

---

# Same-sex sexual behaviour: US frequency estimates from survey data with simultaneous misreporting and non-response

Nathan Berg<sup>a,\*</sup> and Donald Lien<sup>b</sup>

<sup>a</sup>*School of Social Sciences, University of Texas at Dallas, Richardson, Texas, USA*

<sup>b</sup>*University of Texas at San Antonio, San Antonio, Texas, USA*

---

Survey-based research concerning sexual behaviour almost inevitably confronts the simultaneous problems of misreporting and non-response. These problems lead to disparities among estimates of the number and characteristics of those who engage in same-sex sexual behaviour. This paper proposes a statistical model to consistently estimate the frequency of same-sex sexual behaviour in the presence of non-ignorable misreporting and non-response. The model is fitted using 1991–2000 General Social Survey data. Frequency estimates corrected for simultaneous misreporting and non-response are reported. According to the model, 7.1% of US males and 4.1% of females – 15.8 million individuals – are not exclusively heterosexual. Allowing for misreporting and non-response increases the estimated same-sex frequency by more than four million. The model reveals new patterns between misreporting and non-response probabilities and standard demographic variables such as age and income.

## I. Introduction

Do we know what fraction of people engage in same-sex sexual behaviour? Existing estimates in the US span a perplexingly large range, from 1% to 10% of the adult population and beyond (Lauman *et al.*, 1994; Michael *et al.*, 1994; Badgett, 1995; Murray, 1999). Unfortunately, divergent definitions of sexual orientation and disagreements about the meaning

of the label gay (Lauman *et al.*, 1994; Michael *et al.*, 1994; Black *et al.*, 2000; Murray, 1999), result in empirical treatments that do not precisely track the population of interest in this paper – Americans who have had sex with at least one same-sex partner within the last five years.

The population of those who are not exclusively heterosexual plays an important role in a number of policy questions. Forecasts of the spread of AIDS

\*Corresponding author. E-mail: nberg@utdallas.edu

56 and cost-benefit analyses of initiatives to prevent  
 57 the spread of sexually transmitted disease require as  
 58 inputs the size of high-risk populations such as sexu-  
 59 ally active men with same-sex partners (Bloom and  
 60 Glied, 1992; Thomas, 2001). Tabulations of projected  
 61 costs and benefits for legislative proposals to protect  
 62 non-heterosexuals against workplace discrimination,  
 63 such as the Employment Non-Discrimination Act  
 64 which was proposed in the US Congress multiple  
 65 times before being defeated in 1996 (Badgett, 2001),  
 66 require size estimates of the potentially affected popu-  
 67 lation. And anti-sodomy laws, traditionally used to  
 68 prosecute same-sex sexual behaviour (e.g., *Lawrence*  
 69 *and Garner v Texas*), and movements to rescind those  
 70 laws provide yet another example in which quantifi-  
 71 cation of the affected population requires estimates  
 72 for which there exists little reliable data.

73 Related, yet distinct, is the more narrowly  
 74 defined gay population, definitions of which often  
 75 require exclusively homosexual behaviour or public  
 76 self-identification of a gay identity (i.e., being out).  
 77 Previous studies have suggested that sexual orienta-  
 78 tion both at the individual and macro levels are  
 79 indeed economically consequential. Florida (2002)  
 80 argues that cities' gay populations help stimulate  
 81 economic growth. Marketers, religious activists,  
 82 and gay rights advocates who debate the degree  
 83 to which personal income and sexual orientation  
 84 are correlated often agree that gays are economically  
 85 distinct as consumers and workers (Allegretto and  
 86 Arthur, 2001; Arabsheibani *et al.*, 2001; Berg and  
 87 Lien, 2002). Demographers with opposing views  
 88 about the extent and causes of geographic concen-  
 89 tration of gays in particular cities (D'Emilio, 1989;  
 90 Chauncey, 1994; Black *et al.*, 2002) seem to agree  
 91 that differential responses to economic incentives  
 92 by sexual orientation play a role. Those studies  
 93 also inspire questions about whether unconventional  
 94 definitions used to measure gay identity combined  
 95 with systematic misreporting and non-response  
 96 may generate spurious spatial effects. In debates  
 97 surrounding issues such as workplace discrimination  
 98 (Plug and Berkhout, 2004; Black *et al.*, forthcoming),  
 99 gay marriage (Alm *et al.*, 2000), and regulations  
 100 concerning child adoption (Collum, 1993), divergent  
 101 estimates of the gay population's size and  
 102 demographic characteristics clearly contribute to  
 103 persistence of disagreement over policy.

104 Unfortunately, the difficulty of accurately measur-  
 105 ing incidence rates of sexual behaviour and  
 106 sexual orientation raises doubt about the validity  
 107 of existing empirical characterizations of populations  
 108 defined by sexual behaviour and orientation  
 109

(Lauman, *et al.*, 1994; Badgett, 1997). There is a  
 well-known tendency for survey questions about sexual  
 behaviour to elicit non-random patterns of mis-  
 classification and non-response (Pearl and Fairley,  
 1985; Kupek, 1998; Marquie and Baracat, 2000). It  
 is strongly suspected that naive estimates (e.g., those  
 that handle non-response by invoking an unverified  
 missing-at-random assumption) systematically under-  
 count same-sex and gay populations because they  
 have special incentives to misreport and non-respond  
 in many survey settings. Similarly, it is suspected  
 that the empirical distribution of demographic vari-  
 ables among self-reported homosexuals is distorted  
 because those variables are correlated with propensi-  
 ties to misreport and non-respond.

The goal of this paper is to develop a probability  
 model that simultaneously deals with misreporting  
 and non-response and is capable of producing super-  
 ior estimates of the size and characteristics of  
 the same-sex-partner population. The parametric  
 probability model introduced here encompasses  
 the missing-at-random and missing-completely-at-  
 random hypotheses as testable parameter restrictions.  
 In addition to providing improved point estimates  
 of the incidence of non-heterosexuality, the model  
 yields explicit expressions for misreporting and  
 non-response probabilities as functions of observable  
 individual characteristics such as income, age, resi-  
 dential city size, marital and parental status.

The plan of the paper is as follows. Section II  
 reviews related methodological studies of non-  
 response and misreporting. Section III specifies the  
 statistical model and shows how to derive estimates  
 from it. Section IV describes the data, reports  
 estimated probabilities of misreporting and non-  
 response, and provides revised estimates of the  
 non-heterosexual population's frequency and size.  
 Section V concludes with a discussion of the main  
 results and their implications.

## II. Methodological Background

Within the non-response literature, Little and Rubin  
 (1987), Rubin (1987), and Little (1993) distinguish  
 between selection- and pattern-mixture approaches.  
 This paper's model follows the pattern-mixture  
 approach, which assumes that the decisions of respon-  
 dents and non-responders arise from completely  
 different conditional distributions. The argument  
 in favour of this approach is the intuitive appeal  
 of the notion that responders and non-responders  
 are two wholly different groups with covariates that

111 follow entirely different joint distributions. Ekholm  
 112 and Skinner (1998), Forster and Smith (1998),  
 113 and Lee and Marsh (2000) provide applications in  
 114 which the pattern-mixture approach yields clear  
 115 advantages. In contrast, the selection approach  
 116 assumes that a single regression model applies to  
 117 the entire (hypothetically complete) data set and  
 118 appends to it additional equations intended to  
 119 capture the process of selection, most often a prob-  
 120 ability model of the chance of being missing from  
 121 the sample. Applications that successfully adopt  
 122 the selection approach include Heckman (1979),  
 123 Lien and Rearden (1990), Stasny (1991), Conaway  
 124 (1992), Lipsitz *et al.* (1994), Baker (1995), Roy and  
 125 Lin (2002), and Nandram and Choi (2002).

126 In the misclassification literature, economists  
 127 and statisticians have demonstrated that it is not  
 128 necessary to directly observe misreporting in order  
 129 to estimate its frequency (Hausman *et al.*, 1998;  
 130 Black *et al.*, 2000a). This is counterintuitive and  
 131 prompts the question of how rates of misreporting  
 132 can be estimated without observing the phenomenon  
 133 directly. The underlying idea is to exploit information  
 134 contained in the right-hand side variables that  
 135 are correlated with non-response and misreporting  
 136 while including information from incomplete obser-  
 137 vations that are typically discarded in the estimation  
 138 of naive models.

