

# THE STATA JOURNAL

## Editor

H. Joseph Newton  
Department of Statistics  
Texas A&M University  
College Station, Texas 77843  
979-845-8817; fax 979-845-6077  
jnewton@stata-journal.com

## Editor

Nicholas J. Cox  
Department of Geography  
Durham University  
South Road  
Durham City DH1 3LE UK  
n.j.cox@stata-journal.com

## Associate Editors

Christopher F. Baum  
Boston College

Nathaniel Beck  
New York University

Rino Bellocco  
Karolinska Institutet, Sweden, and  
University of Milano-Bicocca, Italy

Maarten L. Buis  
Vrije Universiteit, Amsterdam

A. Colin Cameron  
University of California–Davis

Mario A. Cleves  
Univ. of Arkansas for Medical Sciences

William D. Dupont  
Vanderbilt University

David Epstein  
Columbia University

Allan Gregory  
Queen's University

James Hardin  
University of South Carolina

Ben Jann  
ETH Zürich, Switzerland

Stephen Jenkins  
University of Essex

Ulrich Kohler  
WZB, Berlin

Frauke Kreuter  
University of Maryland–College Park

Jens Lauritsen  
Odense University Hospital

Stanley Lemeshow  
Ohio State University

J. Scott Long  
Indiana University

Thomas Lumley  
University of Washington–Seattle

Roger Newson  
Imperial College, London

Austin Nichols  
Urban Institute, Washington DC

Marcello Pagano  
Harvard School of Public Health

Sophia Rabe-Hesketh  
University of California–Berkeley

J. Patrick Royston  
MRC Clinical Trials Unit, London

Philip Ryan  
University of Adelaide

Mark E. Schaffer  
Heriot-Watt University, Edinburgh

Jeroen Weesie  
Utrecht University

Nicholas J. G. Winter  
University of Virginia

Jeffrey Wooldridge  
Michigan State University

**Stata Press Editorial Manager**

**Stata Press Copy Editors**

Lisa Gilmore

Jennifer Neve and Deirdre Patterson

The *Stata Journal* publishes reviewed papers together with shorter notes or comments, regular columns, book reviews, and other material of interest to Stata users. Examples of the types of papers include 1) expository papers that link the use of Stata commands or programs to associated principles, such as those that will serve as tutorials for users first encountering a new field of statistics or a major new technique; 2) papers that go “beyond the Stata manual” in explaining key features or uses of Stata that are of interest to intermediate or advanced users of Stata; 3) papers that discuss new commands or Stata programs of interest either to a wide spectrum of users (e.g., in data management or graphics) or to some large segment of Stata users (e.g., in survey statistics, survival analysis, panel analysis, or limited dependent variable modeling); 4) papers analyzing the statistical properties of new or existing estimators and tests in Stata; 5) papers that could be of interest or usefulness to researchers, especially in fields that are of practical importance but are not often included in texts or other journals, such as the use of Stata in managing datasets, especially large datasets, with advice from hard-won experience; and 6) papers of interest to those who teach, including Stata with topics such as extended examples of techniques and interpretation of results, simulations of statistical concepts, and overviews of subject areas.

For more information on the *Stata Journal*, including information for authors, see the web page

<http://www.stata-journal.com>

The *Stata Journal* is indexed and abstracted in the following:

- CompuMath Citation Index<sup>®</sup>
- RePEc: Research Papers in Economics
- Science Citation Index Expanded (also known as SciSearch<sup>®</sup>)

**Copyright Statement:** The *Stata Journal* and the contents of the supporting files (programs, datasets, and help files) are copyright © by StataCorp LP. The contents of the supporting files (programs, datasets, and help files) may be copied or reproduced by any means whatsoever, in whole or in part, as long as any copy or reproduction includes attribution to both (1) the author and (2) the *Stata Journal*.

The articles appearing in the *Stata Journal* may be copied or reproduced as printed copies, in whole or in part, as long as any copy or reproduction includes attribution to both (1) the author and (2) the *Stata Journal*.

Written permission must be obtained from StataCorp if you wish to make electronic copies of the insertions. This precludes placing electronic copies of the *Stata Journal*, in whole or in part, on publicly accessible web sites, file servers, or other locations where the copy may be accessed by anyone other than the subscriber.

Users of any of the software, ideas, data, or other materials published in the *Stata Journal* or the supporting files understand that such use is made without warranty of any kind, by either the *Stata Journal*, the author, or StataCorp. In particular, there is no warranty of fitness of purpose or merchantability, nor for special, incidental, or consequential damages such as loss of profits. The purpose of the *Stata Journal* is to promote free communication among Stata users.

The *Stata Journal*, electronic version (ISSN 1536-8734) is a publication of Stata Press. Stata and Mata are registered trademarks of StataCorp LP.

# Methods for estimating adjusted risk ratios

Peter Cummings  
Department of Epidemiology  
School of Public Health and  
Harborview Injury Prevention and Research Center  
University of Washington  
Seattle, WA  
peterc@u.washington.edu

**Abstract.** The risk ratio can be a useful statistic for summarizing the results of cross-sectional, cohort, and randomized trial studies. I discuss several methods for estimating adjusted risk ratios and show how they can be executed in Stata, including 1) Mantel–Haenszel and inverse-variance stratified methods; 2) generalized linear regression with a log link and binomial distribution; 3) generalized linear regression with a log link, normal distribution, and robust variance estimator; 4) Poisson regression with a robust variance estimator; 5) Cox proportional hazards regression with a robust variance estimator; 6) standardized risk ratios from logistic, probit, complementary log-log, and log-log regression; and 7) a substitution method. Advantages and drawbacks are noted for some methods.

**Keywords:** st0162, risk ratio, odds ratio

## 1 Introduction

The case–control study design is typically (but not always) used when outcomes are rare in the population from which study subjects are sampled. In 1951, Cornfield noted that when outcomes are sufficiently rare, the odds ratio from a case–control study will approximate the population risk ratio for the association of an exposure with a disease outcome. It was later realized that if controls are sampled as each case arises in time, the odds ratio will estimate the incidence-rate ratio even when outcomes are common (Greenland and Thomas 1982; Rodrigues and Kirkwood 1990; Rothman, Greenland, and Lash 2008, 113–114). In Stata, case–control data can be analyzed using Mantel–Haenszel stratified methods (`cc`, `tabodds`, `mhodds`), logistic regression (`logistic`), or conditional logistic regression (`clogit`) to estimate adjusted odds ratios that usually can be interpreted either as risk ratios (when outcomes are rare) or incidence-rate ratios (when incidence density sampling is used).

