



^b
**UNIVERSITÄT
BERN**

Faculty of Business, Economics and
Social Sciences

Department of Social Sciences

University of Bern Social Sciences Working Paper No. 45

Marginal Odds Ratios: What They Are, How to Compute Them, and Why Sociologists Might Want to Use Them

Kristian Bernt Karlson and Ben Jann

January 31, 2023

<http://ideas.repec.org/p/bss/wpaper/45.html>
<http://econpapers.repec.org/paper/bsswpaper/45.htm>

Marginal Odds Ratios: What They Are, How to Compute Them, and Why Sociologists Might Want to Use Them

Kristian Bernt Karlson*

Ben Jann**

*Department of Sociology, University of Copenhagen, Oester Farimagsgade 5, Building 16, DK-1353 Copenhagen K, Denmark, email: kbk@soc.ku.dk.

**Institute of Sociology, University of Bern, Fabrikstrasse 8, 3012 Bern, Switzerland, email: ben.jann@unibe.ch

This version: January 31, 2023

Word count (all text incl. notes and references): 6,243

Tables: 2

Figures: 1

Appendices: 1

Keywords: odds ratio, logit, logistic model, regression, marginal effects, average marginal effects, confounding, mediation

Acknowledgments: We thank the following for invaluable comments and feedback on the work presented in this paper: Tim Liao, Mike Hout, Rudolf Farys, and Jesper Fels Birkelund, as well as participants at the Hans Schadee Research Methods Center Seminar on November 3, 2022, at Trento University, the Seminar on Analytical Sociology on November 14-17, 2022, at Venice International University, and the 2022 Swiss Stata Meeting on November 18, 2022, at University of Bern. For Kristian Bernt Karlson, the research leading to the results presented in this article has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement no. 851293).

Abstract

As sociologists are increasingly turning away from using odds ratios, reporting average marginal effects is becoming more popular. We aim to restore the use of odds ratios in sociological research by introducing marginal odds ratios. Unlike conventional odds ratios, marginal odds ratios are not affected by omitted covariates in arbitrary ways. Marginal odds ratios thus behave like average marginal effects but retain the relative effect interpretation of the odds ratio. We argue that marginal odds ratios are well suited for much sociological inquiry and should be reported as a complement to the reporting of average marginal effects. We define marginal odds ratios in terms of potential outcomes, show their close relationship to average marginal effects, and discuss their potential advantages over conventional odds ratios. We also briefly discuss how to estimate marginal odds ratios and present examples comparing marginal odds ratios to conventional odds ratios and average marginal effects.

Introduction

Logit models and their ensuing odds ratios form the backbone of much sociological research. Despite their prominence, recent methodological research has brought to sociologists' attention some serious problems in using and interpreting odds ratios (Allison 1999; Mood 2010; Breen, Karlson, and Holm 2018; Bloome and Ang 2022). These problems are rooted in a peculiar property of the logit model: The magnitude of its coefficients changes even if one controls for a third variable that is uncorrelated with the predictor of interest; a property known as noncollapsibility, rescaling, or sensitivity to unobserved heterogeneity. Although sociologists have responded to this challenge in different ways, the reporting of (average) marginal effects implied by a logit model or obtained from a linear probability model is now recommended in the methods literature (Breen, Karlson, and Holm 2018; Mize 2019; Mize, Doan, and Long 2019; Long and Mustillo 2021). Marginal effects are not arbitrarily affected by the error term and yield readily interpretable effects on the probability scale, which to many is more intuitive than a ratio between odds (Cramer 2007; Norton and Dowd 2018).

Marginal effects are also beginning to replace odds ratios as a preferred effect metric in substantive research. This change in practice becomes clear when one considers papers published in the *American Sociological Review* between 2010 and 2021. Upon conducting a search on the ASR website, we found that the term "marginal effect" appeared in 11 papers of which the vast majority (nine) were published between 2016 and 2021. Similarly, "linear probability model" appeared in 16 papers of which the vast majority (13) were published between 2016 and 2021. In contrast, "odds ratio" appeared in 41 papers of which only a minority (nine) were published between 2016 and 2021. Although marginal effects are gaining popularity over odds ratios, they do not necessarily align with much sociological research in which relative inequality is a key concept (e.g., in stratification research, political sociology, medical sociology, or demography). Indeed, many sociologists still prefer the odds ratio

precisely because it is a relative measure, and because it is insensitive to the marginal distribution of the dependent variable (Mare 1981; Erikson and Goldthorpe 1992). In contrast, the magnitude of a marginal effect depends on the distribution of the binary outcome (Mare 1981:76; Holm, Ejrnæs, and Karlson 2015), a property that makes it difficult to directly compare effect sizes among, say, populations with different overall outcome rates.

In this paper, we aim to restore the odds ratio as a relevant effect metric in sociological research by introducing what we term a “marginal odds ratio.” This effect metric has properties similar to the properties of marginal effects, including being unaffected by noncollapsibility, but it retains the relative effect interpretation. It thus presents itself as a viable alternative or complement to the reporting of marginal effects.¹ Drawing on work in statistics and epidemiology on this topic (e.g., Zhang 2008; Pang, Kaufman, and Platt 2016; Daniel, Zhang, and Farewell 2021), we first define the marginal odds ratios in the potential outcomes framework.² This framework makes clear the marginal odds ratio estimand and shows its close relationship to the average marginal effect estimand. We then explain how the marginal odds ratio should be interpreted as a population-averaged effect, and how this interpretation differs from the conventional odds ratio typically obtained from a logit model that has a conditional interpretation. We then go on to briefly outline how to estimate the marginal odds ratio using counterfactual predictions from a logit model. We also present two examples demonstrating the versatility of the marginal odds ratio. In a companion technical paper, we give a thorough technical introduction to estimation approaches and introduce software that makes the estimation of marginal odds ratios straightforward (Jann and Karlson 2023).

