

Review

Estimating and testing the effect of allelic recombination on the correlation between genotypic relatives

Oyeka I. C. A¹., Oyeka C. A.² and Uzuke C. A.^{1*}

¹Department of Statistics Nnamdi Azikiwe University, Awka.

²Department of Applied Microbiology and Brewing Nnamdi Azikiwe University Awka.

Accepted 21 January, 2011

This paper provides estimates of the correlation between genotypic relatives and the effect of allelic recombination on the correlation assuming random mating. It is shown that the correlation is a non negative quantity and that allelic recombination has the effect of reducing total variation and doubling the correlation between genotypic relatives with respect to measurements on the character of interest. The significance of the correlation coefficient as well as the fitted regression model was obtained using Analysis of Variance method.

Key words: Allele, genotype, regression, correlation, F-ratio, analysis of variance.

INTRODUCTION

Genetic recombination is an effective means of combining one individual trait of two parents, permitting the comparison of one expression of a character with another expression of the same traits (Burns, 1976, Alberts, 1994, Maloy S, 1994).

Although much work has been done on the correlation between relatives for various physical characteristics starting with the pioneering work by Fisher (1918), very little has been written on the effect of genetic recombination on these correlations (Ewens, 1979, Oyeka and Oyeka, 1988). These writers also failed to provide a statistic for testing the significance of the estimated correlation coefficient. In this paper, the work by (Oyeka and Oyeka 1988) is modified to include a test statistic for the estimated correlation coefficient and to know if the hypothesised model fits.

CORRELATION BETWEEN RELATIVES

We will assume that a certain population has a gene locus with possible alleles A and a. We also assume that the probability of occurrence of allele A in the population

of interest is p and that the corresponding probability of the allele a is q = 1-p. We further assume that a certain characteristic or factor V of the population of interest is completely determined by the genotype at the locus in such a way that the individuals of genotype AA have a value or measurement of (V = x) for the character or factor of interest; all individuals of genotype Aa have a measurement of (V = y) and all individuals of genotype aa have measurement (V = z). We finally assume that our population obeys the Hardy-Weinberg law of random mating (Stein, 1943; Clavel et al., 1989).

Assuming n-pairs of relatives are studied, let R₁ and R₂ be a pair of relatives whom we know for sure have at least one allele in common. Then under the law of random mating, the occurrence of the allele A and a in the genotypes are independent. Hence the probabilities of occurrence of AA, Aa, or aA and aa are, respectively,

$$pp = p^2$$

$$p(1-p) + (1-p)p = 2p(1-p)$$

and

$$(1-p)(1-p) = (1-p)^2$$

*Corresponding author. E-mail: chinwe_uzuke@yahoo.com.

We first derive an estimate of the correlation between the

measurements on the character of genotypic relatives. To do this we first find the conditional probabilities that R_2 say, is of a certain genotype given the genotype of R_1 and then proceed to derive the joint probability distribution for the genotypes of R_1 and R_2 .

Now if R_1 is of genotype AA, then R_2 must have allele A in common with R_1 . Also since the occurrence of the second allele in R_2 is independent of the occurrence of the known allele A, the second allele in R_2 is either A with probability p or a with probability $q = 1-p$.

Now, since it is assumed that the relatives R_1 and R_2 must have at least one allele in common, the relative R_2 cannot be of genotype aa if R_1 is of genotype AA.

Hence, the required conditional probabilities are:

$$P(R_2 = AA / R_1 = AA) = p$$

$$P(R_2 = Aa / R_1 = AA) = 1 - p$$

$$P(R_2 = aa / R_1 = AA) = 0$$

If now R_1 is of genotype Aa, then R_2 must have either the allele A or the allele a in common with R_1 . Since the second allele occurs independently of the first allele in R_2 , the second allele is either A with probability p or a with probability $(1-p)$. Hence if R_1 is of genotype Aa, then R_2 is of genotype AA with probability p ; of genotype Aa (or aA) with probability $1-p + p = 1$ and of genotype aa with probability $1-p$.