Cross-sectional, cohort, and randomized controlled trial designs with binary outcomes can often be summarized by estimating odds ratios or risk ratios. If the study outcome is sufficiently rare among exposed and unexposed study subjects, the odds ratio for the exposure–outcome association will closely approximate the risk ratio. But if the outcome is common and the risk ratio is not close to 1, the odds ratio will be further from 1 compared with the risk ratio. Even if the outcome is rare in the entire sample, if an adjustment is made for other variables, then the adjusted odds ratio will

be further from 1 than the adjusted risk ratio if the outcome is common in adjustment variable subgroups that contribute a noteworthy portion of the outcomes (Greenland 1987).

When summary odds and risk ratios differ, there is debate regarding which is preferable. Some have argued that odds ratios are preferred because they are symmetric with regard to the outcome definition (Walter 1998; Olkin 1998; Senn 1999; Newman 2001, 35–40; Cook 2002). Furthermore, when outcomes are common, a constant (homogeneous) adjusted odds ratio for all subjects may be more plausible than a constant risk ratio (Levin 1991; Senn 1998; Cook 2002).

Some who favor risk ratios feel they are more easily understood by physicians (Sackett, Deeks, and Altman 1996). Others have noted that risk ratios have a desirable feature called collapsibility; in the absence of confounding, a weighted average of stratum-specific risk ratios will equal the ratio from one  $2 \times 2$  table of the pooled (collapsed) counts from the stratum-specific tables (Miettinen and Cook 1981; Greenland 1987, 1991b; Greenland, Robins, and Pearl 1999; Newman 2001, 52–55; Rothman, Greenland, and Lash 2008, 62). This means that a crude (unadjusted) risk ratio will not change if we adjust for a variable that is not a confounder. In the absence of confounding, the risk ratio estimates the change in risk, on a ratio scale, for the entire exposed group due to exposure. Because of collapsibility, this risk ratio has a useful interpretation as the ratio change in the average risk in the exposed group due to exposure. It is not the average ratio change in risk (i.e., the average risk ratio) among exposed individuals, except in the unlikely event that the risk ratios for all individuals are the same (Greenland 1987).

Odds ratios lack the property of collapsibility and therefore the interpretation of an odds ratio is more limited; in the absence of confounding, it estimates the change in odds, on a ratio scale, in the exposed group due to exposure. But it does not estimate either the change in the average odds of the exposed due to exposure or the average change in odds (i.e., the average odds ratio) among exposed individuals, not even if all individuals had the same change in odds when exposed (Greenland 1987). The odds ratio will estimate the average change in odds for exposed individuals only if all individual odds ratios are the same and all individual risks without exposure are the same. Except in this unlikely situation, the crude odds ratio will be closer to 1 than the average of stratum-specific or individual odds ratios. Even in the absence of confounding, the adjusted (conditional) odds ratio will be further from 1 than the crude (unadjusted or marginal) odds ratio (Gail et al. 1984; Greenland 1987; Hauck et al. 1998; Steyerberg et al. 2000; Newman 2001, 52–55; Rothman, Greenland, and Lash 2008, 62; Cummings 2009).

For analysts who wish to estimate odds ratios for the association of exposure with disease in a cross-sectional study, cohort study, or randomized trial, the statistical methods in Stata's `cc`, `tabodds`, `mhodds`, and `logistic` commands can be used. If the goal is to estimate risk ratios, these same methods can be used if outcomes are sufficiently rare that odds ratios will closely approximate risk ratios. But if risk ratios are desired when outcomes are common, odds ratio estimates will not suffice. In this article, I describe methods for estimating adjusted risk ratios with confidence intervals (CIs) in Stata.

## 2 Data used to illustrate the methods

I will show how to reproduce the risk-ratio estimates and CIs that [Greenland \(2004a\)](#) gave in a review of risk-ratio estimation. The data (table 1) are from table 5.3 in Newman's (2001, 98 and 126) textbook. Newman described 192 women who were diagnosed with breast cancer in Canada and followed for 5 years; 28% (54/192) of the women died, so the outcome was not rare. Greenland estimated the risk ratio for death at 5 years among women with low estrogen-receptor levels in their breast cancer tissue compared with women who had high receptor levels; these risk ratios were adjusted for cancer stage (I, II, or III) so that women with the same cancer stage were compared.

Table 1. Deaths, total subjects, and risk of death for 192 women with breast cancer followed for 5 years, by stage at diagnosis (I, II, III) and estrogen-receptor-level category (low, high). Also, risk ratios within each cancer stage for death among women with low versus high receptor levels.

Receptor levels	Stage	Died	Total	Risk	Risk ratio for death comparing women with low versus high receptor levels
Low	I	2	12	0.17	1.8
High	I	5	55	0.09	1.0 (reference group)
Low	II	9	22	0.41	1.8
High	II	17	74	0.23	1.0 (reference group)
Low	III	12	14	0.86	1.4
High	III	9	15	0.60	1.0 (reference group)

## 3 Method 1: Mantel–Haenszel and inverse-variance stratified methods

Mantel–Haenszel methods for odds ratios were described in 1959 ([Mantel and Haenszel 1959](#)) and extended to risk ratios in 1981 ([Nurminen 1981](#); [Tarone 1981](#); [Kleinbaum, Kupper, and Morgenstern 1982](#); [Newman 2001](#), 148–149; [Rothman, Greenland, and Lash 2008](#), 274–275). We can estimate the adjusted risk ratio for death associated with low estrogen-receptor levels (the `low` variable) compared with high estrogen-receptor levels by using the `cs` command:

(Continued on next page)

```
. use brcadat
(Breast cancer data)
. cs died low, by(stage) pool
```

Cancer stage	RR	[95% Conf. Interval]		M-H Weight
1	1.833333	.4024728	8.35115	.8955224
2	1.780749	.9269437	3.420991	3.895833
3	1.428571	.8971062	2.274888	4.344828
Crude	2.225806	1.449035	3.418974	
Pooled (direct)	1.554553	1.076372	2.245169	
M-H combined	1.618421	1.093775	2.394719	

```
Test of homogeneity (direct)  chi2(2) = 0.339  Pr>chi2 = 0.8443
Test of homogeneity (M-H)    chi2(2) = 0.385  Pr>chi2 = 0.8251
```

I invoked the `pool` option so that the output shows both the Mantel–Haenszel combined risk ratio and the pooled risk ratio obtained using inverse-variance weights. These methods require that variables be treated as categorical, not continuous.