139 The model used here handles misreporting  
 140 and non-response simultaneously whereas previous  
 141 work usually treats them as separate problems.  
 142 An exception is Wu (2002) who models a complex  
 143 error structure generated by a combination of  
 144 imperfectly-measured patient outcomes and prema-  
 145 ture patient dropouts using data from a medical  
 146 study of an AIDS treatment.

147 Survey design methodologists concerned with  
 148 misreporting and non-response have developed  
 149 techniques tailored to situations where researchers  
 150 can influence sample design or collect additional  
 151 data. For example, Embree and Whitehead (1991)  
 152 compare self-reported alcohol consumption with  
 153 aggregate state alcohol sales statistics to produce  
 154 quantitative corrections for raw survey frequencies.  
 155 Referred to as cross validation, this technique has  
 156 revealed significant bias in self-reported church-  
 157 going, charitable giving, and high-risk sexual contact  
 158 (Kupek, 1998; Turner, 1999; Berg, 2005). By  
 159 conducting multiple surveys of the same population  
 160 using different survey instruments, useful predictions  
 161 of the chances of misreporting and non-response  
 162 can be demonstrated using straightforward techni-  
 163 ques (Whitehead *et al.*, 1993). In contrast to  
 164  
 165

situations in which researchers plan for mis-  
 reporting and non-response at the stage of survey  
 design and have access to multiple sources of data,  
 the present model is specialized for secondary  
 analysis of a single set of survey data without  
 requiring further data collection.

### III. The Model

It is assumed that every individual can be categorized  
 as either heterosexual or non-heterosexual. Non-  
 heterosexuals are defined as those who have had at  
 least one same-sex sexual partner, while the hetero-  
 sexual category includes everyone else, including  
 those who have had no sex partners of either gender.  
 The rationale for this simplifying categorization  
 scheme is to focus information in the data on  
 the quantity of interest, the incidence of same-sex  
 sexual behaviour.

The probability that a randomly selected indivi-  
 dual's true behavioural state is non-heterosexual is  
 represented by the function  $g(\phi'x)$ , which is assumed  
 to depend on a  $K \times 1$  vector of covariates  $x$ , not  
 including a constant, and a vector of parameters  $\phi$ .  
 The probability function  $g$  is assumed to be smooth  
 and non-linear in the linear index  $\phi'x$ , as is the  
 case with a logistic or normal pdf. If same-sex  
 behaviour were easy to observe, standard statistical  
 techniques could be used to estimate  $\phi$  and the  
 empirical relationship between observable traits ( $x$ )  
 and sexual behaviour could be established.

The gender(s) of an individual's sexual partners  
 are difficult to observe directly, however. Surveys  
 asking respondents for such information yield  
 observations of a different variable, self-reported  
 sexual behaviour, taking on one of three possible  
 values: 'heterosexual,' 'non-heterosexual' or  
 'no response.' Rather than assuming that survey  
 responses map transparently into true behavioural  
 states, the model allocates positive probability  
 to all six pairs of the two true behavioural states  
 and three self-reported sexual orientations.

By assumption, misreporting and non-response  
 probabilities for individuals whose true behav-  
 ioural state is non-heterosexual depend on  $x$ . The  
 probability  $m(\beta'x)$  represents the chance that a  
 non-heterosexual with demographic characteristics  
 $x$  misreports his or her behaviour as exclusively  
 heterosexual. The non-heterosexual individual's  
 chance of non-response is represented by the  
 probability function  $n(\gamma'x)$ . The vectors of slope  
 parameters,  $\beta$  and  $\gamma$ , transform  $x$  into two linear

166 indexes whose weights indicate the relative  
167 importance of different covariates.

168 The probability of non-response conditional on  
169 an exclusively heterosexual true behavioural state is  
170 given by the function  $r(\rho'x)$ . Non-response among  
171 heterosexuals also depends on  $x$ . The linear weights  
172 that enter the probability function are not con-  
173 strained to coincide with those of non-heterosexuals,  
174 however. Similar to  $g$ , the functions  $m$ ,  $n$  and  $r$  are  
175 pdfs on linear indexes in  $x$ .<sup>1</sup>

176 Economic and non-economic incentives motivating  
177 the hypothesized relationship between  $x$  and prob-  
178 abilities of misreporting and non-response draw  
179 upon a rich, although partially anecdotal, empirical  
180 backdrop. Horrific stories of violence against homo-  
181 sexuals (e.g., the murders of Brandon Teena in Falls  
182 City, Nebraska [1993], Matthew Shepard in Laramie,  
183 Wyoming [1998], and Danny Overstreet in Roanoke,  
184 Virginia [2000]) illustrate a basic motive – averting  
185 physical threat – to be less than fully open about  
186 same-sex sexual activity. Regarding possible labour  
187 market incentives, legal precedent is noted in some  
188 states supporting employers who refuse to hire homo-  
189 sexuals (*England v the City of Dallas*). Moreover,  
190 outspoken criticism of homosexuality by prominent  
191 political leaders (e.g., Pat Buchanan's speech at the  
192 1992 Republican National Convention) suggests links  
193 between characteristics such as geographic location  
194 and socioeconomic status and non-heterosexuals'  
195 propensities to misreport and non-respond.

196 There is comparatively little theoretical support for  
197 asserting a stable relationship between heterosexual  
198 misreporting and survey respondents' demographic  
199 characteristics. Heterosexual misreporting occurs  
200 when an individual who has never had a same-sex  
201 partner incorrectly reports same-sex partners in his  
202 or her sexual history. Heterosexual misreporting is  
203 assumed to occur for highly idiosyncratic reasons  
204 and, therefore, the model assumes independence  
205 between  $x$  and the heterosexual misreporting  
206 frequency  $M$ .

207 Heterosexual non-response is different. There is,  
208 for example, prior evidence that older heterosexuals  
209 and other demographically defined subsets are more  
210 likely to refuse to answer survey questions about  
211 sexual behaviour (Kupek, 1998). Thus, the hetero-  
212 sexual non-response probability  $r(\rho'x)$  is specified as  
213 a function of  $x$ .

214 In this paper individual  $i$ 's true behavioural  
215 state is represented with the symbol  $G_i \in \{\text{hetero},$

non-hetero} and self-reported sexual behaviour as  
 $Y_i \in \{\text{report hetero}, \text{report non-hetero}, \text{no response}\}$ .  
The joint pdf of true states and self-reports can be  
expressed as:

	report		
	report hetero	non-hetero	no response
true hetero	$(1 - g(\phi'x))$ $\times (1 - M$ $- r(\rho'x))$	$(1 - g(\phi'x))M$	$(1 - g(\phi'x))$ $\times r(\rho'x)$
true non-hetero	$g(\phi'x)m(\beta'x)$	$g(\phi'x)(1 - n(\gamma'x)$ $- m(\beta'x))$	$g(\phi'x)n(\gamma'x)$

The marginal probability of truly having at least  
one same-sex partner is obtained by summing hori-  
zontally. The marginal probabilities of reporting  
heterosexual, reporting non-heterosexual, and non-  
responding, are obtained by summing vertically.

Additional structure must be imposed on the func-  
tions  $g$ ,  $m$ ,  $n$  and  $r$  in order for the slope parameters  
 $\phi$ ,  $\beta$ ,  $\gamma$  and  $\rho$  to be identified. To see what can  
go wrong without additional constraints, consider  
the neighbourhood in parameter space about the  
point  $\phi = \beta = \gamma = \rho = 0$ . This point corresponds  
to the situation in which  $x$  is completely uninfor-  
mative in estimating misreporting and non-response  
probabilities, in which case variation in  $x$  is irrelevant  
and there are only two independent pieces of infor-  
mation, the number of self-reported heterosexuals  $N_s$   
and the number of self-reported non-heterosexuals  
 $N_g$ . (Denoting the number of non-responses as  $N_n$   
and the overall sample size as  $N$ , the equation  
 $N = N_s + N_g + N_n$  holds trivially.) The problem is  
that the model requires five parameter estimates,  
the constants  $M$ ,  $g(0)$ ,  $m(0)$ ,  $n(0)$  and  $r(0)$ , using  
only two observable pieces of information,  $N_s$  and  
 $N_g$ . The model is under-identified. Referring to  
the point in parameter space at which all slope  
parameters are zero as the case in which  $x$  is  
completely uninformative, the following restrictions  
are imposed:

**Assumption 1:** If  $x$  is completely uninformative,  
rates of non-response among non-heterosexuals and  
heterosexuals are equal:  $n(0) = r(0)$ .