¹ Being a ratio, the marginal odds ratio can take on values between zero and infinity, and a value of one means that there is no effect. To obtain a symmetric measure (with zero corresponding to a null effect), one could take the log. Although sociologists employ both the odds ratio and the log odds ratio, we mainly focus on the former in this paper as it is more straightforward to interpret in empirical work.

² Sociologists have also recently begun discussing marginal odds ratios (see Erikson et al. 2005; Breen, Karlson, and Holm 2018:46; Kuha and Mills 2020:521-522; Karlson, Popham, and Holm 2021). Moreover, there is a well-established literature in statistics on this topic for clustered or “multilevel” data (see, e.g., Zeger, Liang, and Albert 1988; Agresti 2002). We also draw on these literatures.

Marginal Odds Ratios

We define the marginal odds ratio using potential outcomes notation (Rubin 1974). This notation makes clear the estimand and its close relationship to the average marginal effect. In our exposition, we exclusively focus on binary treatments and refer readers interested in the extension to continuous treatments to our technical paper (Jann and Karlson 2023). Let Y_t be the potential outcome of an individual receiving either the treatment ($T = 1$) or the control ($T = 0$). Comparing Y_t for the treated (Y_1) or untreated (Y_0) is informative about the effect of T on Y . Scholars are often interested in the average treatment effect, which is defined as $E[Y_1] - E[Y_0]$, i.e., the difference in the expectation over each potential outcome. For binary outcomes, the expectation equals the probability of success, meaning that the average treatment effect equals an average marginal effect defined as

$$\text{AME} = \Pr[Y_1 = 1] - \Pr[Y_0 = 1] \quad (1)$$

The AME is the success probability difference if everyone was treated relative to if everyone was untreated. In a similar vein, we define the *marginal odds ratio* as

$$\text{MOR} = \frac{\text{odds}(\Pr[Y_1 = 1])}{\text{odds}(\Pr[Y_0 = 1])} \quad (2)$$

where $\text{odds}(p)$ stands for $p/(1 - p)$. This odds ratio is the ratio of the odds of success if everyone was treated relative to the odds of success if everyone was untreated.³ The estimands in Equations (1) and (2) involve the same counterfactual quantities but the AME is a probability difference whereas the MOR is a ratio between odds.

The estimands in Equations (1) and (2) can also be expressed as depending on other variables, which we denote \mathbf{X} . For example, applied researchers will often be interested in adjusting for a set of additional covariates if they want to control for potential confounding or are interested in effects for different subpopulations. If we assume that \mathbf{X} has a given

³ This definition only holds under the SUTVA assumption.

distribution in the population, then we can express the conditional success probability given $\mathbf{X} = \mathbf{x}$ as $\Pr(Y_t = 1|\mathbf{X} = \mathbf{x}) = E[Y_t|\mathbf{X} = \mathbf{x}]$. By the law of iterated expectations, we can write the unconditional success probability as

$$\Pr(Y_t = 1) = E_{\mathbf{X}}[\Pr(Y_t = 1|\mathbf{X} = \mathbf{x})] \quad (3)$$

where $E_{\mathbf{X}}$ is the expectation over the distribution of \mathbf{X} . Thus, we can rewrite Equations (1) and (2) as

$$\text{AME} = E_{\mathbf{X}}[\Pr(Y_1 = 1|\mathbf{X} = \mathbf{x})] - E_{\mathbf{X}}[\Pr(Y_0 = 1|\mathbf{X} = \mathbf{x})] \quad (4)$$

$$\text{MOR} = \frac{\text{odds}\{E_{\mathbf{X}}[\Pr(Y_1 = 1|\mathbf{X} = \mathbf{x})]\}}{\text{odds}\{E_{\mathbf{X}}[\Pr(Y_0 = 1|\mathbf{X} = \mathbf{x})]\}} \quad (5)$$

We term the expression in Equation (5) the “adjusted marginal odds ratio”, although it is the same estimand as the marginal odds ratio given in Equation (2). The expression in Equation (5) is useful when estimating marginal odds ratios in substantive research using observational data where confounding is ubiquitous. Given that we refer to Equation (5) as an adjusted MOR, we refer to the Equation (2) as a gross or unadjusted MOR. We may think about them as effects controlling and not controlling for additional and potentially confounding covariates (Karlson, Popham, and Holm 2021). We later give an example showing the difference between the two.

Relationship to the Logit Model

Although we have defined the marginal odds ratio estimand in terms of potential outcomes, sociologists usually obtain odds ratios from a logistic response model. To show the relationship between the two, we first write the unconditional logistic model as

$$\Pr(Y_t = 1) = \text{logit}(\alpha + \delta t) \quad (6)$$

where $\text{logit}(z)$ stands for $\exp(z)/[1 - \exp(z)]$. In this model, the exponent to the treatment logit coefficient, $\exp(\delta)$, has a marginal odds ratio interpretation. However, once we condition on other covariates, \mathbf{X} , the interpretation of the odds ratio changes to a conditional one. To see this, assume that we include \mathbf{X} in the regression equation,

$$\Pr(Y_t = 1|\mathbf{X} = \mathbf{x}) = \text{logit}(\tilde{\alpha} + \tilde{\delta}t + \mathbf{x}\beta) \quad (7)$$

