# The Aggregate Impact of Explanatory Variables in Logit and Linear Probability Models

Charles E. Denk; Steven E. Finkel

# The Aggregate Impact of Explanatory Variables in Logit and Linear Probability Models*

Charles E. Denk, *Department of Sociology, University of Virginia*
Steven E. Finkel, *Department of Government and Foreign Affairs, University of Virginia*

This paper presents methods for computing aggregate change in probabilities of a binary dependent variable from changes in distributions of independent variables in logit and linear probability models. We develop a measure for the logit model based on a Taylor series polynomial expansion that solves the problems inherent in the nonlinearity and nonadditivity of the logit specification. The method can be used to make out-of-sample predictions based on real or hypothetical changes in one or more independent variables and may also be used to assess the relative "importance" of different independent variables by computing the change in dependent probabilities accounted for by each variable. The measure is in the same spirit as Achen's (1982) "level importance" measure for linear models and thus fills an important gap in logit regression analysis. We show, on the basis of simulations and controlled validation in an empirical example, that the aggregate logit impact measure can produce numerical results that differ substantially from the equivalent measure for linear probability models. We provide guidance for future research on the detailed application of the logit method and the criteria for choice of the logit versus the linear aggregate impact measures.

## Introduction

Social scientists may have at least two distinct explanatory goals in assessing the effects of a set of independent variables on a dependent variable. One goal is to explain why some units show higher scores on the dependent variable than others, that is, to explain variation on the dependent variable among the individual units that compose the sample or population in question. Another goal is to assess the independent variables' effects on the aggregate value of the dependent variable, that is, to explain why the dependent variable overall is high or low, based on the distributions of the independent variables among individual units in the sample. Achen (1982, 68–77) refers to these goals, respectively, as explaining the dependent variable's "dispersion" and explaining its "level." For example, the researcher may wish to know not only which variables explain why some individuals are politically sophisticated and some are not but also the extent to which each variable contributes to differences in the overall amount of political sophistication seen across particular samples.

The analytic distinction between these two goals is plainly evident in the

context of electoral behavior. As Markus (1988, 142) notes, the research goal may be either "to account for election outcomes or to explain individual vote decisions. Elections are, of course, aggregations of individual votes. However, explaining differences in individual votes requires attention to *interindividual* differences in relevant causal factors, whereas accounting for variance in election outcomes depends upon *interelection* changes in the distribution of those causal factors." Following this logic, Markus (1988) argues that although voters' perceptions of their personal financial situation play some role in explaining individual votes, the variation in this factor across elections is so small that it does not account for interelection changes in voting outcomes. Several other studies also distinguish between accounting for variation in individual votes and in electoral outcomes (Finkel n.d.; Miller and Shanks 1982; Rosenstone 1983) and conclude that different variables may be important for each level of explanation.

As yet, however, procedures for "explaining" the aggregate level of a dependent variable exist only for the linear model (Achen 1982; Markus 1988, 144; Stover 1987), using a logic we shall describe below. Consequently, researchers are led to employ the linear model because of its statistical convenience, even where a dichotomous dependent variable implies that the underlying behavioral assumptions of the linear model are likely to be false. In this paper, we present procedures for accounting for the aggregate level of a dichotomous dependent variable using the increasingly popular logistic regression model. We show that calculating the "level importance" of a variable, in a sense comparable to that concept for linear models, is possible within the logit context. Further, under certain conditions the results for logit models may differ substantially from calculations made using the linear probability model and least squares estimation. Given the statistical and theoretical superiority of the logistic model over the linear specification in many applications, we believe that the procedures outlined here offer an improved means for accounting for the level of a categorical dependent variable and for assessing the "level importance" of a set of explanatory variables.

In the following sections, we review the basic idea of explaining the level of a dichotomous dependent variable in the context of the linear probability model and least squares regression. We then present our derivation of a comparable measure for the logit model. We follow this with a brief numerical simulation demonstrating the sensitivity of level importance calculations to several important aspects of the independent variable's distribution. We then demonstrate the use of the logit measure in an analysis of presidential voting in the 1972–76 National Election Studies (NES) panel and compare its results with those obtained from its linear least squares counterpart. We conclude with some practical advice on using the logit and least squares procedure in empirical analyses in various contexts.

## Linear Probability Models and Aggregate Change

Probability models describe the distribution of a categorical (typically binary) outcome variable conditioned on some set of explanatory variables across individuals. We denote the *predicted* probabilities for a single outcome at the individual level as:

$$P_i = \Pr [Y_i = y \,|\mathbf{x}_i].$$

The two principal regression-like contenders are the dummy dependent or linear regression model:

$$P_i = \mathbf{b}'\mathbf{x}_i. \tag{1}$$

And the log-odds or logit or logistic regression model:[1]

$$\ln\!\left(\frac{P_i}{1 - P_i}\right) = \mathbf{b}'\mathbf{x}_i. \tag{2}$$

In each model $P_i$ denotes the probability that a categorical dependent variable, $Y$, for individual $i$, attains a specific value $y$; $\mathbf{x}$ is a vector of explanatory variables; and $\mathbf{b}$ is a vector of effect coefficients. The inner product $\mathbf{b}'\mathbf{x}_i$ denotes, for each individual, the sum of products of paired regression variables and coefficients, implicitly including an intercept term.

The impact of a change in a specific $x_{ij}$ on $P_i$ for an individual $i$ can be defined for any model via the derivative $dP_i/dx_{ij}$, which depends on the model form.[2] For the linear model:

$$dP_i/dx_{ij} = b_j \tag{3}$$

and for the logit model:

$$dP_i/dx_{ij} = b_j P_i(1 - P_i). \tag{4}$$

Note that the impact of a change in $x_j$ for the linear model is a constant, while for the logit model it is proportional to $P_i(1 - P_i)$, which is the variance of the individual $Y_i$. This variance is at a maximum when $P_i = .5$, which leads to the characterization of the logit and similar models as "tipping models," denoting the existence of a range of values where the probability is most variable, and hence explanatory factors have maximum impact. Since the effects of independent variables on the probability of a given outcome are often assumed to be

---

[1]Later we shall show how the probit model, also popular in this context, may be approached in the same general way as we present for the logit model.