## 4 Method 2: Generalized linear regression with a log link and binomial distribution

Estimation of risk ratios using a generalized linear model with a log link and binomial distribution was proposed in 1986 (Wacholder 1986). This approach has been described in several articles (Robbins, Chao, and Fonseca 2002; McNutt et al. 2003; Barros and Hirakata 2003); it has been called log-binomial (Blizzard and Hosmer 2006) or binomial log-linear regression (Greenland 2004a). This approach can be implemented in Stata:

```
. glm died low stage2 stage3, family(binomial) link(log) eform difficult
Iteration 0:  log likelihood = -154.81266   (not concave)
Iteration 1:  log likelihood = -99.23223
Iteration 2:  log likelihood = -94.764125
Iteration 3:  log likelihood = -93.93508
Iteration 4:  log likelihood = -93.050613
Iteration 5:  log likelihood = -92.930394
Iteration 6:  log likelihood = -92.927072
Iteration 7:  log likelihood = -92.927069

Generalized linear models               No. of obs   =       192
Optimization      : ML                 Residual df   =       188
                                           Scale parameter =       1
Deviance          = 185.8541388         (1/df) Deviance = .9885858
Pearson           = 190.2212968         (1/df) Pearson  = 1.011815
Variance function: V(u) = u*(1-u)      [Bernoulli]
Link function     : g(u) = ln(u)       [Log]

                                           AIC           = 1.009657
                                           BIC           = -802.555

Log likelihood    = -92.92706941
```

died	OIM					
	Risk Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.558321	.3148624	2.20	0.028	1.048745	2.315497
stage2	2.538159	.9991488	2.37	0.018	1.17339	5.490288
stage3	5.868042	2.273768	4.57	0.000	2.745787	12.54064

Above Stata reported that for iteration 0, the likelihood region was not concave. When the command is run without the `difficult` option, Stata 10.0 will repeatedly report a not-concave region and fail to converge. The `difficult` option changed Stata's convergence algorithm and solved the problem in this example, but that option may not always work. The convergence problem arose because among women with Stage III cancer and low estrogen-receptor levels, the risk of death was close to 1:  $12/14 = 0.86$ . When the risk is close to 1 in a stratum of the data, maximum-likelihood convergence may fail. This problem has been discussed in several articles (Carter, Lipsitz, and Tilley 2005; Blizzard and Hosmer 2006; Lumley, Kronmal, and Ma 2006; Localio, Margolis, and Berlin 2007).

Wacholder (1986; Lumley, Kronmal, and Ma 2006) described a method that modified the convergence by truncating estimated risks to values slightly greater than 0 and less than 1. This is implemented in Stata's `binreg` command:

(Continued on next page)

```

. xi: binreg died low i.stage, rr
i.stage      _Istage_1-3      (naturally coded; _Istage_1 omitted)
Iteration 1:  deviance = 309.6253
Iteration 2:  deviance = 190.79
Iteration 3:  deviance = 185.9416
Iteration 4:  deviance = 185.8543
Iteration 5:  deviance = 185.8541
Iteration 6:  deviance = 185.8541

Generalized linear models      No. of obs      =      192
Optimization      : MQL Fisher scoring      Residual df      =      188
                    (IRLS EIM)              Scale parameter =      1
Deviance          = 185.8541388              (1/df) Deviance = .9885858
Pearson           = 190.2164034              (1/df) Pearson  = 1.011789
Variance function: V(u) = u*(1-u)           [Bernoulli]
Link function     : g(u) = ln(u)             [Log]
BIC               = -802.555

```

	EIM					
died	Risk Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.558326	.3067419	2.25	0.024	1.059515	2.291972
_Istage_2	2.538158	.9976133	2.37	0.018	1.174781	5.483782
_Istage_3	5.868047	2.259727	4.60	0.000	2.758699	12.48196

The risk ratios and standard errors estimated by `glm, family(binomial)` `link(log)` and by `binreg, rr` are similar, but not identical, because the convergence methods differ. The default for `glm` is maximum likelihood and the default for `binreg` is iterated, reweighted least squares; by changing the default optimization, either command can produce the estimates obtained by the other, provided that both methods achieve convergence. Because `binreg` constrains risk estimates to be greater than 0 and less than 1, it may converge when maximum likelihood will not. Sometimes both methods will fail to converge.

Another approach to convergence difficulty is to make several copies of the original data and append them into one data file (Deddens, Petersen, and Lei 2003; Petersen and Deddens 2006; Deddens and Petersen 2008). Then make one more copy, but recode all the 0 outcomes to 1 and all the 1 outcomes to 0 in that copy, and append this recoded copy to all the other copies. Then analyze all these data together; including one set of data with reversed outcomes may help the maximum-likelihood algorithm converge. If the number of copies is sufficiently large, the risk-ratio estimates will approximate those from maximum-likelihood methods. Because a set of records larger than the original data is used, corrections must be made to the standard errors and CIs. This extra step of correcting the standard errors can be avoided by using just two copies of the data, one with recoded outcomes, with appropriate weights (Lumley, Kronmal, and Ma 2006). Below I used importance weights of 0.999 and 0.001, and produced the risk ratio I would get by analyzing 999 copies of the original data with just 1 copy of recoded data. The `difficult` option was still required.



```

. generate iweight = .999
. append using brcadat
(label noyes already defined)
. recode died 0=1 1=0 if iweight==.
(died: 192 changes made)
. replace iweight = .001 if iweight==.
(192 real changes made)
. glm died low stage2 stage3 [iweight=iweight], family(binomial) link(log)
> eform difficult

Iteration 0:  log likelihood = -155.22946   (not concave)
Iteration 1:  log likelihood = -99.401459
Iteration 2:  log likelihood = -94.953215
Iteration 3:  log likelihood = -94.13882
Iteration 4:  log likelihood = -93.23838
Iteration 5:  log likelihood = -93.113209
Iteration 6:  log likelihood = -93.109614
Iteration 7:  log likelihood = -93.109611

Generalized linear models                               No. of obs   =       384
Optimization      : ML                               Residual df   =       380
                                                         Scale parameter =         1
                                                         (1/df) Deviance =   .4900506
                                                         (1/df) Pearson  =   .5006674
Deviance          =  186.2192211                               [Bernoulli]
Pearson           =  190.253604                               [Log]
Variance function: V(u) = u*(1-u)
Link function     : g(u) = ln(u)
                                                         AIC           =   .5057792
Log likelihood    = -93.10961053                             BIC           = -2075.025

```

died	OIM					
	Risk Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.556724	.3143539	2.19	0.028	1.047915	2.312583
stage2	2.523499	.9897355	2.36	0.018	1.169918	5.443157
stage3	5.823722	2.248418	4.56	0.000	2.732558	12.41172

## 5 Method 3: Generalized linear regression with a log link, normal distribution, and robust variance estimator

Convergence problems for generalized linear regression with a log link can also be resolved by using a Gaussian (normal) distribution (Lumley, Kronmal, and Ma 2006). The resulting standard errors may be too big or small, but a robust variance estimator will correct the standard errors. Below convergence was achieved without the `difficult` option:

