

**Web-based Supplementary Materials for “Estimating Disease Prevalence Using Relatives of Case and Control Probands” by Javaras, Laird, Hudson, & Ripley**

**Web Appendix A. Proof that  $\hat{\pi}$  is approximately unbiased at the first-order for  $\pi$**

The overall population prevalence is defined as  $\pi \equiv f(Y_{ij} = 1)$ , where individual  $ij$  is randomly selected from the population of interest. Assumption (i) about the availability of relatives allows us to expand  $f(Y_{ij} = 1)$  as follows

$$\begin{aligned}\pi &\equiv f(Y_{ij} = 1) \\ &= f(Y_{ij} = 1|Y_{ij'} = 1)f(Y_{ij'} = 1) + f(Y_{ij} = 1|Y_{ij'} = 0)f(Y_{ij'} = 0),\end{aligned}\quad (\text{A.1})$$

where individual  $ij'$  is randomly selected from among  $Y_{ij}$ 's relatives with disease status  $Y_{ij'}$ . (In the remainder of this proof, we will assume that  $j' \neq j$ .) We can rewrite the above equation as

$$\pi = f(Y_{ij} = 1|Y_{ij'} = 1)\pi + f(Y_{ij} = 1|Y_{ij'} = 0)(1 - \pi),$$

which can be rearranged to give

$$\pi = \frac{\pi_U}{1 - \pi_A + \pi_U},\quad (\text{A.2})$$

where  $\pi_U \equiv f(Y_{ij} = 1|Y_{ij'} = 0)$  and  $\pi_A \equiv f(Y_{ij} = 1|Y_{ij'} = 1)$ . The parameters  $\pi_A$  and  $\pi_U$  can be defined in terms of the finite population:

$$\begin{aligned}\pi_A &\equiv \frac{f(Y_{ij} = 1|Y_{ij'} = 1)}{\sum_{i=1}^F \sum_{j=1}^{N_i} \sum_{j' \neq j} \mathbf{I}(Y_{ij} = 1) \mathbf{I}(Y_{ij'} = 1)} \\ &= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} \sum_{j' \neq j} \mathbf{I}(Y_{ij} = 1) \mathbf{I}(Y_{ij'} = 1)}{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1)(N_i^A - 1)} \\ &= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} (\mathbf{I}(Y_{ij} = 0)N_i^A + \mathbf{I}(Y_{ij} = 1)(N_i^A - 1))}{\sum_{i=1}^F (N_i^A - 1)N_i^A};\end{aligned}\quad (\text{A.3})$$

where  $N_i^A = \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1)$ , the number of affected members in family  $i$ , and

$$\begin{aligned}
\pi_U &\equiv f(Y_{ij} = 1 | Y_{ij'} = 0) \\
&= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} \sum_{j' \neq j} \mathbf{I}(Y_{ij} = 1) \mathbf{I}(Y_{ij'} = 0)}{\sum_{i=1}^F \sum_{j=1}^{N_i} \sum_{j' \neq j} \mathbf{I}(Y_{ij'} = 0)} \\
&= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1) N_i^U}{\sum_{i=1}^F \sum_{j=1}^{N_i} (\mathbf{I}(Y_{ij} = 0)(N_i^U - 1) + \mathbf{I}(Y_{ij} = 1)N_i^U)} \\
&= \frac{\sum_{i=1}^F N_i^A N_i^U}{\sum_{i=1}^F (N_i - 1)N_i^U}, \tag{A.4}
\end{aligned}$$

where  $N_i^U = \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 0)$ , the number of unaffected members in family  $i$ .

Now, recall that

$$\hat{\pi} = \frac{p_U}{1 - p_A + p_U}, \tag{A.5}$$

where  $p_A$  is the proportion of case probands' relatives who are affected,

$$p_A = \frac{\sum_{i^*=1}^{F_A} \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(Y_{i^*j^*} = 1)}{\sum_{i^*=1}^{F_A} \sum_{j^*=2}^{n_{i^*}} 1};$$

and  $p_U$  is the proportion of control probands' relatives who are affected,

$$p_U = \frac{\sum_{i^*=F_A+1}^{F_A+F_U} \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(Y_{i^*j^*} = 1)}{\sum_{i^*=F_A+1}^{F_A+F_U} \sum_{j^*=2}^{n_{i^*}} 1}.$$

The estimator  $\hat{\pi}$  in (A.5) can be approximated by a second-order Taylor expansion

around  $\mathbb{E}p_U$  and  $\mathbb{E}p_A$ , the expected values of  $p_U$  and  $p_A$ , respectively:

$$\begin{aligned}\widehat{\pi} &\approx \frac{\mathbb{E}p_U}{1 - \mathbb{E}p_A + \mathbb{E}p_U} + (p_U - \mathbb{E}p_U) \left. \frac{\partial \widehat{\pi}}{\partial p_U} \right|_{\mathbb{E}p_U, \mathbb{E}p_A} + (p_A - \mathbb{E}p_A) \left. \frac{\partial \widehat{\pi}}{\partial p_A} \right|_{\mathbb{E}p_U, \mathbb{E}p_A} \\ &\quad + \frac{1}{2} (p_U - \mathbb{E}p_U)^2 \left. \frac{\partial^2 \widehat{\pi}}{\partial p_U^2} \right|_{\mathbb{E}p_U, \mathbb{E}p_A} + \frac{1}{2} (p_A - \mathbb{E}p_A)^2 \left. \frac{\partial^2 \widehat{\pi}}{\partial p_A^2} \right|_{\mathbb{E}p_U, \mathbb{E}p_A} \\ &\quad + (p_U - \mathbb{E}p_U) (p_A - \mathbb{E}p_A) \left. \frac{\partial^2 \widehat{\pi}}{\partial p_U \partial p_A} \right|_{\mathbb{E}p_U, \mathbb{E}p_A}\end{aligned}$$

Inserting expressions for the derivatives and then taking the expectation of both sides of the above equation yields

$$\begin{aligned}E\widehat{\pi} &\approx \frac{\mathbb{E}p_U}{1 - \mathbb{E}p_A + \mathbb{E}p_U} - \frac{\text{Var}(p_U)(1 - \mathbb{E}p_A)}{(1 - \mathbb{E}p_A + \mathbb{E}p_U)^3} + \frac{\text{Var}(p_A)\mathbb{E}p_U}{(1 - \mathbb{E}p_A + \mathbb{E}p_U)^3} \quad (\text{A.6}) \\ &\quad + \frac{\text{Cov}(p_U, p_A)(1 - \mathbb{E}p_A - \mathbb{E}p_U)}{(1 - \mathbb{E}p_A + \mathbb{E}p_U)^3}\end{aligned}$$