[2]The derivative $dP/dX$ may be interpreted as the proportional change or slope $\Delta P/\Delta X$ when $\Delta X$ is very small. Nonlinear models have varying slopes, estimates of which also vary with the precise value of $\Delta X$. Hence, calculus offers the only precise treatment.

dependent on the value of $P_i$, the logit (or probit) model is a more accurate depiction of the underlying process than the linear specification (Aldrich and Nelson 1984; Hanushek and Jackson 1977).[3] In voting analyses, the logit specification better captures the behavior of "swing voters," that is, those who are approximately equally likely to vote either way and therefore exhibit more variance in their choices.

Accounting for the level of a dependent variable in the linear model is a relatively straightforward exercise that exploits a mathematical identity relating the means of all variables (Achen 1982, 68–77):

$$\bar{P} = \frac{1}{n} \sum_{i=1}^{n} P_i = \frac{1}{n} \sum_{i=1}^{n} \left( b_0 + \sum_j b_j x_{ij} \right) = b_0 + \sum_j b_j \bar{x}_j. \tag{5}$$

Thus, the aggregate probability of the dependent variable being equal to one is a function of the intercept term and each independent variable's mean weighted by its unstandardized regression coefficient. The contribution of each independent variable in producing the "level" of the dependent variable observed in the sample is simply $b_j \bar{x}_j$.

However, the intrinsic meaning of each variable's contribution, computed in this way, depends critically on each variable's scale. For example, if an independent variable's scale range is changed from [0, 10] to [−5, 5], its mean would change, while the regression coefficient would remain constant. Thus, each independent variable's contribution to the dependent variable's level directly depends on the apparently nonsubstantive decision of scale location. Consequently, as Achen's discussion makes clear, any meaningful "level importance" calculation implicitly invokes some form of cross-sample or hypothetical comparison, for example: "why this Republican did so much better *than others before him*" or "How much difference does [a newspaper's reporting] bias make *(compared to fair reporting)*?" (Achen 1982, 72, emphasis added). Thus, a precise accounting of level importance involves relating aggregate changes between a specific sample and an explicit baseline, either observed or hypothetical.[4]

In the linear model, the levels of the dependent variable in the observed and baseline distributions (with the baseline denoted $\bar{P}^0$ and $\bar{x}_j^0$) are related by:

$$\bar{P} = \bar{P}^0 + \sum_j b_j (\bar{x}_j - \bar{x}_j^0)$$

---

[3] The linear probability model also violates other OLS assumptions, namely, normality and equal variability of residuals at all points (homoscedasticity). Weighted least squares may be used to estimate more efficiently the parameters of the linear probability model, but the model's specification of constant effects for explanatory variables remains a fundamental problem.

[4] In this context, standardizing the independent variables worsens the situation by further obscuring the substantive reference point of the scale.

and the deviation of a specific sample from the baseline:

$$\Delta \bar{P} = \sum_j b_j(\bar{x}_j - \bar{x}_j^0) = \sum_j b_j \Delta \bar{x}_j \tag{6}$$

where $\Delta \bar{x}_j$ expresses the (intratemporal or intertemporal) deviation of $x_j$ from the baseline distribution.

Studies that make use of this logic for both continuous variable and linear probability models are widespread in electoral and other social science research, beginning at least with Stokes's (1966) seminal analysis of the importance of cross-election changes in the electorate's partisan attitudes in explaining aggregate shifts in the presidential vote. More recently Smith (1989) uses a model explaining an individual's political sophistication to simulate what the level of the electorate's overall sophistication would be if education levels were higher than they are now, or if political interest were at a higher level. In this case, the baseline comparison is the sample for which he has relevant data to estimate the initial regression model. Markus (1982) accounts for the "importance" of a set of independent variables on the level of the presidential vote in 1980 in a slightly different way. By measuring each variable on a scale such that zero represents theoretical "neutrality," Markus can show each independent variable's effect on the overall vote *compared with* a hypothetical "neutral" electorate. For example, the public's low rating of Jimmy Carter's performance in office had an overall pro-Republican effect of 6.5 percentage points, compared with an electorate that would have been neutral in their ratings of the incumbent's job performance (Markus 1982, 559). Finkel (n.d.) uses this same procedure to account for the impact of change in attitudes on the vote over the course of the 1980 campaign. In a later article, Markus (1988) pools several election data sets and notes that, given his individual-level model of the vote, "shifting from a distribution of perceived financial well-being such as occurred in 1976 (a relatively bad year) to a distribution such as that for 1956 (a relatively good one) would increase the incumbent's vote share by only 1.3%." Lewis-Beck (1988, 85–87) uses these procedures in a cross-national study to estimate the aggregate vote shifts that would occur in several West European democracies if 20% more of the electorate (from a 1984 baseline) believed that the future performance of the government will "worsen" the national economy.

In all of these examples, the relative simplicity of the results is a direct consequence of the assumption of the linear model that the effects of explanatory variables are uniform throughout their range at the individual level. Given this assumption, it is clear from equation (6) that the aggregate impact of each variable depends only on the change in its mean from sample to sample. This assumption, and the resultant ease of the mathematical calculations, undoubtedly account for the use of the linear model in contexts where the researcher wishes to discuss level importance. Yet in most of these examples, a nonlinear

specification such as the logit model would have been considered theoretically more plausible at the individual level. We now turn to the development of a similar aggregate level impact strategy for the logit model.

## Aggregate Change under the Logit Specification

Because of the nonlinearity of the logit model, applying the same logic of aggregating individual changes in equations (5) and (6) fails to produce a result with the desirable properties of the linear model. For change in the distribution of any single $x_j$:

$$\Delta \bar{P} = \frac{1}{n}\sum_i \Delta P_i = \frac{1}{n}\sum_i \left( \frac{1}{1 + e^{-\mathbf{b}'(\mathbf{x}_i + \Delta \mathbf{x}_i)}} - \frac{1}{1 + e^{-\mathbf{b}'\mathbf{x}_i}} \right). \tag{7}$$

This equation has two apparent shortcomings. First, computation of the change in the aggregate probability cannot be obtained directly from changes in the means of the independent variables (Markus 1988, note 8). This implies that more information must be employed to weight changes in independent variables. Second, and more important, in this form the change in the aggregate probability cannot be decomposed into individual, additive components for each explanatory variable. This nonadditivity is inherent in nonlinear models.[5] However, we can estimate additive components using a standard tool from multivariate calculus: the Taylor series approximation for the individual change in $P_i$.