In this model,

$$\text{COR}_{\mathbf{x}} = \frac{\text{odds}[\Pr(Y_1 = 1|\mathbf{X} = \mathbf{x})]}{\text{odds}[\Pr(Y_0 = 1|\mathbf{X} = \mathbf{x})]} = \exp(\tilde{\delta}) \quad (8)$$

is a *conditional* odds ratio (COR), where conditional refers to the effects operating at the subgroup level defined by the covariates in \mathbf{X} . The conditional odds ratio differs from the marginal counterpart adjusting for \mathbf{X} , given by

$$\text{MOR}_{\mathbf{x}} = \frac{\text{odds}\{E_{\mathbf{X}}[\Pr(Y_1 = 1|\mathbf{X} = \mathbf{x})]\}}{\text{odds}\{E_{\mathbf{X}}[\Pr(Y_0 = 1|\mathbf{X} = \mathbf{x})]\}} = \frac{\text{odds}\{E_{\mathbf{X}}[\text{logit}(\tilde{\alpha} + \tilde{\delta} + \mathbf{x}\beta)]\}}{\text{odds}\{E_{\mathbf{X}}[\text{logit}(\tilde{\alpha} + \mathbf{x}\beta)]\}} \quad (9)$$

which has a population-averaged interpretation, i.e., the average population response to changing treatment status (Zeger, Liang, Albert 1988:1050). The COR in Equation (8) is the response of the subgroup defined by the covariates in \mathbf{X} (and the COR is assumed to be constant across those groups as the covariates enter additively on the logit scale). In other words, the COR in Equation (8) and the MOR in Equation (9) refer to different estimands, have different interpretations, and cannot be directly compared (Pang, Kaufman, and Platt 2016; Breen, Karlson, and Holm 2018; Daniel, Zhang, and Farewell 2021; Schuster et al. 2021).

To provide further intuition about the difference between the substantive interpretations of the two estimands, we find it instructive to compare them to the distinction between conditional and unconditional quantile regression (Firpo, Fortin, and Lemieux 2009; Killewald and Bearak 2014). In quantile regression, we are typically interested in modelling a given percentile in an outcome distribution as function of a treatment variable and some potential confounders, thus echoing the setup we have described for the odds ratios. In conditional quantile regression, the percentiles in the *conditional* outcome distribution are modelled (i.e., the percentiles in the residual distribution after netting out the effects of the confounders). Thus, it recovers the treatment effect on a percentile in the outcome distribution among individuals with similar values on the confounders. In unconditional quantile regression, the percentiles in

the *unconditional* outcome distribution are modelled, i.e., how the treatment affects the percentile in the overall or marginal outcome distribution, but still controlling for the potential effects of confounders (via their correlation with the treatment variable).

Like marginal and conditional odds ratios, the two quantile estimands have different interpretations. For example, for the conditional one, researchers would be interested in whether a labor market program affects the 25th earnings percentile within groups with different levels of schooling. For the unconditional one, researchers would be interested in whether the labor market program affects the 25th earnings percentile in the overall earnings distribution, controlling for the possibility that schooling and participation in the labor market program may be correlated (and thus schooling will act as a confounder). If we extend these descriptions to odds ratios, then imagine that we replace the earnings outcome with the binary outcome of obtaining a job (and consider the odds ratio as the preferred effect metric). The conditional odds ratio would then capture the labor market program effect among participants within groups with different levels of schooling, whereas the marginal odds ratio would capture the labor market effect on average in the population, accounting for any sorting into the labor market program on schooling.⁴

Marginal and conditional odds ratios are equally valid estimands and their respective uses should depend on the research question. However, from a mathematical perspective, the difference between them arises from what statisticians call noncollapsibility: “Noncollapsibility of the OR derives from the fact that when the expected probability of outcome is modeled as a nonlinear function of the exposure, the marginal effect cannot be expressed as a weighted average of the conditional effects” (Pang, Kaufman, and Platt

⁴ In light of this comparison to quantile regression, an alternative term to “marginal” would be “unconditional.” However, we adopt the former term because this terminology is already established in the literature (Stampf et al. 2010).

2016:1926). Indeed, the key difference between the two estimands is whether the averaging occurs on the log odds scale or on the probability scale.

The mathematical relationship between the two estimands (noncollapsibility) is well-described in the methods literature. From this literature, we highlight two key points. First, from Equations (8) and (9), we see that the MOR will differ from the COR whenever $\beta \neq 0$ (i.e., if there are other relevant predictors apart from the treatment variable). Moreover, the MOR will be attenuated relative to the COR. For example, if there is just a single covariate X , the relationship between the two can be approximated by

$$\ln \text{MOR}_X = \frac{\ln \text{COR}_X}{\sqrt{1 + 0.35\beta^2 \text{var}(X)}} \quad (10)$$

(Zeger, Liang, and Albert 1988:1054).⁵ Whenever $\beta = 0$, the COR collapses to the MOR. Whenever $\beta \neq 0$, i.e., if the adjusting covariate has a non-zero effect on the outcome, the COR will be larger than the MOR *even if* X is not confounding the treatment effect.⁶ Moreover, the attenuation of the MOR relative to the COR depends on the magnitude of β and the dispersion in X . In the example we later provide, we demonstrate how the attenuating effect of noncollapsibility operates.

Second, while there is only one MOR, there are in principle an infinite number of CORs. The interpretation of the conditional odds ratio will depend on the covariates included in the regression equation as it refers to effects specific to subgroups defined by those covariates. Each of these CORs is not directly comparable to the other CORs. In practical terms, whenever researchers successively add variables to a logit regression equation—a widespread practice in sociological research—the COR estimand changes and so coefficients of the treatment of

⁵ This approximation assumes that X is normally distributed, and the number 0.35 is the approximation of $[16\sqrt{3}/(15\pi)]^2$ from the expression derived in Zeger, Liang, and Albert (1988:1054). We obtain a similar approximation if we formulate the logit model in terms of an underlying latent variable model (see Breen, Karlson, and Holm 2018).

⁶ This situation is sometimes referred to as rescaling bias (Karlson, Holm, and Breen 2012; Breen, Karlson, and Holm 2013).

interest are not directly comparable. Karlson, Holm, and Breen (2012) suggested solving this issue by holding one underlying COR estimand constant (what they refer to as the full model including all covariates) and then changing the set of conditioning control variables using residualized predictors (what they refer to as a reduced model). For this reason, the method by Karlson, Holm, and Breen (2012) recovers a COR estimand (Karlson, Popham, and Holm 2021).

Why Should Sociologists Use Marginal Odds Ratios?