In order to determine the bias in the leading term on the right-hand side of (A.6), we must derive expressions for the bias of  $p_A$  and  $p_U$  as estimators for  $\pi_A$  and  $\pi_U$ , respectively. Beginning with the former, we introduce indicators in order to rewrite  $p_A$  as a sum over every member of every family in the population, except for one affected member of each family who is arbitrarily designated as the (case) proband:

$$p_A = \frac{\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A \mathbb{I}(Y_{ij} = 1)}{\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A}$$

where  $J_i^* = \{j : 1 \leq j \leq N_i \text{ and } r_i(j) \neq 1\}$ , a set containing the indices of the non-proband members of family  $i$  (after ascertainment); and  $\delta_{ij}^A$  equals 1 if family member  $ij$  is sampled as part of a case-ascertained family and equals 0 otherwise. The indicator  $\delta_{ij}^A$  will depend on  $\delta_i^A$ , the case-ascertainment indicator for family  $i$ , which equals 1 if family  $i$  is ascertained via an affected proband and 0 otherwise (with the constraint that  $\sum_{i=1}^F \delta_i^A = F_A$ ). If  $\delta_i^A = 0$ , then  $\delta_{ij}^A = 0$  by definition, but if  $\delta_i^A = 1$ , then  $\delta_{ij}^A$  can equal 0 or 1.

To obtain an expression for the bias of  $p_A$  as an estimator for  $\pi_A$ , we employ the strategy of Hartley and Ross (1954) for determining the bias of a ratio estimator. This strategy begins by expanding the covariance between  $p_A$  and its denominator:

$$\begin{aligned}
\text{Cov}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right) &= \text{E}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A \mathbf{I}(Y_{ij} = 1)\right) - \text{E}p_A \cdot \text{E}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right) \\
&= \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \text{E}(\delta_{ij}^A \mathbf{I}(Y_{ij} = 1)) - \text{E}p_A \cdot \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \text{E}(\delta_{ij}^A) \\
&= \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \text{E}\left(\text{E}(\delta_{ij}^A | \delta_i^A)\right) \mathbf{I}(Y_{ij} = 1) - \text{E}p_A \cdot \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \text{E}\left(\text{E}(\delta_{ij}^A | \delta_i^A)\right).
\end{aligned}$$

Under Assumption (v), the probability that relative  $ij$  is sampled is a constant referred to as  $s$ . Using this fact to replace  $\text{E}(\delta_{ij}^A | \delta_i^A)$  in the last line of the above equation yields

$$\begin{aligned}
\text{Cov}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right) &= \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \text{E}(s\delta_i^A + 0(1 - \delta_i^A)) \mathbf{I}(Y_{ij} = 1) \\
&\quad - \text{E}p_A \cdot \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \text{E}(s\delta_i^A + 0(1 - \delta_i^A)) \\
&= s \sum_{i=1}^F \text{E}(\delta_i^A) \sum_{\{j: j \in J_i^*\}} \mathbf{I}(Y_{ij} = 1) - \text{E}p_A \cdot s \sum_{i=1}^F \text{E}(\delta_i^A) \sum_{\{j: j \in J_i^*\}} 1 \\
&= s \sum_{i=1}^F \text{E}(\delta_i^A) (N_i^A - 1) - \text{E}p_A \cdot s \sum_{i=1}^F \text{E}(\delta_i^A) (N_i - 1) \\
&= s F_A s_A \sum_{i=1}^F N_i^A (N_i^A - 1) - \text{E}p_A \cdot s F_A s_A \sum_{i=1}^F N_i^A (N_i - 1)
\end{aligned}$$

where the third line follows because  $J_i^*$  does not include one affected member of family  $i$  who is designated as the case proband; and where the fourth line follows because, under Assumption (iii) about proband selection and Assumption (iv) about single ascertainment,  $\text{E}(\delta_i^A)$  can be rewritten as  $N_i^A F_A s_A$ , where  $s_A = 1 / \sum_{i=1}^F N_i^A$ , the sampling fraction for

case probands. The last line above can be rearranged to give

$$E p_A - \frac{\sum_{i=1}^F N_i^A (N_i^A - 1)}{F \sum_{i=1}^F N_i^A (N_i - 1)} = - \frac{\text{Cov}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)}{s F_A s_A \sum_{i=1}^F N_i^A (N_i - 1)}.$$

Since the second term on the lefthand side of the above equation is just  $\pi_A$  as written in (A.3), the bias of  $p_A$  can be written as

$$E p_A - \pi_A = - \frac{\text{Cor}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right) \cdot \text{SD}(p_A) \cdot \text{SD}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)}{s F_A s_A \sum_{i=1}^F N_i^A (N_i - 1)}. \quad (\text{A.7})$$

Since the denominator on the righthand side of the above equation equals the expectation of  $\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A$ , we can rewrite (A.7) as

$$E p_A - \pi_A = - \text{Cor}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right) \cdot \text{SD}(p_A) \cdot \text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right), \quad (\text{A.8})$$

where  $\text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)$ , the coefficient of variation, is defined as the ratio of the standard deviation of  $\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A$  to the mean of  $\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A$ .

We examine the magnitude of each of the three multiplicands on the righthand side of (A.8). First,  $\text{SD}(p_A)$  must be less than 0.5 because  $p_A$  is a proportion. Second, it is difficult (if not impossible) to construct a population where  $\text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)$  is larger than 2, as would be expected for a quantity that is effectively the sum of binary variables (albeit non-identical, non-independent ones). Third, under Assumption (ii), which states that family size and disease status are uncorrelated,  $\text{Cor}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)$  is negligible. This is true because Assumption (ii) ensures that the average size of the case-ascertained families included in the study will have a negligible correlation with the proportion of the relatives in those families who have the disease. Putting these three facts together, we see

that the bias of  $p_A$  is negligible when Assumption (ii) holds.

To illustrate the importance of Assumption (ii) in guaranteeing that  $p_A$  is approximately unbiased for  $\pi_A$ , we examine the bias of  $p_A$  in two fictional populations. The first is the same population used in the simulation experiments in Section 4, where the lifetime prevalence of disease was set to 11.5% for females from all families and 5.9% for males from all families. The second is a population created in an identical fashion, except that the prevalence of disease was set to 16.5% for females and 7.9% for males from families with three or fewer members, and to 6.5% for females and 3.9% for males from families with four or more members. Note that the overall prevalence of disease is equal in both populations, but in the second population, a disproportionately large number of the diseased individuals belong to small families. For each population, we sampled 1,000 datasets, each consisting of 64 case probands ( $F_A = 64$ ) and all of their relatives ( $s = 1$ ). We calculated  $p_A$  and  $\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A$  for each dataset and then used the resulting values to estimate the three multiplicands in (A.8) for that population. Web Table 1 presents the values of the three multiplicands in the two populations, as well as the bias of  $p_A$  in percentage terms. The middle column of Web Table 1 reveals that  $\text{Cor}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)$  and the percentage bias of  $p_A$  are negligible in the first population, where Assumption (ii) holds. Comparing the middle column to the right-most column in Web Table 1 reveals that the percentage bias of  $p_A$  is approximately 100 times larger for the second population, where Assumption (ii) does not hold. This increase in bias is due to the larger value of  $\text{Cor}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)$ .