**Assumption 2:** If  $x$  is completely uninformative,  
rates of misreporting among non-heterosexuals and  
heterosexuals are equal:  $m(0) = M$ .

**Assumption 3:** If  $x$  is completely uninformative,  
the marginal probability of non-response is equal

1It is possible to let  $g$ ,  $m$ ,  $n$ , and  $r$  depend on different sets of covariates. This can be implemented by concatenating all independent variables into the vector  $x$  and imposing zero restrictions on particular elements of  $\phi$ ,  $\beta$ ,  $\gamma$  and  $\rho$ . Since there is no reason to rule out possible connections between  $x$  and reporting probabilities *a priori*, and because identification of the model (discussed subsequently) does not require it, no such restrictions are imposed.

221 to the empirical non-response rate:  $(1 - g(0))r(0) +$   
 222  $g(0)n(0) = (N_n/N)$ .

223 **Assumption 4:** If  $x$  is completely uninformative,  
 224 the marginal non-heterosexual probability is  
 225 equal to the empirical frequency of self-reported  
 226 non-heterosexuality among those who are not  
 227 non-responders:  $(1 - g(0))M + g(0)(1 - m(0) - n(0)) =$   
 228  $(N_g/(N_g + N_s))$ .

229 In assessing how reasonable these assumptions are,  
 230 it is important to point out what they allow and pre-  
 231 cisely what they rule out. A key feature of the model  
 232 is that misreporting and non-response probability  
 233 functions may differ according to whether an indivi-  
 234 dual's true behavioural state is heterosexual or other-  
 235 wise. Assumptions 1 and 2 require that whenever the  
 236 covariates  $x$  contain no information about misreport-  
 237 ing and non-response, then no difference across  
 238 unobserved underlying behavioural states can be  
 239 claimed. This reflects an agnostic prior that begins  
 240 search in parameter space using symmetric guesses  
 241 based on unconditional frequencies. As long as  $x$   
 242 helps predict misreporting and non-response, then  
 243 misreporting and non-response probabilities are free  
 244 to vary across heterosexual and non-heterosexual  
 245 behavioural states. However, when  $x$  contains no  
 246 information, the default misreporting and non-  
 247 response probabilities are symmetric. Thus, the  
 248 model allows the data to decide the extent to which  
 249 there is variation across types without building in  
 250 differential rates of misreporting and non-response.  
 251

252 Similarly, Assumption 3 centres estimated non-  
 253 response probabilities on the empirical non-response  
 254 rate while allowing the data to guide search elsewhere  
 255 using a likelihood criterion. If  $x$  provides no basis for  
 256 adjusting an individual's estimated probability  
 257 of non-response by the likelihood criterion, then  
 258 non-response estimates do not move away from  
 259 the unconditional frequency estimates which are  
 260 based on face-value interpretations of the data.  
 261 Assumption 4 establishes an analogous prior for the  
 262 probability that an individual's true behavioural state  
 263 is non-heterosexual. Guesses are centred at face-value  
 264 empirical frequencies, based on self-reported sexual  
 265 behaviour, and adjusted away from those priors  
 266 only when  $x$  is informative in the sense that at least  
 267 one of the slope parameters in the non-heterosexual  
 268 probability function is non-zero.

269 Assumptions 1–4 imply four functional relation-  
 270 ships between the constants  $g(0)$ ,  $m(0)$ ,  $n(0)$ ,  $r(0)$   
 271 and the parameter  $M$  conditional on the observed  
 272 number of self-reported heterosexuals and non-  
 273 heterosexuals,  $N_s$  and  $N_g$ . This suggests the possibi-  
 274 lity of line search on  $M$  and, at each prospective  
 275 value, maximum likelihood estimation of the

slope parameters. The likelihood function then deci-  
 des which value of  $M$  and corresponding MLE slope  
 estimates are best. This technique is used to derive the  
 estimates reported below. Additional details of the  
 estimation procedure are relegated to the Appendix.  
 The explicit functional dependence of  $g(0)$ ,  $m(0)$ ,  
 $n(0)$ , and  $r(0)$  on  $M$ ,  $N_s$  and  $N_g$  is presented there,  
 along with a modified logistic specification of the  
 probability functions  $g$ ,  $m$ ,  $n$  and  $r$ .

To develop the likelihood function, it is convenient  
 to represent the three possible realizations of self-  
 reported sexual behaviour as a set of three indicator  
 variables:

$$z_{i0} = \begin{cases} 1 & \text{if } i \text{ reports heterosexual} \\ 0 & \text{otherwise} \end{cases}$$

$$z_{i1} = \begin{cases} 1 & \text{if } i \text{ reports non-heterosexual} \\ 0 & \text{otherwise} \end{cases}$$

and

$$z_{i2} \equiv 1 - z_{i0} - z_{i1} = \begin{cases} 1 & \text{if } i \text{ refuses to respond} \\ 0 & \text{otherwise} \end{cases}$$

The likelihood function is:

$$L(\phi, \beta, \gamma, \rho, M|x)$$

$$= \prod_{i=1}^N [(1 - g(\phi'x_i))(1 - M - r(\rho'x_i))$$

$$+ g(\phi'x_i)m(\beta'x_i)]^{z_{i0}} \times [(1 - g(\phi'x_i))M + g(\phi'x_i)$$

$$\times (1 - m(\beta'x_i) - n(\gamma'x_i))]^{z_{i1}} \times [(1 - g(\phi'x_i))r(\rho'x_i)$$

$$+ g(\phi'x_i)n(\gamma'x_i)]^{z_{i2}}. \quad (1)$$

Although it is not explicit in Equation 1, the func-  
 tions  $g$  and  $m$  depend on  $M$  through the constants  
 $g(0)$  and  $m(0)$ .

The strategy for numerically optimizing Equation 1  
 combines line search on  $M$  with steepest descent and  
 Newton–Raphson refinement on the  $4K \times 1$  vector  
 $[\phi' \beta' \gamma' \rho']'$ . For each proposed value of  $M$ , a  
 numerical solution is computed to solve first order  
 conditions (Equations A24–A27 in the Appendix).  
 Among all pairs of proposed  $M$  values and associated  
 slope estimates, the pair that maximizes Equation 1  
 provides a bias-corrected estimate of misreporting  
 and non-response probabilities that vary according  
 to  $x$ . Taking the average estimated probability or  
 the estimated probability evaluated at the average  
 value of  $x$ , an estimated rate of same-sex sexual  
 behaviour follows. Multiplying the estimated rate  
 times the total adult US population provides an  
 estimate of the number of Americans who engage  
 in same-sex sexual behaviour. Standard errors and  
 $t$  statistics for these quantities are computed based  
 on Taylor expansions of the non-linear probability

276 functions with respect to  $[\phi' \beta' \gamma' \rho']'$  and  $M$  and the  
 277 expressions for asymptotic variance they yield.

278 At the point in parameter space where  $x$  is  
 279 uninformative, the model is, by construction, locally  
 280 identified as a result of assumptions 1–4. At other  
 281 points in parameter space, non-linearity of the  
 282 probability functions succeeds in locally identifying  
 283 the model. This can be verified by checking the  
 284 rank of the second derivative matrix of the log-  
 285 likelihood function, an expression for which appears  
 286 in the Appendix. After unsuccessfully attempting  
 287 to prove global identification, numerical rank tests  
 288 at a variety of points in parameter space were  
 289 settled for, including all maximum likelihood  
 290 estimates reported here.

#### 293 IV. Data and Estimation

294  
 295 The model is estimated using GSS data from 1991  
 296 through 2000.<sup>2</sup> Out of a potential pool of 10 458  
 297 survey responses in the GSS from 1991–2000, over  
 298 20% contain one or more missing items among the  
 299 variables in the model. The variables most frequently  
 300 missing are those that code sexual behaviour variable  
 301 and income.