(Continued on next page)

```

. use brccat, clear
(Breast cancer data)

. glm died low stage2 stage3, family(gaussian) link(log) eform robust
Iteration 0:  log pseudolikelihood = -169.42788
Iteration 1:  log pseudolikelihood = -119.04978
Iteration 2:  log pseudolikelihood = -98.365962
Iteration 3:  log pseudolikelihood = -94.332246
Iteration 4:  log pseudolikelihood = -94.242364
Iteration 5:  log pseudolikelihood = -94.242364

Generalized linear models      No. of obs      =      192
Optimization      : ML              Residual df      =      188
                                   Scale parameter = .1595924
Deviance          = 30.00336652      (1/df) Deviance = .1595924
Pearson           = 30.00336652      (1/df) Pearson  = .1595924

Variance function: V(u) = 1          [Gaussian]
Link function     : g(u) = ln(u)     [Log]

                                   AIC          = 1.023358
                                   BIC          = -958.4058

Log pseudolikelihood = -94.24236355

```

died	exp(b)	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.553274	.3193155	2.14	0.032	1.038154	2.323992
stage2	2.524622	1.005045	2.33	0.020	1.157006	5.508803
stage3	5.86819	2.312817	4.49	0.000	2.710329	12.70534

## 6 Method 4: Poisson regression with a robust variance estimator

Poisson regression is a generalized linear model with a log link and a Poisson distribution. When the outcome is binary, the exponentiated coefficients are risk ratios instead of incidence-rate ratios (Gourieroux, Monfort, and Trognon 1984a,b; Lloyd 1999, 85–86; Wooldridge 2002, 648–649; Greenland 2004a; Zou 2004; Carter, Lipsitz, and Tilley 2005). Methods that rely on the Poisson distribution assume that the mean count and its variance are equal. In the breast cancer data, the mean count of deaths per woman was  $54/192 = 0.28125$ . If the variance of the mean count is also 0.28125, then the standard error of the mean count is the square root of the variance divided by the square root of the number of women ( $192$ ) = 0.0382733. This is indeed the standard error that Stata reports using the `ci`, `poisson` command or using `lincom` after the `poisson` command:

```

. ci died, poisson

```

Variable	Exposure	Mean	Std. Err.	— Poisson — [95% Conf. Interval]	Exact —
died	192	.28125	.0382733	.2112837	.3669702

```
. poisson died, nolog
Poisson regression
Log likelihood = -122.49961
Number of obs = 192
LR chi2(0) = 0.00
Prob > chi2 = .
Pseudo R2 = 0.0000
```

died	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_cons	-1.268511	.1360828	-9.32	0.000	-1.535229	-1.001794

```
. lincom _cons, irr
( 1) [died]_cons = 0
```

died	IRR	Std. Err.	z	P> z	[95% Conf. Interval]	
(1)	.28125	.0382733	-9.32	0.000	.2154064	.3672201

If the deaths were from a Poisson distribution, women would have nonnegative integer counts of 0, 1, 2, 3, ..., or more deaths. The data cannot be Poisson, because no woman dies more than once. The data are from a binomial distribution, and the binomial variance is assumed to be the proportion that died multiplied by 1 minus that proportion. The standard error of the mean proportion is the square root of the variance divided by the square root of the number of women, which is 0.0324477. This is the standard error reported by `ci`, `binomial`:

```
. ci died, binomial
```

Variable	Obs	Mean	Std. Err.	— Binomial Exact — [95% Conf. Interval]	
died	192	.28125	.0324477	.2188833	.3505085

If we use Poisson methods for these binomial data, the standard error for the outcome proportion (risk) is too large: 0.03827 instead of 0.03245. As the outcome becomes less common, the Poisson standard error will converge toward the binomial standard error (Armitage, Berry, and Matthews 2002, 71–76). But in the breast cancer data, use of Poisson regression to estimate risk ratios will produce standard errors, *p*-values, and CIs that are too large:

```
. poisson died low stage2 stage3, irr nolog
Poisson regression
Log likelihood = -109.14601
Number of obs = 192
LR chi2(3) = 26.71
Prob > chi2 = 0.0000
Pseudo R2 = 0.1090
```

died	IRR	Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.630775	.4688634	1.70	0.089	.9282513	2.864987
stage2	2.520742	1.074375	2.17	0.030	1.093288	5.811955
stage3	5.913372	2.645148	3.97	0.000	2.460814	14.20992

Above, the 95% CI for the `low` variable is wide, 0.93 to 2.86, compared with the CIs from other methods. We can obtain standard errors and CIs that are approximately correct by using a robust variance estimator, which can relax the assumption that the data are from a Poisson distribution (Wooldridge 2002, 650–651; Greenland 2004a; Zou 2004; Carter, Lipsitz, and Tilley 2005). The robust variance estimator is sometimes called the Huber, White, Huber–White, sandwich, or survey estimator, as well as other names (Hardin and Hilbe 2007, 35–36). In Stata, we can invoke this estimator with the `vce(robust)` option and the CI for the `low` variable becomes narrower, 1.07 to 2.48:

```
. poisson died low stage2 stage3, irr nolog vce(robust)
Poisson regression                                Number of obs   =       192
                                                    Wald chi2(3)     =       53.61
                                                    Prob > chi2      =       0.0000
Log pseudolikelihood = -109.14601                Pseudo R2       =       0.1090
```

died	IRR	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.630775	.3480542	2.29	0.022	1.073305	2.477792
stage2	2.520742	.9937819	2.35	0.019	1.16399	5.458932
stage3	5.913372	2.28568	4.60	0.000	2.772187	12.61386

## 7 Method 5: Cox proportional hazards regression with a robust variance estimator

Cox proportional hazards regression can estimate risk ratios if we set the follow-up time to 1, or any quantity that is the same for all subjects, and use the Breslow method to break ties. The robust variance estimator should be used, because otherwise the standard errors will be too large:

```
. generate byte time = 1
. stset time, failure(died) noshow
      failure event:  died != 0 & died < .
obs. time interval:  (0, time]
exit on or before:  failure
```

---

```
192  total obs.
   0  exclusions
```

---

```
192  obs. remaining, representing
   54  failures in single record/single failure data
192  total analysis time at risk, at risk from t =      0
      earliest observed entry t =      0
      last observed exit t =      1
```

```
. stcox low stage2 stage3, hr breslow vce(robust) nolog
Cox regression -- Breslow method for ties
No. of subjects      =          192          Number of obs   =          192
No. of failures      =           54
Time at risk        =          192
Log pseudolikelihood = -270.55115          Wald chi2(3)      =          53.61
                                          Prob > chi2       =          0.0000
```

_t	Haz. Ratio	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.630775	.3480542	2.29	0.022	1.073305	2.477792
stage2	2.520742	.9937819	2.35	0.019	1.16399	5.458932
stage3	5.913372	2.28568	4.60	0.000	2.772187	12.61386

The results above reproduce exactly the results from Poisson regression with the robust variance estimator; the Poisson and Cox methods are identical when implemented in this way. Options other than the Breslow method for dealing with tied survival times will produce risk-ratio estimates for exposure to a low estrogen-receptor-level tumor that are too large: 1) the `efron` option produces a risk ratio of 1.91, 2) the `exactm` option yields 2.04, and 3) the `exactp` option risk ratio is 2.49.

