

APPENDIX / AVAILABLE FROM THE AUTHORS

It is assumed that the study population is a stratified population with a total of J strata ($j = 1, 2, \dots, J$). In each stratum, the occurrence of disease follows a log-linear risk model, that is, for a subject in the j th stratum, $\log(\text{risk}) = \alpha_j + \mathbf{x}_{(\text{gene, exposure, sex})}^t \boldsymbol{\beta}$, where α_j represents the log background risk for the j th stratum, $\mathbf{x}_{(\text{gene, exposure, sex})}$ is a vector of data or codes regarding the gene, the environmental exposure, the sex, and any possible cross-product terms between them, for the subject under concern (\mathbf{x}^t is the transpose of \mathbf{x}), and $\boldsymbol{\beta}$ is a vector of parameters of present interest (to be estimated).

Supposed that a sample of n ($i=1, \dots, n$) case-spouse pairs has been recruited. Let T_i (T'_i) represent the stratum to which the i th proband(his/her spouse) belongs. Let $G_i = A_{1i}A_{2i}$ ($G'_i = A_{3i}A_{4i}$) represent the genotype for the i th proband(his/her spouse) (A_{1i}, A_{2i}, A_{3i} , and A_{4i} are 'alleles'). Let D_i represent the event that the i th proband is a case (i.e., he/she contracts the disease). Let E_i (S_i) represent the environmental exposure(sex) of the i th proband, and E'_i (S'_i), the environmental exposure(sex) of his/her spouse. Further, we let U_i represent the set that contains as its two elements the genotype of the i th proband and the genotype of the i th spouse, that is, $U_i = \{A_{1i}A_{2i}, A_{3i}A_{4i}\}$. And we let V_i be the set containing as its elements the four alleles in the i th case-spouse pair, that is, $V_i = \{A_{1i}, A_{2i}, A_{3i}, A_{4i}\}$.

Conditioned on U_i, D_i, E_i, E'_i, S_i , and S'_i , the probability that the i th proband has genotype of $A_{1i}A_{2i}$ and the i th spouse has genotype of $A_{3i}A_{4i}$ is denoted as Q_i . With elementary algebra, this conditional probability is (the index i was suppressed for simplicity):

$$Q = \Pr(G = A_1A_2, G' = A_3A_4 | U, D, E, E', S, S')$$

$$= \sum_{j=1}^J \Pr(T = j | D, E, E', S, S') \cdot \frac{\Pr(G = A_1A_2, G' = A_3A_4, D, E, E', S, S', T = j)}{\left\{ \Pr(G = A_1A_2, G' = A_3A_4, D, E, E', S, S', T = j) \right.}$$

$$\left. + \Pr(G = A_3A_4, G' = A_1A_2, D, E, E', S, S', T = j) \right\}}$$

The numerator can be expressed as the product of three terms:

$$\begin{aligned} \Pr(G = A_1A_2, G' = A_3A_4, D, E, E', S, S', T = j) &= \Pr(D|G = A_1A_2, G' = A_3A_4, E, E', S, S', T = j) \\ &\quad \cdot \Pr(G = A_1A_2, G' = A_3A_4 | E, E', S, S', T = j) \\ &\quad \cdot \Pr(E, E', S, S', T = j). \end{aligned}$$

With the assumed disease model, the first term of the product is

$$\begin{aligned} \Pr(D|G = A_1A_2, G' = A_3A_4, E, E', S, S', T = j) \\ &= \Pr(D|G = A_1A_2, E, S, T = j) \\ &= \exp[\alpha_j + \mathbf{x}_{(A_1A_2, E, S)}^t \boldsymbol{\beta}]. \end{aligned}$$

The second term of the product can be shown to be (the parenthesis shows the assumption used in each step of derivation; A1: mating is restricted to subjects in the same stratum; A2: mating is independent to genotypes, for males and females in each and every stratum; A3: environmental exposures are independent to genotypes, for males and females in each and every stratum; A4: the genotype frequencies for males are equal to the corresponding frequencies for females, in each and every stratum):

$$\begin{aligned} \Pr(G = A_1A_2, G' = A_3A_4 | E, E', S, S', T = j) \\ &= \Pr(G = A_1A_2, G' = A_3A_4 | E, E', S, S', T = T' = j) \quad (\text{A1}) \\ &= \Pr(G = A_1A_2 | E, S, T = j) \cdot \Pr(G' = A_3A_4 | E', S', T' = j) \quad (\text{A2}) \\ &= \Pr(G = A_1A_2 | S, T = j) \cdot \Pr(G' = A_3A_4 | S', T' = j) \quad (\text{A3}) \\ &= \Pr(G = A_1A_2 | T = j) \cdot \Pr(G = A_3A_4 | T = j). \quad (\text{A4}) \end{aligned}$$