Hence,

$$P(R_2 = AA / R_1 = Aa) = p$$

$$P(R_2 = Aa / R_1 = Aa) = 1$$

$$P(R_2 = aa / R_1 = Aa) = 1 - p$$

Similarly,

$$P(R_2 = AA / R_1 = aa) = 0$$

$$P(R_2 = Aa / R_1 = aa) = p$$

$$P(R_2 = aa / R_1 = aa) = 1 - p$$

Now to find the joint probability distribution of R_1 and R_2 , we apply the multiplication law of probability (Miller J 1996), which states that for any two events x and y ,

$$P(X = x, Y = y) = P(Y = y / X = x)P(X = x)$$

Hence,

$$P(R_1 = AA, R_2 = AA) = P(R_2 = AA / R_1 = AA)P(R_1 = AA) = p(p.p) = p^3$$

since the alleles in the genotype occur independently.

Similarly,

$$P(R_1 = AA, R_2 = Aa) = P(R_2 = Aa / R_1 = AA)P(R_1 = AA) = (1-p)(p.p) = p^2(1-p).$$

Other probabilities are:

$$P(R_1 = AA, R_2 = aa) = P(R_2 = aa / R_1 = AA)P(R_1 = AA) = (0)(p^2) = 0$$

$$P(R_1 = Aa, R_2 = AA) = P(R_2 = AA / R_1 = Aa)P(R_1 = Aa) = p \times p (1-p) = p^2 (1-p)$$

$$P(R_1 = Aa, R_2 = Aa) = P(R_2 = Aa / R_1 = Aa)P(R_1 = Aa) = 1 \times p (1-p) = p (1-p)$$

$$P(R_1 = Aa, R_2 = aa) = P(R_2 = aa / R_1 = Aa)P(R_1 = Aa) = (1-p) (p (1-p)) = p(1-p)^2$$

$$P(R_1 = aa, R_2 = AA) = P(R_2 = AA / R_1 = aa)P(R_1 = aa) = 0 \times (1-p)^2 = 0$$

$$P(R_1 = aa, R_2 = Aa) = P(R_2 = Aa / R_1 = aa) \times P(R_1 = aa) = p \times (1-p)^2 = p(1-p)^2$$

$$P(R_1 = aa, R_2 = aa) = P(R_2 = aa / R_1 = aa) \times P(R_1 = aa) = (1-p)(1-p)^2 = (1-p)^3$$

These calculations yield the results of Table 1 which shows the joint probability distribution of R_1 and R_2 , the marginal probability distribution and the corresponding measurements on the character of interest in the population.

Hence from the Table 1,

$$E(R_1) = p^2x + 2p(1-p)y + (1-p)^2z = m$$

And

$$E(R_2) = p^2x + 2p(1-p)y + (1-p)^2z = m \quad (1)$$

Table 1. Joint probability distribution of R₁ and R₂.

Genotype for relative R ₁	Measurement V	Genotype for relative R ₂			Marginal probability P(R ₁)
		AA	Aa	aa	
		X	Y	Z	
AA	X	P ³	P ² (1-p)	0	P ²
Aa	Y	P ² (1-p)	P(1-p)	P(1-p) ²	2P(1-p)
Aa	Z	0	P(1-p) ²	(1-p) ³	(1-p) ²
Marginal probability P(R ₂)		P ²	2P(1-p)	(1-p) ²	1

Therefore, the expectation of R₁ is equal to that of R₂. The corresponding variance S_t² on the measurement R₁ is given as:

$$\begin{aligned}
 V(R_1) &= E(R_1^2) - (E(R_1))^2 \\
 &= P^4X^2 - 2p(1-p)Y^2 + (1-p)^2Z^2 - (P^2X + 2P(1-p)Y + (1-p)^2Z)^2 \\
 &= P^2(1-P^2)X^2 + 2P(1-p)(1-2P-2P^2)Y^2 + (1-p)^2(2P(1-p)Z^2 - 2P^2(1-p^2)XZ - 4P(1-p)^3ZY - 4P^3(1-p)XY)
 \end{aligned}$$