302 Respondents who non-respond to items pertaining  
 303 to sexual behaviour while providing valid responses  
 304 to other items are, of course, included in the sample.  
 305 Survey respondents who non-respond on other items  
 306 aside from those concerning sexual behaviour (giving  
 307 rise to incomplete observations of the independent  
 308 variables  $x$ ) do, however, create a potential selection  
 309 problem. The variables in  $x$  correspond to survey  
 310 items that seek to elicit standard, non-stigmatizing  
 311 demographic information. Therefore, no attempt is  
 312 made to augment what is already a highly parameter-  
 313 ized model by accounting for other possible forms  
 314 of selection bias.

315 After dropping partial responses with missing items  
 316 other than self-reported sexual behaviour, the sample  
 317 size drops to 8446, with 4263 women and 4183 men.  
 318 Respondents with missing self-reported sexual behav-  
 319 iour but otherwise complete item responses remain in  
 320 the sample and are coded with a dependent variable  
 321 distinguishing non-response,  $(z_{0i}, z_{1i}, z_{2i}) = (0, 0, 1)$ ,

322  
 323 <sup>2</sup>It is possible to extend the sample backward in time to 1972. There are at least two reasons for restricting the sample period  
 324 to 1991–2000, however. First, the income variable in the GSS is coded by income brackets and those brackets (both the  
 325 thresholds and the number of categories) have changed periodically from 1972 through 2000. Inflation and changes in the real  
 326 income distribution make it difficult to link bracketed income data over long periods of time. Choosing a sample that begins  
 327 in 1991 avoids all but one of the bracket changes, which occurred between GSS coding schemes in 1996 and 1998. Income data  
 328 are adjusted in the analysis by a simple linear transformation of the initial 21 categories onto the range of the later  
 329 23 categories. Shifting attitudes toward homosexuality (Newport, 2001) provide a second rationale for sampling over  
 330 relatively short time frames. Changes in the anticipated consequences and self-identification habits of non-heterosexuals  
 331 would result in unstable or time-dependent functions  $m$ ,  $n$  and  $r$ .

from the other two possible self-reports of hetero-  
 sexual  $(z_{0i}, z_{1i}, z_{2i}) = (1, 0, 0)$  and non-heterosexual  
 $(z_{0i}, z_{1i}, z_{2i}) = (0, 1, 0)$ . Only 157 women (3.9% exclud-  
 ing non-responders) self-report non-heterosexual, and  
 214 (5.0% of all women in the sample) non-respond.  
 A total of 176 men (4.4% excluding non-responders)  
 self-report as non-heterosexual and 180 (4.3% all men  
 in the sample) non-respond.

Table 1 presents mean values of the independent  
 variables  $x$  broken out by reporting decision and  
 own gender. Standard errors appear in parentheses  
 below each mean. According to Table 1, men who  
 report having same-sex partners earn less than  
 self-reported heterosexual men. In contrast, self-  
 reported non-heterosexual women earn more than  
 self-reported heterosexual women. The average self-  
 reported non-heterosexual is younger than the sample  
 average and lives in a larger city than the average  
 respondent does. Self-reported non-heterosexuals  
 are less likely to have children, less likely to be  
 married, and are better educated (as measured by  
 the number of degrees held, where 0 indicates no  
 high school, 1 indicates completion of high school,  
 2 indicates completion of junior college, 3 indicates  
 college graduate status, and 4 indicates the comple-  
 tion of at least one graduate degree). Non-responders  
 earn less than average, with non-whites dispro-  
 portionately represented among them. The average  
 non-responder is less well educated and is signifi-  
 cantly older than average.

Table 2 reports model estimates based on separate  
 male and female samples. Raw parameter estimates  
 and estimated changes in probabilities are reported  
 expressed in percentage points as the effect of a  
 one-unit change in  $x$  on the chances, respectively,  
 of misreporting and non-response. Raw parameter  
 estimates and transformed probability effects give  
 somewhat different impressions about the relative  
 importance of different regressors. The magnitudes  
 of the effects are variable, and parameters with  
 small  $t$  statistics sometimes have large and statisti-  
 cally significant probability effects.

Estimates under the headings  $\Delta p$  in Table 2  
 are defined in such a way so as to compare the  
 probabilities  $m, n, g$  and  $r$  for two individuals who  
 have opposite values of one characteristic but are

Table 1. Sample means among self-reported heterosexuals, non-heterosexuals and non-responders (GSS 1991–2000)

Variables	Men				Women			
	All	Report hetero	Report non-hetero	Non-respond	All	Report hetero	Report non-hetero	Non-respond
Income <sup>a</sup>	14.92 (0.08)	15.01 (0.08)	13.91 (0.40)	14.07 (0.42)	11.91 (0.09)	11.91 (0.09)	12.47 (0.46)	11.52 (0.41)
Parent	0.66 (0.01)	0.67 (0.01)	0.32 (0.04)	0.65 (0.04)	0.73 (0.01)	0.74 (0.01)	0.50 (0.04)	0.73 (0.03)
White	0.85 (0.01)	0.85 (0.01)	0.80 (0.03)	0.81 (0.03)	0.80 (0.01)	0.81 (0.01)	0.83 (0.03)	0.66 (0.03)
Married	0.58 (0.01)	0.60 (0.01)	0.23 (0.03)	0.60 (0.04)	0.56 (0.01)	0.57 (0.01)	0.31 (0.04)	0.55 (0.03)
Degrees <sup>b</sup>	1.64 (0.02)	1.64 (0.02)	1.85 (0.09)	1.47 (0.09)	1.62 (0.02)	1.63 (0.02)	1.69 (0.10)	1.40 (0.07)
Age	40.19 (0.19)	40.07 (0.20)	39.01 (0.85)	43.78 (0.98)	39.49 (0.18)	39.30 (0.18)	36.18 (0.84)	45.33 (1.05)
City size <sup>c</sup>	0.34 (0.02)	0.32 (0.02)	0.85 (0.14)	0.33 (0.09)	0.36 (0.02)	0.34 (0.02)	0.58 (0.13)	0.43 (0.10)
<i>N</i>	4183	3827	176	180	4263	3892	157	214

## Notes:

<sup>a</sup> Income is measured categorically on a 21-category scale in which the top bracket is US \$75 000 and above, for 1991, 1993, 1994 and 1996. In 1998 and 2000, income is measured on a 23-category scale in which the top bracket is US \$110 000 and above. Surveys were not conducted in 1992, 1995, 1997 and 1999. The categorical variable from the earlier years was transformed to the 23-category scale by assigning top bracket individuals from earlier years to the top bracket for later years, and by adjusting midpoints for inflation (using the CPI-U index), recording the adjusted values on a 23-point scale. The average male's income (in the 13–15 range) corresponds to US \$20 000–US \$30 000 in 1998 dollars. The average female income (in the 11–12 range) corresponds to US \$15 000–US \$20,000.

<sup>b</sup> The variable Degrees is a count variable ranging from zero to four that indicates the number of degrees each respondent has earned.

<sup>c</sup> The variable City size is adjusted to cover the unit interval. City size in its unadjusted form ranges from 1000 to 7.3 million people. The overall average adjusted city size of 0.35 corresponds to a city population of 2.6 million. The average self-reported homosexual male's city size of 0.85 corresponds to a city population 6.2 million. The median city size for all males is approximately 28 000, and 76 000 among self-reported gay males.

otherwise average. Denote the arithmetic mean of  $x$  with its  $k$ th component replaced by the maximum sample value  $\max(x_k)$  as  $\bar{x}_k^{\max}$ . Similarly, denote average  $x$  with its  $k$ th component replaced by the minimum sample value as  $\bar{x}_k^{\min}$ .<sup>3</sup> For a generic probability function  $p \in \{m(\cdot), n(\cdot), g(\cdot), r(\cdot)\}$ , and generic slope parameter  $\eta \in \{\beta, \gamma, \phi, \rho\}$ , define

$$\Delta p_k \equiv p(\eta' \bar{x}_k^{\max}) - p(\eta' \bar{x}_k^{\min}) \quad (2)$$

The estimated asymptotic standard error  $se_{\Delta p_k}$  used in computing the  $t$  statistics labeled  $t_{\Delta p}$  in Table 2 are computed as the square root of:

$$\begin{aligned} & \left[ \partial p(\eta' \bar{x}_k^{\max}) / \partial \eta - \partial p(\eta' \bar{x}_k^{\min}) / \partial \eta \right]' \text{var}(\hat{\eta}) \\ & \times \left[ \partial p(\eta' \bar{x}_k^{\max}) / \partial \eta - \partial p(\eta' \bar{x}_k^{\min}) / \partial \eta \right], \quad (3) \end{aligned}$$

<sup>3</sup> To improve the numerical stability of computer routines used in likelihood maximization, all independent variables were scaled to range over the unit interval. Because most are 0/1 variables, rescaling changed only three variables: degrees, age and city size. Thus,  $\bar{x}_k^{\max}$  and  $\bar{x}_k^{\min}$  are equal to  $\bar{x}$  with the  $k$ th component replaced by 1 or 0.

where the theoretical value of the matrix  $(\hat{\eta})$  is replaced by an outer product estimator.

