

Supplementary material for the article

ESTIMATING EFFECT SIZES IN GENOME-WIDE ASSOCIATION STUDIES

József Bukszár and Edwin J. C. G. van den Oord

Derivation of (5):

In our article ([1]) we showed that an asymptotic approximation to Pearson's statistic on a $2 \times v$ contingency table is given by

$$Z_1^2 + \dots + Z_v^2, \quad (1)$$

where Z_1, \dots, Z_v are independent normal random variables with

$$E(Z_i) = \mathbf{a}_i \mu \quad \text{and} \quad \text{Var}(Z_i) = \lambda_i, \quad (2)$$

where

$$\mu^T = \left(\frac{(p_1 - q_1) \sqrt{\gamma \delta n}}{\sqrt{\gamma p_1 + \delta q_1}}, \dots, \frac{(p_m - q_m) \sqrt{\gamma \delta n}}{\sqrt{\gamma p_m + \delta q_m}} \right),$$

and \mathbf{a}_i , $i = 1, \dots, v$, are the unit length eigenvectors of matrix J with corresponding eigenvalues λ_i , and the entries of J are given by

$$J_{ij} = \begin{cases} -\frac{1}{\sqrt{p_i \gamma + q_i \delta} \sqrt{p_j \gamma + q_j \delta}} \left[\left(1 + \frac{(p_i - q_i) \delta}{2(p_i \gamma + q_i \delta)} \right) \left(1 + \frac{(p_j - q_j) \delta}{2(p_j \gamma + q_j \delta)} \right) \gamma q_i q_j + \right. \\ \left. \left(1 - \frac{(p_i - q_i) \gamma}{2(p_i \gamma + q_i \delta)} \right) \left(1 - \frac{(p_j - q_j) \gamma}{2(p_j \gamma + q_j \delta)} \right) \delta p_i p_j \right] & \text{if } i \neq j \\ \frac{1}{p_i \gamma + q_i \delta} \left[\left(1 + \frac{(p_i - q_i) \delta}{2(p_i \gamma + q_i \delta)} \right)^2 \gamma q_i (1 - q_i) + \left(1 - \frac{(p_i - q_i) \gamma}{2(p_i \gamma + q_i \delta)} \right)^2 \delta p_i (1 - p_i) \right] & \text{if } i = j. \end{cases} \quad (3)$$

It can be verified that under the usual null hypothesis, i.e. when $p_i = q_i$ for all $i = 1, \dots, v$, the sum in (1) follows the chi-square distribution with $v - 1$ degree of freedom. In fact, under the null hypothesis $\mu = \mathbf{0}$ and J reduces to J^0 with $J_{ij}^0 = -\sqrt{p_i p_j}$ if $i \neq j$ and $J_{ij}^0 = 1 - p_i$ if $i = j$. Thus, J^0 has two eigenvalues, 1 with multiplicity $v - 1$, and 0 with multiplicity 1, which implies that $v - 1$ of the Z 's in (1) has standard normal distribution and the remaining one is 0 with probability 1.

The approximation in (1) is the asymptotic equivalent ([1]), which is at least as accurate as the central chi-square with $v - 1$ degrees of freedom for approximating the distribution of Pearson's statistic under the null hypothesis. However, there are too many parameters involved in J . Therefore, we now give another approximation that involves just one parameter and will subsequently be used to estimate p_0 . Note that all small fractions in parentheses in (3) are close to 0. By deleting these fractions, we obtain matrix G whose entries are

$$G_{ij} = \begin{cases} -\frac{1}{\sqrt{p_i \gamma + q_i \delta} \sqrt{p_j \gamma + q_j \delta}} [\gamma q_i q_j + \delta p_i p_j] & \text{if } i \neq j \\ \frac{1}{p_i \gamma + q_i \delta} [\gamma q_i (1 - q_i) + \delta p_i (1 - p_i)] & \text{if } i = j. \end{cases} \quad (4)$$