And that of R₂ is V(R₂)= E(R₂²) - (E(R₂))²

$$\begin{aligned}
 &= P^2X^2 - 2p(1-p)Y^2 + (1-p)^2Z^2 - (P^2X + 2P(1-p)Y + (1-p)^2Z)^2 \\
 &= P^2(1-P^2)X^2 + 2P(1-p)(1-2P-2P^2)Y^2 + (1-p)^2(2P(1-p)Z^2 - 2P^2(1-p^2)XZ - 4P(1-p)^3ZY - 4P^3(1-p)XY) \\
 &= V(R_1)
 \end{aligned}$$

Hence, the variance of the measurement on R₁ which is the same as the measurement on R₂ is given as

Equation 1.

$$\begin{aligned}
 S_t^2 &= p^2x^2 + 2p(1-p)y^2 + (1-p)^2z^2 - m^2 \quad (2) \\
 \text{Where } m &= E(R_1) = p^2x + 2p(1-p)y + (1-p)^2z \text{ from}
 \end{aligned}$$

The covariance between R₁ and R₂ is also calculated in the usual way from the table as:

$$S_{12} = \text{Cov}(R_1, R_2) = E(R_1, R_2) - E(R_1)E(R_2) \quad (\text{Uche, 2004}).$$

$$S_{12} = p^3x^2 + p^2(1-p)xy + p^2(1-p)xy + p(1-p)y^2 + p(1-p)^2yz + p(1-p)^2yz + (1-p)^3z^2 - m^2$$

Where m = E(R₁) = E(R₂)

$$\begin{aligned}
 &= p^3x^2 + 2p^2(1-p)xy + p(1-p)y^2 + 2p(1-p)^2yz + (1-p)^3z^2 - m^2 \\
 &= p^3x^2 + 2p^2(1-p)xy + p(1-p)y^2 + 2p(1-p)^2yz + (1-p)^3z^2 - (p^2x + 2p(1-p)y + (1-p)^2z)^2
 \end{aligned}$$

$$\begin{aligned}
 S_{12} &= p(1-p)((p^2x^2 + 2p(1-2p)xy + (1-2p)^2y^2 - 2p(1-p)xz - 2(1-p)(1-2p)yz + (1-p)^2z^2) \\
 S_{12} &= p(1-p)[px + (1-2p)y - (1-p)z]^2 \quad (3)
 \end{aligned}$$

The correlation, r, between R₁ and R₂ is found by dividing Equation (3) by Equation (2) since the variance of R₁ is the same as that of R₂.

$$r = \frac{\text{cov}(R_1, R_2)}{\text{var}(R_1)} \text{ , since } \text{var}(R_1) = \text{var}(R_2)$$

Thus,

Then,

$$r = \frac{S_{12}}{S_e^2} \quad (4)$$

Note that since $p \geq 0$.

The covariance, S_{12} and hence the correlation, r , is a non negative quantity, and for $0 < p < 1$, has a value zero only when

$$p = \frac{x-y}{x+z-2y} \quad (5)$$

Provided x and z are both greater than y or x and z are both less than y .

EFFECT OF ALLELIC RECOMBINATION

Let us now examine what would happen to the measurements of the character of interest and hence the correlation if we recombine the alleles by replacing one allele of a genotype by another allele. Specifically, and without loss of generality, suppose we replace an allele A by an allele a in a genotype determining a certain character in an individual, by this allelic replacement model, the original individual must possess an A allele and hence must be of genotype AA or of genotype Aa .