In general, the Taylor series can be used to approximate any nonlinear function as a polynomial based on partial derivatives.[6] The general Taylor series formula for the approximate change in $P_i$ for change in a given $x_{ij}$, which we denote $\Delta P_{ij}^*$ is:

---

[5] Attempts to calculate theoretical effects at the individual level from the logit model suffer from similar problems. A common tactic is to compute a base probability from the logit model using the mean of each variable:

$$P^0 = \frac{1}{1 + e^{-\mathbf{b}'\bar{\mathbf{x}}}}$$

and then a series of changes in that probability by adding one unit to each variable in succession:

$$\Delta P = \frac{1}{1 + e^{-(\mathbf{b}'\bar{\mathbf{x}} + b_j \Delta x_j)}} - P^0.$$

Although this "difference model" computes changes in probability that are *exact* for the given choices of initial probability and change in each independent variable, they will not add correctly to the change in probability that occurs when several or all variables are changed simultaneously. Moreover, since this method is based on hypothetical choices of initial probability and amount of change in the independent variable, it does not yield an empirical, sample–specific average, as our measure (below) does.

[6] An exposition of the Taylor series can be found in introductory calculus texts such as Saltz (1977).

$$\Delta P_{ij}^* = \sum_{k=1}^{m} \frac{1}{k!} \frac{dP_i^{(k)}}{dx_{ij}} (\Delta x_{ij})^k \tag{8}$$

where the parenthesized superscript $(k)$ indicates the $k$th partial derivative. For example, a first order expansion $(m = 1)$ of the logit model would use the first derivative for the logit function, defined in equation (4), to produce a linear approximation of the change in $P_i$ at the individual level:

$$\Delta P_{ij}^* = \frac{dP_i}{dx_{ij}}\Delta x_{ij} = b_j P_i(1 - P_i)\Delta x_{ij}. \tag{9}$$

These individual approximations would then be averaged over the sample:

$$\Delta \bar{P}_j^* = \frac{1}{n} \sum_i \Delta P_{ij} = b_j \frac{1}{n} \sum_i P_i(1 - P_i)\Delta x_{ij}. \tag{10}$$

A second-order expansion $(m = 2)$ adds a quadratic term with a second derivative for greater accuracy in approximating the change in $P_i$:

$$\Delta P_{ij}^* = b_j P_i(1 - P_i)\Delta x_{ij} + \frac{1}{2} b_j^2 P_i(1 - P_i)(1 - 2P_i)\Delta x_{ij}^2, \tag{11}$$

which would then be averaged to:

$$\Delta \bar{P}_j^* = b_j\left(\frac{1}{n} \sum_i P_i(1 - P_i)\Delta x_{ij}\right) \\ + \frac{1}{2} b_j^2 \left(\frac{1}{n} \sum_i P_i(1 - P_i)(1 - 2P_i)\Delta x_{ij}^2\right). \tag{12}$$

This expansion could in theory contain an infinite number of terms, but two are adequate for our purposes.[7] We refer to the quantity in equation (12) as our *aggregate logit impact measure*. The impact of changes in several variables would be the sum of the estimated changes for each variable:

$$\Delta \bar{P} = \sum_j \Delta \bar{P}_j^*. \tag{13}$$

This equation is the logit equivalent of equation (6) for linear models and provides an accurate decomposition of equation (7) into effects of individual variables for the logit model.

Equations (9)–(12) express an important point about the nonlinear nature of the logit and similar models. Since the impact of an independent variable's change at the individual level varies, the aggregate or average impact of a shift in the distribution of that variable under *any* nonlinear model will depend on more than the change in its mean. The aggregate impact depends on three separable

---

[7]We shall consider the issue of numerical accuracy in a later section of this paper.

components: (1) the individual-level impact of that variable on the outcome ($b_j$ as interpreted within the given functional form); (2) the distribution of the individual probabilities ($P_i$) prior to any shift (e.g., large or small concentration of "swing" voters with $P_i$ near .5); and (3) the exact distribution of individual changes in the explanatory variable ($\Delta x_{ij}$). Thus, the aggregate impact of change in a nonlinear model is specific to the distribution of each independent variable in the sample in a way that is impossible to capture by the linear model. We illustrate this point in the following section.

## A Numerical Simulation

An extended numerical example will help to illustrate the potential variation in estimates of our aggregate logit impact measure in specific samples. Assume that a particular coefficient in the logit model is $b_j = .5$. Then the aggregate logit impact in equation (12) depends on the distribution of individual probabilities and the distribution of change in the explanatory variable. In our example, we shall compare four different distributions for the individual $P_i$ and three patterns of change in the explanatory variable. We use the highly flexible beta distributional family to generate differently shaped distributions for individual probabilities with the proper zero-one range (Mood, Graybill, and Boes 1974). Our comparison choices for the distributions of initial probabilities were:

    a.  bell-shaped, with mean .5 and standard deviation .2;
    b.  uniform, with mean .5 and standard deviation .289;
    c.  U-shaped, with mean .5 and standard deviation .4;
    d.  skewed, with mean .33 and standard deviation .2.

In a voting context, one might characterize the bell-shaped distribution as representing a competitive distribution of voters and the U-shaped distribution as polarized. Diagrams of these distributions are shown in Figure 1.

We chose three very simple patterns of change in the explanatory variable ($\Delta x_{ij}$):

    1.  all cases increase by exactly one unit;
    2.  increase proportional to $P_i$, averaging one unit;
    3.  increase largest when $P_i = .5$, declining proportionately with distance from .5 in either direction and averaging one unit.

