electorate.

King DE: Method identical to Goodman double-equation method except that King's ecological inference method is used for each of the two equations instead of ecological regression.

King DE 2: Method estimated in the identical manner to King DE, except that calculations for nonvoters are conducted on a precinct-specific basis, as opposed to an aggregate estimate for Hispanic and non-Hispanic, and the overall estimate is calculated as a voter-weighted average of the precinct values.

Tables shows results of the five different methodologies as applied to precinct-level data for the 2001 Hahn-Villaraigosa contest. In addition, we present the results of the *Los Angeles Times* exit poll of election-day voters in the final row. In other studies, Gelman et al. (2001) provided a similar table comparing ecological vote estimates and exit-poll results (see their Table 2, p. 112). We follow a similar approach in this article and also replicate the scatterplot charts provided by Gelman et al. (see Figures 1 to 4). The scatterplots suggest a high degree of homogeneity among the heavily Latino precincts, with consistency

Table 1
Difference in VAP, CVAP, and Registration, Latino and White, City of Los Angeles, 2000

	White	Latino	White	Latino
Adult population	100	100	927,714	1,110,330
Noncitizen	-9	-59	82,637	649,655
Adult eligible	91	41	845,077	460,675
Nonregistered	-14	-12	126,762	137,709
Registered voters	77	29	718,315	322,966

Sources: U.S. census 2000 Summary Files 1 to 4 and Los Angeles County Registrar of Voter Records, 2000.

Note: CVAP = citizen voting-age population; VAP = voting-age population.

in vote choice and turnout and their residuals. In particular, the confluence of heavily Latino precincts around the zero point in the residual plot (Figure 2) suggests a very good model fit. Across all four plots, the results in Figures 1 to 4 demonstrate that the models fit the data very well.

For data runs based on VAP, we find a discrepancy of .13 (.75 vs. .88) between the results of Goodman's single-equation ecological regression and the exit-poll estimate of the proportion of Hispanics who voted for Villaraigosa and a discrepancy of .15 between the results of King's basic single-equation ecological inference program and the exit-poll data. The corresponding estimate from double regression is off by .11. For data runs based on registration data, we find a discrepancy of .12 (.99 vs. .88) between the results of Goodman's ecological regression and the exit-poll estimate of the proportion of Hispanics who voted for Villaraigosa and a discrepancy of .10 between the results of King's basic ecological inference program and the exit-poll data. In contrast, the corresponding estimate from double regression is off by only .02. When we turn to VAP-based estimates of the proportion of non-Hispanics who voted for Villaraigosa, we find discrepancies of .03, .03, and .02 for Goodman, single-equation ecological inference, and double regression, respectively. Similarly, we find discrepancies of .03, .04, and .05, respectively, for estimates of the proportion of non-Hispanics who voted for Villaraigosa on the basis of the registration data. Thus, when using similar data, on VAP (or on registration) with one exception, the first three of our methods give very similar answers. In the one exception to the finding that the three methods give similar results, that for the estimates of Latino voting patterns derived from registration data, we find that the