Hence, if the replacement is being made in an AA

individual, the resulting effect on the measurement on the character of interest is to reduce y by x ; that is

$$y - x,$$

while if the allelic replacement is being made on an Aa individual the effect would be

$$z - y.$$

Interest is now on finding the differential effect of this allelic recombination on the measurement of the character concerned and its significance. We propose to do this using the method of least squares. Let us, therefore, find the best estimates, in the least square sense, of the parameters μ , α and β that would minimize the expected sum of squared deviations of x , y , and z from $\mu + 2\alpha$, $\mu + \alpha + \beta$, and $\mu + 2\beta$, respectively, assuming random mating and subject to the constraint

$$p\alpha + (1-p)\beta = 0 \quad (6)$$

Where α = the effect of allele A on the character of interest and β = the effect of allele a on the character of interest.

The expected sum of squared deviations of observed from their true values of the measurements using the marginal probability distribution of Table 1.

$$S^2_e = p^2(x - \mu - 2\alpha)^2 + 2p(1-p)(y - \mu - \alpha - \beta)^2 + (1-p)^2(z - \mu - 2\beta)^2 \quad (7)$$

Differentiating S^2_e with respect to μ and minimizing we have

$$\frac{dS^2_e}{d\mu} = 2p^2(x - \mu - 2\alpha) + 4p(1-p)(y - \mu - \alpha - \beta) + 2(1-p)^2(z - \mu - 2\beta) = 0.$$

When solved replacing μ , α , and β by their estimates $\hat{\mu}$, $\hat{\alpha}$ and $\hat{\beta}$, respectively, and simplifying yields

$$p^2x + 2p(1-p)y + (1-p)^2z = p^2(\hat{\mu} + 2\hat{\alpha}) + 2p(1-p)(\hat{\mu} + \hat{\alpha} + \hat{\beta}) + (1-p)^2(\hat{\mu} + 2\hat{\beta})$$

The right hand side of the above equation reduces to

$$p^2\hat{\mu} + 2p^2\hat{\alpha} + 2p\hat{\mu} + 2p\hat{\alpha} + 2p\hat{\beta} - 2p^2\hat{\mu} - 2p^2\hat{\alpha} - 2p^2\hat{\beta} + \hat{\mu} + 2\hat{\alpha} - 2p\hat{\mu} - 4p\hat{\alpha} + p^2\hat{\mu} + 2p^2\hat{\beta}$$

$$= 2[p\hat{\alpha} + (1-p)\hat{\beta}] + \hat{\mu} \\ = \hat{\mu}.$$

Since from equation (6)

$$p\alpha + (1-p)\beta = 0.$$

Hence we have that

$$\hat{\mu} = m = p^2x + 2p(1-p)y + (1-p)^2z \quad (8)$$

Also, differentiating S^2_e with respect to α , yields

$$\frac{dS^2_e}{d\alpha} = -4p^2(x - \mu - 2\alpha) - 4p(1-p)(y - \mu - \alpha - \beta) = 0 \\ = -4[p^2(x - \mu - 2\alpha) + p(1-p)(y - \mu - \alpha - \beta)] = 0$$

Therefore,

$$p^2(x - \mu - 2\alpha) + p(1 - p)(y - \mu - \alpha - \beta) = 0$$

Replacing the values of α and β by their estimates $\hat{\alpha}$ and $\hat{\beta}$ and solving, we have:

$$p^2(x - m) + p(1 - p)(y - m) = p[p\hat{\alpha} + (1 - p)\hat{\beta} + \hat{\alpha}]$$

$$\begin{aligned} \frac{dS^2_e}{d\beta} &= -2p(1 - p)(y - \mu - \alpha - \beta) - 4p(1 - p)^2(2 - \mu - 2\beta) \\ &= -2p[(1 - p)(y - \mu - \alpha - \beta) + 2(1 - p)^2(z - \mu - 2\beta)] = 0 \end{aligned}$$