Note that all change models are parameterized to produce a mean change of one unit in the underlying explanatory variable.[8] For each distribution a–d, we take the 100 percentile points as simulated observations,[9] then impute changes in the explanatory factor according to each pattern 1–3, and finally compute the

---

    [8]In the second change model, $\Delta x_{ij} = 2 \times P_i/a$, where $a$ is a normalizing factor to force $\Delta \bar{x}_j = 1$. In the third model, $\Delta x_{ij} = (2 - 4 \times |P_i - .5|)/a$, again with a normalizing factor.
    [9]This is a deterministic, rather than a Monte Carlo, simulation.

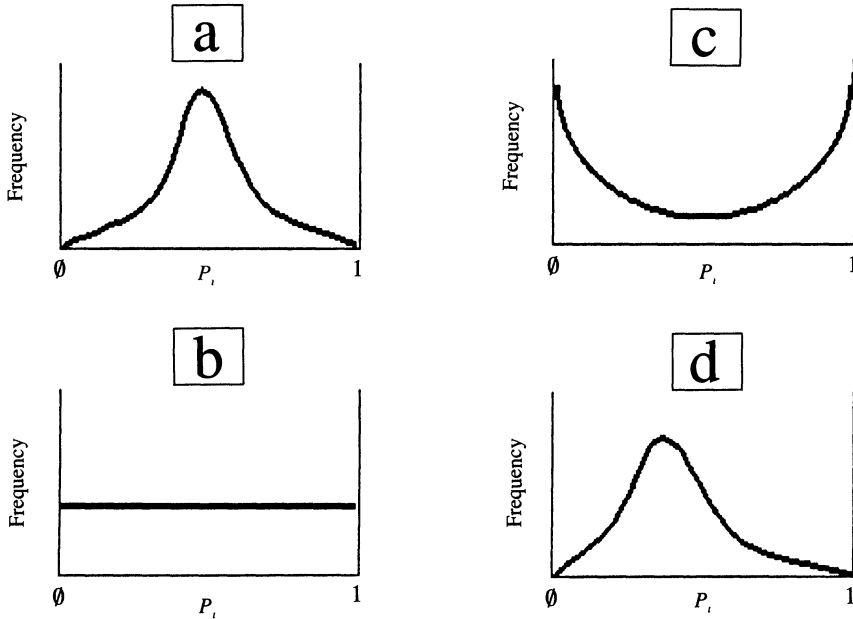**Figure 1. Hypothetical Beta Distributions Describing Prior Probabilities**



**Table 1. Values of Aggregate Logit Impact Measure under Varying
Hypothetical Conditions
(Logit Coefficient $b = 0.5$)**

|  | Change Model | | |
|---|---|---|---|
|  | (1) | (2) | (3) |
|  | Uniform | Proportional | Declining |
| Initial Distribution | Change | to $P_i$ | from $P_i = .5$ |
| (a) bell-shaped ($\mu = .5, \sigma = .2$) | .105 | .099 | .112 |
| (b) uniform ($\mu = .5, \sigma = .3$) | .083 | .075 | .104 |
| (c) U-shaped ($\mu = .5, \sigma = .4$) | .045 | .039 | .096 |
| (d) skewed ($\mu = .33, \sigma = .2$) | .096 | .103 | .110 |

summation in equation (12) for each combination. The results of our compari-
sons appear in Table 1.

We remind the reader that all entries in the last three columns of Table 1 are
estimates of the aggregate change in the outcome probability for an average
change of one unit in the explanatory variable, and that all are premised on the
assumption that the individual-level coefficient from the logit model is the same

in each case ($b_j$ = .5). Given this, the range of estimates is quite impressive, from a low of .039 (the combination of distribution c and change model 2) to a high of .112 (distribution a and change model 3). That is, under these varying conditions the effect of a unit mean change in the explanatory variable raises the aggregate percentage for some outcomes from between 3.9% and 11.2%—a factor of almost three. These are by no means the limits of potential variation. The variations are due to the assumption of the logit model that all individuals are not equally changeable (i.e., those whose probabilities are near .5 are more volatile than those near the extremes). Thus, the extent of aggregate impact depends on whether the initial distribution has a large number of "variable" individuals and on how the changes in the independent variable are distributed among individuals. For example, the column labeled 1 in Table 1 presents the aggregate impact of a uniform change in an explanatory variable in each of the four initial probability distributions. As can be seen, the aggregate effect of .105 in the bell-shaped distribution, which has the highest concentration of individuals near .5 in initial probability, is larger than in the uniform distribution (.083), which is in turn larger than in the U-shaped distribution of initial probabilities (.045). The impact for the skewed distribution is .096, falling as one would guess between the bell-shaped and uniform conditions. These patterns confirm the intuitive hypothesis in voting behavior studies that the heavier the concentration of "swing voters," the greater the impact that change in an explanatory variable should have on the aggregate outcome.

In column 2 in Table 1, change in the explanatory factor is not uniform, but proportional to the initial probability, and hence asymmetric. Since the first three distributions of initial probabilities are symmetric, however, the aggregate logit impacts are only slightly different from those for uniform change (combinations a.2–c.2), usually smaller. Unlike the case for the symmetric distributions, the impact for the skewed distribution (combination d.2) is largest among all distributions for this type of change.

The last column in Table 1 describes the situation where change in the explanatory variable is concentrated toward the middle of the initial probability distribution. For each distribution, the aggregate impact is larger than in previous columns. This is because change in the explanatory factor is concentrated among the more volatile cases near $P_i$ = .5. Otherwise, effects remain in roughly the same rank order as in previous columns.

## A Predictive Evaluation: The 1972 and 1976 Presidential Elections

Having derived our aggregate logit impact measure and demonstrated its sensitivity to various distributional factors, we now turn to a substantive application. We use the National Election Study (NES) panel data from the 1972 and 1976 U.S. presidential elections to illustrate the computation and interpretation of our measure and to compare it with the parallel measure for the linear

probability model in equation (6). A panel data set is especially useful for the purpose of validation because, unlike the usual cross-sectional situation, all the information needed for equation (12), including the actual changes in each independent variable, is known. Thus, we may compare the predictions of the linear and logit model measures against each other and against the actual aggregate change across the panel waves.