The top panel of Table 2 shows that moving from the lowest to the highest income category increases a behaviourally homosexual man's misreporting probability by 42 percentage points. Among behaviourally homosexual men, the model suggests that whites are significantly more likely to misreport, as are young males and those who live in small towns. Among female non-heterosexuals, only city size appears to have noticeable effects, with lower chances of misreporting in large cities.

The non-response probability for female non-heterosexuals is generally more sensitive to  $x$  than for males. Moving from the lowest to the highest income bracket, or from the highest to the lowest age group (72 to 18), reduces the average

Table 2. Estimated parameters and changes in probability

	Men ( $N=4183$ )				Women ( $N=4263$ )			
	$\theta$	$t$	$\Delta p$	$t_{\Delta p}$	$\theta$	$t$	$\Delta p$	$t_{\Delta p}$
Effect on Pr (Misclassification Homosexual), $m(\beta'x)$								
Income	15.60	0.6	0.42	0.5	-83.67	-0.1	0.00	-0.1
Parent	9.37	0.5	0.05	0.2	76.47	0.1	0.00	1.1
White	18.79	0.7	0.47	144.0	-42.37	-0.1	-0.00	-0.1
Married	5.05	0.4	0.00	0.1	97.45	0.1	0.00	0.2
Degrees	-12.94	-0.6	-0.16	-0.1	-9.73	-0.1	-0.00	-0.1
Age	-34.03	-0.6	-0.48	-1283.8	84.45	0.1	0.00	0.1
City size	-50.98	-0.5	-0.50	-2577.7	-95.77	-0.1	-0.46	-31.1
Effect on Pr (Nonresponse Homosexual), $n(\gamma'x)$								
Income	0.63	0.2	0.07	0.2	-11.31	-0.4	-0.49	-3.8
Parent	0.61	0.3	0.06	0.3	-2.04	-0.2	-0.19	-0.4
White	-1.49	-0.5	-0.11	-0.5	-23.25	-0.4	-0.47	-34.3
Married	2.11	0.6	0.22	0.8	2.90	0.3	0.20	0.0
Degrees	1.88	0.7	0.17	0.6	6.97	0.4	0.46	1.5
Age	4.80	0.8	0.37	1.2	53.18	0.4	0.50	1338.2
City size	-8.64	-0.5	-0.40	-2.0	-63.88	-0.3	-0.42	-0.8
Effect on Pr (True behavioral state is Homosexual), $g(\phi'x)$								
Income	0.44	0.9	0.02	0.9	-0.16	-0.3	0.00	-0.4
Parent	-1.35	-3.9	-0.07	-3.1	-0.66	-1.9	-0.02	-1.9
White	0.15	0.5	0.01	0.5	-0.11	-0.3	-0.00	-0.5
Married	-1.51	-3.4	-0.08	-3.6	-1.49	-2.5	-0.05	-3.4
Degrees	0.63	1.5	0.03	1.4	0.24	0.6	0.01	0.6
Age	3.25	4.0	0.21	2.3	3.48	3.1	0.15	2.1
City size	0.09	0.2	0.00	0.2	-0.02	0.0	-0.01	-0.0
Effect on Pr (Nonresponse Heterosexual), $r(\rho'x)$								
Income	-1.20	-1.0	-0.04	-0.7	-0.59	-0.7	-0.02	-0.3
Parent	-0.19	-0.3	0.00	-0.1	-0.47	-0.9	-0.02	-0.2
White	-0.06	-0.1	-0.00	-0.0	-1.12	-2.8	-0.05	-1.1
Married	0.42	0.5	0.01	0.2	-0.04	-0.1	-0.00	-0.2
Degrees	-1.79	-1.3	-0.04	-0.9	-1.02	-1.2	-0.03	-0.3
Age	2.02	1.2	0.06	0.3	3.38	2.6	0.15	0.6
City size	0.41	0.4	0.01	0.1	0.50	0.6	0.02	0.1
Estimated Pr (Misreport Heterosexual)								
$M$	0.013	1.7			0.016	1.4		
Pseudo $R^2$	0.08				0.07			
Log likli.	-1364.4				-1414.2			
LL ratio	240.19				197.17			

non-heterosexual woman's non-response probability by more than 45 percentage points. There is also a strong education effect, where non-heterosexual females with more degrees are more likely to non-respond. Among non-heterosexual males, the magnitudes of the effects indicate suggest that only city size has much of an effect on the probability of non-response, a decreasing function of city size.

The third panel of Table 2 presents estimated relationships between the components of  $x$  and the probability of being behaviourally non-heterosexual controlling for misreporting and non-response. As one might expect, being a parent or married in a heterosexual marriage decreases the chance of homosexual behaviour. Recall, however, from Table 1 that

more than 30% of self-reported non-heterosexuals (50% among female non-heterosexuals) are parents, and more than 20% are married. Therefore, it is incorrect to assume that same-sex sexual behaviour and being a parent, or being in a heterosexual marriage, are mutually exclusive. Also noteworthy is that, after correcting for reporting bias, non-heterosexuals do not appear more numerous in large cities. Knowing that an individual lives in a large city appears to increase the chance that he or she will truthfully disclose non-heterosexuality conditional on non-heterosexual behaviour. But there is no evidence that big city residents are actually more likely to engage in same-sex behaviour, contrary to numerous claims in the demographic literature that



441 rely on survey frequencies while making no allowance  
 442 for misreporting and non-response.

443 The fourth panel of Table 2 applies to hetero-  
 444 sexuals and the chance of non-response. Age appears  
 445 to be the single most important factor, with older  
 446 heterosexuals disproportionately refusing to answer  
 447 questions about sex partners and the gender of  
 448 those partners. Small magnitudes and low *t* statistics  
 449 among the other variables suggest that non-response  
 450 among heterosexuals is difficult to predict. The final  
 451 estimate presented in Table 2 is *M*, the probability  
 452 that heterosexuals misreport, which turns out to be  
 453 less than 2% for both men and women.

454 Table 3 presents likelihood ratio test statistics  
 455 for several additional parameter restrictions corre-  
 456 sponding to various notions of random misreporting  
 457 and non-response. The first row of Table 3 suggests  
 458 that *x* is a useful predictor for misreporting among  
 459 non-heterosexuals. The second row tests the missin-  
 460 g-at-random assumption (i.e., the hypothesis that  
 461 non-response mechanisms among heterosexuals and  
 462 non-heterosexuals depend on *x* in the same way)  
 463 leading to rejection. Reported in the next row is  
 464 that the missing-completely-at-random hypothesis  
 465 (i.e., the proposition that non-response is not a  
 466 function of *x*) is rejected as well. The last two rows  
 467 demonstrate that non-heterosexual behavioural  
 468 parameters are jointly significant (at least in the  
 469 statistical sense) as is the model as a whole.