Therefore, we see that the numerator in Q is

$$\begin{aligned} \Pr(G = A_1A_2, G' = A_3A_4, D, E, E', S, S', T = j) &= \exp[\alpha_j + \mathbf{x}_{(A_1A_2, E, S)}^t \boldsymbol{\beta}] \\ &\quad \cdot \Pr(G = A_1A_2 | T = j) \cdot \Pr(G = A_3A_4 | T = j) \\ &\quad \cdot \Pr(E, E', S, S', T = j). \end{aligned}$$

Similarly, we can show that the second term of the denominator in Q is

$$\begin{aligned} \Pr(G = A_3A_4, G' = A_1A_2, D, E, E', S, S', T = j) &= \exp[\alpha_j + \mathbf{x}_{(A_3A_4, E, S)}^t \boldsymbol{\beta}] \\ &\quad \cdot \Pr(G = A_3A_4 | T = j) \cdot \Pr(G = A_1A_2 | T = j) \\ &\quad \cdot \Pr(E, E', S, S', T = j). \end{aligned}$$

Thus, we have (with the index i denoting the i th case-spouse pair)

$$\begin{aligned}
Q_i &= \sum_{j=1}^J \Pr(T_i = j | D_i, E_i, E'_i, S_i, S'_i) \cdot \frac{\exp[\mathbf{x}_{(A_{1i}A_{2i}, E_i, S_i)}^t \boldsymbol{\beta}]}{\exp[\mathbf{x}_{(A_{1i}A_{2i}, E_i, S_i)}^t \boldsymbol{\beta}] + \exp[\mathbf{x}_{(A_{3i}A_{4i}, E_i, S_i)}^t \boldsymbol{\beta}]} \\
&= \frac{\exp[\mathbf{x}_{(A_{1i}A_{2i}, E_i, S_i)}^t \boldsymbol{\beta}]}{\exp[\mathbf{x}_{(A_{1i}A_{2i}, E_i, S_i)}^t \boldsymbol{\beta}] + \exp[\mathbf{x}_{(A_{3i}A_{4i}, E_i, S_i)}^t \boldsymbol{\beta}]}.
\end{aligned}$$

And the conditional likelihood function for the case-spouse data (1:1 case-counterfactual-control analysis) is

$$L_{1:1} = \prod_{i=1}^n Q_i.$$

Let $F_i (F'_i)$ represent the allele that the i th proband(spouse) inherited from his/her father, and $M_i (M'_i)$, the allele that the i th proband(spouse) inherited from his/her mother. Conditioned on $V_i, D_i, E_i, E'_i, S_i,$ and S'_i , the probability that the i th proband has genotype of $A_{1i}A_{2i}$ and the i th spouse has genotype of $A_{3i}A_{4i}$ is denoted as R_i . With elementary algebra, this conditional probability is (the index i was suppressed):

$$\begin{aligned}
R &= \Pr(G = A_1A_2, G' = A_3A_4 | V, D, E, E', S, S') \\
&= \sum_{j=1}^J \Pr(T = j | D, E, E', S, S') \\
&\quad \cdot \frac{\Pr(G = A_1A_2, G' = A_3A_4, D, E, E', S, S', T = j)}{\sum_h \sum_{k \neq h} \sum_{t \neq k} \sum_{\substack{u \neq t \\ t \neq h \\ u \neq k \\ u \neq h}} \Pr(F = A_h, M = A_k, F' = A_t, M' = A_u, D, E, E', S, S', T = j)} \\
&= \sum_{j=1}^J \Pr(T = j | D, E, E', S, S') \\
&\quad \cdot \left\{ \frac{\Pr(G = A_1A_2, G' = A_3A_4, D, E, E', S, S', T = j)}{\sum_h \sum_{k \neq h} \sum_{t \neq k} \sum_{\substack{u \neq t \\ t \neq h \\ u \neq k \\ u \neq h}} \Pr(D | F = A_h, M = A_k, F' = A_t, M' = A_u, E, E', S, S', T = j)} \right. \\
&\quad \cdot \Pr(F = A_h, M = A_k, F' = A_t, M' = A_u | E, E', S, S', T = j) \\
&\quad \left. \cdot \Pr(E, E', S, S', T = j) \right\}
\end{aligned}$$