## 8 Method 6: Regression-based standardized risk ratios

Any regression model for binomial outcomes can estimate the probability (risk) of death for women with high and low estrogen-receptor-level tumors within each cancer stage. With this information, we can estimate the average risk of death that would be expected if all 192 women had low estrogen-receptor-level tumors and the distribution of cancer stages observed in the data. This estimate is said to be standardized to the distribution of the other variables, cancer stage in this example, in the regression model (Lane and Nelder 1982; Flanders and Rhodes 1987; Joffe and Greenland 1995; Greenland 1991a, 2004a; Localio, Margolis, and Berlin 2007; Rothman, Greenland, and Lash 2008, 442–446). The average risk can also be estimated assuming that all 192 women had a high estrogen-receptor-level tumor. The ratio of the low estrogen-receptor-level average risk divided by the high estrogen-receptor-level average risk can be calculated and the standard error for this risk ratio can be estimated using the delta method (Casella and Berger 2002, 240–245). The word “standardized” is used just as for standardized mortality rates or any statistic standardized to a given population distribution. We can estimate this standardized risk ratio by using logistic regression:

(Continued on next page)

```
. use brcadat, clear
(Breast cancer data)

. logistic died low stage2 stage3, nolog

Logistic regression
```

Number of obs	=	192
LR chi2(3)	=	42.27
Prob > chi2	=	0.0000
Pseudo R2	=	0.1853

```
Log likelihood = -92.939847
```

died	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
low	2.508065	.9916923	2.33	0.020	1.155507	5.443836
stage2	3.109772	1.44851	2.44	0.015	1.248087	7.748406
stage3	18.8389	11.03231	5.01	0.000	5.978344	59.36498

```
. #delimit ;
delimiter now ;
. predictnl lnrr =
> ln(
> sum(1/
> (1+exp(-(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3+_b[low]))))
> /
> sum(1/
> (1+exp(-(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3))))))
> , se(lnrr_se);

. #delimit cr
delimiter now cr
. scalar rr = exp(lnrr[_N])
. scalar upper = exp(lnrr[_N] + invnormal(1-.05/2)*lnrr_se[_N])
. scalar lower = exp(lnrr[_N] - invnormal(1-.05/2)*lnrr_se[_N])
. display "Risk ratio = " rr " 95% CI = " lower " , " upper
Risk ratio = 1.6755988 95% CI = 1.0935713, 2.5673969
```

The adjusted odds ratio for death among women with a low estrogen-receptor level–tumor, compared with women with a high estrogen-receptor–level tumor, was 2.5. Because the outcome of death was common, this odds ratio does not closely approximate the risk ratio.

Above I used `predictnl` to estimate the  $\ln$  of the risk ratio (`lnrr` variable); this command can estimate nonlinear comparisons from regression coefficients. The `se` option estimated the standard error for the  $\ln$  risk ratio using the delta method. To make the output less cluttered, I used `delimit` to change how Stata recognizes the end of a command line. To estimate the risk or probability of death, I used the expression  $1/\{1 + \exp(-\text{linear predictor})\}$ . The first sum used by `predictnl` is for risk estimates if all women had a low estrogen-receptor–level tumor, because the  $\ln$  odds term for low receptor status, `_b[low]`, is included in the sum, regardless of each woman’s actual receptor status. In this first sum, the regression coefficients (which are  $\ln$  odds estimates) for `_Istage_2` and `_Istage_3` were both multiplied by each woman’s observed cancer stage, thereby standardizing the estimate to the observed distribution of cancer stage. Stata’s `sum()` function is the running sum from the first record to the last, so the sum in the last record of the data is the sum of all the estimated risks if all 192 women had a low estrogen-receptor–level tumor but each had her observed cancer stage. The second sum

in the expression, after the line with only a division sign, is the sum of all estimated risks if all women had a high estrogen-receptor-level tumor, again standardized to the observed cancer-stage distribution. The first sum is divided by the second and the  $\ln$  taken of this ratio so that the  $\ln$  of the risk ratio is estimated. I then estimated the risk ratio, which is  $\exp(\ln rr)$ , and the 95% upper and lower confidence limits for the risk ratio, and used the `display` command to show these results for the last record by using the subscript `[_N]`: risk ratio = 1.7, 95% CI is [1.1, 2.6].

We can use simpler commands to estimate the risk ratio, but they do not provide a CI. Still, these commands show how the risks and risk ratio may be estimated and are shown below to clarify how Stata is using the regression estimates. After fitting the logistic model, the commands are

```
. replace low=0
(48 real changes made)
. predict risk0
(option pr assumed; Pr(died))
. summ risk0, meanonly
. local avrisk0 = r(mean)
. replace low=1
(192 real changes made)
. predict risk1
(option pr assumed; Pr(died))
. summ risk1, meanonly
. local avrisk1 = r(mean)
. local rr = `avrisk1' / `avrisk0'
. display "Risk1 = " `avrisk1' " Risk0 = " `avrisk0' " Risk ratio = " `rr'
Risk1 = .40087948 Risk0 = .23924549 Risk ratio = 1.6755989
```

Risks for binomial outcomes can also be estimated after probit regression. In the probit model, the outcome risk estimate applies the cumulative standard normal distribution function (`normal()`) to the linear predictor, instead of the  $\ln$  odds function used in logistic regression:

```
. use brcadat, clear
(Breast cancer data)
. probit died low stage2 stage3, nolog

Probit regression                                Number of obs   =       192
                                                LR chi2(3)      =       42.21
                                                Prob > chi2     =       0.0000
Log likelihood = -92.968357                    Pseudo R2      =       0.1850
```

died	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
low	.5386148	.2343501	2.30	0.022	.079297	.9979327
stage2	.6290485	.2503224	2.51	0.012	.1384256	1.119671
stage3	1.739085	.3302966	5.27	0.000	1.091715	2.386454
_cons	-1.376363	.2165793	-6.36	0.000	-1.800851	-.9518752

```

. #delimit ;
delimiter now ;
. predictnl lnrr = ln(
>   sum(normal(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3+_b[low]))
>   /
>   sum(normal(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3)))
>   , se(lnrr_se);

. #delimit cr
delimiter now cr
. scalar rr = exp(lnrr[_N])
. scalar upper = exp(lnrr[_N] + invnormal(1-.05/2)*lnrr_se[_N])
. scalar lower = exp(lnrr[_N] - invnormal(1-.05/2)*lnrr_se[_N])
. display "Risk ratio = " rr " 95% CI = " lower " , " upper
Risk ratio = 1.6751332 95% CI = 1.0913484, 2.5711965

```

Nelder (2001) has suggested that when a dichotomous outcome is common, the complementary log-log regression model may fit the data well:

```

. use brccat, clear
(Breast cancer data)

. cloglog died low stage2 stage3, nolog
Complementary log-log regression