Figure 1
Vote for Villaraigosa by Proportion Latino in Precinct

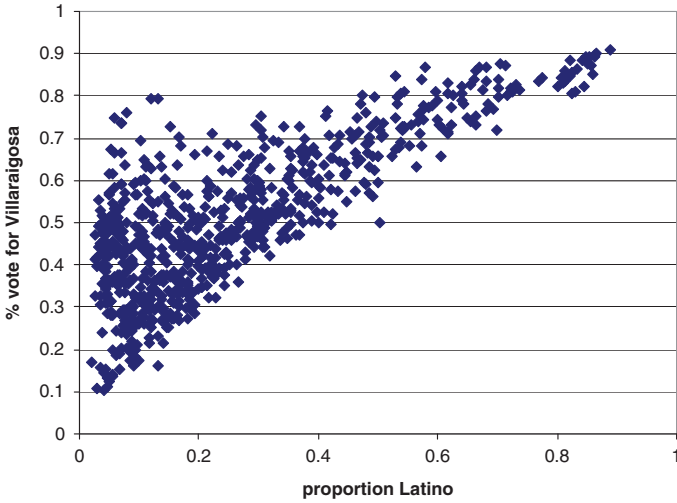


Figure 2
Residuals of Villaraigosa Vote by Proportion Latino in Precinct

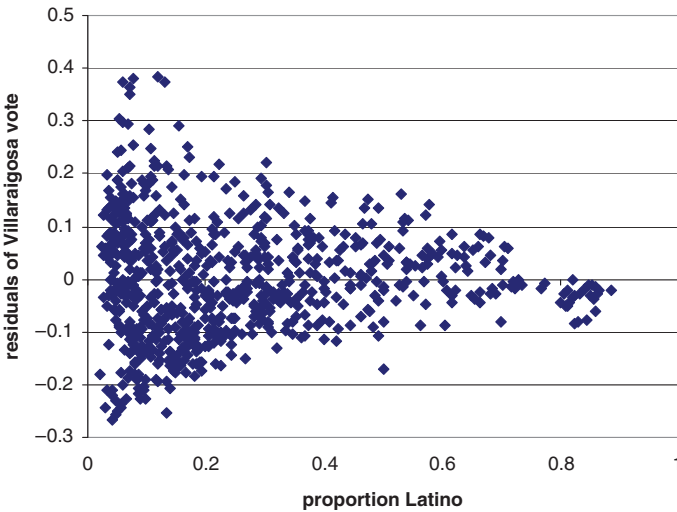


Figure 3
Voter Turnout by Proportion Latino in Precinct

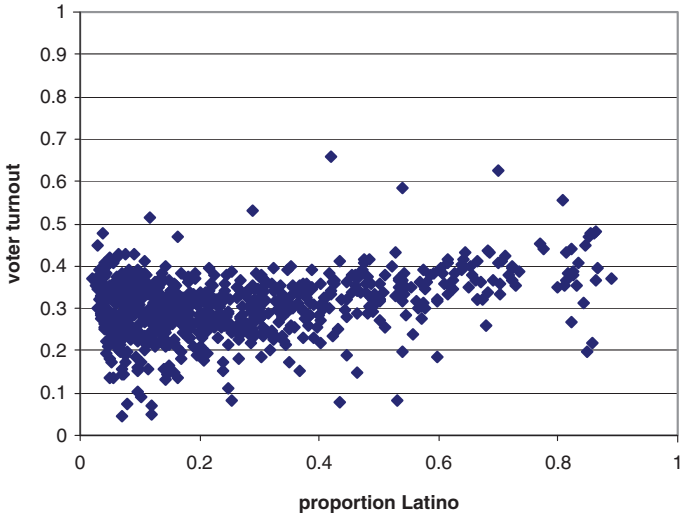
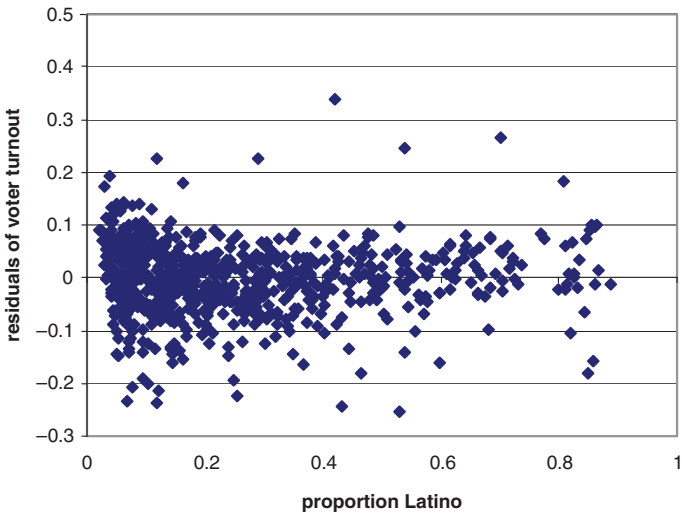


Figure 4
Residuals of Voter Turnout by Proportion Latino in Precinct



double-regression method gives us the most accurate results by a considerable margin.

Now we turn to the results of the two variants of the double-equation approach that make use of King's ecological inference methods instead of simple regression. Contrary to what Zax (2002) would seem to imply, for both VAP and registration data, and for both estimates of Latino and non-Latino voting patterns, the two-equation techniques using King's basic ecological inference software (DE and DE 2) and the double-regression approach described by Grofman and Migalski (1988) give virtually identical results.¹⁰

Another obviously important point about the results shown in Table 2 is that in 9 of 10 of the possible comparisons shown in Table 2, the VAP-derived estimates are not as accurate as the registration-derived estimates. As academics who have served as expert witnesses in voting rights cases involving Latino voting patterns, we are not surprised that for Hispanics in Los Angeles, where there is a dramatic drop-off between their proportion of the VAP and their proportion of the citizen VAP, and a further substantial drop-off in their proportion of the registered voter population, and a further drop-off in their proportion of the actual voting electorate, we generally get better estimates using estimated Spanish surname registration levels than we do using VAP data. For example, for every 100 Latino adults (VAP) in Los Angeles, approximately 59 are not citizens, compared with only 9 of 100 Whites who are noncitizens. Furthermore, among the adult eligible population, 30 percent of Latinos are not registered to vote, compared with just 15 percent of Whites (see Table 1). Thus, it is very important to have as accurate an estimate as possible of the independent variable: percentage Hispanic. The registration data are closer to the actual election-day turnout proportions than are the VAP data.