$$p(1 - p)(y - m) + 2(1 - p)^2(z - m) = (1 - p)\beta$$

Replacing β by its estimate $\hat{\beta}$ and solving we obtain

$$\begin{aligned} \hat{\beta} &= p(y - m) + 2(1 - p)(z - m) \\ \hat{\beta} - \hat{\alpha} &= [p(y - m) + (1 - p)(z - m)] - [p(x - m) + (1 - p)(y - m)] \\ &= [py - pm + z - pz - m + pm] - [px - pm + y - py - m + pm] \\ &= py - pm + z - pz - m + pm - px + pm - y + py + m - pm \\ &= py - px + py - pz + z - y \end{aligned} \tag{10}$$

$$= p(y - x) + (1 - p)(z - y) \tag{11}$$

Substituting equations (8) – (10) into equation (7), we obtain, after simplification, the minimum value of the sum of squared deviation, S^2_e , or the so-called “error sum of squares” in regression analysis parlance, as

$$S^2_e = p^2(1 - p)^2[x + z - 2y]^2 \tag{12}$$

For $0 < p < 1$, S^2_e assumes the value of zero only if

$$S^2_e = S^2_f - S^2_g = 2p(1 - p)(px + (1 - 2p)y + (1 - p)z)^2 \tag{14}$$

S^2_e is similar to a regression sum of squares and may be interpreted as that part of the total variation S^2_f in the measurement of the characters of interest attributable to the average effects of α and β of the alleles A and a. We note from equation (4) that recombination has the effect of reducing the value of the total variation S^2_f by S^2_g , the residual variance. Also, it can be seen from Equation (3) that:

$$= p(0 + \hat{\alpha}) = p\hat{\alpha}$$

Hence, we then have that

$$\hat{\alpha} = p(x - m) + (1 - p)(y - m) \tag{9}$$

Differentiating S^2_e with respect to β we have:

Hence, the differential effect of replacing the allele A by the allele a on the measurement of the character concerned is estimated using Equations (9) and (10) as:

$$y = \frac{1}{2}(x + z) \tag{13}$$

That is if the measurement on individuals who are heterozygous on the character of interest is equal to half the sum of the measurements on homozygous individuals. The difference between the total sum of square S^2_f of equation (2) and S^2_e of equation 12 namely, the sum of squares due to regression model that is that part of the total variation S^2_g accounted for or explained by the regression model.

$$S^2_g = 2S_{12} \tag{15}$$

Thus the allelic recombination doubles the covariance between genotypic relatives.

Now the proportion of total variance, S^2_g explained by α and β is

$$Q = \frac{S^2_g}{S^2_f} = R^2 \tag{16}$$

Table 2. Analysis of variance table for the hypothesised regression model.

Source of variation	Sum of square	Degrees of freedom	Mean square	F-ratio
Regression	$SSR = S_{\hat{y}}^2 = 2p(1-p)(px + (1-2p)y + (1-p)z)^2$	2	$MSR = \frac{SSR}{2} = \frac{S_{\hat{y}}^2}{2}$	$F = \frac{MSR}{MSE} = \frac{(n-3)S_{\hat{y}}^2}{2S_{\hat{e}}^2} = \frac{(n-3)(2p(1-p)(px + (1-2p)y + (1-p)z)^2)}{2p^2(1-p)^2[x+z-2y]^2}$
Error	$SSE = S_{\hat{e}}^2 = p^2(1-p)^2[x+z-2y]^2$	n-3	$MSE = \frac{SSE}{n-3} = \frac{S_{\hat{e}}^2}{n-3}$	
Total	$SST = S_{\hat{y}}^2 = p^2x^2 + 2p(1-p)y^2 + (1-p)^2z^2 - m^2$	n-1		

Where R^2 is the usual coefficient of determination in regression parlance.

Or equivalently from equation (4) we have that

$$Q = \frac{2S_{\hat{y}}^2}{S_{\hat{e}}^2} = 2r \tag{17}$$

In other words, the proportion of total variance in the measurements on the character of interest that is accounted for by the effect of manipulation of the alleles is equal to twice the correlation between genotypic relatives in the absence of allelic recombination.