We can use the approach employed above to obtain an analogous expression for the bias of  $p_U$  as an estimator for  $\pi_U$

$$E p_U - \pi_U = -\text{Cor}\left(p_U, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^U\right) \cdot \text{SD}(p_U) \cdot \text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^U\right) \quad (\text{A.9})$$

where  $\delta_{ij}^U$  equals 1 if relative  $ij$  is sampled as part of a control-ascertained family and equals 0 otherwise. Since the same arguments made for the multiplicands in (A.8) also apply to the multiplicands in (A.9), we see that the bias of  $p_U$  will be negligible when Assumption (ii) holds. For both populations described above, Web Table 2 presents values of the three multiplicands in (A.9) and the percentage bias of  $p_U$ , calculated from 1,000 datasets containing  $F_U = 58$  control probands and all of their relatives ( $s = 1$ ). The bias of

**Web Table 1**

*Components of the bias of  $p_A$  when Assumption (ii) does and does not hold*

| Term   | Value<br>(for $F_A = 64$ and $s = 1$ ) |  |
|--|--|--|
|  | Population 1                           | Population 2                               |
|  | $\text{Cor}(N_i, N_i^A/N_i) \approx 0$ | $\text{Cor}(N_i, N_i^A/N_i) \approx -0.19$ |
| $\text{Cor}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)$ | 0.001459                               | -0.137009                                  |
| $\text{SD}(p_A)$   | 0.033318                               | 0.038691                                   |
| $\text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right)$       | 0.058140                               | 0.054310                                   |
| Percentage bias of $p_A^\dagger$   | -0.0017%                               | 0.173%                                     |

† Percentage bias of  $p_A$  equals  $100 \cdot (Ep_A - \pi_A)/\pi_A$ , where  $\pi_A \approx 0.16$ ; and where

$$Ep_A - \pi_A = \text{Cor}\left(p_A, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right) \cdot \text{SD}(p_A) \cdot \text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^A\right).$$

$p_U$  is negligible for the first population, but is again approximately 100 times larger for the second population because the correlation between  $p_U$  and its denominator is larger. Thus, the bias of  $p_U$  is negligible when Assumption (ii) holds.

The fact that  $p_A$  and  $p_U$  have negligible bias under the assumptions enumerated in Section 2 implies that the bias in the leading term of (A.6) is negligible. Thus, to a first-degree approximation,  $\hat{\pi}$  is an unbiased estimator for  $\pi$ .

We now turn to the bias introduced by the second-order terms in (A.6). First, note that the final second-order term in (A.6) introduces no bias because  $\text{Cov}(p_A, p_U) = 0$ . Next, note that the numerators of the first two second-order terms in (A.6) can be re-written as  $\left[Ep_U(1 - Ep_U)/\sum_{i^*=F_A+1}^{F_A+F_U} (n_{i^*} - 1)\right] (1 - Ep_A)$  and  $\left[Ep_A(1 - Ep_A)/\sum_{i^*=1}^{F_A} (n_{i^*} - 1)\right] Ep_U$ , respectively. We can ignore the summations in the denominators of these terms because they are approximately equal under assumption (ii) that disease status is uncorrelated with family size and under the assumption that  $F_A \approx F_U$ . Now, the only difference between the

**Web Table 2***Components of the bias of  $p_U$  when Assumption (ii) does and does not hold*

| Term   | Value<br>(for $F_U = 58$ and $s = 1$ ) |  |
|--|--|--|
|  | Population 1                           | Population 2                               |
|  | $\text{Cor}(N_i, N_i^A/N_i) \approx 0$ | $\text{Cor}(N_i, N_i^A/N_i) \approx -0.19$ |
| $\text{Cor}\left(p_U, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^U\right)$ | -0.008865                              | -0.098235                                  |
| $\text{SD}(p_U)$   | 0.025768                               | 0.020937                                   |
| $\text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^U\right)$       | 0.061170                               | 0.060793                                   |
| Percentage bias of $p_U$ <sup>†</sup>  | 0.0018%                                | 0.2281%                                    |

<sup>†</sup> Percentage bias of  $p_U$  equals  $100 \cdot (Ep_U - \pi_U)/\pi_U$ , where  $\pi_U \approx 0.055$ ; and where

$$Ep_U - \pi_U = \text{Cor}\left(p_U, \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^U\right) \cdot \text{SD}(p_U) \cdot \text{CV}\left(\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^U\right).$$

two terms is that  $(1 - Ep_A)Ep_U$  is multiplied by  $(1 - Ep_U)$  in the first term and by  $Ep_A$  in the second term. Thus, if we assume that  $E(1 - p_U) > Ep_A$ , then the first second-order term, which has a negative sign in front of it, is larger in magnitude than the second second-order term, which has a positive sign in front of it. As a result, the bias introduced through the second-order terms in (A.6) will be non-positive if  $E(1 - p_U) > Ep_A$  and  $F_A \approx F_U$ . However, the results of the simulation experiments in Section 4, where  $\hat{\pi}$  underestimates  $\pi$  by only a very small amount, suggest that the bias introduced through the second- (and higher-) order terms is very small in practice. ■



**Web Appendix B. Proof that  $\hat{\pi}^x$  is only slightly biased at the first-order for  $\pi^x$**

We define  $\pi^x \equiv f(Y_{ij} = 1|X_{ij} = x)$ , where individual  $ij$  is randomly selected from among the members of the population with  $X_{ij} = x$ . Assumption (i) allows us to expand  $f(Y_{ij} = 1|X_{ij} = x)$  as

$$\begin{aligned}\pi^x &= f(Y_{ij} = 1|Y_{ij'} = 1, X_{ij} = x)f(Y_{ij'} = 1|X_{ij} = x) \\ &\quad + f(Y_{ij} = 1|Y_{ij'} = 0, X_{ij} = x)f(Y_{ij'} = 0|X_{ij} = x),\end{aligned}\tag{B.1}$$

where individual  $ij'$  is randomly selected from among  $Y_{ij}$ 's relatives with disease status  $Y_{ij'}$ . (In the remainder of this proof, we will assume that  $j' \neq j$ .) Under Assumption (vii) about the independence of probands' disease status and relatives' covariates,