```

Number of obs	=	192
Zero outcomes	=	138
Nonzero outcomes	=	54
LR chi2(3)	=	42.60
Prob > chi2	=	0.0000

```

Log likelihood = -92.771237

```

	died	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
	low	.7138795	.2936374	2.43	0.015	.1383607	1.289398
	stage2	1.022028	.426798	2.39	0.017	.1855198	1.858537
	stage3	2.302261	.4522776	5.09	0.000	1.415813	3.188709
	_cons	-2.369669	.3882655	-6.10	0.000	-3.130655	-1.608683

```

. #delimit ;
delimiter now ;
. predictnl lnrr = ln(
>   sum(1-exp(-exp(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3+_b[low])))
>   /
>   sum(1-exp(-exp(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3))))
>   , se(lnrr_se);

. #delimit cr
delimiter now cr
. scalar rr = exp(lnrr[_N])
. scalar upper = exp(lnrr[_N] + invnormal(1-.05/2)*lnrr_se[_N])
. scalar lower = exp(lnrr[_N] - invnormal(1-.05/2)*lnrr_se[_N])
. display "Risk ratio = " rr " 95% CI = " lower " , " upper
Risk ratio = 1.6652749 95% CI = 1.1009763, 2.5188013

```

Hardin and Hilbe (2007, 147) note that if most subjects either have or do not have the outcome, the complementary log-log and log-log models may fit better than logistic or probit models. We can fit the log-log model using the `glm` command:



```

. use brccat, clear
(Breast cancer data)

. glm died low stage2 stage3, nolog family(bin) link(loglog) eform

Generalized linear models                                No. of obs      =       192
Optimization      : ML                                Residual df    =       188
                                                         Scale parameter =         1
Deviance          = 186.5949538                        (1/df) Deviance = .9925263
Pearson          = 192.8460376                        (1/df) Pearson  = 1.025777
Variance function: V(u) = u*(1-u)                    [Bernoulli]
Link function     : g(u) = -ln(-ln(u))                [Log-log]
                                                         AIC            = 1.013515
Log likelihood    = -93.2974769                       BIC            = -801.8142

```

died	exp(b)	OIM Std. Err.	z	P> z	[95% Conf. Interval]	
low	1.642276	.3938706	2.07	0.039	1.026362	2.627796
stage2	1.695259	.3494023	2.56	0.010	1.131876	2.539063
stage3	6.133979	2.415435	4.61	0.000	2.835022	13.27174

```

. #delimit ;
delimiter now ;
. predictnl lnrr = ln(
> sum(exp(-exp(-(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3+_b[low]))))
> /
> sum(exp(-exp(-(_b[_cons]+_b[stage2]*stage2+_b[stage3]*stage3))))
> , se(lnrr_se);

. #delimit cr
delimiter now cr
. scalar rr = exp(lnrr[_N])

. scalar upper = exp(lnrr[_N] + invnormal(1-.05/2)*lnrr_se[_N])
. scalar lower = exp(lnrr[_N] - invnormal(1-.05/2)*lnrr_se[_N])

. display "Risk ratio = " rr " 95% CI = " lower " , " upper
Risk ratio = 1.6312705 95% CI = 1.0545183, 2.5234682

```

## 9 Method 7: A substitution method

A crude (unadjusted) odds ratio can be converted to a risk ratio: crude risk ratio = (crude odds ratio)/{(1-Po)+(Po×crude odds ratio)}, where Po is the proportion of all unexposed subjects who had the outcome in the data. [Zhang and Yu \(1998\)](#) suggested that, in a cohort study, one can use this same formula to convert an adjusted odds ratio to an adjusted risk ratio. This substitution method was described by [Holland \(1989\)](#), who used it to estimate an adjusted risk difference from a Mantel-Haenszel summary odds ratio. [Greenland and Holland \(1991\)](#) reported that this will produce ratio estimates biased away from 1 when outcomes are common and risk among those not exposed varies substantially. The bias occurs because the summary odds ratio is not a weighted average of stratum-specific odds ratios; odds ratios lack the property of collapsibility ([Greenland 1987](#)). The method can be implemented in Stata:

```
. use brcadat, clear
(Breast cancer data)

. summ died if low==0, meanonly
. local p0 = r(mean)
. display `p0'
.21527778

. logistic died low stage2 stage3, nolog
Logistic regression
```

Number of obs	=	192
LR chi2(3)	=	42.27
Prob > chi2	=	0.0000
Pseudo R2	=	0.1853

```
Log likelihood = -92.939847
```

died	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
low	2.508065	.9916923	2.33	0.020	1.155507	5.443836
stage2	3.109772	1.44851	2.44	0.015	1.248087	7.748406
stage3	18.8389	11.03231	5.01	0.000	5.978344	59.36498

```
. scalar rr = exp(_b[low])/[(1-`p0')+(`p0'*exp(_b[low]))]
. scalar lower = exp(_b[low]-invnormal(1-.05/2)*_se[low])/[(1-`p0')+
> (`p0'*exp(_b[low]-invnormal(1-.05/2)*_se[low]))]
. scalar upper = exp(_b[low]+invnormal(1-.05/2)*_se[low])/[(1-`p0')+
> (`p0'*exp(_b[low]+invnormal(1-.05/2)*_se[low]))]
. display "Risk ratio = " rr " 95% CI = " lower " , " upper
Risk ratio = 1.8933751 95% CI = 1.1180767, 2.7822097
```

## 10 Bootstrap CIs

In some examples above, approximately correct CIs were obtained using robust or delta methods. Bootstrap methods can also be used for CIs. Here are commands to estimate bootstrap CIs for the risk ratio by using a logistic model:

```
. use brcadat, clear
(Breast cancer data)

. program stlogit, rclass
1.      version 10
2.      logistic died low stage2 stage3, nolog
3.      preserve
4.      replace low=0
5.      predict risk0
6.      summ risk0, meanonly
7.      scalar avrisk0 = r(mean)
8.      replace low=1
9.      predict risk1
10.     summ risk1, meanonly
11.     scalar avrisk1 = r(mean)
12.     return scalar lnrr = ln(avrisk1/avrisk0)
13.     restore
14. end

. set seed 93514
```

```
. bootstrap lnrr=r(lnrr), saving(bsanrr3a, replace) reps(400001) nowarn nodots:
> stlogit
(output omitted)
. estat bootstrap, all eform

Bootstrap results                                Number of obs      =       192
                                                Replications       =    400001
```

```
command: stlogit
lnrr: r(lnrr)
```

	Observed exp(b)	Bias	Bootstrap Std. Err.	[95% Conf. Interval]		
lnrr	1.6755989	-.0049751	.37575458	1.079661	2.600476	(N)
				1.068839	2.587837	(P)
				1.071832	2.594932	(BC)

```
(N) normal confidence interval
(P) percentile confidence interval
(BC) bias-corrected confidence interval
```

Stata's `bootstrap` command simplifies the task of estimating bootstrap CIs by using four methods: 1) normal, 2) percentile, 3) bias corrected, and 4) bias corrected and accelerated. Other methods are available ([Carpenter and Bithell 2000](#)). For the risk ratios estimated in this article, the choice among Stata's four methods makes little difference. But for some epidemiologic data, the normal and percentile methods should be used with caution because they may have substantial coverage error ([Efron and Tibshirani 1993](#); [Carpenter and Bithell 2000](#); [Greenland 2004b](#)).