In the panel sample, 68.8% of respondents reported voting for Nixon in 1972, while 54.4% reported voting for Ford in 1976.[10] Our task is to account for as much as possible of this 14.4% change in the aggregate level of the Republican vote and to determine which variables contribute most to this overall net change. For illustrative purposes, we simplify our explanatory model to include only three variables prominent in many electoral analyses: the respondent's party identification, presidential job approval (of Nixon in 1972 and Ford in 1976), and the difference in the respondent's "feeling thermometer" ratings of the Republican and Democratic candidates in both election years.[11] All variables were coded so that higher numbers indicated pro-Republican attitudes. Party identification ranged from zero for "strong Democrat" to six for "strong Republican"; presidential approval was a dichotomous variable coded as zero for "disapproval" and one for "approval"; and the thermometer score difference ranged from $-100$ to 100 in strength of feelings toward the Republican candidate relative to his Democratic opponent.

Table 2 shows the means for each variable in 1972 and 1976, and the difference between these figures represents the average change for the variable between the two elections. As can be seen, party identification in 1972 averaged just over three on the seven-point scale and became slightly more Democratic in 1976, while the two other variables were strongly pro-Republican in 1972 and changed more significantly in the Democratic direction by 1976.

What was the contribution of each variable to the overall change in the Republican share of the vote between 1972 and 1976? To answer this question, we first estimate both a linear probability and a logit model to predict the probability of an individual voting Republican in the 1972 panel wave. We then use equations (6) and (12) to calculate, for the linear and logit models respectively, the estimated aggregate impact of the changes in each independent variable observed between the two elections. The results of these calculations are shown in Table 3.

The OLS regression coefficients for the linear probability specification

---

[10]The 54.4% Ford vote is greater than the 50.2% vote reported in the entire 1972–76 panel data set. The discrepancies are due to differential rates of panel attrition and to exclusion from the analysis of all respondents who were missing on any of our explanatory variables.

[11]The variables in the 1972–76 NES panel study (ICPSR 7010) were: party identification— v140 (1972) and v3174 (1976); presidential job approval—v221 (1972) and v3135 (1976); thermometer scores—Nixon, v255; McGovern, v254; Ford, v3299; Carter, v3298.

### Table 2. Changes in Variables Predicting Presidential Vote, 1972–76

| Variable | 1972 Mean | 1976 Mean | Change |
|---|---|---|---|
| Vote | .688 | .544 | − .144 |
| Party ID | 3.009 | 2.976 | − .033 |
| Incumbent approval | .756 | .691 | − .065 |
| Thermometer difference | 22.385 | 3.865 | − 18.52 |

*Source:* NES 1972–76 Panel Study.

### Table 3. OLS and Logit Impact Estimates from the 1972–76 Panel

| Variable | OLS *b* | OLS Impact | Logit *b* | Logit Impact |
|---|---|---|---|---|
| Intercept | .287 | — | − 1.265 | — |
| Party ID | .040 | − .001 | .453 | − .006 |
| Approval | .228 | − .015 | .491 | − .004 |
| Thermometer difference | .0049 | − .091 | .071 | − .101 |
| $R^2$ | .62 | | .70[a] | |
| Explained change | − .107 | | − .111 | |
| Residual change | − .037 | | − .033 | |
| Total change | − .144 | | − .144 | |
| $N = 340$ | | | | |

[a]The squared correlation between the observed (dummy) outcome and the predicted probability.

*Source:* NES 1972–76 Panel Study.

appear in the first column of the table. Following equation (6), we multiply the coefficient of each independent variable by the corresponding change in its mean to obtain its estimated aggregate impact on the probability of the vote for the Republican; these estimates appear in the second column. Thus, in the linear model, party identification contributed − .001 to the − .144 Republican loss between 1972 and 1976 (multiplying the coefficient .040 by the mean change of − .033). Presidential approval contributed − .015; and the thermometer score differences, − .091. Summing these figures yields a total predicted change of − .107, or a 10.7 percentage point decline in the Republican vote. The prediction differs from the actual change by 3.7 percentage points, indicating the change between the two elections attributable to omitted variables and other sources of specification error.

The next two columns of Table 3 are the corresponding figures for the logit model. The logistic regression coefficients for the 1972 election are found in the third column. To compute the aggregate logit impact measure for each variable

according to equation (12), we utilize the logit coefficients, the estimated individual probabilities from the 1972 model, and the individual 1972–76 change scores. The SPSS commands we used to perform these calculations are shown in Appendix A. The logit impact estimates for each independent variable are shown in the fourth column of Table 3.[12]. According to these calculations, changes in party identification contribute $-.006$ ($-0.6\%$) to the vote shift; changes in presidential approval contribute $-.004$; and changes in the thermometer difference measure contribute $-.101$ to the vote shift. The changes in the three variables together yield a total of $-.111$ in the overall predicted difference in the probability of voting Republican in 1976 compared with 1972. This prediction is slightly closer to the actual $-.144$ value than the OLS estimates. More important, the total predicted change calculated from the individual logit estimates exactly reproduces the results of applying equation (7) to calculate the aggregate effect of changing all variables simultaneously, or equivalently, the results of applying the 1972 logit model directly to each individual in the 1976 wave of data and computing the mean change in vote probabilities. Thus, our measure fully captures the *predictable* portion of change in the mean dependent probability.

The choice of specification has a more significant effect on the assessment of different variables' *relative* contributions to change in the level of the dependent variable. While both models indicate that changes in thermometer difference scores accounted for the largest portion of change in the dependent variable, the logit specification indicates that party identification has six times greater impact than in the linear model, while presidential approval has less than one-third the overall impact it does in the linear specification. There is no simple explanation for these differences; we found a quite complicated relationship between predicted prior probabilities and changes in each independent variable at the individual level. We can only reiterate the need for an estimation strategy that does not assume that only a variable's *mean* change is important for predicting its impact on aggregate change in the dependent variable.

### Constructing Predicted Impacts in Practice

In this section, we consider how to estimate aggregate logit impacts in practice, depending on the amount of information available to the researcher. In general, making assessments of aggregate change requires three steps:

1. *Estimate the logit model to obtain coefficients and predicted prior probabilities.* The predictive power of model should be as high as possible, since reasonable predicted probabilities are crucial.