$$\pi^x = f(Y_{ij} = 1|Y_{ij'} = 1, X_{ij} = x)f(Y_{ij'} = 1) + f(Y_{ij} = 1|Y_{ij'} = 0, X_{ij} = x)f(Y_{ij'} = 0),$$

follows from (B.1). We can rewrite the preceding equation as

$$\pi^x = \pi_A^x \pi + \pi_U^x (1 - \pi),\tag{B.2}$$

where  $\pi_A^x \equiv f(Y_{ij} = 1|Y_{ij'} = 1, X_{ij} = x)$  and  $\pi_U^x \equiv f(Y_{ij} = 1|Y_{ij'} = 0, X_{ij} = x)$ . The parameters  $\pi_A^x$  and  $\pi_U^x$  can be defined in terms of the finite population:

$$\begin{aligned}\pi_A^x &= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1)\mathbf{I}(X_{ij} = x) \sum_{j' \neq j} \mathbf{I}(Y_{ij'} = 1)}{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(X_{ij} = x) \sum_{j' \neq j} \mathbf{I}(Y_{ij'} = 1)} \\ &= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1)\mathbf{I}(X_{ij} = x)(N_i^A - 1)}{\sum_{i=1}^F \sum_{j=1}^{N_i} (\mathbf{I}(X_{ij} = x)\mathbf{I}(Y_{ij} = 1)(N_i^A - 1) + \mathbf{I}(X_{ij} = x)\mathbf{I}(Y_{ij} = 0)N_i^A)} \\ &= \frac{\sum_{i=1}^F N_i^{Ax}(N_i^A - 1)}{\sum_{i=1}^F (N_i^x N_i^A - N_i^{Ax})},\end{aligned}\tag{B.3}$$

and

$$\begin{aligned}
\pi_U^x &= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1) \mathbf{I}(X_{ij} = x) \sum_{j' \neq j} \mathbf{I}(Y_{ij'} = 0)}{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(X_{ij} = x) \sum_{j' \neq j} \mathbf{I}(Y_{ij'} = 0)} \\
&= \frac{\sum_{i=1}^F \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1) \mathbf{I}(X_{ij} = x) N_i^U}{\sum_{i=1}^F \sum_{j=1}^{N_i} (\mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) N_i^U + \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 0) (N_i^U - 1))} \\
&= \frac{\sum_{i=1}^F N_i^{Ax} N_i^U}{\sum_{i=1}^F (N_i^x N_i^U - N_i^{Ux})} \tag{B.4}
\end{aligned}$$

where  $N_i^A = \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1)$ ;  $N_i^U = \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 0)$ ;  $N_i^x = \sum_{j=1}^{N_i} \mathbf{I}(X_{ij} = x)$ ;  $N_i^{Ax} = \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 1) \mathbf{I}(X_{ij} = x)$ ; and  $N_i^{Ux} = \sum_{j=1}^{N_i} \mathbf{I}(Y_{ij} = 0) \mathbf{I}(X_{ij} = x)$ .

Now, recall that

$$\hat{\pi}^x = p_A^x \hat{\pi} + p_U^x (1 - \hat{\pi}), \tag{B.5}$$

where

$$p_A^x = \frac{\sum_{i^*=1}^{F_A} \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} = x) \mathbf{I}(Y_{i^*j^*} = 1)}{\sum_{i^*=1}^{F_A} \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} = x)}$$

and

$$p_U^x = \frac{\sum_{i^*=F_A+1}^{F_A+F_U} \sum_{i^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} = x) \mathbf{I}(Y_{i^*j^*} = 1)}{\sum_{i^*=F_A+1}^{F_A+F_U} \sum_{i^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} = x)}.$$

The estimator  $\hat{\pi}^x$  in (B.5) can be approximated by a second-order Taylor expansion

around  $\mathbb{E}\widehat{\pi}$ ,  $\mathbb{E}p_U^x$  and  $\mathbb{E}p_A^x$ :

$$\begin{aligned}
\widehat{\pi}^x &\approx (\mathbb{E}p_A^x \mathbb{E}\widehat{\pi} + \mathbb{E}p_U^x (1 - \mathbb{E}\widehat{\pi})) \\
&+ (\widehat{\pi} - \mathbb{E}\widehat{\pi}) \left. \frac{\partial \widehat{\pi}^x}{\partial \widehat{\pi}} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} + (p_U^x - \mathbb{E}p_U^x) \left. \frac{\partial \widehat{\pi}^x}{\partial p_U^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} + (p_A^x - \mathbb{E}p_A^x) \left. \frac{\partial \widehat{\pi}^x}{\partial p_A^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} \\
&+ \frac{1}{2} (\widehat{\pi} - \mathbb{E}\widehat{\pi})^2 \left. \frac{\partial^2 \widehat{\pi}^x}{\partial \widehat{\pi}^2} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} + \frac{1}{2} (p_U^x - \mathbb{E}p_U^x)^2 \left. \frac{\partial^2 \widehat{\pi}^x}{\partial p_U^{x2}} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} \\
&+ \frac{1}{2} (p_A^x - \mathbb{E}p_A^x)^2 \left. \frac{\partial^2 \widehat{\pi}^x}{\partial p_A^{x2}} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} \\
&+ (\widehat{\pi} - \mathbb{E}\widehat{\pi}) (p_A^x - \mathbb{E}p_A^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial \widehat{\pi} \partial p_A^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} + (\widehat{\pi} - \mathbb{E}\widehat{\pi}) (p_U^x - \mathbb{E}p_U^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial \widehat{\pi} \partial p_U^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} \\
&+ (p_U^x - \mathbb{E}p_U^x) (p_A^x - \mathbb{E}p_A^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial p_U^x \partial p_A^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x}.
\end{aligned}$$

Taking the expectation of both sides of the above equation yields

$$\begin{aligned}
\mathbb{E}\widehat{\pi}^x &\approx (\mathbb{E}p_A^x \mathbb{E}\widehat{\pi} + \mathbb{E}p_U^x (1 - \mathbb{E}\widehat{\pi})) \tag{B.6} \\
&+ \frac{1}{2} \text{Var}(\widehat{\pi}) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial \widehat{\pi}^2} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} + \frac{1}{2} \text{Var}(p_U^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial p_U^{x2}} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} + \frac{1}{2} \text{Var}(p_A^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial p_A^{x2}} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} \\
&+ \text{Cov}(\widehat{\pi}, p_A^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial \widehat{\pi} \partial p_A^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} + \text{Cov}(\widehat{\pi}, p_U^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial \widehat{\pi} \partial p_U^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x} \\
&+ \text{Cov}(p_U^x, p_A^x) \left. \frac{\partial^2 \widehat{\pi}^x}{\partial p_U^x \partial p_A^x} \right|_{\mathbb{E}\widehat{\pi}, \mathbb{E}p_U^x, \mathbb{E}p_A^x}.
\end{aligned}$$