Still, even using VAP as the basis for our calculations, we do not do a bad job in using ecological methods to estimate the parameters of Latino and non-Latino voting, once we recognize that the basic legal issue being addressed is simply whether there are differences between Latino and non-Latino voting patterns in this contest such that a majority of Latinos support the Latino candidate and a majority of non-Latinos do not (Grofman et al. 1992). Using all five techniques, and including the VAP-derived data, the range of our estimates argue that support for Villaraigosa was somewhere from 73 percent to 99 percent or higher among Latinos and from 26 percent to 35 percent among non-Latinos. The exit-poll data (88 percent and 31 percent, respectively) are toward the middle of these ranges. Of course, we

Table 2
Comparisons of the Accuracy of Five Ecological Methods of Estimating
Racial Bloc Voting Using Real Election Data, City of Los Angeles
Mayoral Election, 2001

Estimate	Estimates Using VAP		Estimates Using Voter Registration	
	Latino	Non-Latino	Latino	Non-Latino
Goodman (single equation)	.75 (.012)	.34 (.009)	≥.99 (.026)	.34 (.006)
King EI (single equation)	.73 (.002)	.34 (.001)	.98 (.003)	.35 (.001)
Goodman DE (aggregate turnout)	.77	.30	.91	.26
King DE (aggregate turnout)	.76	.27	.91	.28
King DE 2 (precinct turnout)	.78	.27	.91	.30
Election-day exit poll (minus absentee)	.88	.31	.88	.31

Sources: Los Angeles City Clerk, Statement of Votes Cast, for polling-place voters, 2001. Los Angeles County Registrar of Voters Database, merged with U.S. Census Bureau Spanish surname list.

Note: VAP = voting-age population; Goodman = standard single-equation estimate of Goodman (1953, 1959) ecological regression; King E-I = standard single-equation estimate of King's (1997) ecological inference using King's EZI software; Goodman DE = Goodman double-equation ecological regression (Grofman et al. 1985), intended to allow for differential turnout rates for minority and nonminority populations; King DE = method identical to Goodman double-equation method except that King's ecological inference method is used for each of the two equations instead of ecological regression; King DE 2 = method estimated in the identical manner to King DE, except that calculations for nonvoters are conducted on a precinct-specific basis, as opposed to an aggregate estimate for Hispanic and non-Hispanic, and the overall estimate is calculated as a voter-weighted average of the precinct values.

cannot be sure that the legally irrelevant differences in estimation across methods we find for this one election will always be found, but between us, we have examined hundreds of contests and found that differences across methods tend to be minor when looking at African American voting patterns, and also minor for Hispanic voting patterns, at least as long as these are estimated from Spanish surname registration or turnout data.¹¹

Errors in Formulating Equations to Calculate Double-Regression Parameters When Voters May Vote for More Than a Single Candidate

We are grateful to Zax (2002) for exposing a mistake in the statement of the identity of equation 5 in Grofman and Migalski (1988). This equation is used to expand the double-equation approach to the case in which voters can vote for more than a single candidate (as in multimember districts or at-large elections not involving numbered places). However, we find his apparent disinterest in a constructive response as how to fix the error to which he calls attention somewhat puzzling, because the error turns out to be relatively minor, requiring the replacement of \bar{n} (the mean number of ballots casts per voter) with n (the maximum number of ballots each voter is eligible to cast).

Moreover, and even more important, although this error does effect equations 5 through 11 in Grofman and Migalski (1988), it turns out not to matter for the final estimating equations used in the article, equations 12 and 13. That is because \bar{n} in the original version of equation 5 and n in the respecified version of that equation may be taken as scaling constants (one specified by the voting system and the other derived directly from the aggregate data), and the nature of the double-equation estimation in equations 12 and 13 involves a division in which this scaling constant (being found in both equations) simply cancels out.¹² Thus, even though there was a mistake in the specification of some of the earlier equations, it did not affect the calculations of the actual parameters of interest reported by Grofman and Migalski. And in any case, we note that the model proposed in equation 5 in that article only applies to the case in which voters have multiple *votes* to cast and thus does not apply to most of the applications to date of the double ecological regression methodology, that in expert witness testimony in voting rights cases involving *single* offices involving two candidates (or, if more than two candidates for the same single seat, in which the votes for candidates of the same race can be combined so as to treat the election as if it were a two-candidate contest).