---

[12]Appendix A shows commands only for the party ID and incumbent approval variables. The individual changes in the thermometer scores were so large in some cases that we employed a piecewise procedure in SAS to arrive at the predicted changes in probabilities. We explain this problem and the piecewise procedure more fully in the next section.

2. *Measure changes in the explanatory factors of interest, or select a model for imputing such changes.* Imputation may employ one of the models we used in our hypothetical example, or may imply change only in certain subgroups of a sample.

3. *Compute the impacts as in equation (12).* The ingredients for the calculation are the logit coefficients, predicted initial probabilities and individual changes. The computer commands in Appendix A first compute terms for individual cases as in equation (11); a descriptive procedure is then used to produce the required average impact for the entire sample.

We constructed the preceding example, with its use of data from a second panel wave, in order to *validate* our assessment of impacts against a future aggregate outcome. When actual change data are available, a precise accounting of absolute and relative aggregate impacts may be obtained. In applications using only cross-sectional data, empirical information is available for step 1, but not for step 2. In that case aggregate impact assessments require a model that matches (imputes) individual-level changes in explanatory factors to prior probabilities estimated within the cross-section. A generalization of the method we used in the numerical simulations is contained in the following model:

$$\Delta \hat{x}_i = a_0 + a_1 P_i + a_2 D_i + a_3 P_i D_i, \tag{14}$$

where $P_i$ is an individual's initial probability as before, and $D_i$ is a dummy indicator for membership in a particular subgroup. The choice of particular values of the coefficients will determine the mean change in the given explanatory variable. The coefficient $a_0$ denotes uniform change across the entire sample; $a_1$ denotes change proportional to initial probabilities across all subgroups; $a_2$ denotes uniform change within the subgroup only; and $a_3$ denotes potential interaction between initial probability and subgroup membership. For example, if an explanatory variable changed in the aggregate by five units, a uniform change model would specify $a_0 = 5$ and all other coefficients equal to zero. A five-unit change restricted only to Democratic identifiers, for example, would imply $a_0 = 0$ and $a_2 = 5$. Of course, the imputation model may contain more variables and greater complexity of functional form if desired.[13]
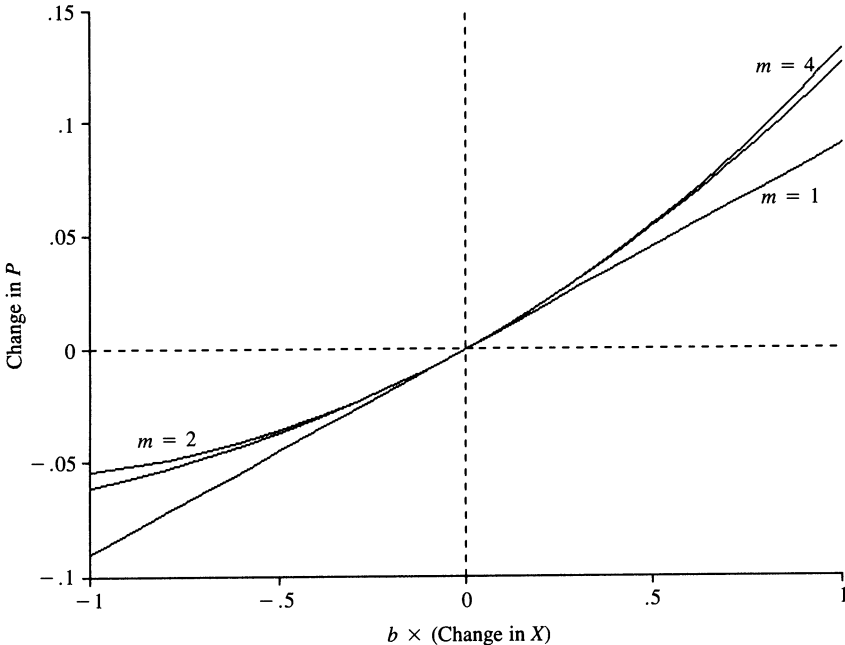
Whenever possible, the researcher should choose a theoretically informed model of change in the independent variable. For example, certain macrolevel events such as war may exert a relatively uniform effect on changes in individual's ratings of presidential performance ($a_0 \neq 0$). In other instances, change

---

[13]Predictions of impact may be calculated in a very rough fashion without *any* individual-level information, by combining published logit coefficients with actual or hypothetical information about both the distributions of prior probabilities (step 1) and changes in independent variables (step 2) gleaned from other sources. The quality of the information that is utilized must be fairly high to make such predictions meaningful.

related to prior probability ($a_1 \neq 0$) is frequently plausible, as when changes in independent variables are related to prior ideological and policy preferences. For example, successful management of the economy may bring an overall increase of approval for an incumbent candidate, but individual increases may be strongly related to prior evaluations and partisan dispositions that condition original propensity to vote for the incumbent. In a more extreme case, *all* of the increases in approval might be hypothesized to occur among individuals who identify with the incumbent's party, or who otherwise have a strong propensity to favor the incumbent ($a_2 \neq 0$). Our general point is that impact assessments should be structured in accordance with the type of variable involved and the most appropriate theoretical model of change.

One final point concerning the accuracy of equation (12) applies to all of the situations we have described in this section. The approximation for change in each individual $P_i$, and hence for the aggregate change, depends on $\Delta x_{ij}$ being "small" relative to the magnitude of the logit coefficient $b_j$. In practice this amounts to the restriction that $b_j \Delta x_{ij}$ be no more than .5. We illustrate this with Figure 2. In this figure, the change in $P_i$ is predicted as a continuous function of $b_j$ times the change in $x_{ij}$, when $P_i$ is equal to 0.1, for the first-, second-, and fourth-order Taylor series approximation ($m = 1$, $2$, and $4$, respectively). We

Figure 2. Predicted Change in Pi for Taylor Series
Approximation of Order m; P_i = .1



$b \times$ (Change in X)

can see that the second- and fourth-order approximations are virtually equal in the range $-.5$ to $.5$. This empirical convergence of the approximation is not, however, formal convergence in the sense that each successive term in the approximation series declines in magnitude for any $\Delta x_{ij}$; for larger changes, the curves begin to diverge rapidly, and the Taylor approximation is only valid as an infinite series. When change in the independent variable exceeds this range restriction, we recommend computing the change in successive pieces, by changing $x_{ij}$ by small increments (at most $.5/b_j$), and cumulating the corresponding increments in $P_i$, recomputing $P_i$ each time using the logit equation (2). We were required to do this in our computations for the impact of thermometer difference scores. The programming algorithm we used in SAS to produce our estimate is presented in Appendix B.