We focus now on the leading term of the expectation of the Taylor expansion in (B.6):  $\mathbb{E}p_A^x \mathbb{E}\widehat{\pi} + \mathbb{E}p_U^x (1 - \mathbb{E}\widehat{\pi})$ . We have already shown in Web Appendix A that, under conditions (i)-(v),  $\widehat{\pi}$  has a very small negative bias as an estimator for  $\pi$ . To derive the bias in  $p_A^x$  and  $p_U^x$ , we introduce indicators in order to rewrite them as

$$p_A^x = \frac{\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax} \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax^c} \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1)}{\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax} \mathbf{I}(X_{ij} = x) + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax^c} \mathbf{I}(X_{ij} = x)}$$

where  $\delta_{ij}^{Ax}$  (or  $\delta_{ij}^{Ax^c}$ ) equals 1 if family member  $ij$  is sampled as part of a family ascertained through a case proband with covariate value  $x$  (or covariate value in the complement of  $x$ ) and equals 0 otherwise; and

$$p_U^x = \frac{\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ux} \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ux^c} \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1)}{\sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ux} \mathbf{I}(X_{ij} = x) + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ux^c} \mathbf{I}(X_{ij} = x)}$$

where  $\delta_{ij}^{Ux}$  (or  $\delta_{ij}^{Ux^c}$ ) equals 1 if family member  $ij$  is sampled as part of a family ascertained through a control proband with covariate value  $x$  (or covariate value in the complement of  $x$ ) and equals 0 otherwise. For the sake of brevity, we will refer to the denominators of  $p_A^x$  and  $p_U^x$  as  $d_A^x$  and  $d_U^x$ , respectively.

To derive the bias of  $p_A^x$  as an estimator for  $\pi_A^x$ , we use the same Hartley-Ross (1954) approach employed for  $p_A$  in Web Appendix A:

$$\begin{aligned} \text{Cov}(p_A^x, d_A^x) &= \mathbf{E} \left( \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax} \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax^c} \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) \right) \\ &\quad - \mathbf{E} p_A^x \cdot \mathbf{E} \left( \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax} \mathbf{I}(X_{ij} = x) + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \delta_{ij}^{Ax^c} \mathbf{I}(X_{ij} = x) \right) \\ &= \left[ \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E} \left( \mathbf{E}(\delta_{ij}^{Ax} | \delta_i^{Ax}) \right) \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) \right. \\ &\quad \left. + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E} \left( \mathbf{E}(\delta_{ij}^{Ax^c} | \delta_i^{Ax^c}) \right) \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) \right] \\ &\quad - \mathbf{E} p_A^x \cdot \left[ \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E} \left( \mathbf{E}(\delta_{ij}^{Ax} | \delta_i^{Ax}) \right) \mathbf{I}(X_{ij} = x) + \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E} \left( \mathbf{E}(\delta_{ij}^{Ax^c} | \delta_i^{Ax^c}) \right) \mathbf{I}(X_{ij} = x) \right] \end{aligned}$$

where  $\delta_i^{Ax}$  equals 1 if family  $i$  is ascertained via an affected proband with covariate value  $x$  and 0 otherwise, and where  $\delta_i^{Ax^c}$  equals 1 if family  $i$  is ascertained via an affected proband with covariate value in the complement of  $x$  and 0 otherwise. Invoking Assumption (v),

the preceding line reduces to

$$\begin{aligned}
\text{Cov}(p_A^x, d_A^x) &= \left[ s \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E}(\delta_i^{Ax}) \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) + s \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E}(\delta_i^{Ax^c}) \mathbf{I}(X_{ij} = x) \mathbf{I}(Y_{ij} = 1) \right. \\
&\quad \left. - \mathbf{E}p_A^x \cdot \left[ s \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E}(\delta_i^{Ax}) \mathbf{I}(X_{ij} = x) + s \sum_{i=1}^F \sum_{\{j: j \in J_i^*\}} \mathbf{E}(\delta_i^{Ax^c}) \mathbf{I}(X_{ij} = x) \right] \right] \\
&= \left[ s \sum_{i=1}^F \mathbf{E}(\delta_i^{Ax}) (N_i^{Ax} - 1) + s \sum_{i=1}^F \mathbf{E}(\delta_i^{Ax^c}) N_i^{Ax} \right] \\
&\quad - \mathbf{E}p_A^x \cdot \left[ s \sum_{i=1}^F \mathbf{E}(\delta_i^{Ax}) (N_i^x - 1) + s \sum_{i=1}^F \mathbf{E}(\delta_i^{Ax^c}) N_i^x \right] \\
&= \left[ s F_A s_A \sum_{i=1}^F N_i^{Ax} (N_i^{Ax} - 1) + s F_A s_A \sum_{i=1}^F N_i^{Ax^c} N_i^{Ax} \right] \\
&\quad - \mathbf{E}p_A^x \cdot \left[ s F_A s_A \sum_{i=1}^F N_i^{Ax} (N_i^x - 1) + s F_A s_A \sum_{i=1}^F N_i^{Ax^c} N_i^x \right]
\end{aligned}$$

where  $N_i^{Ax^c} = \sum_{j=1}^{N_i} \mathbf{I}(X_{ij} \neq x) \mathbf{I}(Y_{ij} = 1)$ . Note that the last expression above follows from the second-to-last expression above under Assumption (iii) about proband selection and Assumption (iv) about single ascertainment. The last expression above can be rewritten as

$$\text{Cov}(p_A^x, d_A^x) = s F_A s_A \sum_{i=1}^F N_i^{Ax} (N_i^A - 1) - \mathbf{E}p_A^x \cdot s F_A s_A \sum_{i=1}^F (N_i^A N_i^x - N_i^{Ax}),$$

which, when rearranged and combined with (B.3), gives

$$\begin{aligned}
\mathbf{E}p_A^x - \pi_A^x &= - \frac{\text{Cov}(p_A^x, d_A^x)}{s F_A s_A \sum_{i=1}^F (N_i^A N_i^x - N_i^{Ax})} \\
&= - \text{Cor}(p_A^x, d_A^x) \cdot \text{SD}(p_A^x) \cdot \text{CV}(d_A^x). \tag{B.7}
\end{aligned}$$

We can use the Hartley-Ross (1954) approach to obtain an analogous expression for the

bias of  $p_U^x$  as an estimator for  $\pi_U^x$ :

$$\begin{aligned}
\mathbb{E}p_U^x - \pi_U^x &= -\frac{\text{Cov}(p_U^x, d_U^x)}{s F_U s_U \sum_{i=1}^F (N_i^U N_i^x - N_i^{Ux})} \\
&= -\text{Cor}(p_U^x, d_U^x) \cdot \text{SD}(p_U^x) \cdot \text{CV}(d_U^x). \tag{B.8}
\end{aligned}$$