## 11 Risk-ratio methods for matched data

Adjusted risk ratios for matched data can be estimated using conditional Poisson regression, which Stata implements in the `xtpoisson`, `fe` command. I have previously reviewed the analysis of matched cohort data in the *Stata Journal* ([Cummings and McKnight 2004](#)) and elsewhere ([Cummings, McKnight, and Weiss 2003](#); [Cummings, McKnight, and Greenland 2003](#)).

## 12 Summary

When the risk-ratio estimates in this article are rounded to one decimal, nearly all the methods produced estimates of 1.6 or 1.7 (table 2). They also differed little with regard to estimated CIs: the 95% lower bound was 1.0 or 1.1 and the upper bound, 2.3 to 2.6.

The risk ratio that stands out as different came from the substitution method: risk ratio = 1.9 and 95% CI is [1.1, 3.1]. The substitution method has nothing to recommend it; it will usually produce estimates biased away from 1 when outcomes are common, and in Stata it offers little advantage in terms of simplicity. Stata users who wish to estimate an adjusted risk ratio have better methods that they can use, all of which are fairly easy to implement.

Table 2. Risk-ratio estimates for death within 5 years among 192 women with breast cancer, comparing women with low estrogen-receptor-level tumors with women with high estrogen-receptor-level tumors, adjusted for cancer stage at diagnosis. Results are shown using the methods described in this article.

Method	Risk ratio	95% CI	Bootstrap 95% CI†	Akaike information criteria‡
1. Stratified methods				
Mantel-Haenszel	1.62	1.09, 2.39	1.07, 2.48	...
Inverse-variance weights	1.55	1.08, 2.25	...	...
2. Generalized linear regression with log link and binomial distribution				
Maximum likelihood	1.56	1.05, 2.32	...	193.9
Wacholder's truncated method	1.56	1.06, 2.29	1.03, 2.44	...
Copy method	1.56	1.05, 2.31	...	...
3. Generalized linear regression with log link, Gaussian distribution, robust variance estimator	1.55	1.04, 2.32	1.04, 2.43	196.5
4. Poisson regression with robust variance estimator	1.63	1.07, 2.48	1.06, 2.55	226.3
5. Cox proportional hazards with robust variance estimator	1.63	1.07, 2.48	1.06, 2.55	547.1
6. Regression-based standardized risk ratios				
Logistic	1.68	1.09, 2.57	1.07, 2.59	193.9
Probit	1.68	1.09, 2.57	1.07, 2.59	193.9
Complementary log-log	1.67	1.10, 2.52	1.08, 2.56	193.5
Log-log	1.63	1.05, 2.52	1.03, 2.54	194.6
7. Substitution method	1.89	1.12, 2.78	1.08, 3.11	...

† Bias-corrected bootstrap CIs based upon 400001 replications. Not estimated for the inverse-variance stratified method because the `cs` command does not return the pooled risk ratio from this method. Not estimated for the maximum-likelihood version of the generalized linear model with a log link and binomial distribution because convergence failed in many bootstrap samples. Convergence also failed in 100 bootstrap samples (0.025%) using Wacholder's truncated method (`binreg`), 10 samples (0.0025%) using the generalized linear model with a log link and Gaussian distribution, and 5 samples (0.00125%) using the complementary log-log method.

‡ Akaike information criteria statistic for models fit using maximum likelihood. The statistic compares the fitted model with a model that has only the outcome variable. In Stata, smaller Akaike information criteria statistics indicate better fit.

## 13 Acknowledgment

This work was supported by grant R49/CE000197-04 from the Centers for Disease Control and Prevention, Atlanta, GA.

## 14 References

- Armitage, P., G. Berry, and J. N. S. Matthews. 2002. *Statistical Methods in Medical Research*. 4th ed. Oxford: Blackwell.
- Barros, A. J. D., and V. N. Hirakata. 2003. Alternatives for logistic regression in cross-sectional studies: An empirical comparison of models that directly estimate the prevalence ratio. *BioMed Central Medical Research Methodology* 3: 21.  
<http://www.biomedcentral.com/1471-2288/3/21>.
- Blizzard, L., and D. W. Hosmer. 2006. Parameter estimation and goodness-of-fit in log binomial regression. *Biometrical Journal* 48: 5–22.
- Carpenter, J., and J. Bithell. 2000. Bootstrap confidence intervals: When, which, what? A practical guide for medical statisticians. *Statistics in Medicine* 19: 1141–1164.
- Carter, R. E., S. R. Lipsitz, and B. C. Tilley. 2005. Quasi-likelihood estimation for relative risk regression models. *Biostatistics* 6: 39–44.
- Casella, G., and R. L. Berger. 2002. *Statistical Inference*. 2nd ed. Pacific Grove, CA: Duxbury.
- Cook, T. D. 2002. Advanced statistics: Up with odds ratios! A case for odds ratios when outcomes are common. *Academic Emergency Medicine* 9: 1430–1434.
- Cornfield, J. 1951. A method of estimating comparative rates from clinical data: Applications to cancer of the lung, breast, and cervix. *Journal of the National Cancer Institute* 11: 1269–1275.
- Cummings, P. 2009. The relative merits of risk ratios and odds ratios. *Archives of Pediatrics and Adolescent Medicine* 163: 438–445.
- Cummings, P., and B. McKnight. 2004. Analysis of matched cohort data. *Stata Journal* 4: 274–281.
- Cummings, P., B. McKnight, and S. Greenland. 2003. Matched cohort methods for injury research. *Epidemiologic Reviews* 25: 43–50.
- Cummings, P., B. McKnight, and N. S. Weiss. 2003. Matched-pair cohort methods in traffic crash research. *Accident Analysis and Prevention* 35: 131–141.
- Deddens, J. A., and M. R. Petersen. 2008. Approaches for estimating prevalence ratios. *Occupational and Environmental Medicine* 65: 501–506.