Standard Errors for Double Regression

We turn now to the third issue raised by Zax (2002), the problem of inaccurate standard errors. Once again, Zax is better at criticizing than suggesting improvements. While he asserts that the variance-covariance

method we offer to calculate standard errors for the double-regression technique is inappropriate because of potential problems such as correlated error, Zax provides no formula for how standard errors should be calculated for the double-equation ecological regression approach. Perhaps even more important, neither Zax (2002) nor Zax (2005) offers any empirical evidence that, for the kinds of real-world data sets at issue in voting rights litigation, the failure to take heteroscedasticity or correlated errors into account substantially biases error estimates derived from variance and covariance calculations in a fashion that would affect legal conclusions drawn from the election data.¹³

Equally important, Zax's (2002) claim that unlike what is the case with double regression, "trustworthy standard errors can be accomplished almost automatically through King's (1997) maximum likelihood model of random coefficients" (p. 83) simply is not accurate, because using King's approach in the presence of error in the independent variable would seem to require a multiequation approach, and calculating standard errors when King's approach is applied twice (in interrelated equations of any sort) requires additional complex calculations involving variances and covariances in exactly the same way that calculating errors for double regression on the basis of Goodman's ecological regression approach does, or a process of estimating those standard errors through simulation (cf. Herron and Shotts 2003a, 2003b).

Discussion

Zax (2002) is somewhat helpful to students interested in problems of reliable ecological inference by catching a mistake in the statement of formulas in Grofman and Migalski (1988) used to extend the double-regression approach to the case in which voters have more than one vote to cast, but he does not realize that the mistake in question is not only minor but also that it does not apply to the equations they use to perform the actual calculations reported in that article. Zax is also correct about the inappropriateness of SUR as a validity check on double-regression results, but he neglects the use of SUR as a reliability check on the variance-covariance-based standard error calculations in Grofman and Migalski's article. On the other hand, Zax's most important claim, that a double-equation approach is obviated by reliance on King's ecological inference methods, is at best overstated and at worst dead wrong, and it contradicts the views of King (1997:71-72) himself.

While we do in fact share with Zax (2002) the view that King's ecological inference approach is superior, all else equal, to ecological regression approaches (see Grofman 2000), it is important to recognize that in practice, in applications of bloc voting analysis in situations in which voting is highly racially polarized, the two families of methods will tend to give very similar results. Contrary to what Zax would lead us to expect, we see from Table 2 that for the Villaraigosa-Hahn contest, the estimates derived from King's method and from ecological inference are essentially identical when we compare the single-equation variants of these methods with each other, and they are essentially identical as well when we compare, one to the other, the double-equation variants of each model that we used. Furthermore, where there are substantial differences found between single-equation and double-equation approaches, it is the double-equation approach that appears best. Thus, contrary to Zax, if we are to apply King's ecological inferences approach when we have substantial misspecification in the independent variable due to substantial differences in turnout across racial groups, we need to correct for such misspecification using a double-equation approach.¹⁴

The finding that for real-world data in situations in which voting is highly polarized, single-equation ecological regression and single-equation ecological inference give essentially identical or near identical estimates has been corroborated numerous times by experts who have testified in voting-rights cases who have applied both King's and Goodman's approaches to the same data (Richard Engstrom, personal communication, March 2006). Unfortunately, such comparisons are generally only reported in trial transcripts, or only available in expert witness reports, rather than being published in the academic literature. However, several recent published works (Bullock and Gaddie (2005, especially Table 11; Greiner 2007, especially Table 5; see also Bullock and Gaddie 2006) show how remarkably near to identical various ecological methods are when they are applied to real data of the kind (biracial contests) studied in voting-rights litigation and how, for these data sets, for legal purposes, it would not matter as to which estimate is used to determine if Whites and Blacks vote differently.