We show that the approximation based on matrix G rather than J has the form given in (??) if $\gamma = \delta$. If $\gamma = \delta$, then the entries of G equal

$$G_{ij} = \begin{cases} -\frac{1}{\sqrt{p_i+q_i}\sqrt{p_j+q_j}} [q_i q_j + p_i p_j] & \text{if } i \neq j \\ \frac{1}{p_i+q_i} [q_i (1-q_i) + p_i (1-p_i)] & \text{if } i = j. \end{cases} \quad (5)$$

First we show that the eigenvalues of G are 1 with multiplicity $m-2$, 0 with multiplicity 1 and $2 \sum_{i=1}^m \frac{p_i q_i}{p_i+q_i}$ with multiplicity 1. Matrix G can be written in the form

$$G = I - D (\mathbf{p}\mathbf{p}^T + \mathbf{q}\mathbf{q}^T) D,$$

where $D = \text{diag} \left(\frac{1}{\sqrt{p_i+q_i}} \right)$, $\mathbf{p} = (p_1, \dots, p_v)^T$ and $\mathbf{q} = (q_1, \dots, q_v)$. First we verify that vector $\mathbf{x} = D^{-1}\mathbf{1} = (\sqrt{p_1+q_1}, \dots, \sqrt{p_m+q_m})$ is an eigenvector of G with eigenvalue 0 by

$$G\mathbf{x} = \mathbf{x} - D (\mathbf{p}\mathbf{p}^T + \mathbf{q}\mathbf{q}^T) \mathbf{1} = \mathbf{x} - D (\mathbf{p} (\mathbf{p}^T \mathbf{1}) + \mathbf{q} (\mathbf{q}^T \mathbf{1})) = \mathbf{x} - D (\mathbf{p} + \mathbf{q}) = \mathbf{x} - \mathbf{x} = \mathbf{0},$$

where $\mathbf{1}$ is a v -dimensional vector with components 1. Furthermore, if $\mathbf{x} = D^{-1}\mathbf{z}$, where \mathbf{z} is orthogonal to both \mathbf{p} and \mathbf{q} , i.e. $\mathbf{p}^T \mathbf{z} = \mathbf{q}^T \mathbf{z} = 0$, then \mathbf{x} is an eigenvector of G with eigenvalue 1, because

$$G\mathbf{x} = \mathbf{x} - D (\mathbf{p}\mathbf{p}^T + \mathbf{q}\mathbf{q}^T) \mathbf{z} = \mathbf{x} - D (\mathbf{p}\mathbf{p}^T \mathbf{z} + \mathbf{q}\mathbf{q}^T \mathbf{z}) = \mathbf{x}.$$

Consequently, the eigenvectors of G with eigenvalue 1 span a $(v-2)$ -dimensional eigenspace if $\mathbf{p} \neq \mathbf{q}$, and a $(v-1)$ -dimensional eigenspace if $\mathbf{p} = \mathbf{q}$. Thus, if $\mathbf{p} = \mathbf{q}$, then there is no other eigenvector of G , and if $\mathbf{p} \neq \mathbf{q}$, then there is one more left, particularly $D(\mathbf{p} - \mathbf{q})$. This is verified by

$$\begin{aligned} \{I - D (\mathbf{p}\mathbf{p}^T + \mathbf{q}\mathbf{q}^T) D\} D(\mathbf{p} - \mathbf{q}) &= D \{(\mathbf{p} - \mathbf{q}) - (\mathbf{p}\mathbf{p}^T + \mathbf{q}\mathbf{q}^T) D^2 (\mathbf{p} - \mathbf{q})\}^* \\ &= D \{(\mathbf{p} - \mathbf{q}) + (\mathbf{p} - \mathbf{q}) \mathbf{q}^T D^2 (\mathbf{p} - \mathbf{q})\} = D(\mathbf{p} - \mathbf{q}) \{1 + \mathbf{q}^T D^2 (\mathbf{p} - \mathbf{q})\}, \end{aligned} \quad (6)$$

where at $*$ we used that $\mathbf{p}^T D^2 (\mathbf{p} - \mathbf{q}) = -\mathbf{q}^T D^2 (\mathbf{p} - \mathbf{q})$, which holds because of

$$\mathbf{p}^T D^2 (\mathbf{p} - \mathbf{q}) + \mathbf{q}^T D^2 (\mathbf{p} - \mathbf{q}) = (\mathbf{p} + \mathbf{q})^T D^2 (\mathbf{p} - \mathbf{q}) = \mathbf{1}^T (\mathbf{p} - \mathbf{q}) = 0.$$