The numerator of R is the same as that of Q and has been previously shown to be (with assumptions A1~A4)

$$\begin{aligned} \Pr(G = A_1A_2, G' = A_3A_4, D, E, E', S, S', T = j) &= \exp[\alpha_j + \mathbf{x}_{(A_1A_2, E, S)}^t \boldsymbol{\beta}] \\ &\cdot \Pr(G = A_1A_2 | T = j) \cdot \Pr(G = A_3A_4 | T = j) \\ &\cdot \Pr(E, E', S, S', T = j). \end{aligned}$$

With the A5 assumption that there is no imprinting effect for the gene under study, the first term of the product in the denominator of R is

$$\begin{aligned} \Pr(D|F = A_h, M = A_k, F' = A_t, M' = A_u, E, E', S, S', T = j) \\ = \Pr(D|F = A_h, M = A_k, E, S, T = j) \\ = \exp[\alpha_j + \mathbf{x}_{(A_hA_k, E, S)}^t \boldsymbol{\beta}]. \end{aligned}$$

With the A1~A4 assumptions, the second term of the product in the denominator of R is

$$\begin{aligned} \Pr(F = A_h, M = A_k, F' = A_t, M' = A_u | E, E', S, S', T = j) \\ = \Pr(F = A_h, M = A_k | T = j) \cdot \Pr(F = A_t, M = A_u | T = j). \end{aligned}$$

Therefore,

$$\begin{aligned} R &= \sum_{j=1}^J \Pr(T = j | D, E, E', S, S') \\ &\cdot \frac{\exp[\mathbf{x}_{(A_1A_2, E, S)}^t \boldsymbol{\beta}] \cdot \Pr(G = A_1A_2 | T = j) \cdot \Pr(G = A_3A_4 | T = j)}{\sum_h \sum_{k \neq h} \sum_{\substack{t \neq k \\ t \neq h}} \sum_{\substack{u \neq t \\ u \neq h}} \exp[\mathbf{x}_{(A_hA_k, E, S)}^t \boldsymbol{\beta}] \cdot \Pr(F = A_h, M = A_k | T = j) \cdot \Pr(F = A_t, M = A_u | T = j)} \\ &= \sum_{j=1}^J \Pr(T = j | D, E, E', S, S') \\ &\cdot \frac{\left\{ \exp[\mathbf{x}_{(A_1A_2, E, S)}^t \boldsymbol{\beta}] \cdot [\Pr(F = A_1, M = A_2 | T = j) + \Pr(F = A_2, M = A_1 | T = j)] \right\}}{\sum_h \sum_{k \neq h} \sum_{\substack{t \neq k \\ t \neq h}} \sum_{\substack{u \neq t \\ u \neq h}} \exp[\mathbf{x}_{(A_hA_k, E, S)}^t \boldsymbol{\beta}] \cdot \Pr(F = A_h, M = A_k | T = j) \cdot \Pr(F = A_t, M = A_u | T = j)}. \end{aligned}$$

With the A6 assumption that each and every stratum is in Hardy-Weinberg equilibrium, we have

$$\Pr(F = x, M = y | T = j) = \Pr(F = x | T = j) \cdot \Pr(M = y | T = j),$$

for arbitrary alleles x and y and arbitrary stratum j . Thus we see that (with the index i

denoting the i th case-spouse pair):

$$\begin{aligned}
R_i &= \sum_{j=1}^J \Pr(T_i = j | D_i, E_i, E'_i, S_i, S'_i) \cdot \frac{4 \cdot \exp[\mathbf{x}_{(A_{1i}, A_{2i}, E_i, S_i)}^t \boldsymbol{\beta}]}{2 \cdot \sum_h \sum_{k \neq h} \exp[\mathbf{x}_{(A_{hi}, A_{ki}, E_i, S_i)}^t \boldsymbol{\beta}]} \\
&= \frac{\exp[\mathbf{x}_{(A_{1i}, A_{2i}, E_i, S_i)}^t \boldsymbol{\beta}]}{\sum_{h=1}^3 \sum_{k=h+1}^4 \exp[\mathbf{x}_{(A_{hi}, A_{ki}, E_i, S_i)}^t \boldsymbol{\beta}]} .
\end{aligned}$$

And the conditional likelihood function for the case-spouse data (1:5 case-counterfactual-controls analysis) is

$$L_{1:5} = \prod_{i=1}^n R_i .$$