There are fewer examples of work comparing various types of ecological inference in situations involving candidacies of different races in which the "true" (average) individual-level voting parameters are known from other data sources. Two special issues of *Historical Methods* do contain articles that directly bear on comparisons between aggregate-level inference and real data on elections involving candidates of more than one race, and our findings are generally consistent with those reported in

articles by Kousser (2001), Phelan (2001), and Lewis (2001). Similarly, Liu (2007) shows that for a biracial contest in New Orleans in which racial sign-in data were available, different ecological methods gave very similar results in estimating minority and nonminority turnout levels and did a good job in reliably estimating turnout.¹⁵

We would emphasize, however, that the ability of any form of ecological inference to recover individual-level parameters from aggregated data is in part contingent on the nature of the patterns in that data. The clearer the “signal” in the data, and the more substantial is the range of variation in the independent variable(s), the easier it is for that signal to be correctly detected by ecological inference methods, even with a limited number of data points. Still, when we compare methods, for RBV analyses in biracial or biethnic contests, using real data, if there are very substantial turnout or eligibility differences between the groups under study that are not already largely being taken into account in the independent variable we use as our race proxy, then we should to be using some *two-equation* method to double-check our results.¹⁶

In this connection, it is important to note that, as mentioned earlier, there is an alternative two-equation technique to the double-equation model presented here. In that model, we first estimate turnout through a single-equation method (whether Goodman or EI, with the latter preferred) and then run a second single-equation version of Goodman or EI (with the latter preferred) in which we substitute in our new turnout estimates for the less accurate independent variable that was the best we previously had to work with. Using EI (and a model within the EZI program called EI2), this approach has been used by Liu (2001; cf. Cho and Gaines 2004; Voss 2004) and some expert witnesses in voting rights cases (Lisa Handley, personal communication, August 9, 2007). We have run some preliminary comparisons of the double-equation approaches presented in this article and in our earlier work with this alternative two-equation approach and find, for highly polarized elections in jurisdictions where there are many racially homogeneous or near homogeneous precincts, the answers from the various double-equation approaches tend to be virtually identical.¹⁷

We have generated King EI2 (double-equation) estimates of polarization for the Los Angeles data on the Villaraigosa-Hahn contest.¹⁸ These are as follows: estimated Latino vote for Villaraigosa = .90 ($SE = .0811$), and estimated non-Latino vote for Villaraigosa = .26 ($SE = .1770$). The double-equation method used by Grofman and Migalski (1988) gives us virtually identical estimates of .91 and .26 for these two parameters. In

sum, the choice of double-equation method does not affect our conclusion as to whether voting in this contest is polarized along Hispanic versus non-Hispanic lines, and because the exit-poll data gave us estimates of .88 and .31 for these same parameters, we can be confident that our ecological results are pretty much “on the money.”¹⁹

Notes

1. These experts now either use King’s approach exclusively, instead of Goodman-like ecological regression, or use ecological inference as a check on results generated by ecological regression methods.

2. Zax (2005) added a fifth criticism, the claim that double regression produces asymptotically biased estimates. That argument takes us into issues beyond the scope of this note. Suffice it to remark that whatever problems this causes are (1) likely to be minor in real-world RBV data and (2) shared by any double-equation variant of King’s ecological inference approach.

3. The estimating equations in this double-equation method are distinct from those in the double-equation method in Grofman and Migalski (1988). We will have more to say about in the concluding section of this article.

4. Hahn won the general election on June 5, 2001, with 54 percent of the vote, defeating Villaraigosa by 7 percentage points.

5. We use *Latino* and *Hispanic* as synonyms in this article.

6. We have election data from two main sources: the Los Angeles County Registrar of Voters database and the Los Angeles City Clerk’s Statement of Votes Cast. The unit of aggregation is the voting tabulation unit (here called precincts), and our data include a complete enumeration of registered voters and votes cast in each of the 1,730 precincts in Los Angeles.

7. We have Hispanic VAP in each of the Los Angeles precincts (ca. 2000) from the U.S. census.

8. The Spanish surname list is based on the 1990 census and is constructed by tabulating the responses to the Hispanic-origin question. Each surname is categorized by the percentage of individuals who identified themselves as “Hispanic.” Although the use of this instrument results in a modest underestimate, given the presence of Latinos with non-Hispanic surnames, the U.S. Census Bureau estimates that this captures 93.6% of all Hispanics, and fewer than 5% of those identified are false. For a full explanation on the methodology of the list, see Word and Perkins (1996). While there are both Type I and Type II errors possible in such matching, unpublished work done by the demographer William O’Hare (in conjunction with one of the present authors) suggests that for Los Angeles, the errors are rather minor. In any case, we can compare results using this estimate of the size of the Latino population to the exit-poll data.