The equality in (6) also shows that the eigenvalue corresponding to eigenvector $D(\mathbf{p} - \mathbf{q})$ is

$$1 + \mathbf{q}^T D^2 (\mathbf{p} - \mathbf{q}) = 1 + \sum_j \frac{q_j (p_j - q_j)}{p_j + q_j} = \sum_j \frac{q_j (p_j + q_j)}{p_j + q_j} + \sum_j \frac{q_j (p_j - q_j)}{p_j + q_j} = 2 \sum_j \frac{p_j q_j}{p_j + q_j}.$$

Finally, note that

$$D(\mathbf{p} - \mathbf{q}) = \left(\frac{p_1 - q_1}{\sqrt{p_1 + q_1}}, \dots, \frac{p_v - q_v}{\sqrt{p_v + q_v}} \right)^T = \sqrt{\frac{2}{n}} \mu.$$

Since G is symmetric, all eigenvectors corresponding to eigenvalue 1 are orthogonal to eigenvector $D(\mathbf{p} - \mathbf{q})$, and hence to μ .

The approximation to Pearson's statistic based on matrix G rather than J has the form

$$U_1^2 + \dots + U_v^2,$$

where U_1, \dots, U_v are independent normal random variables with $E(U_i) = \mathbf{b}_i \mu$ and $\text{Var}(U_i) = \varepsilon_i$, and \mathbf{b}_i , $i = 1, \dots, v$, are the unit length eigenvectors of G with corresponding eigenvalues ε_i . As we have shown above the eigenvalues of G are $\varepsilon_1 = \dots = \varepsilon_{v-2} = 1$, $\varepsilon_{v-1} = 2 \sum_{j=1}^v \frac{p_j q_j}{p_j + q_j}$, $\varepsilon_v = 0$ and $\mathbf{b}_1, \dots, \mathbf{b}_{v-2}$ are orthogonal to

μ , hence $U_v = 0$ and $U_1^2 + \dots + U_{v-2}^2$ is a central chi-square random variable with $v - 2$ degrees of freedom. Furthermore, since

$$\mathbf{b}_{v-1}^T = \frac{1}{\sqrt{\sum_{j=1}^v \frac{(p_j - q_j)^2}{p_j + q_j}}} \left(\frac{p_1 - q_1}{\sqrt{p_1 + q_1}}, \dots, \frac{p_v - q_v}{\sqrt{p_v + q_v}} \right),$$

we have

$$E(U_{v-1}) = \mathbf{b}_{v-1}^T \mu = \sqrt{\frac{n}{2} \sum_{j=1}^v \frac{(p_j - q_j)^2}{p_j + q_j}} = \sqrt{n} \Delta,$$

where

$$\Delta = \sqrt{\frac{1}{2} \sum_{j=1}^v \frac{(p_j - q_j)^2}{p_j + q_j}}.$$

Since

$$1 - \Delta^2 = \frac{1}{2} \sum_{j=1}^v \frac{(p_j + q_j)^2}{p_j + q_j} - \frac{1}{2} \sum_{j=1}^v \frac{(p_j - q_j)^2}{p_j + q_j} = 2 \sum_{j=1}^v \frac{p_j q_j}{p_j + q_j},$$

we have

$$\text{Var}(U_{v-1}) = 1 - \Delta^2.$$

We obtained that the approximation to Pearson's statistic based on G is the one given in (5) if $\gamma = \delta$ in the manuscript.

References

- [1] Bukszár J, Van den Oord EJCG. Accurate and efficient power calculations for $2 \times m$ tables in unmatched case-control designs. *Statistics in Medicine*. 2006; **25**:2632-2646.