### Discussion: OLS and Logit (and Probit) Compared

One of the major points of our analysis of the logit model is the high level of specificity necessary for comparisons between baseline and new samples. If the correct underlying model for individual outcomes is in fact linear, all aspects of the distribution of the individual $P_i$ and changes in $x_{ij}$, except their means, are irrelevant, as demonstrated in equation (6). If the correct underlying model for individuals is instead the logit model, aggregate impact is not indifferent to these distributions, as we have shown in equation (12) and in the simulations. Much research, such as that cited earlier, has made deliberate use of the linear probability model as a simplifying approximation to an underlying logit model in order to make estimates of aggregate impact. We then should consider under what conditions the linear estimate would essentially duplicate the logit estimate.

One obvious situation where the two methods converge in prediction of aggregate impact is when the linear and logit models are indistinguishable at the individual level. This occurs when all individual dependent probabilities are concentrated in the 30%–70% range, and hence the deviations from linearity predicted by the logit model are trivial. For cases where the two models do differ substantially at the individual level, we explore other possible situations of convergence at the aggregate level by way of the uniform change model. This model provides a bridge to the usual definition of the OLS slope as the *average* effect of changing $x_j$ *by one unit*. Inserting a constant change of $\Delta \bar{x}_j$ for each individual into equation (12) implies that:

$$\Delta \bar{P}_j^* = b_j \frac{1}{n} \sum_i P_i(1 - P_i) \Delta \bar{x}_j$$

$$+ \frac{1}{2} (b_j)^2 \frac{1}{n} \sum_i P_i(1 - P_i)(1 - 2P_i) (\Delta \bar{x}_j)^2. \tag{15}$$

But since $\Delta \bar{x}_j$ is a constant, it can be moved outside both summations. Without this source of variation, the second summation is equal to zero whenever the $P_i$

are distributed symmetrically around .5. In that case, the previous equation reduces to:

$$\Delta \bar{P}_j^* = \Delta \bar{x}_j \frac{1}{n} \sum_i b_j P_i (1 - P_i) = \Delta \bar{x}_j \frac{1}{n} \sum_i \frac{dP_i}{dx_{ij}}. \tag{16}$$

The last term in this equation is the average linear effect of a small increase in $x_{ij}$ for the logit model, which would be the interpretation of the slope from the linear probability model. So under these conditions, the linear estimate of aggregate impact would be an adequate approximation of the logit measure. We reiterate, however, that this connection between the linear impact and the logit estimate for uniform change relies on three distinct restrictions: (1) a symmetric distribution of initial $P_i$; (2) averaging the nonlinear effects of $x_j$ over the *observed* sample of $P_i$; and (3) uniform change in $\Delta x_{ij}$. Each of these restrictions implies a different cost in substituting a linear approximation for the logit model.

The condition of symmetric $P_i$ can, of course, be verified in any particular analysis. The linear impact estimator can still come close to the logit estimator if the deviations from symmetry are minor, or if a different change model somehow compensates by chance. Both our simulations and our example, however, indicate that the differences may indeed be substantial. We know of no method to determine the degree of closeness between the two estimators a priori, but it seems reasonable to assume that whatever calculations would need to be done would be at least as involved as simply calculating the logit estimates from equation (12) directly. The sample specificity implied by condition (2) reiterates the well-known point that the linear approximation to the logit model depends critically on the underlying variability of individual probabilities (Hanushek and Jackson 1977, 185). Hence, a linear model from one symmetric distribution is unlikely to generalize accurately to another, for example, from a unimodal to a more polarized distribution of vote propensity. Aggregate impact assessments derived from linear probability models inherit this sample specificity directly.

We noted in the previous section that the uniform change assumption is not the only, nor necessarily the most plausible, change model from which to consider aggregate cross-sample comparisons. In most empirical situations, uniform change will be the exception rather than the rule. Our analysis of the 1972–76 NES panel data, for example, showed large variations in individual change in the independent variables. In addition, the changes in presidential popularity and candidate thermometer ratings were related to the individual's probability of voting Republican in the first wave. To the extent that individual changes differ widely and are correlated with initial probabilities, then the uniform change model, and hence the linear impact estimates, becomes less and less relevant.

Finally, we address the somewhat similar probit model. Our use of the Taylor series approximation is generally appropriate for any nonlinear model. For the probit model:

$$P_t = \Phi(\mathbf{b}'\mathbf{x}_t) \tag{17}$$

where $\phi(z)$ is the standard normal cumulative distribution function. The equivalent to equation (11), giving individual contributions to the aggregate impact, is:

$$\Delta P_{ij}^* = b_j\Phi(z_i)\Delta x_{ij} - \frac{1}{2}b_j^2 z_i \Phi(z_i)\Delta x_{ij}^2 \tag{18}$$

where $\Phi(z)$ is the standard normal probability *density* function and $z_i = \mathbf{b}'\mathbf{x}_i$. The logit and probit models are known to be highly similar except for the scale of parameters (Aldrich and Nelson 1984, 30–47), and we do not expect substantial differences in impact estimates between the two.

## Conclusion

We hope that readers of this paper will take away three points. First, aggregate impact analysis offers another means of assessing the "importance" of explanatory variables, aside from examination of unstandardized or standardized regression coefficients (Achen 1982; Lewis-Beck and Mohr 1976; Stover 1987). Researchers often are interested in how aggregate outcomes are affected by changes in the distributions of individual independent variables. We argue here that such questions imply the comparison of an observed sample with an alternative sample, either hypothetical or actual.