470 Having shown that misreporting and non-  
 471 response mechanisms are predictable in the sense  
 472  
 473

that they arise from a model whose parameters  
 (and therefore derivatives with respect to *x*)  
 are non-zero, Table 4 presents corrected frequency  
 estimates for same-sex behaviour in the USA.  
 According to the model, 7.1% of men have had  
 at least one same-sex sexual partner, significantly  
 more than the self-reported rate of 4.4%. The cor-  
 rected same-sex-behaviour frequency among women  
 of 4.1% is not much higher than the self-reported  
 rate of 3.9%. Nevertheless, rates of misreporting  
 among non-heterosexuals men and women are  
 significantly above zero, and their rates of non-  
 response are considerably higher than heterosexuals.  
 This contrasts sharply with naive frequency  
 estimates based on face-value interpretations of  
 self-reported data. The bias correction technique  
 does not build in or presuppose positive amounts  
 of misreporting and non-response or any differences  
 between heterosexuals and homosexuals. The model  
 is symmetric in that it allows for heterosexuals  
 who self-report as non-heterosexuals (measured  
 by the frequency *M*). Because there are no  
 restrictions on the slope parameters requiring  
 that non-response among heterosexuals be less  
 than among non-heterosexuals, the model could in  
 theory produce a lower same-sex frequency than  
 is estimated naively from self reports. Therefore,  
 the higher estimated number of individuals engag-  
 ing in same-sex behaviour reflects information  
 present in the data rather than priors built into  
 the model.

474 **Table 3. Likelihood ratio test statistics for misreporting/missing at random hypotheses**

477 Hypothesis		476 Men		476 Women		477 df
		477 Lik. ratio	477 <i>p</i> -value	477 Lik. ratio	477 <i>p</i> -value	
478 Misreporting is random w.r.t. <i>x</i>	$\beta = 0$	27.3	0.0003	23.5	0.0014	7
479 Missing at random	$\gamma = \rho$	29.9	0.0001	21.6	0.0030	7
480 Missing completely at random	$\gamma = \rho = 0$	40.8	0.0002	30.3	0.0070	14
481 Misrep. and nonresp. independent of <i>x</i>	$\beta = \gamma = \rho = 0$	70.0	0.0000	127.1	0.0000	21
482 Sexual orientation independent of <i>x</i>	$\phi = 0$	144.6	0.0000	71.4	0.0000	7
483 All slopes zero	$\beta = \gamma = \phi = \rho = 0$	240.2	0.0000	197.2	0.0000	28

486 **Table 4. Estimated homosexual, misreporting and nonresponse frequencies**

488 Bias-corrected 489 frequency: 490 estimates	488 Men			488 Women			488 Naive formula
	489 $\Sigma_i p(n'_i x_i)/n$	489 se	489 Naive	489 $\Sigma_i p(n'_i x_i)/n$	489 se	489 Naive	
491 Homosexual ( <i>g</i> )	0.071	0.014	0.043	0.041	0.009	0.039	$N_g/(N_s + N_g)$
492 Misreport Homosexual ( <i>m</i> )	0.41	0.002	0	0.39	0.001	0	–
493 Nonrespond Homosexual ( <i>n</i> )	0.33	0.233	0.043	0.21	0.174	0.05	$N_n/N$
494 Nonrespond Heterosexual ( <i>r</i> )	0.03	0.051	0.043	0.04	0.061	0.05	$N_n/N$
495 Misreport Heterosexual ( <i>M</i> )	0.013	0.008	0	0.016	0.012	0	–

## V. Conclusion

Previous attempts at measuring the incidence of non-heterosexual behaviour have produced highly variable results. The range has been between 1% and 10% according to Badgett (1995), although larger estimates are not unheard of (Kinsey *et al.*, 1948; Sorenson, 1973; Dover, 1978). None of the previous attempts deals effectively with the problems of misreporting and non-response. Against this background of inconsistent empirical findings and questions of methodological reliability, the model in this paper attempts to simultaneously deal with misreporting and non-response within a unified maximum-likelihood framework.

According to the misreporting- and non-response-corrected model, 7.1% of men and 4.1% of women have had at least one same-sex sexual partner. That translates into 15.8 million non-heterosexual Americans: 5.9 ( $\pm 1.3$ ) million women and 10.0 ( $\pm 1.9$ ) million men based on total US Census Bureau figures from July 2001. More than four million non-heterosexuals go uncounted using the naive model.

The estimates reported above demonstrate interesting patterns that condition probabilities of misreporting and non-response. Among non-heterosexual males, misreporting is most likely for those who have high income, are young, and live in small cities. Among non-heterosexual females, non-response is most likely among those who have low income, are older, and hold more academic degrees. Older heterosexuals are disproportionately represented among non-responders. These relationships suggest that face-value interpretations of survey data concerning sexual behaviour distort the true joint distribution of behavioural variables and their covariates. The data reject formal tests of both strong and weak forms of the missing-at-random hypothesis.

Four assumptions were required to estimate the model. The assumptions effectively require the model to revert to face-value frequencies whenever the independent variables fail to be informative. Near the point in parameter space where the hypothesized covariates  $x$  do not influence misreporting and non-response probabilities, the assumptions force the model to produce estimated same-sex frequencies equal to empirical face-value same-sex frequencies. Similarly, the assumptions dictate that predicted non-response probabilities coincide with observed non-response rates whenever the coefficients on  $x$  are zero. This symmetry allows the data to decide whether conditioning on  $x$  changes the likelihood that an individual will misreport or non-respond.

Whether or not to adjust away from face-value interpretations based on self-reported frequencies and, if so, in which direction, depends on the data. Adjustment can occur in either direction.

There are no parameter restrictions preventing the model from estimating lower same-sex frequencies than those derived from naive, or face-value, computations. Thus, the paper's main finding that, after accounting for misreporting and non-response, there are four million more behaviourally non-heterosexual individuals in the USA – fully one-third more than would be estimated otherwise – reflects information derived solely from the data rather than priors built into the model. By simultaneously dealing with misreporting and non-response within an otherwise conventional maximum likelihood framework, the estimators enjoy consistency and asymptotic efficiency, conditional on correct model specification. A similar modelling approach could be applied in other settings where misreporting and non-response are suspected to be significant problems.

## Acknowledgements

The authors thank Brian Bucks, Ted Juhl, Joseph Sicilian, and Wim Vijverberg for helpful comments.

## References

- Allegretto, S. and Arthur, M. (2001) An empirical analysis of homosexual/heterosexual male earnings differentials: unmarried and unequal?, *Industrial and Labor Relations Review*, **54**, 631–46.
- Alm, J., Badgett, M. V. L. and Whittington, L. A. (2000) Wedding bell blues: the income tax consequences of legalizing same-sex marriage, *National Tax Journal*, **53**, 201–14.
- Arabsheibani, G. R., Marin, A. and Wadsworth, J. (2001) Gays' pay in the UK, Working Paper, University of Wales, Aberystwyth.
- Badgett, L. M. V. (1995) The wage effects of sexual orientation discrimination, *Industrial Labor Relations Review*, **48**, 726–39.
- Badgett, L. M. V. (1997) Beyond biased samples: Challenging the myths on the economic status of lesbians and gay men, in *Homo Economics: Capitalism, Community and Lesbian and Gay Life* (Eds) A. Gluckman and B. Reed, Routledge, London, pp. 65–72.
- Badgett, L. M. V. (2001) *Money, Myths and Change: The Economic Lives of Lesbians and Gay Men*, University of Chicago Press, Chicago.
- Baker, S. G. (1995) Marginal regression for repeated binary data with outcome subject to non-ignorable non-response, *Biometrics*, **51**, 1042–52.
- Berg, N. (2005) Non-response bias, in *Encyclopedia of Social Measurement*, vol. 2 (Ed.) K. Kempf-Leonard, Academic Press, London, pp. 865–73.