- Deddens, J. A., M. R. Petersen, and X. Lei. 2003. Estimation of prevalence ratios when PROC GENMOD does not converge. SAS Users Group International Proceedings. <http://www2.sas.com/proceedings/sugi28/270-28.pdf>.
- Efron, B., and R. J. Tibshirani. 1993. *An Introduction to the Bootstrap*. New York: Chapman & Hall/CRC.
- Flanders, W. D., and P. H. Rhodes. 1987. Large sample confidence intervals for regression standardized risks, risk ratios, and risk differences. *Journal of Chronic Diseases* 40: 697–704.
- Gail, M. H., S. Wieand, and S. Piantadosi. 1984. Biased estimates of treatment effect in randomized experiments with nonlinear regressions and omitted covariates. *Biometrika* 71: 431–444.
- Gourieroux, C., A. Monfort, and A. Trognon. 1984a. Pseudo maximum likelihood methods: Theory. *Econometrica* 52: 681–700.
- . 1984b. Pseudo maximum likelihood methods: Applications to Poisson models. *Econometrica* 52: 701–720.
- Greenland, S. 1987. Interpretation and choice of effect measures in epidemiologic analyses. *American Journal of Epidemiology* 125: 761–768.
- . 1991a. Estimating standardized parameters from generalized linear models. *Statistics in Medicine* 10: 1069–1074.
- . 1991b. Letter to the editor: The author replies. *American Journal of Epidemiology* 133: 964–965.
- . 2004a. Model-based estimation of relative risks and other epidemiologic measures in studies of common outcomes and in case-control studies. *American Journal of Epidemiology* 160: 301–305.
- . 2004b. Interval estimation by simulation as an alternative to and extension of confidence intervals. *International Journal of Epidemiology* 33: 1389–1397.
- Greenland, S., and P. W. Holland. 1991. Estimating standardized risk differences from odds ratios. *Biometrics* 47: 319–322.
- Greenland, S., J. M. Robins, and J. Pearl. 1999. Confounding and collapsibility in causal inference. *Statistical Science* 14: 29–46.
- Greenland, S., and D. C. Thomas. 1982. On the need for the rare disease assumption in case-control studies. *American Journal of Epidemiology* 116: 547–553.
- Hardin, J. W., and J. M. Hilbe. 2007. *Generalized Linear Models and Extensions*. 2nd ed. College Station, TX: Stata Press.

- Hauck, W. W., S. Anderson, and S. M. Marcus. 1998. Should we adjust for covariates in nonlinear regression analyses of randomized trials? *Controlled Clinical Trials* 19: 249–256.
- Holland, P. W. 1989. A note on the covariance of the Mantel–Haenszel log-odds-ratio estimator and the sample marginal rates. *Biometrics* 45: 1009–1016.
- Joffe, M. M., and S. Greenland. 1995. Standardized estimates from categorical regression models. *Statistics in Medicine* 14: 2131–2141.
- Kleinbaum, D. G., L. L. Kupper, and H. Morgenstern. 1982. *Epidemiologic Research: Principles and Quantitative Methods (Industrial Health and Safety)*. New York: Nosstrand Reinhold.
- Lane, P. W., and J. A. Nelder. 1982. Analysis of covariance and standardization as instances of prediction. *Biometrics* 38: 613–621.
- Levin, B. 1991. Letter to the editor. *American Journal of Epidemiology* 133: 963–964.
- Lloyd, C. J. 1999. *Statistical Analysis of Categorical Data*. New York: Wiley.
- Localio, A. R., D. J. Margolis, and J. A. Berlin. 2007. Relative risks and confidence intervals were easily computed indirectly from multivariable logistic regression. *Journal of Clinical Epidemiology* 60: 874–882.
- Lumley, T., R. Kronmal, and S. Ma. 2006. Relative risk regression in medical research: Models, contrasts, estimators, and algorithms. Working Paper 293, UW Biostatistics Working Paper Series. <http://www.bepress.com/uwbiostat/paper293>.
- Mantel, N., and W. Haenszel. 1959. Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute* 22: 719–748.
- McNutt, L.-A., C. Wu, X. Xue, and J. P. Hafner. 2003. Estimating the relative risk in cohort studies and clinical trials of common outcomes. *American Journal of Epidemiology* 157: 940–943.
- Miettinen, O. S., and E. F. Cook. 1981. Confounding: Essence and detection. *American Journal of Epidemiology* 114: 593–603.
- Nelder, J. A. 2001. Letter to the editor. *Statistics in Medicine* 20: 2205.
- Newman, S. C. 2001. *Biostatistical Methods in Epidemiology*. New York: Wiley.
- Nurminen, M. 1981. Asymptotic efficiency of general noniterative estimators of common relative risk. *Biometrika* 68: 525–530.
- Olkin, I. 1998. Letter to the editor. *Evidence-Based Medicine* 3: 71.
- Petersen, M. R., and J. A. Deddens. 2006. Letter to the editor. *American Journal of Epidemiology* 163: 1158–1159.

- Robbins, A. S., S. Y. Chao, and V. P. Fonseca. 2002. What's the relative risk? A method to directly estimate risk ratios in cohort studies of common outcomes. *Annals of Epidemiology* 12: 452–454.
- Rodrigues, L., and B. R. Kirkwood. 1990. Case-control designs in the study of common diseases: Updates on the demise of the rare disease assumption and the choice of sampling scheme for controls. *International Journal of Epidemiology* 19: 205–213.
- Rothman, K. J., S. Greenland, and T. L. Lash. 2008. *Modern Epidemiology*. 3rd ed. Philadelphia: Lippincott Williams & Wilkins.
- Sackett, D. L., J. J. Deeks, and D. G. Altman. 1996. Down with odds ratios! *Evidence-Based Medicine* 1: 164–166.
- Senn, S. 1998. Letter to the editor. *Evidence-Based Medicine* 3: 71.
- . 1999. Rare distinction and common fallacy [letter].  
<http://www.bmj.com/cgi/eletters/317/7168/1318#3089>.
- Steyerberg, E. W., P. M. Bossuyt, and K. L. Lee. 2000. Clinical trials in acute myocardial infarction: Should we adjust for baseline characteristics? *American Heart Journal* 139: 745–751.
- Tarone, R. E. 1981. On summary estimators of relative risk. *Journal of Chronic Diseases* 34: 463–468.
- Wacholder, S. 1986. Binomial regression in GLIM: Estimating risk ratios and risk differences. *American Journal of Epidemiology* 123: 174–184.
- Walter, S. 1998. Letter to the editor. *Evidence-Based Medicine* 3: 71.
- Wooldridge, J. M. 2002. *Econometric Analysis of Cross Section and Panel Data*. Cambridge, MA: MIT Press.
- Zhang, J., and K. F. Yu. 1998. What's the relative risk? A method of correcting the odds ratio in cohort studies of common outcomes. *Journal of the American Medical Association* 280: 1690–1691.
- Zou, G. 2004. A modified Poisson regression approach to prospective studies with binary data. *American Journal of Epidemiology* 159: 702–706.

#### About the author

Peter Cummings is a professor in the Department of Epidemiology in the School of Public Health and a core faculty member at the Harborview Injury Research and Prevention Research Center, University of Washington, Seattle, WA.