Second, aggregate impact estimates based on a Taylor series expansion of the logit model are available and intuitively meaningful, if more complicated to calculate than their linear counterparts. In order to estimate the impact of explanatory variables on aggregate shifts in the dependent variable in logit models, the researcher must take into account the variation in individual change in an independent variable and the relationship between change and initial probabilities. If the procedure is applied to cross-sectional data, the method can be used to make out-of-sample predictions of the aggregate shift in the dependent variable by specifying a theoretically meaningful pattern of individual-level change in an independent variable. It may then be used to assess the relative importance of different independent variables by estimating the change in dependent probabilities accounted for by each variable. If this procedure is applied to panel data, all the necessary information on change in independent variables is observed.

Third, researchers should not generally employ linear probability models to approximate aggregate impact in contexts when logit models are substantively more sensible at the individual level. In those cases, the linear model can be expected to approximate the logit estimates accurately when individual probabilities are symmetrically distributed *and* individual change in independent variables is uniform across the distribution of initial probabilities. We think this combination of circumstances will be rare in practice.

While our exposition has focused on electoral behavior, this type of analysis

may be applied whenever aggregate outcomes are of substantive interest, as in social demography or cross-national comparative analysis. Our approach can also be expanded beyond cross-sectional and two-wave panel analyses. Aggregate impact analysis may be even more fruitfully exploited in multiwave pooled cross-section and other longitudinal designs, where data are available for many intersample comparisons. Development of procedures for such applications should be a focus of future research.

*Manuscript submitted 14 May 1991*
*Final manuscript received 2 December 1991*

## APPENDIX A
### SPSS Algorithm for Computing Aggregate Logit Impact Measure

We assume that the data contain predictor variables P, A, and T, change scores CH_P and CH_A. Probabilities computed from logit model are saved as PRED_Y. Aggregate impacts for party identification (P) and approval (A) are calculated below. A more accurate algorithm is needed for the impact of thermometer differences; that algorithm is given in Appendix B (see text).

```
logistic regression VOTE72 with P A T /save pred (PRED_Y).
* insert logit coefficients into the following equations.
compute P_IMP = .453*PRED_Y*(1-PRED_Y)*CH_P
    + .5*((.453)**2)*PRED_Y*(1-PRED_Y)*(1-2*PRED_Y)*(CH_P**2).
compute A_IMP = .491*PRED_Y*(1-PRED_Y)*CH_A
    + .5*((.491)**2)*PRED_Y*(1-PRED_Y)*(1-2*PRED_Y)*(CH_A**2).
descriptives P_IMP A_IMP.
```

## APPENDIX B
### SAS Algorithm for Computing Piecewise Changes in Probabilities

As in Appendix A, we assume that the data contain predictor variables P, A, and T, change score CH_T. Probabilities computed from logit model have already been saved in this data set, as PRED_Y. In this algorithm, thermometer differences (T) are allowed to change only a little at a time. New probabilities are computed incrementally until the full change due to CH_T is accumulated.

```
data logit2; set nes.logit1;

* logit coefficients;
b_0 = - 1.265; b_p = .453; b_a = .491; b_t = .071;

* piecewise loop for thermometer differences—individual cases;
cdx = 0;                        * cdx = cumulative change in x;
cdp = 0;                        * cdp = cumulative change in probability;

newp = pred_y;                  * set dependent probability = initial probability;

do until (cdx = ch_t);          * repeat loop until x changes entire amount;
   dx = min(ch_t-cdx,.5/b_t);   * change x by at most .5/b;
```

```
cdx = cdx + dx;                          * add current change to cumulative change;
dp = b__t*newp*(1-newp)*dx;              * first term of equation (11);
dp = dp + b__t*dp*(1-2*newp)*dx/2;       * add second term of equation (11);
cdp = cdp + dp;                          * add this increment to cumulative change in prob-
                                           ability;

z = b__0 + b__p*p + b__a*a +             * update logit with new value of x;
    b__t*(t + cdx);
newp = 1/(1 + exp(-z));                  * update probability;
end; * end of loop;

t__imp = cdp;                            * impact is cumulative change over all increments;
run;

* compute mean of individual changes;
proc means mean; var t__imp; run;
```

## REFERENCES

Achen, Christopher H. 1982. *Interpreting and Using Regression*. Beverly Hills: Sage.

Aldrich, John H., and Forrest D. Nelson. 1984. *Linear Probability, Logit and Probit Models*. Beverly Hills: Sage.

Finkel, Steven E. N.d. "Reexamining the 'Minimal Effects' Model in Recent Presidential Elections." *Journal of Politics*. Forthcoming.

Hanushek, Eric A., and John E. Jackson. 1977. *Statistical Methods for Social Scientists*. New York: Academic Press.

Lewis-Beck, Michael S. 1988. *Economics and Elections: The Major Western Democracies*. Ann Arbor: University of Michigan Press.

Lewis-Beck, Michael S., and Lawrence B. Mohr. 1976. "Evaluating the Effects of Independent Variables." *Political Methodology* 3:27–47.

Markus, Gregory B. 1982. "Political Attitudes during an Election Year: A Report on the 1980 NES Panel Study." *American Political Science Review* 76:538–60.

———. 1988. "The Impact of Personal and National Economic Conditions on the Presidential Vote: A Pooled Cross-sectional Analysis." *American Journal of Political Science* 32:137–54.

Miller, Warren E., and J. Merrill Shanks. 1982. "Policy Directions and Presidential Leadership: Alternative Interpretations of the 1980 Presidential Election." *British Journal of Political Science* 12:299–356.

Mood, Alexander M., Franklin A. Graybill, and Duane C. Boes. 1974. *Introduction to the Theory of Statistics*. 3d ed. New York: McGraw-Hill.

Rosenstone, Steven J. 1983. *Forecasting Presidential Elections*. New Haven: Yale University Press.

Saltz, Daniel. 1977. *A Short Calculus: An Applied Approach*. 2d ed. Santa Monica: Goodyear.

Smith, Eric R. A. N. 1989. *The Unchanging American Voter*. Berkeley: University of California Press.

Stokes, Donald E. 1966. "Some Dynamic Elements of Contests for the Presidency." *American Political Science Review* 60:19–28.

Stover, Robert V. 1987. "The Coefficient of Aggregate Impact: A Method of Accounting for Aggregate Change." *Social Science Journal* 24:127–38.