

Expectation of R^2 with stratified phenotypes and genotypes

September 3, 2015

Let $Y_i, i = 1, \dots, n$ be the stratification-adjusted phenotype centered around 0 and $X_i, i = 1, \dots, n$ be the stratification-adjusted genotype centered around 0.

We have

$$\begin{aligned} E[X_i] &= 0 \forall i \\ E[Y_i] &= 0 \forall i \\ Var[X_i] &= \sigma_{X,i}^2 \\ Var[Y_i] &= \sigma_{Y,i}^2 \end{aligned}$$

The Pearson correlation coefficient is

$$\begin{aligned} \hat{\rho} &= \frac{\sum_{i=1}^n (X_i - \bar{X}) (Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}} \\ \hat{\rho} &= \frac{\sum_{i=1}^n X_i Y_i}{\sqrt{\sum_{i=1}^n (X_i)^2} \sqrt{\sum_{i=1}^n (Y_i)^2}} \\ R^2 &= \frac{(\sum_{i=1}^n X_i Y_i)^2}{\sum_{i=1}^n (X_i)^2 \sum_{i=1}^n (Y_i)^2} \\ &= \frac{(\sum_{i=1}^n X_i Y_i)^2}{\sum_{i=1}^n Var(X_i) \sum_{i=1}^n Var(Y_i)} \\ E(R^2) &= E \left[\frac{(\sum_{i=1}^n X_i Y_i)^2}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \right] \\ &= \frac{E[(\sum_{i=1}^n X_i Y_i)^2]}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \\ &= \frac{Var[\sum_{i=1}^n X_i Y_i] + E[\sum_{i=1}^n X_i Y_i]^2}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \\ &= \frac{Var[\sum_{i=1}^n X_i Y_i]}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \end{aligned}$$

$$\begin{aligned}
&= \frac{\sum_{i=1}^n \text{Var}(X_i Y_i)}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \\
&= \frac{\sum_{i=1}^n [E(X_i^2) E(Y_i^2) - E(X_i)^2 E(Y_i)^2]}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \\
&= \frac{\sum_{i=1}^n [E(X_i^2) E(Y_i^2)]}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \\
&= \frac{\sum_{i=1}^n [\text{Var}(X_i) \text{Var}(Y_i)]}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2} \\
&= \frac{\sum_{i=1}^n [\sigma_{X,i}^2 \sigma_{Y,i}^2]}{\sum_{i=1}^n \sigma_{X,i}^2 \sum_{i=1}^n \sigma_{Y,i}^2}
\end{aligned}$$

For a binary phenotype and a diploid genotype we can express this as

$$E(R^2) = \frac{\sum_{i=1}^n [2\hat{p}_{X,i}(1-\hat{p}_{X,i})\hat{p}_{Y,i}(1-\hat{p}_{Y,i})]}{\sum_{i=1}^n 2\hat{p}_{X,i}(1-\hat{p}_{X,i}) \sum_{i=1}^n \hat{p}_{Y,i}(1-\hat{p}_{Y,i})}$$