9. However, considerable care must be used in making comparisons with these exit-poll data. Our ecological methods only make use of data on polling-place votes cast in the 2001 mayoral election and exclude votes cast in absentia. Although the *Los Angeles Times* exit poll only interviewed polling-place voters on election day and also contained no data on the choices

of absentee voters, the *Times* adjusted its final published estimates of voting by race and Hispanic ethnicity using a reweighting of the exit-poll data to ensure that its announced exit-poll results were consistent with the overall election outcome. That is, the available poll data are only for polling-place voters. Thus, we cannot directly compare our ecological results about bloc voting with the corresponding figures published in the *Times*'s postelection news stories, because the latter incorporate absentee voter preferences. To match our Los Angeles city election data to the exit poll conducted by the *Times*, we had to go back to the raw precinct-level data from the *Times* poll and retabulate vote preference by race and ethnicity for polling-place voters. To improve comparability, we also chose to limit ourselves to precinct-level data for just polling-place voters (excluding absentee voters). Our analyses of the raw data suggest that 88% of Latinos voted for Villaraigosa, while the published *Times* figure, which attempts to account for absentee preferences, is 82%. Of course, even an exit poll is far from perfect, but then no method is. Still, an exit poll provides a useful baseline against which to compare the results of aggregate-level ecological methods using precinct-level data, even though we recognize that sometimes, survey results may actually be less accurate than those derived from ecological methods (when the surveys are not fully random and/or there are interviewer effects or response bias effects in the data; see especially Liu 2007).

10. For some estimates, one or both of the King-based methods is marginally superior to double regression, and for other estimates, double regression comes closer to the exit-poll data, but the differences are in all cases trivial and, when averaged over the two parameters being estimated, essentially nonexistent (see Table 1).

11. In this article, in Table 2, we have simply omitted standard errors for the three double-equation approaches, because presenting the standard errors for each equation singly can be misleading. For present purposes, the proof is in the pudding, that is, in the direct comparisons of ecological estimates with exit-poll data shown in Table 2. It is apparent, however, given misspecification in the independent variable, that the standard errors reported for both the single-equation ecological regressions and ecological inference method are too low. It is straightforward to show that the presence of such errors creates nonlinearity in the data.

12. A correct respecification of equations 5 through 11 is available on request from the authors.

13. When Zax (2002) attacks Grofman and Migalski's (1988) use of SUR, he is technically correct in that SUR is intended for use in interrelated equations in which there are different independent variables, but SUR can still be used to check the *reliability* of the calculations in Grofman and Migalski (1988), even if not their *validity*. In that context, we would note that except for one trivial case of rounding error and one missing negative sign that was inadvertently omitted from the table, the computer-generated SUR estimates in Grofman and Migalski are in fact identical to those of the authors' spreadsheet-based calculations using variances and covariances. Moreover, after we correct a typo, a missing negative sign, all of the covariances shown in their Table 2 (p. 449) are negative and are similar in magnitude to the variances. Moreover, Zax (2002) does not mention that Grofman and Migalski look for the kind of contextual effects that might bias estimates of standard errors by running polynomial regressions and find no substantial improvement in fit (cf. Owen and Grofman 1997). Because the standard errors for each of the two individual equations in the various double-equation approaches in table can be misleading, as noted earlier, we have omitted them from Table 2.

14. There is one further complication that we neglect, namely, the fact that there is a substantial African American population in the Los Angeles electorate, and thus, seeking to

distinguish voting patterns of non-Hispanic Whites from those of non-Hispanic Blacks might allow us some improvement in our estimates of Hispanic and non-Hispanic voting patterns. However, standard techniques for multibloc ecological inference do not really exist, although ideas along this line are discussed by King (1997:chap. 15) and in subsequent work (see review in Greiner 2007). In any case, this issue takes us well beyond the limited scope of this essay. Moreover, in the Los Angeles mayoral election we report on here, we know from exit-poll and other data that Blacks, like non-Hispanic Whites, were largely supporting Hahn. Also, there is considerable residential segregation in Los Angeles for both Hispanics and Blacks, so the kinds of inference problems involving multibloc electorates referred to by Greiner (2007), who analyzed data on an election in Georgia in which there were no heavily homogeneous Hispanic precincts, are not really of moment for the analyses presented here.