- 551 Berg, N. and Lien, D. (2002) Measuring the effect of sexual  
552 orientation on income: evidence of discrimination?,  
553 *Contemporary Economic Policy*, **20**, 394–414.
- 554 Black, D. A., Berger, M. C. and Scott, F. A. (2000a)  
555 Bounding parameter estimates with nonclassical  
556 measurement error, *Journal of the American  
557 Statistical Association*, **95**, 739–48.
- 558 Black, D., Gates, C., Sanders, S. and Taylor, L. (2000b)  
559 Demographics of the gay and lesbian population in  
560 the United States: evidence from available systematic  
561 data sources, *Demography*, **37**, 139–54.
- 562 Black, D., Gates, C., Sanders, S. and Taylor, L. (2002)  
563 Why do gay men live in San Francisco?, *Journal of  
564 Urban Economics*, **51**, 54–76.
- 565 Black, D. A., Makar, H. R., Sanders, S. G. and Taylor, L.  
566 (forthcoming) The effects of sexual orientation on  
567 earnings, *Industrial Labor Relations Review*.
- 568 Bloom, D. E. and Glied, S. (1992) Projecting the number of  
569 new AIDS cases in the United States, *International  
570 Journal of Forecasting*, **8**, 339–65.
- 571 Chauncey, G. (1994) *Gay New York: Gender, Urban Culture  
572 and the Making of the Gay Male World, 1890–1940*,  
573 Basic, New York.
- 574 Collum, C. S. (1993) Co-parent adoptions: lesbian and  
575 gay parenting, *Trial*, **29**, 28.
- 576 Conaway, M. R. (1992) The analysis of repeated categorical  
577 measurements subject to nonignorable non-response,  
578 *Journal of the American Statistical Association*, **87**,  
579 817–24.
- 580 D’Emilio, J. (1989) Gay politics and community in  
581 San Francisco since World War II, in *Hidden From  
582 History: Reclaiming the Gay and Lesbian Past*  
583 (Eds) M. B. Duberman *et al.*, New American  
584 Library, New York, pp. 456–73.
- 585 Dover, K. J. (1978) *Greek Homosexuality*, Harvard  
586 University Press, Cambridge.
- 587 Ekholm, A. and Skinner, C. (1998) The muscatine  
588 children’s obesity data reanalyzed using pattern  
589 mixture models, *Applied Statistics*, **47**, 251–63.
- 590 Embree, B. G. and Whitehead, P. C. (1991) Validity  
591 and reliability of self-reported drinking behaviour:  
592 dealing with the problem of response bias, *Journal of  
593 Studies on Alcohol*, **54**, 334–44.
- 594 Florida, R. (2002) *The Rise of the Creative Class*, Basic  
595 Books, New York.
- 596 Forster, J. J. and Smith, P. W. F. (1998) Model-based  
597 inference for categorical survey data subject to  
598 non-ignorable non-response, *Journal of the Royal  
599 Statistical Society Ser. B*, **60**, 57–70.
- 600 Hausman, J. A., Abrevaya, J. and Scott-Morton, F. M.  
601 (1998) Misclassification of the dependent variable in  
602 a discrete-response setting, *Journal of Econometrics*,  
603 **87**, 239–69.
- 604 Heckman, J. J. (1979) Sample selection bias as a specifica-  
605 tion error, *Econometrica*, **47**, 153–61.
- 606 Kinsey, A. C., Pomeroy W. B. and Martin, C. E. (1948)  
607 *Sexual Behaviour in the Human Male*, Saunders,  
608 Philadelphia.
- 609 Kupek, E. (1998) Determinants of item non-response  
610 in a large national sex survey, *Archives of Sexual  
611 Behaviour*, **27**, 581–89.
- 612 Laumann, E. O., Gagnon, J. H., Michael, R. T.  
613 and Michaels, S. (1994) *The Social Organization  
614 of Sexuality: Sexual Practices in the United States*,  
615 University of Chicago Press, Chicago.
- 616 Lee, B. and Marsh, L. C. (2000) Sample selection bias  
617 correction for missing response observations, *Oxford  
618 Bulletin of Economics and Statistics*, **62**(2), 305–22.
- 619 Lien, D. and Rearden, D. (1990) Missing measurements  
620 in discrete response models, *Economics Letters*, **32**,  
621 231–35.
- 622 Lipsitz, S. R., Laird, N. M. and Harrington, D. P. (1994)  
623 Weighted least squares analysis of repeated categorical  
624 measurement with outcomes subject to non-response,  
625 *Biometrics*, **50**, 11–24.
- 626 Little, R. J. (1993) Pattern-mixture models for multivariate  
627 incomplete data, *Journal of the American Statistical  
628 Association*, **88**, 125–34.
- 629 Little, R. J. A. and Rubin, D. B. (1987) *Statistical Analysis  
630 With Missing Data*, Wiley, New York.
- 631 Marquie, J. C. and Baracat, B. (2000) Effects of age, educa-  
632 tion and sex on response bias in a recognition task,  
633 *Journal of Gerontology*, **55B**(5), 266–72.
- 634 Michael, R. T., Gagnon, J. H., Laumann, E. O. and  
635 Kolata, G. (1994) *Sex in America: A Definitive Study*,  
636 University of Chicago Press, Chicago.
- 637 Murray, S. O. (1999) *Homosexualities*, University of  
638 Chicago Press, Chicago.
- 639 Nandram, B. and Choi, J. W. (2002) Hierarchical Bayesian  
640 non-response models for binary data from small areas  
641 with uncertainty about ignorability, *Journal of the  
642 American Statistical Association*, **97**, 381–388.
- 643 Newport, F. (2001) American attitudes toward homosexu-  
644 ality continue to become more tolerant, *The Gallup  
645 Organization*, < [http://www.gallup.com/subscription/  
646 ?m=f&c\\_id=10680](http://www.gallup.com/subscription/?m=f&c_id=10680) >, June, 4.
- 647 Pearl, D. K. and Fairley, D. (1985) Testing for the potential  
648 for non-response bias in sample surveys, *Public  
649 Opinion Quarterly*, **49**(4), 553–60.
- 650 Plug, E. and Berkhout, P. (2004) Effects of sexual  
651 preferences on earnings in the Netherlands, *Journal  
652 of Population Economics*, **17**, 117–31.
- 653 Rubin, D. B. (1987) *Multiple Imputation for Non-response  
654 in Surveys*, Wiley, New York.
- 655 Roy, J. and Lin, X. (2002) Analysis of multivariate  
656 longitudinal outcomes with nonignorable dropouts  
657 and missing covariates: changes in methadone treat-  
658 ment practices, *Journal of the American Statistical  
659 Association*, **97**, 49–52.
- 660 Sorenson, R. C. (1973) *Adolescent Sexuality in  
661 Contemporary America*, World Publishing, New York.
- 662 Stasny, E. A. (1991) Hierarchical models for the probabili-  
663 ties of a survey classification and non-response: an  
664 example from the national crime survey, *Journal of  
665 the American Statistical Association*, **86**, 296–303.
- 666 Thomas, R. (2001) Estimated population mixing by  
667 country and risk cohort for the HIV/AIDS epidemic  
668 in Western Europe, *Journal of Geographical Systems*,  
669 **3**, 283–301.
- 670 Turner, H. A. (1999) Participation bias in AIDS-related  
671 telephone surveys: results from the national AIDS  
672 behavioural survey (NABS) non-response study,  
673 *The Journal of Sex Research*, **36**, 52–66.
- 674 Whitehead, J. C., Groothuis, P. A. and Blomquist, G. C.  
675 (1993) Testing for non-response and sample selection  
676 bias in contingent valuation: analysis of a combination  
677 phone/mail survey, *Economics Letters*, **41**, 215–20.
- 678 Wu, L. (2002) A joint model for nonlinear mixed-effects  
679 models with censoring and covariates measured with  
680 error, with application to AIDS studies, *Journal of the  
681 American Statistical Association*, **97**, 955–64.

3

4

## Appendix

The empirical results reported in the body of the paper are derived using modified logistic functional forms as follows:

$$g(\phi'x) = (1 + a_g e^{-\phi'x})^{-1} \quad (\text{A1})$$

$$m(\beta'x) = 0.5(1 + a_m e^{-\beta'x})^{-1} \quad (\text{A2})$$

$$n(\gamma'x) = 0.5(1 + a_n e^{-\gamma'x})^{-1} \quad (\text{A3})$$

$$r(\rho'x) = 0.5(1 + a_r e^{-\rho'x})^{-1}, \quad (\text{A4})$$

where the constants  $a_g$ ,  $a_m$ ,  $a_n$  and  $a_r$  are given by the formulas:

$$a_g = \frac{1 - 2M - N_n/N}{N_g/(N_s + N_g) - M} - 1 \quad (\text{A5})$$

$$a_m = 1/M - 1 \quad (\text{A6})$$

$$a_n = N/N_n - 1 \quad (\text{A7})$$

$$a_r = N/N_n - 1. \quad (\text{A8})$$

The presence of the term '0.5' arises from requiring the probabilities to lie in the unit interval, explained as follows. The expression  $1 - m(\beta'x) - n(\gamma'x)$  is the probability of self-reporting non-heterosexual conditional on non-heterosexual behaviour. Because it is a probability, it must be bounded in the unit interval. The idea that either misreporting or non-response would exceed 50% so wildly contradicts the priors, that is applied a parameterization that bounds  $m$  and  $n$  from above by 0.50. The upper bound of 0.5 is also imposed on the heterosexual refusal probability  $r$  out of symmetry considerations, so that the non-response parameters  $\rho$  and  $\gamma$  can be more easily compared.