Because MOR and COR both are valid estimands, there is no *a priori* argument for choosing one over the other. However, although we agree with the general point that the choice of estimand should depend on the research question, for most practical purposes we find that the MOR estimand is superior to COR estimands. We highlight four reasons. First, the MOR has an interpretation equivalent to an average marginal effect on the probability scale: It is a population-average effect focusing on the average “population response” to a treatment of interest. Given the increasing reporting of average marginal effects in sociological research, the MOR presents itself as a notable alternative or complement to the reporting of AMEs. Second, because MORs are unaffected by noncollapsibility, they can be used for comparing coefficients from same-sample models including different covariates (i.e., for mediation analyses or effect decompositions). Third, MORs are straightforward to compare across different studies or populations as their magnitude does not depend in arbitrary ways on the conditioning set (i.e., set of control variables). Fourth, because many COR estimands exist (depending on the conditioning set), but only one MOR estimand, researchers are free from presenting arguments for why a specific COR estimand is more interesting than another. Such arguments would require highly developed theoretical frameworks which are rare in sociology.

Estimating Marginal Odds Ratios

Marginal odds ratios can be estimated in different ways. In a technical companion paper (Jann and Karlson 2023), we review these approaches and show which estimands each estimation technique recovers. Here, we present an estimation approach based on counterfactual predictions (known as G-computation) and focus on binary treatments for making the exposition as accessible as possible (Robins 1974). This approach compares counterfactual predictions from a (typically parametric) model involving the treatment and conditioning covariates (Zhang 2008). The approach proceeds in four steps:

- (1) Regress Y on T and \mathbf{X} using a logit model (or, in principle, any other model).
- (2) Generate two sets of predictions of the success probability for each observation in the data, one setting everyone to be treated, $T = 1$, and one setting everyone to be untreated, $T = 0$ (i.e., $\hat{p}_{i,T=1}$ and $\hat{p}_{i,T=0}$, where i indexes observations).
- (3) Average each set of predictions to obtain the marginal or population-averaged success probabilities if treated ($\bar{p}_{T=1}$) and if untreated ($\bar{p}_{T=0}$), respectively.
- (4) To obtain the marginal odds ratio, plug in the average marginal predictions from step 3 into the formula for the marginal odds ratio,

$$\widehat{\text{MOR}} = \frac{\text{odds}(\bar{p}_{T=1})}{\text{odds}(\bar{p}_{T=0})} \quad (11)$$

G-computation is a straightforward way of obtaining marginal odds ratios. This approach is available in the user-written Stata command *lnmor*, which we present in our companion paper (Jann and Karlson 2023).⁷ It is also worth noting that the average marginal effect can be obtained by the four steps outlined above, except that in the fourth step, one plugs in the average

⁷ Stata command *lnmor* also supports continuous treatments and provides consistent standard errors.

marginal predictions into the estimand formula for AMEs, i.e., $\bar{p}_{T=1} - \bar{p}_{T=0}$. From this estimation perspective, the close relationship between AME and MOR also becomes apparent.

Examples

Academic Ability and Intergenerational College Mobility

Stratification scholars are interested in quantifying the extent to which academic abilities explain or mediate family background inequalities in educational attainment (Boudon 1974; Erikson et al. 2005; Jackson 2013). Gaps by family background in educational attainment that operate independently of demonstrated academic abilities are theorized to represent the “secondary effects” of family background, i.e., how class origin-based aspirations, preferences, and outlooks feed into educational decisions over and above those difference that come about through unequal skill levels. In our example, we examine the extent to which academic ability (as measured by a cognitive test) accounts for the gap in college attainment between children born to parents with and without a college degree. To fully illustrate the difference between MOR and COR, we conduct this analysis on representative samples from the United States and Denmark. Comparing the United States and Denmark is substantively interesting because, for the birth cohorts we analyze here (born in the mid-1950s through the mid-1960s), Denmark was a more educationally mobile country than the United States (Landersø and Karlson 2021).

For the United States, we analyze data from the National Longitudinal Survey of Youth 1979, which follows a national probability sample of children aged 14 through 21 in 1979 (Bureau of Labor Statistics 2019). For Denmark, we examine data from the Danish National Longitudinal Survey of Youth, which follows a national-probability sample of 7-graders in 1968/1969 (Hansen 1995). Both datasets provide information on parental college attainment,

respondent college attainment (as adults), and a standardized cognitive ability test.⁸ Our analytical strategy is straightforward: We compare the unadjusted or gross gap by parental college attainment to the one adjusted for academic ability.

Table 1 shows the results. We find that the unadjusted or gross marginal odds ratio is about twice as large in the United States (7.6) as in Denmark (3.8), meaning that Denmark is significantly more educationally mobile (row 1). This estimate has a marginal or population-averaged interpretation as the “population response in child college attainment” to changing from non-college to college-educated parents. Adjusting for academic ability (row 2), we find that the marginal odds ratio reduces to 2.5 in both countries. We interpret this adjusted marginal odds ratio as the impact, *on average in each population*, of parental college attainment on child college attainment, accounting for the unequal distribution of academic abilities across family background. Because the adjusted MOR reduces to the same number (2.5), it means that the “secondary effects” of family background are of similar magnitude in the two countries. Moreover, because the unadjusted MOR is much larger in the U.S. than in Denmark, it means that academic ability “mediates” a significantly larger portion of the gap by parental college attainment in college completion in the U.S. than in Denmark.⁹

[TABLE 1 ABOUT HERE]

In contrast to the adjusted marginal odds ratio, the conditional counterpart in row 3 is 3.4 for the United States and 2.9 for Denmark. Thus, had we been using this adjusted COR for comparing the two countries, we would have concluded that, net of academic ability, Denmark is a (albeit only slightly) more educationally mobile country. The adjusted COR has an interpretation that is different from the marginal counterpart: It is the odds ratio for groups with

⁸ In the replication package for this article, we provide the Stata code used for generating the results in this analysis, including the recoding of variables. NLSY79 is available from the Bureau of Labor Statistics; DLSY is available from the Danish National Archives. The final NLSY sample is 10,068; the final DLSY sample is 2,185.