We can then use the same arguments made in Web Appendix A to establish that the right-hand sides of (B.7) and (B.8) will be negligible when Assumption (ii) holds. Thus, under Assumption (ii),  $p_A^x$  and  $p_U^x$  are approximately unbiased estimators for  $\pi_A^x$  and  $\pi_U^x$ , respectively. Using our previous finding that  $\hat{\pi}$  slightly underestimate  $\pi$ , along with the fact that  $p_A^x$  will typically exceed  $p_U^x$  for diseases that aggregate in families, we see that the leading term in (B.6) underestimates  $\pi_x$  slightly. Thus, to a first-degree approximation,  $\hat{\pi}^x$  is a slightly downwardly biased estimator for  $\pi_x$ . However, the results of the simulation experiments in Section 4 suggest that the bias introduced by the leading term and also the higher order terms in (B.6) is very small. ■

### Web Appendix C. Standard Errors and Confidence Intervals for $\hat{\pi}$ and $\hat{\pi}^x$

The delta method can be used to obtain approximate standard errors for  $\hat{\pi}$  and  $\hat{\pi}^x$ . The approximate standard error for  $\hat{\pi}$  is

$$\text{se}(\hat{\pi}) = \pi(1 - \pi) \sqrt{\frac{\pi_A}{d_A(1 - \pi_A)} \left[ 1 + \frac{2\rho_A}{d_A} \sum_{i^*=1}^{F_A} \binom{n_{i^*} - 1}{2} \right] + \frac{(1 - \pi_U)}{d_U \pi_U} \left[ 1 + \frac{2\rho_U}{d_U} \sum_{i^*=F_A+1}^{F_A+F_U} \binom{n_{i^*} - 1}{2} \right]}, \quad (\text{C.1})$$

where  $d_A = \sum_{i^*=1}^{F_A} \sum_{j^*=2}^{n_{i^*}} 1$  and  $d_U = \sum_{i^*=F_A+1}^{F_A+F_U} \sum_{j^*=2}^{n_{i^*}} 1$ ;  $\rho_A = \text{Cor}(Y_{i^*j^*}, Y_{i^*j^{*'}})$  for  $i^* = 1, \dots, F_A$ ,  $j^* > 1$ ,  $j^{*'} > 1$ , and  $j^* \neq j^{*'}$ ; and  $\rho_U = \text{Cor}(Y_{i^*j^*}, Y_{i^*j^{*'}})$  for  $i^* = F_A + 1, \dots, F_A + F_U$ ,  $j^* > 1$ ,  $j^{*'} > 1$ , and  $j^* \neq j^{*'}$ . Note that the more the disease aggregates in families, the larger  $\rho_A$  and  $\rho_U$  will be and therefore the larger the standard error for  $\hat{\pi}$  will be.

The approximate standard error for  $\hat{\pi}^x$  is

$$\text{se}(\hat{\pi}^x) = \sqrt{a \Sigma a^T}, \quad (\text{C.2})$$

where

$$a = \left[ \pi \quad (1 - \pi) \quad \frac{(\pi_A^x - \pi_U^x)\pi^2}{\pi_U} \quad \frac{(\pi_A^x - \pi_U^x)(1 - \pi)^2}{1 - \pi_A} \right]; \quad (\text{C.2.a})$$

and

$$\Sigma = \begin{bmatrix} \sigma_{1,1} & 0 & \sigma_{1,3} & 0 \\ 0 & \sigma_{2,2} & 0 & \sigma_{2,4} \\ \sigma_{1,3} & 0 & \sigma_{3,3} & 0 \\ 0 & \sigma_{2,4} & 0 & \sigma_{4,4} \end{bmatrix} \quad (\text{C.2.b})$$

with

$$\sigma_{1,1} = \frac{\pi_A^x(1 - \pi_A^x)}{d_A^x} \left[ 1 + \frac{2\rho_A^x}{d_A^x} \sum_{i^*=1}^{F_A} \binom{n_{i^*}^x}{2} \right],$$

$$\sigma_{2,2} = \frac{\pi_U^x(1 - \pi_U^x)}{d_U^x} \left[ 1 + \frac{2\rho_U^x}{d_U^x} \sum_{i^*=F_A+1}^{F_A+F_U} \binom{n_{i^*}^x}{2} \right],$$

$$\sigma_{3,3} = \frac{\pi_A(1 - \pi_A)}{d_A} \left[ 1 + \frac{2\rho_A}{d_A} \sum_{i^*=1}^{F_A} \binom{n_{i^*} - 1}{2} \right],$$

$$\sigma_{4,4} = \frac{\pi_U(1 - \pi_U)}{d_U} \left[ 1 + \frac{2\rho_U}{d_U} \sum_{i^*=F_A+1}^{F_A+F_U} \binom{n_{i^*} - 1}{2} \right],$$

$$\sigma_{1,3} = \frac{1}{d_A^x d_A} \left[ \sigma_{1,1} + \rho_A^{x,x^c} \sum_{i^*=1}^{F_A} n_{i^*}^x n_{i^*}^{x^c} \right],$$

and

$$\sigma_{2,4} = \frac{1}{d_U^x d_U} \left[ \sigma_{2,2} + \rho_U^{x,x^c} \sum_{i^*=F_A+1}^{F_A+F_U} n_{i^*}^x n_{i^*}^{x^c} \right].$$

Further,  $d_A^x = \sum_{i^*=1}^{F_A} \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} = x)$  and  $d_U^x = \sum_{i^*=F_A+1}^{F_A+F_U} \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} = x)$ ;  $n_{i^*}^x = \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} = x)$  and  $n_{i^*}^{x^c} = \sum_{j^*=2}^{n_{i^*}} \mathbf{I}(X_{i^*j^*} \neq x)$ ;  $\rho_A^x = \text{Cor}(Y_{i^*j^*}, Y_{i^*j^{*'}})$  for  $i^* = 1, \dots, F_A, j^* > 1, j^{*'} > 1, j^* \neq j^{*'}$ ,  $X_{i^*j^*} = X_{i^*j^{*'}} = x$ ;  $\rho_U^x = \text{Cor}(Y_{i^*j^*}, Y_{i^*j^{*'}})$  for  $i^* = F_A + 1, \dots, F_A + F_U, j^* > 1, j^{*'} > 1, j^* \neq j^{*'}$ ,  $X_{i^*j^*} = X_{i^*j^{*'}} = x$ ;  $\rho_A^{x,x^c} = \text{Cor}(Y_{i^*j^*}, Y_{i^*j^{*'}})$  for  $i^* = 1, \dots, F_A, j^* > 1, j^{*'} > 1, j^* \neq j^{*'}$ ,  $X_{i^*j^*} = x, X_{i^*j^{*'}} \neq x$ ; and  $\rho_U^{x,x^c} = \text{Cor}(Y_{i^*j^*}, Y_{i^*j^{*'}})$  for  $i^* = F_A + 1, \dots, F_A + F_U, j^*, j^{*' > 1, j^* \neq j^{*'}$ ,  $X_{i^*j^*} = x, X_{i^*j^{*'}} \neq x$ .