Alternative specifications that do away with the upper bound restrictions are possible, such as the following:

$$g(\phi'x) = (1 + a_g e^{-\phi'x})^{-1} \quad (\text{A9})$$

$$m(\beta'x, \gamma'x) = a_m e^{\beta'x} / (1 + a_m e^{\beta'x} + a_n e^{\gamma'x})^{-1} \quad (\text{A10})$$

$$n(\beta'x, \gamma'x) = a_n e^{\gamma'x} / (1 + a_m e^{\beta'x} + a_n e^{\gamma'x}) \quad (\text{A11})$$

$$r(\rho'x) = (1 + a_r e^{-\rho'x})^{-1}. \quad (\text{A12})$$

This specification ensures that the probability non-heterosexuals correctly report their sexual behaviour  $1 - m(\cdot) - n(\cdot)$  is contained in the unit interval for any value of  $M$  (in the unit interval) without limiting the range of  $m(\cdot)$  or  $n(\cdot)$ . In practice, this model encounters problems with numerical stability in computing maximum likelihood estimates. Therefore, the

analysis proceeds under the specification given by (A1)–(A3).

A second detail worth clarifying is the role of the constants  $a_g$ ,  $a_m$ ,  $a_n$  and  $a_r$ . They serve to centre the respective probabilities at their corresponding empirical face-value frequencies when the slope parameters are zero, in accordance with assumptions 1–4. Those assumptions result in the equations:

$$g(0) = \left( \frac{N_g}{N_s + N_g} - M \right) / \left( 1 - \frac{N_n}{N} - 2M \right) \quad (\text{A13})$$

$$m(0) = M \quad (\text{A14})$$

$$n(0) = \frac{N_n}{N} \quad (\text{A15})$$

$$r(0) = \frac{N_n}{N} \quad (\text{A16})$$

In finite samples the numerator and denominator in the formula for  $g(0)$ , Equation A30, may be either negative or positive. To guarantee that  $g(0) > 0$ ,  $M$  must be chosen such that one of two inequalities holds:

$$0 < M < \min \left\{ \frac{N_g}{N_s + N_g}, \frac{1}{2} \left( 1 - \frac{N_n}{N} \right) \right\} \quad \text{or} \\ \max \left\{ \frac{N_g}{N_s + N_g}, \frac{1}{2} \left( 1 - \frac{N_n}{N} \right) \right\} < M < 1$$

Asymptotically,  $g(0)$  approaches the population frequency  $\text{plim}[g(\phi'x)] \equiv g$  with respect to sample size  $N$ , as the slope parameters approach zero. Therefore,  $g(0)$  must be positive in the limit. To see this, denote the plims of  $g(0)$ ,  $m(0)$ ,  $n(0)$  and  $r(0)$  as  $g_l$ ,  $m_l$ ,  $n_l$  and  $r_l$ , respectively, and use the calculations

$$\text{plim} \left[ \frac{N_g}{N_s + N_g} \right] = g_l(1 - m_l - n_l) + (1 - g_l)M,$$

and

$$\text{plim} \left[ \frac{N_n}{N} \right] = g_l m_l + (1 - g_l) r_l \quad (\text{A17})$$

to substitute into the right hand side of Equation A20. Thus,

$$\text{plim}[g(0)] = \frac{g_l(1 - m_l - n_l) + (1 - g_l)M - M}{1 - 2M - g_l m_l - (1 - g_l) r_l} \quad (\text{A18})$$

The zero slope restriction together with assumptions 1–4 allows the substitution of  $m_l$  for  $M$  and  $n_l$  for  $r_l$ , implying that

$$\text{plim}[g(0)] = g_l |_{\phi=\beta=\gamma=\rho=0} \quad (\text{A19})$$

Together, the four assumptions in equations (A1)–(A4) imply a functional relationship between

661  $g(0)$ ,  $m(0)$ ,  $n(0)$ , and  $r(0)$  and  $M$ , conditional on the  
 662 observed number of self-reported heterosexuals and  
 663 non-heterosexuals,  $N_s$  and  $N_g$ :

$$664 \quad g(0) = \left( \frac{N_g}{N_s + N_g} - M \right) / \left( 1 - \frac{N_n}{N} - 2M \right) \quad (\text{A20})$$

$$667 \quad m(0) = M \quad (\text{A21})$$

$$669 \quad n(0) = \frac{N_n}{N} \quad (\text{A22})$$

$$672 \quad r(0) = \frac{N_n}{N} \quad (\text{A23})$$

673 The likelihood function for the parameters  
 674  $[\phi' \beta' \gamma' \rho']'$  and  $M$ , Equation 1 in the body of the  
 675 paper, has first order conditions as follows:

$$677 \quad \frac{\partial L}{\partial \phi} = \sum_{i=1}^N \left\{ \frac{z_{0i}}{A_{0i}} [-(1 - M - r_i) + m_i] \right. \\
 678 \quad \quad \quad + \frac{z_{1i}}{A_{1i}} [-M + (1 - m_i - n_i)] \\
 681 \quad \quad \quad \left. + \frac{z_{2i}}{A_{2i}} [-r_i + n_i] \right\} \frac{\partial g_i}{\partial \phi} = 0 \quad (\text{A24})$$

$$685 \quad \frac{\partial L}{\partial \beta} = \sum_{i=1}^N \left\{ \frac{z_{0i}}{A_{0i}} - \frac{z_{1i}}{A_{1i}} \right\} g_i \frac{\partial m_i}{\partial \beta} = 0 \quad (\text{A25})$$

$$688 \quad \frac{\partial L}{\partial \gamma} = \sum_{i=1}^N \left\{ -\frac{z_{1i}}{A_{1i}} + \frac{z_{2i}}{A_{2i}} \right\} g_i \frac{\partial n_i}{\partial \gamma} = 0 \quad (\text{A26})$$

$$692 \quad \frac{\partial L}{\partial \rho} = \sum_{i=1}^N \left\{ -\frac{z_{0i}}{A_{0i}} + \frac{z_{2i}}{A_{2i}} \right\} (1 - g_i) \frac{\partial r_i}{\partial \rho} = 0 \quad (\text{A27})$$

where

$$\frac{\partial g_i}{\partial \phi} = \frac{a_g e^{-\phi' x_i}}{(1 + a_g e^{-\phi' x_i})^2} x_i \quad (\text{A28})$$

$$\frac{\partial m_i}{\partial \beta} = \frac{0.5 a_m e^{-\beta' x_i}}{(1 + a_m e^{-\beta' x_i})^2} x_i \quad (\text{A29})$$

$$\frac{\partial n_i}{\partial \phi} = \frac{0.5 a_n e^{-\gamma' x_i}}{(1 + a_n e^{-\gamma' x_i})^2} x_i \quad (\text{A30})$$

$$\frac{\partial r_i}{\partial \rho} = \frac{0.5 a_r e^{-\rho' x_i}}{(1 + a_r e^{-\rho' x_i})^2} x_i, \quad (\text{A31})$$

and

$$A_{0i} = (1 - g(\phi' x_i))(1 - M - r(\rho' x_i)) \\
 + g(\phi' x_i)m(\beta' x_i) \quad (\text{A32})$$

$$A_{1i} = (1 - g(\phi' x_i))M + g(\phi' x_i) \\
 \times (1 - m(\beta' x_i) - n(\gamma' x_i)) \quad (\text{A33})$$

$$A_{2i} = (1 - g(\phi' x_i))r(\rho' x_i) + g(\phi' x_i)n(\gamma' x_i). \quad (\text{A34})$$

The estimates presented in the paper are solutions in  $[\phi' \beta' \gamma' \rho']'$  and  $M$  to the system of equations above. Expressions for the second derivative matrix from which variance estimators were computed is available from the authors upon request.