⁹ In log odds ratios, academic ability explains 56 and 33 percent of the gap in United States and Denmark, respectively.

similar levels of demonstrated academic ability, and it does not refer to population-level effects. Moreover, the difference between the adjusted MOR and COR points to the attenuating impact of noncollapsibility. This impact is significantly larger in the United States than in Denmark, meaning that academic ability is a much stronger predictor of college attainment in the United States than in Denmark (net of parental college attainment).¹⁰

Using the Karlson-Holm-Breen (KHB) approach, we present the unadjusted COR in row 4 (Karlson, Popham, and Holm 2021).¹¹ The KHB approach holds constant the COR estimand (i.e., the COR within groups of children with different ability levels). We make three observations regarding this estimand. First, had we used this COR for comparing the two countries' gross mobility levels, Denmark would be almost three times as mobile by this measure (compared to twice as mobile with the unadjusted MOR). Had we adjusted for additional covariates (e.g., aspirations) that also have more predictive power in the United States than in Denmark, this ratio would only grow (and the estimand would change). As we stated earlier, CORs are valid estimands but as there are an infinite amount of them (depending on the conditioning set) and sociological theory rarely is sufficiently detailed to make informed choices about which COR is the better one, it is difficult to argue for choosing one over the other. The MOR does not have this property (there is only one estimand) and so appears to be the best choice for any initial comparison of mobility levels in this example.

Second, as is the case with MORs, the unadjusted COR can be compared to the adjusted COR to gauge mediation. Here we find that the percent mediated is virtually identical to that based on MORs, indicating that conclusions about mediation are similar using MORs or the

¹⁰ By the “strength of the predictor,” we refer to $\beta^2 \text{var}(X)$ in Equation 10; that is, both the impact of ability and the dispersion in ability affect the degree of attenuation. Because academic ability is a latent variable, we cannot meaningfully disentangle the two here. Had we controlled for a variable with a natural metric (e.g., parental income or number of books in the home), we could have decomposed the attenuating impact into the contribution of each of these two components.

¹¹ The KHB approach uses residualized control variables to make constant the scales of the coefficients across logit models with different covariates (see Karlson, Holm, and Breen 2012).

KHB approach (Karlson, Popham, and Holm 2021).¹² Third, comparing the unadjusted COR to the unadjusted MOR is informative about the impact of noncollapsibility (as were comparing the adjusted odds ratio counterparts). Here again we find that the bias stemming from noncollapsibility is larger in the United States than in Denmark (resulting from academic ability being a stronger predictor in the United States).

In rows 5 and 6, we also present average marginal effects, i.e., the probability difference estimand [cf. Equations (1) and (4)]. Similar to the odds ratios, we find that the unadjusted AME is larger in the United States (43 percentage points) than in Denmark (31 percentage points). Moreover, adjusting for academic ability significantly reduces the AMEs and, in relative terms, to an extent similar to the reductions seen in log odds ratios.¹³ However, the adjusted AME is now *smaller* in the United States than in Denmark (about 14 percent), pointing to the well-known “flipped-signs phenomenon” of interaction terms depending on the scale of measurement (Bloome and Ang 2022). Still, had we relied exclusively on AMEs, we would have concluded that the “secondary effects” of family background are (slightly) larger in Denmark than in the United States, a conclusion that would run counter to the conclusion based on the MOR (similarity) and, in particular, the COR (opposite country difference).

Trends in the College Gap in Attitudes towards Racial Segregation

Political sociologists are interested in how schooling shapes attitudes. We examine the gap in attitudes towards racial segregation between respondents with and without a four-year college degree. In particular, we study whether this gap has changed over two decades, focusing on the results based on average marginal effects (absolute gaps) and marginal odds ratios (relative gaps) when we also control for a range of other covariates. We examine data from the General

¹² For the CORs reported in Table 1, in log odds ratios academic ability explains 54 and 33 percent of the gap in the United States and Denmark, respectively.

¹³ For the United States, the percent explained is 59 percent; for Denmark, 35 percent.

Social Surveys cumulative file (Smith et al. 2019), focusing here on the years 1976 through 1996 when information on attitudes towards racial segregation was collected. Our outcome variable is the response to a question about whether white people have a right to keep black people out of their neighborhoods if they feel like it (and that black people should respect that right). We collapse the outcome variable into a binary variable indicating agreement (1) or disagreement (0) with the stated opinion. We measure college attainment as having completed at least 16 years of schooling. Moreover, we include additional covariates, including survey year (for studying trends from 1976 through 1996), age, gender, race, marital status, and a generic 7-point political views variable indicating whether the respondent thinks of him- or herself as liberal (1) or as conservative (7). The final sample with valid information on all variables comprises 12,239 respondents.¹⁴

In this example, we are not interested in quantifying the degree of confounding but merely in summarizing the trends in the college gap net of other factors. We specify a logit model in which calendar year is fully interacted with the college dummy and all other covariates (allowing for the effects of the covariates to change over time). We estimate the model specified both as a linear probability model and as a logit model. For all models, we derive average marginal effects and marginal odds ratios evaluated at calendar years 1976, 1981, 1986, 1991, and 1996, and report these implied quantities in Table 2 (estimates of the coefficients in the underlying regression models are available in the Appendix).

[TABLE 2 ABOUT HERE]

The main finding in Table 2 is that the absolute college gap (as measured by average marginal effects) in the attitude towards racial segregation has declined significantly over the 20-year period, whereas the relative gap (as measured by marginal odds ratios) has remained

¹⁴ In the replication package for this article, we provide the final sample of the GSS used in this example and Stata code for generating the results in this analysis, including the recoding of variables.

unchanged. In 1976, the absolute college gap net of the other covariates is around 21 percentage points on average (for both AMEs implied by the linear probability model and the logit model, respectively), suggesting that college educated individuals were much less likely than individuals without a college degree to support racial segregation. In 1996, the absolute gap is reduced to about 7 or 8 percentage points on average, pointing to a decline of around 60 percent in just 20 years. This decrease is also highly statistically significant. By way of contrast, the relative gap implied by the marginal odds ratios is virtually constant (we can detect a minor change in the odds ratio towards 1, but this trend is not statistically significant).