TESTING THE SIGNIFICANCE OF THE FIT OF THE REGRESSION MODEL AND THE CORRELATION BETWEEN GENOTYPIC RELATIVES

One may be interested in testing the hypothesis that the regression model fits. That is, testing whether the differential effects of the allelic recombination on the character of interest are statistically different from zero. To test the null hypothesis H_0 , we may use the Analysis of Variance method (Montgomery and Peck, 1992). The three sums of squares, their associated degrees of freedom, their mean squares, and the resulting F-ratio are summarised in Table 2.

The F-ratio =

$$\frac{(n-3)(2p(1-p)(px + (1-2p)y + (1-p)z)^2)}{2p^2(1-p)^2[x+z-2y]^2} \tag{18}$$

which has an F-distribution with (2, n-3) degrees of freedom and may be compared at an α significance level with tabulated critical F-value to test that the regression model fits. If:

$$F\text{-ratio} > F_{(1-\alpha), (2, n-3)}$$

We reject the null hypothesis of no differential effects of allelic recombination on the genotypic relatives.

We may also wish to test the null hypothesis that allelic recombination has no significant effect on the correlation between genotypic relatives that is that the population correlation coefficient ρ due to allelic recombination of genotypic relatives is zero versus the alternative hypothesis that ρ is different from zero.

Now, note that testing that the regression model fits is equivalent to testing the hypothesis that, in the sampled population, the values of Q which is equal to R^2 is equal to $2r$ not equal to zero.

The significance of the sample estimates of these population parameters is tested using the usual F-test of equation (18). The rejection of the null hypothesis implies that the population values of Q

equals R^2 is not equal to or equivalently, Q equals $2r$ is not equal to zero in the population implying that r is not equal to zero in the population sampled. Hence the usual F-test provides a test statistic for testing the null hypothesis that the correlation between genotypic relatives is zero versus the alternative hypothesis that the correlation between genotypic relatives is different from zero. The null hypothesis is rejected at an appropriately chosen significance level α .

CONCLUSION

This paper provided estimates of the correlation between genotypic relatives both in the presence and in the absence of allelic recombination. It is shown that the correlation between genotypic relatives in the absence of allelic recombination is double the correlation between genotypic relatives in the presence of allelic recombination. The correlation obtained is a non negative quantity and except for trivial cases ($p = 0$ or 1), assuming the value zero only for some critical value p . The significance of the correlation obtained and the regression model fitted are tested using the analysis of variance technique.

REFERENCES

Alberts B (1994). Molecular Biology of Cells (3rd edition)

- Garland Publishing Company.
- Burns G (1976). *The Science of Genetics: An Introduction to Heredity* (3rd edition) Macmillan Publishing Company, New York.
- Clavel F, Hoggan MD, Willey RL, Strelbel K, Martin M A, Rapaske R (1989). Genetic Recombination of Human Immunodeficiency Virus. *J. Virol.*, 63(3): 1455-1459.
- Ewens WJ (1979). *Mathematical Population Genetics*, Springer – Varlang, New York.
- Fisher RA (1918). The Correlation Between Relatives on the Supposition of Mendelian Inheritance. *Transactions Royal Society of Edinburgh*. 52: 399–433.
- Maloy S, Cronan J, Freifelder (1994). *Microbial Genetics* (2nd Edition), Jones and Bartlet MA.
- Miller J (1996). *Statistics for Advanced Level* (2nd Edition) Cambridge University Press.
- Montgomery DC, Peck EA (1992). *Introduction to Linear Regression Analysis*. John Willey and Sons, 2nd Edition.
- Oyeka ICA (2009). Applied Statistical Methods in Sciences Norben Avocation. pp: 135–137.
- Oyeka ICA, Oyeka CA (1988). Estimating the Effect of Allelic Recombination on the Correlation Between Genotypic Relatives. *Nigerian J. Biotechnol.*, Enugu, Nigeria.
- Uche PI (2004). *Probability: Theory and Practice* Longman Nig PLC.