15. Liu does find some cases in which an ecological regression method actually comes slightly closer to the true values than the corresponding estimates derived from ecological inference, but his data generally support agreement among the basic methods. The only real exception to congruity among ecological approaches in estimates of polarization in Liu's data (see his Table 2) comes from Freedman et al.'s "neighborhood model," but this model is quite different from the standard approaches and in our view is simply not very useful for analysis of racial polarization for candidates in biracial contests, except when there is almost perfect segregation, although it will sometimes give estimates for turnout that match more closely with what we get from standard methods (see Grofman 1991; cf. Liu 2007, Table 3). Liu's best estimate of turnout comes from a variant of King's ecological inferences model that includes a covariate, but even this model does not generate large differences from other approaches (e.g., King's basic model estimates Black turnout at 51.8%; the model with one covariate gives us an estimate of 49.1%). Also, as Liu acknowledges, there are no agreed-upon procedures for determining which covariates to introduce.

16. In this context, we would note that Liu (2007) describes his own empirical findings as seemingly inconsistent with the claim made by Zax (2002, 2005) that "there is no use in the double regression approach."

17. It has been our experience that in biracial contests, ecological estimates of turnout by race tend to have higher standard errors than estimates of RBV because there can be considerable within-group variance in turnout behavior (often related to socioeconomic status differences). The EI2 double-equation method involves one stage at which turnout is directly estimated. In contrast, the double-regression methods reported in Table 2 estimate turnout more indirectly.

18. There is also some new work on ecological inference in $r \times c$ tables larger than 2×2 , extending ideas in King (1997:chap. 15). Such methods can be used to deal with abstention (and thus with differential levels of turnout across racial or ethnic groupings) by estimating parameters for a 2×3 table in which the first column is votes for the minority candidate, the second column is votes for the nonminority candidate, and the third column is abstention. However, this methodology is not yet well established and raises issues beyond the scope of this brief essay.

19. We remind the reader that exit-poll results are themselves subject to sampling and other forms of error (such as response bias), so that we would not expect a perfect fit between ecological inference and survey data. Indeed, in the New Orleans election examined by Liu (2007), he is able to compare both survey findings and ecological regression estimates of turnout to actual sign-in data by race, and he actually finds the latter more accurate than the former in estimating turnout by race.

References

- Abrajano, Marissa, Jonathan Nagler, and Michael Alvarez. 2005. "A Natural Experiment of Race-Based and Issue Voting: The 2001 City of Los Angeles Elections." *Political Research Quarterly* 58 (2): 203-18.
- Barreto, Matt A., Mario Villarreal, and Nathan Woods. 2005. "Metropolitan Latino Political Behavior: Turnout and Candidate Preference in Los Angeles." *Journal of Urban Affairs* 27:71-91.
- Bullock, Charles S., III and Ronald Keith Gaddie. 2005. "An Assessment of Voting Rights Progress in Georgia." Available at http://www.aei.org/docLib/20060210_Georgia.pdf.
- . 2006. "An Assessment of Voting Rights Progress in Louisiana." Available at http://www.aei.org/docLib/20060308_Louisiana.pdf.
- Cho, Wendy Tam and Brian J. Gaines. 2004. "The Limits of Ecological Inference: The Case of Split-Ticket Voting." *American Journal of Political Science* 48 (1): 152-71.
- Cho, Wendy Tam, George C. Judge, and Bruce Cain. 2002. "Some Empirical Evidence on the Impact of Measurement Errors in Making Ecological Inferences." Unpublished manuscript.
- Freedman, David, S.P. Klein, J. Sacks, C.A. Smyth and C.G. Everett. 1991. "Ecological Regression and Voting Rights." *Evaluation Review* 15(6): 673-711.
- Gelman, Andrew, Stephen Ansolabehere, Phillip N. Price, David K. Park, and Lorraine C. Minnite. 2001. "Models, Assumptions, and Model Checking in Ecological Regressions." *Journal of the Royal Statistical Society* 164:101-18.
- Goodman, Leo. 1953. "Ecological Regression and the Behavior of Individuals." *American Sociological Review* 18 (6): 663-64.
- . 1959. "Some Alternatives to Ecological Correlation." *American Journal of Sociology* 64:610-25.
- Greiner, D. James. 2007. "Ecological Inference in Voting Rights Disputes: Where Are We Now, and Where Do We Want to Be?" *Jurimetrics* 47 (2): 115-67.
- Grofman, Bernard. 1991. "Statistics Without Substance: A Critique of Freedman et al. and Clark and Morrison." *Evaluation Review* 15 (6): 746-69.
- . 2000. "A Primer of Racial Bloc Voting." Pp. 44-67 in *The Real Y2K Problem: Census 2000 Data and Redistricting Technology*, edited by Nathaniel Persily. New York: Brennan Center for Law and Justice, New York University.
- Grofman, Bernard, Lisa Handley, and Richard Niemi. 1992. *Minority Representation and the Quest for Voting Equality*. New York: Cambridge University Press.
- Grofman, Bernard and Michael Migalski. 1988. "Estimating the Extent of Racially Polarized Voting in Multicandidate Elections." *Sociological Methods & Research* 16 (4): 427-54.
- Grofman, Bernard, Michael Migalski, and Nicholas Noviello. 1985. "The 'Totality of Circumstances' Test in Section 2 of the 1982 Extension of the Voting Rights Act: A Social Science Perspective." *Law and Policy* 7 (2), 209-23.
- Herron, Michael C. and Kenneth W. Shotts. 2003a. "Cross Contamination and EI-R." *Political Analysis* 11:77-85.
- . 2003b. "Using Ecological Inference Point Estimates as Dependent Variables in Second-Stage Linear Regressions." *Political Analysis* 11:44-64.
- King, Gary. 1997. *A Solution to the Ecological Inference Problem*. Princeton, NJ: Princeton University Press.