Thus, while both average marginal effects and marginal odds ratios point to a substantial college divide in attitudes towards racial segregation (with non-college educated being more supportive of this opinion)—even when we control for potentially “confounding” variables—they disagree on the trend in this gap. To see why this is the case, we report in Figure 1 the average marginal predictions from the logit model by college attainment and survey year. From this figure, we can easily see that the absolute gap reduces over time because there is a general decline in support of racial segregation; this decline is steeper among non-college educated in absolute terms because they start at a higher level than the college educated. However, the relative difference does not change much, resulting in constant odds ratios.¹⁵ In conclusion, the support for racial segregation declined steadily over the period in question, resulting in a decline in the absolute college gap, but the relative difference between non-college and college did not change.

[FIGURE 1 ABOUT HERE]

¹⁵ We also calculated risk ratios with this outcome definition (agreeing to the opinion) and the results are identical to those based on odds ratios, indicating that the relative gap has remained constant over time.

Discussion

We have introduced to sociologists the marginal odds ratio, an odds ratio that “behaves like” the increasingly popular average marginal effect (on the probability scale): The marginal odds ratio is unaffected by noncollapsibility, has a population-averaged interpretation, and is “derived from” a given model. We have demonstrated the close relationship between the marginal odds ratio and average marginal effects, and we have outlined why we believe that marginal odds ratios should be preferred over conditional odds ratios in many areas of sociology.

In addition to introducing to sociologists the marginal odds ratio as a complement to the reporting of average marginal effects, our defining the marginal odds ratio in terms of potential outcomes also highlights the crucial distinction between estimands and estimation (Lundberg, Johnson, and Stewart 2021). Many sociologists think of odds ratios as the exponentiated coefficients from logistic response models and have been trained in interpreting these coefficients as if they behave like coefficients from linear regression models. By separating the estimands from their estimation, as we do in this paper, we hope to contribute to sociologists being more precise about the quantities they are interested in estimating.

We have also presented empirical examples to illustrate the uses and interpretation of marginal odds ratios relative to the conditional counterparts or the average marginal effect. Although these examples are stylized, they represent types of analyses that are widespread in mainstream sociology. We show how overall conclusions can depend on the chosen estimand. As all of the estimands are equally valid from a statistical perspective, the choice should depend on the research question. For the examples we provided, the marginal odds ratio appears as an obvious candidate. However, in most applied research, it will be useful to report and interpret estimates of several estimands. In particular, reporting both average marginal effects and

marginal odds ratios could be very informative about absolute and relative differences, even if it results in the “flipped-signs phenomenon” (Bloome and Ang 2022).

For readers interested in an in-depth description of estimation techniques and the estimands they each recover, including user-written Stata software that implements the discussed methods, we refer to our technical companion paper (Jann and Karlson 2023). In the replication package for this paper, we share code and sample data that reproduce the two examples reported earlier. We hope that these tools will urge sociologists to consider using the marginal odds ratio and reporting it as a complement to average marginal effects in substantive research.

References

- Agresti, A. 2002. *Categorical Data Analysis*. New York: John Wiley & Sons.
- Allison, P.D. 1999. "Comparing logit and probit coefficients across groups." *Sociological Methods & Research* 28(2):186-208
- Breen, R., Karlson, K.B. and Holm, A. 2013. "Total, direct, and indirect effects in logit and probit models." *Sociological Methods & Research* 42(2):164-191.
- Breen, R., Karlson, K.B. and Holm, A. 2018. "Interpreting and understanding logits, probits, and other nonlinear probability models." *Annual Review of Sociology* 44:39-54.
- Bloome, D. and Ang, S. 2022. "Is the Effect Larger in Group A or B? It Depends: Understanding Results From Nonlinear Probability Models." *Demography* 59(4):1459-1488.
- Boudon, R. 1974. *Education, opportunity, and social inequality: Changing prospects in western society*. New York: John Wiley and Sons.
- Bureau of Labor Statistics, U.S. Department of Labor. 2019. *National Longitudinal Survey of Youth 1979 cohort, 1979-2016 (rounds 1-27)*. The Ohio State University. Columbus, Ohio: Center for Human Resource Research (CHRR).
- Cramer, J.S. 2007. "Robustness of logit analysis: Unobserved heterogeneity and mis-specified disturbances." *Oxford Bulletin of Economics and Statistics* 69(4):545-555.
- Cummings, P. 2009. "The relative merits of risk ratios and odds ratios." *Archives of Pediatrics & Adolescent Medicine* 163(5):438-445.
- Daniel, R., Zhang, J. and Farewell, D. 2021. "Making apples from oranges: Comparing noncollapsible effect estimators and their standard errors after adjustment for different covariate sets." *Biometrical Journal* 63(3):528-557.
- Erikson, R. and Goldthorpe, J.H. 1992. *The constant flux*. Oxford: Oxford University Press.
- Erikson, R., Goldthorpe, J.H., Jackson, M., Yaish, M. and Cox, D.R. 2005. "On class differentials in educational attainment." *Proceedings of the National Academy of Sciences* 102(27):9730-9733.
- Firpo, S., Fortin, N.M. and Lemieux, T. 2009. "Unconditional quantile regressions." *Econometrica* 77(3):953-973.
- Hansen, E.J. 1995. *En generation blev voksne*. København: SFI.
- Holm, A., Ejrnæs, M. and Karlson, K. 2015. "Comparing linear probability model coefficients across groups." *Quality & Quantity* 49(5):1823-1834.
- Jackson, M. 2013. *Determined to succeed? Performance versus choice in educational attainment*. Stanford: Stanford University Press.
- Jann, B. and Karlson, K.B. 2023. "Estimation of Marginal Odds Ratios." Working Paper available at repec. Web: <https://ideas.repec.org/p/bss/wpaper/44.html>.