In practice, the population parameters in Equations (C.1) and (C.2) are replaced with estimates, which yields estimated standard errors that we refer to as  $\widehat{\text{se}}(\widehat{\pi})$  and  $\widehat{\text{se}}(\widehat{\pi}^x)$ , respectively. Although estimating the parameters  $\pi, \pi_A, \pi_U, \pi_A^x$ , and  $\pi_U^x$  will not require additional calculation because they appear in Equations (1) and (2), the parameters  $\rho_A, \rho_U, \rho_A^x, \rho_U^x, \rho_A^{x,x^c}$ , and  $\rho_U^{x,x^c}$  will need to be estimated. Estimates can be calculated from the data using Pearson correlations for all relevant pairs of relatives or instead can be based on prior information. In situations where the data do not contain sufficient information to estimate the correlations and there is no prior information about them, the value 0.30 can be used as an extremely conservative estimate (by comparison, the correlations are approximately 0.10 for a disease with heritability and prevalence similar to the Austrian case-control family study). In our simulation experiments, we calculated the correlations from the data whenever possible and otherwise set them to 0.30.

The estimated quantities  $\widehat{\pi}$  and  $\widehat{\text{se}}(\widehat{\pi})$  could be used to form a Wald interval for  $\pi$ , which



would take the form  $CI = [\hat{\pi} \pm z_{\alpha/2} \widehat{se}(\hat{\pi})]$ , where  $z_{\alpha/2}$  is the  $\alpha/2$  quantile of the standard normal distribution. However, for small population prevalences ( $\pi < 0.2$ ), Wald intervals do not achieve their nominal coverage level because of the frequent occurrence of  $[0, 0]$  intervals for samples with no diseased relatives. Various adjusted confidence intervals with improved coverage probabilities have been proposed for population proportions, including the Agresti-Coull (1998) interval, which has its roots in the work of Wilson (1927). Simply put, the Agresti-Coull interval improves the Wald interval's coverage properties by smoothing the proportion estimates and the estimated standard errors away from zero. Miao and Gastwirth (2004) have performed simulations to examine the performance of an Agresti-Coull-type interval (and various other intervals) for proportions estimated from moderately-sized samples containing dependent clusters. Since the Agresti-Coull-type interval appears to perform well in the simulations and, further, is easy to compute, we adopt confidence intervals based on the same concept.

For the overall prevalence, the  $100 \cdot (1 - \alpha)\%$  interval takes the form:

$$CI = \tilde{\pi} \pm z_{\alpha/2} \sqrt{\tilde{se}} \quad (C.3)$$

where  $\tilde{\pi}$  and  $\tilde{se}$  are calculated using the formulas for  $\hat{\pi}$  and  $\widehat{se}(\hat{\pi})$ , respectively, with  $p_A$  replaced by

$$\tilde{p}_A = \frac{d_A p_A + 0.5 z_{\alpha/2}^2}{d_A + z_{\alpha/2}^2}$$

and  $p_U$  replaced by

$$\tilde{p}_U = \frac{d_U p_U + 0.5 z_{\alpha/2}^2}{d_U + z_{\alpha/2}^2}.$$

For the stratum-specific prevalence,  $\pi^x$ , the  $100 \cdot (1 - \alpha)\%$  interval takes the form:

$$CI = \tilde{\pi}^x \pm z_{\alpha/2} \sqrt{\tilde{se}^x} \quad (C.4)$$

where  $\tilde{\pi}^x$  and  $\tilde{se}^x$  are calculated using the formulas for  $\hat{\pi}^x$  and  $\widehat{se}(\hat{\pi}^x)$ , respectively, with  $p_A$  replaced by  $\tilde{p}_A$ ,  $p_U$  replaced by  $\tilde{p}_U$ ,  $p_A^x$  replaced by

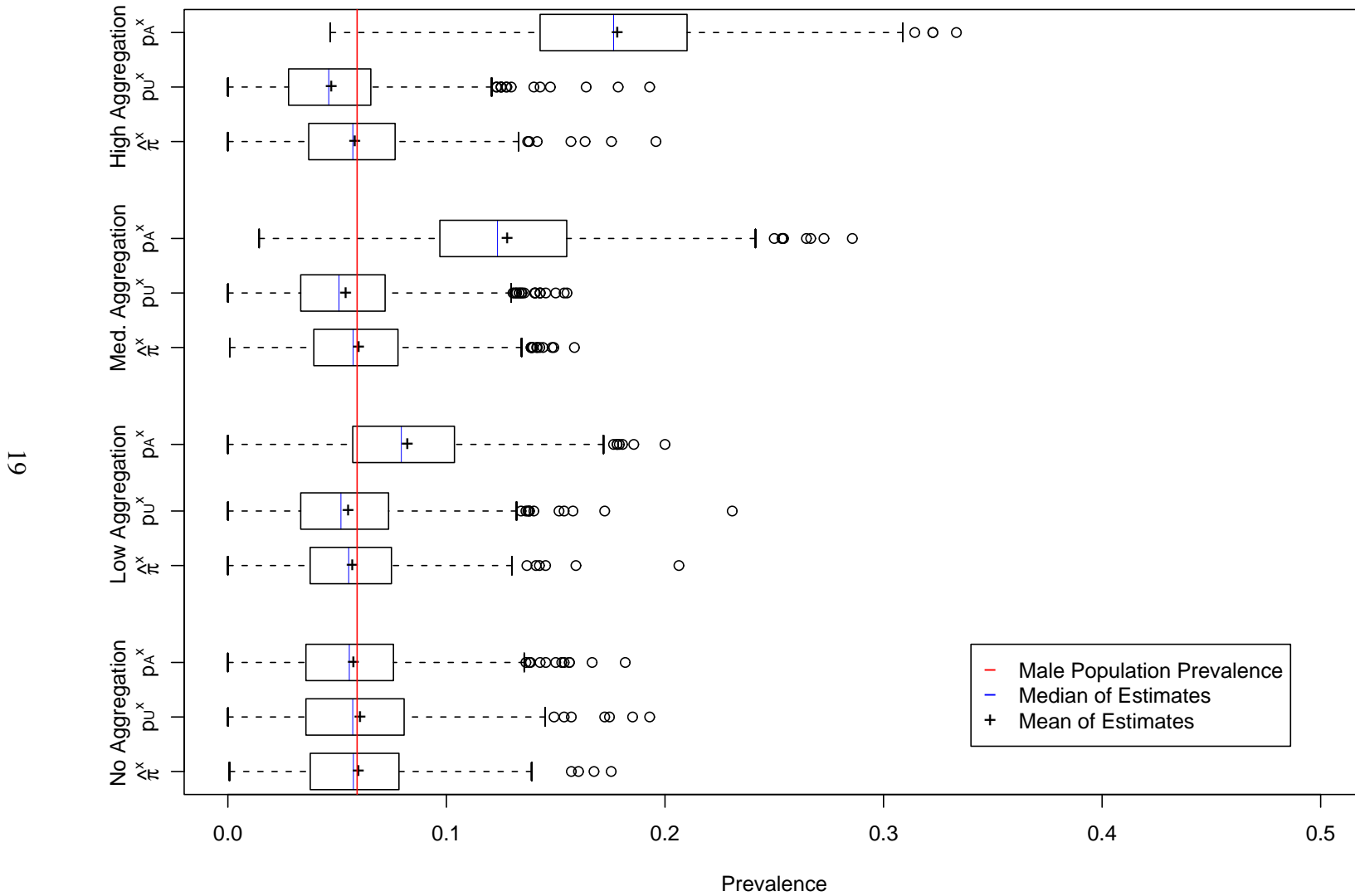
$$\tilde{p}_A^x = \frac{d_A^x p_A^x + 0.5 z_{\alpha/2}^2}{d_A^x + z_{\alpha/2}^2}$$