- Kousser, Morgan J. 1973. "Ecological Regression and Analysis of Past Politics." *Journal of Interdisciplinary History* 4:237-62.
- . 2001. "Ecological Inference From Goodman to King." *Historical Methods* 34:100-26.
- Lewis, Jeffrey B. 2001. "Understanding King's Ecological Inference Model: A Method-of-Moments Approach." *Historical Methods* 34:170-88.
- Liu, Baodong. 2001. "The Positive Effect of Black Density on White Crossover Voting: Reconsidering Social Interaction Theory." *Social Science Quarterly* 82:600-13.
- . 2007. "EI Extended Model and the Fear of Ecological Fallacy." *Sociological Methods & Research* 36:3-25.
- Loewen, James and Bernard Grofman. 1989. "Comment on Recent Developments in Methods Used in Voting Rights Litigation." *Urban Lawyer* 21 (3): 589-604.
- Owen, Guillermo and Bernard Grofman. 1997. "Estimating the Likelihood of Fallacious Ecological Inference: Linear Ecological Regression in the Presence of Context Effects." *Political Geography* 16 (8): 657-90.
- Phelan, Thomas. 2001. "Comparing Individual-Level Returns With Aggregates: A Historical Appraisal of the King Solution." *Historical Methods* 34:127-34.
- Sonenshein, Raphael and Susan Pinkus. 2002. "The Dynamics of Latino Political Incorporation: The 2001 Los Angeles Mayoral Election as Seen in Los Angeles Times Exit Polls." *PS: Political Science & Politics* 35:67-74.
- Voss, Stephen. 2004. "Using Ecological Inference for Contextual Research." Pp. 69-96 in *Ecological Inference: New Methodological Strategies*, edited by Gary King, Ori Rosen, and Martin Tanner. New York: Cambridge University Press.
- Word, David L. and R. Colby Perkins. 1996. "Building a Spanish Surname List for the 1990's—A New Approach to an Old Problem." Technical Working Paper No. 13. Washington: U.S. Census Bureau.
- Zax, Jeffrey S. 2002. "Comment on 'Estimating the Extent of Racially Polarized Voting in Multicandidate Contests' by Bernard Grofman and Michael Migalski." *Sociological Methods & Research* 31:75-86.
- . 2005. "The Statistical Properties and Empirical Performance of Double Regression." *Political Analysis* 13 (1): 57-76.

Bernard Grofman received his B.S. in Mathematics at the University of Chicago in 1966 and his Ph.D. in Political Science at the University of Chicago in 1972. He has been on the faculty of the University of California, Irvine since 1976 and Professor of Political Science since 1980. In 1986 he was a Fellow at the Center for Advanced Study in the Behavioral Sciences at Stanford. In 2000 he was elected President of the Public Choice Society. In 2001 he became a Fellow of the American Academy of Arts and Sciences. In 2007 he became the Director of the Center for the Study of Democracy. In January of 2008 he was selected as the first Jack W. Peltason (Bren Foundation) Endowed Chair. His past research has dealt with legislative representation, electoral rules, redistricting, and voting rights.

Matthew A. Barreto is an assistant professor of Political Science at the University of Washington. He has published in leading political science journals, including the *American Political Science Review*. His work deals primarily with issues of racial and ethnic politics and immigration, as well as political participation more broadly.