- Karlson, K.B., Holm, A. and Breen, R. 2012. "Comparing regression coefficients between same-sample nested models using logit and probit: A new method." *Sociological Methodology* 42(1):286-313.
- Karlson, K.B., Popham, F. and Holm, A. 2021. "Marginal and Conditional Confounding Using Logits." *Sociological Methods & Research*, published in advance online: <https://doi.org/10.1177/004912412199554>.
- Killewald, A. and Bearak, J. 2014. "Is the motherhood penalty larger for low-wage women? A comment on quantile regression." *American Sociological Review* 79(2):350-357.
- Kuha, J. and Mills, C. 2020. "On group comparisons with logistic regression models." *Sociological Methods & Research* 49(2):498-525.
- Karlson, K. and Landersø, R. 2021. "The making and unmaking of opportunity: Educational mobility in 20th century-Denmark." *IZA Discussion Paper Series*, IZA DP No. 14135. Web: <https://www.iza.org/publications/dp/14135/>
- Long, J.S. and Mustillo, S.A. 2021. "Using predictions and marginal effects to compare groups in regression models for binary outcomes." *Sociological Methods & Research* 50(3):1284-1320.
- Lundberg, I., Johnson, R. and Stewart, B.M. 2021. "What is your estimand? Defining the target quantity connects statistical evidence to theory." *American Sociological Review* 86(3):532-565.
- Mare, R.D. 1981. "Change and stability in educational stratification." *American Sociological Review* 46(1):72-87.
- Mize, T.D. 2019. "Best practices for estimating, interpreting, and presenting nonlinear interaction effects." *Sociological Science* 6:81-117.
- Mize, T.D., Doan, L. and Long, J.S. 2019. "A general framework for comparing predictions and marginal effects across models." *Sociological Methodology* 49(1):152-189.
- Mood, C. 2010. "Logistic regression: Why we cannot do what we think we can do, and what we can do about it." *European Sociological Review* 26(1):67-82.
- Norton, E. and Dowd, B.E. 2018. "Log Odds and the Interpretation of Logit Models." *Health Services Research* 52(2):859-878.
- Pang, M., Kaufman, J.S. and Platt, R.W. 2016. "Studying noncollapsibility of the odds ratio with marginal structural and logistic regression models." *Statistical Methods in Medical Research* 25(5):1925-1937.
- Robins J. 1986. "A new approach to causal inference in mortality studies with a sustained exposure period — application to control of the healthy worker survivor effect." *Mathematical Modelling* 7(9–12):1393–1512.
- Rubin, D.B. 1974. "Estimating causal effects of treatments in randomized and nonrandomized studies." *Journal of Educational Psychology* 66(5):688-701.

Smith, T.W., Davern, M., Freese, J., and Morgan, S.L.. 2019. *General Social Surveys, 1972-2018 [machine-readable data file] /Principal Investigator, Smith, Tom W.; Co-Principal Investigators, Michael Davern, Jeremy Freese and Stephen L. Morgan; Sponsored by National Science Foundation*. Chicago: NORC.

Schuster, N.A., Twisk, J.W., Ter Riet, G., Heymans, M.W. and Rijnhart, J.J. 2021. "Noncollapsibility and its role in quantifying confounding bias in logistic regression." *BMC Medical Research Methodology* 21(1):1-9.

Stampf, S., Graf, E., Schmoor, C. and Schumacher, M. 2010. "Estimators and confidence intervals for the marginal odds ratio using logistic regression and propensity score stratification." *Statistics in Medicine* 29(7-8):760-769.

Zeger, S.L., Liang, K.Y. and Albert, P.S. 1988. "Models for longitudinal data: a generalized estimating equation approach." *Biometrics* 44(4):1049-1060.

Zhang, Z. 2008. "Estimating a marginal causal odds ratio subject to confounding." *Communications in Statistics - Theory and Methods* 38(3):309-321.

Tables and Figures

Table 1. Odds ratios and average marginal effects of parental college attainment gap in college attainment unadjusted and adjusted for academic ability. The United States and Denmark. Standard errors in parentheses.

	USA (N = 10,068)	DNK (N = 2,185)	USA/DNK
1: MOR: Unadjusted	7.7 (0.46)	3.8 (0.55)	2.03*
2: MOR: Adjusted	2.5 (0.13)	2.5 (0.33)	1.00
3: COR: Adjusted	3.4 (0.23)	2.9 (0.45)	1.17
4: COR: Unadjusted (knb)	14.2 (1.05)	4.9 (0.78)	2.90*
5: AME: Unadjusted	0.43 (0.01)	0.31 (0.03)	1.38*
6: AME: Adjusted	0.17 (0.01)	0.20 (0.20)	0.86