and  $p_U^x$  replaced by

$$\tilde{p}_U^x = \frac{d_U^x p_U^x + 0.5 z_{\alpha/2}^2}{d_U^x + z_{\alpha/2}^2}.$$

Of course, if the lower (upper) bound of the confidence interval in either (C.3) or (C.4) turns out to be less (greater) than 0 (1), then it should be replaced with 0 (1).

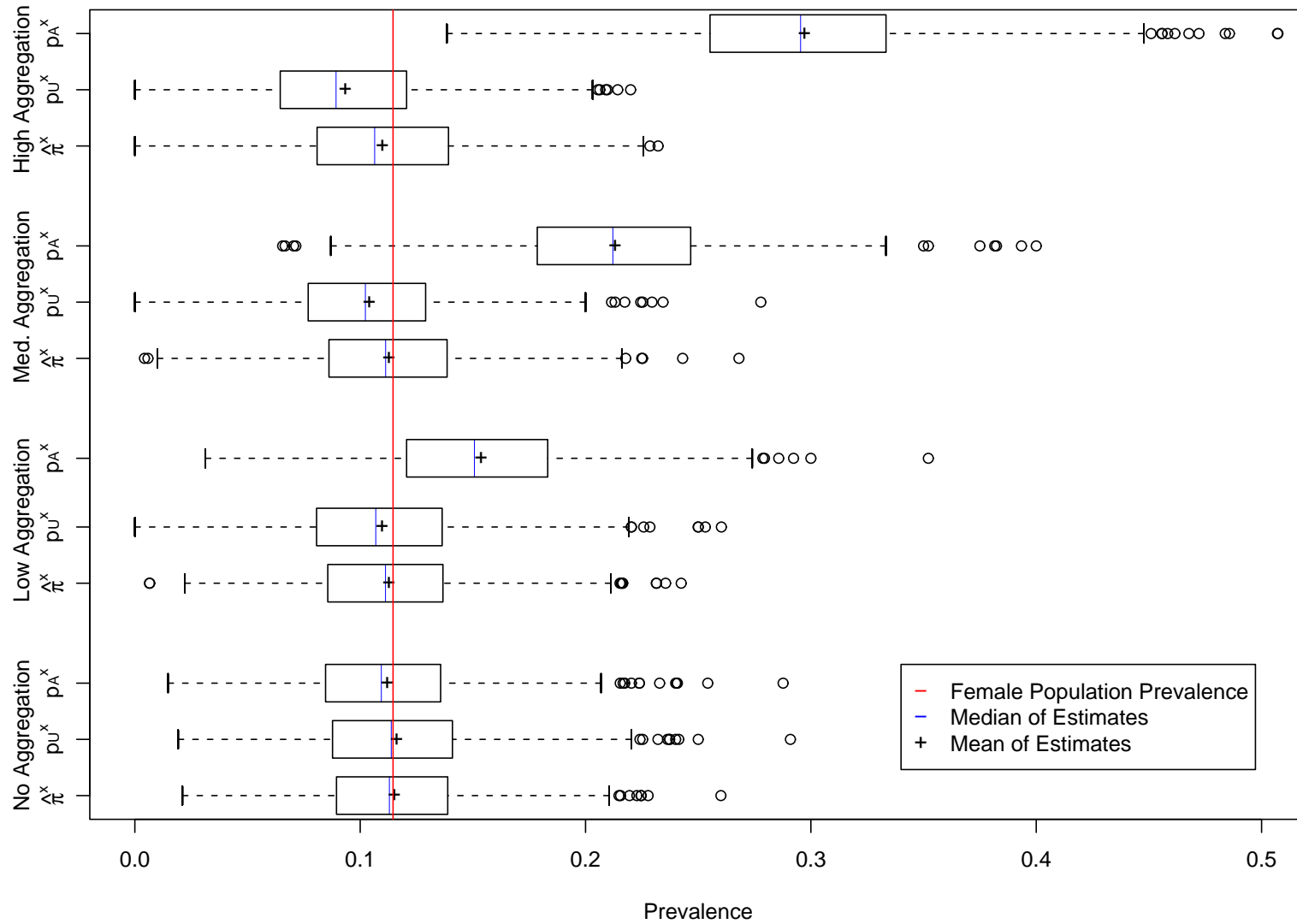
Boxplots of Male  $\hat{\pi}^x$ ,  $p_U^x$ , and  $p_A^x$  for 1000 Samples from Simulated Populations with Different Familial Aggregations  
(Only Assumption iv. Violated)



**Web Figure 1:** Boxplots of male  $\hat{\pi}^x$ ,  $p_U^x$ , and  $p_A^x$  values calculated for 1000 samples drawn from four different populations with varying degrees of disease familiarity.

Boxplots of Female  $\hat{\pi}^x$ ,  $p_U^x$ , and  $p_A^x$  for 1000 Samples from Simulated Populations with Different Familial Aggregations  
(Only Assumption iv. Violated)

20



**Web Figure 2:** Boxplots of female  $\hat{\pi}^x$ ,  $p_U^x$ , and  $p_A^x$  values calculated for 1000 samples drawn from four different populations with varying degrees of disease familiarity.