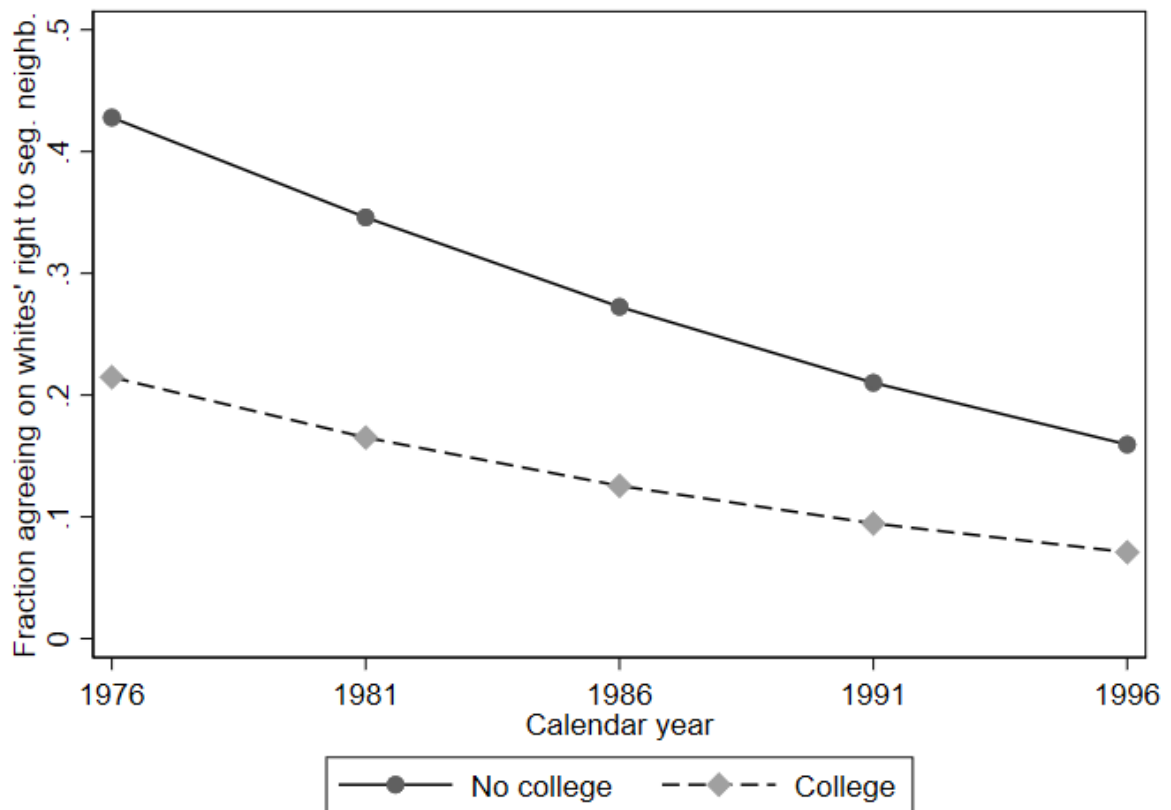
Note: MOR is marginal odds ratio; COR is conditional odds ratio; AME is average marginal effect; knb is the Karlson-Holm-Breen decomposition method (using orthogonalized predictors). US data are from the NLSY79; the Danish data are from the Danish Longitudinal Survey of Youth. * indicates that the country difference in log odds ratios is statistically significant at a five percent level.

Table 2. Average marginal effects and marginal odds ratios of the college gap in the attitude toward racial segregation in 1976, 1981, 1986, 1991, and 1996. Standard errors in parenthesis.

	AME_{LPM}	AME_{LOGIT}	MOR	lnMOR
1976	-0.215 (0.019)	-0.213 (0.021)	0.366 (0.043)	-1.006 (0.118)
1981	-0.180 (0.013)	-0.181 (0.012)	0.374 (0.030)	-0.983 (0.079)
1986	-0.146 (0.009)	-0.147 (0.008)	0.383 (0.025)	-0.959 (0.065)
1991	-0.111 (0.011)	-0.115 (0.009)	0.393 (0.035)	-0.934 (0.089)
1996	-0.076 (0.017)	-0.088 (0.010)	0.403 (0.053)	-0.909 (0.130)
<i>1976–1996 difference</i>	0.138 (0.031)	0.125 (0.028)	-	0.097 (0.211)
<i>1976–1996 prop. reduction</i>	64.4% (9.6)	58.5% (7.7)	-	9.7% (20.0)

Note: MOR is marginal odds ratio; LPM is linear probability model; AME is average marginal effect. Estimates are adjusted for gender, race, age, marital status, and overall political view. Data are from General Social Surveys Cumulative File, N = 12,239.

Figure 1. Trends among college and non-college educated in the attitude towards whites' right to live segregated from blacks, 1976–1996. Average marginal predictions.



Note: Estimates are adjusted for gender, race, age, marital status, and overall political view. Data are from General Social Surveys Cumulative File, N = 12,239.

Appendix A

Appendix Table A1. Linear probability and logit models of the college gap in the attitude toward racial segregation, 1976–1996.

	Linear Probability Model		Logit Model	
	B	SE	B	SE
year	-0.0128	0.0034	-0.1113	0.0208
college	-0.2146	0.0193	-1.0500	0.1222
age	0.0043	0.0006	0.0172	0.0033
<i>race (ref. white)</i>				
black	-0.2556	0.0296	-0.1341	0.2203
other	0.0114	0.0637	0.0864	0.3619
<i>marital (ref. married)</i>				
widowed	0.0702	0.0274	0.2756	0.1410
divorced	-0.0002	0.0245	0.0163	0.1356
separated	0.0029	0.0422	-0.1625	0.2417
never married	0.0175	0.0289	0.0770	0.1666
polviews	0.0160	0.0057	0.0756	0.0317
female	0.0018	0.0152	0.0101	0.0824
year*college	0.0069	0.0015	0.0056	0.0109
year*age	0.0000	0.0001	0.0007	0.0003
<i>year*race</i>				
year*black	0.0096	0.0024	0.0261	0.0192
year*other	0.0002	0.0046	0.0014	0.0283
<i>year*marital</i>				
year*widowed	0.0008	0.0024	0.0110	0.0130
year*divorced	0.0007	0.0019	0.0042	0.0117
year*separated	0.0019	0.0035	0.0186	0.0215
year*never married	0.0005	0.0023	0.0038	0.0144
year*polviews	-0.0004	0.0005	0.0003	0.0029
year*female	-0.0007	0.0013	-0.0057	0.0077
constant	0.1669	0.0401	-1.3402	0.2206

Note: Data are from General Social Surveys Cumulative File, N = 12,239. The variable year is centered around 1976 to facilitate interpretation of the main effects and interaction terms.