

Caitlin McHugh
Oct 2014

To investigate association testing on the X chromosome, I implemented simulation studies using various numbers of genotypes that are generated using the pedigree shown in Figure 1.

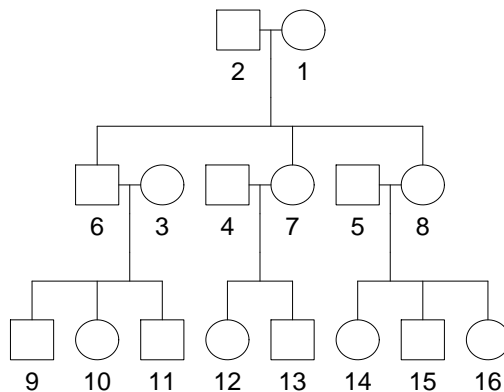


Figure 1: The 16-person pedigree used for the simulations.

The full model assumed when testing for association on X chromosome SNPs (and one of the models used for simulating the X chromosome quantitative phenotype) is

$$y = \beta_0 + \beta_1 \text{SNP}_x + g_A + g_X + \epsilon \quad (1)$$

$$g_A \sim MVN(0, \sigma_A^2 \Phi_A) \quad (2)$$

$$g_X \sim MVN(0, \sigma_X^2 \Phi_X) \quad (3)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2) \quad (4)$$

where SNP_x is the vector of genotypes on the X chromosome SNP that is being tested for association, Φ_A is the matrix of kinship coefficients as measured on the autosomes and Φ_X is the matrix of X chromosome specific kinship coefficients. The male X chromosome genotypes are coded as 0, 2 and the female genotypes are coded as 0, 1, 2.

We calculate the variance for a given individual i to be the sum of the variance of the SNP being tested, and the variances due to X chromosome, autosomal and other effects as

$$\text{var}(y_i) = \beta_1^2 \text{var}(\text{SNP}_x) + \sigma_X^2 + \sigma_A^2 + \sigma_\epsilon^2 \quad (5)$$

The variance of an X chromosome SNP can be calculated conditionally based on

whether the sample is female or male.

$$\mathbb{E}(\text{SNP}_x^F) = 2p^2 + 2p(1 - p) = 2p \quad (6)$$

$$\mathbb{E}(\text{SNP}_x^M) = 2p \quad (7)$$

when the male genotypes are coded as 0, 2 and the female genotypes are coded as 0, 1, 2. Then, we find that

$$\text{var}(\text{SNP}_x^F) = \mathbb{E}((\text{SNP}_x^F)^2) - \mathbb{E}^2(\text{SNP}_x^F) \quad (8)$$

$$= 4p^2 + 2p(1 - p) - (2p)^2 \quad (9)$$

$$= 2p(1 - p) \quad (10)$$

$$\text{var}(\text{SNP}_x^M) = \mathbb{E}((\text{SNP}_x^M)^2) - \mathbb{E}^2(\text{SNP}_x^M) \quad (11)$$

$$= 4p - (2p)^2 \quad (12)$$

$$= 4p(1 - p) \quad (13)$$

To calculate the covariance of genotypes between a pair of individuals, we must consider their sex. In what follows, I am denoting the X chromosome kinship value between a pair of individuals as Φ_X . Technically, this should be written as $\Phi_{X,ij}$ where ij indexes the individuals i and j for which the X chromosome kinship value represents. First, we calculate the covariance for a SNP between a pair of males as

$$\text{cov}(\text{SNP}_x^M, \text{SNP}_x^M) = \mathbb{E}(\text{SNP}_x^M \text{SNP}_x^M) - \mathbb{E}^2(\text{SNP}_x^M) \quad (14)$$

$$= \mathbb{E}(\text{SNP}_x^M \text{SNP}_x^M | \text{IBD}) \mathbb{P}(\text{IBD}) \quad (15)$$

$$+ \mathbb{E}(\text{SNP}_x^M \text{SNP}_x^M | \text{no IBD}) \mathbb{P}(\text{no IBD}) - (2p)^2 \quad (16)$$

$$= \mathbb{E}(\text{SNP}_x^M \text{SNP}_x^M | \text{IBD}) \Phi_X \quad (17)$$

$$+ \mathbb{E}(\text{SNP}_x^M \text{SNP}_x^M | \text{no IBD}) (1 - \Phi_X) - 4p^2 \quad (18)$$

$$= 4p\Phi_X + 4p^2(1 - \Phi_X) - 4p^2 \quad (19)$$

$$= 4p(1 - p)\Phi_X \quad (20)$$

Next we consider the covariance between a pair of female genotypes

$$\text{cov}(\text{SNP}_x^F, \text{SNP}_x^F) = \mathbb{E}(\text{SNP}_x^F \text{SNP}_x^F) - \mathbb{E}^2(\text{SNP}_x^F) \quad (21)$$

$$= \mathbb{E}(\text{SNP}_x^F \text{SNP}_x^F | \text{IBD}) \mathbb{P}(\text{IBD}) \quad (22)$$

$$+ \mathbb{E}(\text{SNP}_x^F \text{SNP}_x^F | \text{no IBD}) \mathbb{P}(\text{no IBD}) - (2p)^2 \quad (23)$$

$$= (2(2p(1 - p)) + 4p^2)\Phi_X + 4p^2(1 - \Phi_X) - 4p^2 \quad (24)$$

$$= 4p(1 - p)\Phi_X \quad (25)$$

Finally, we calculate the covariance between a pair of genotypes where one is a female

and one is a male

$$\text{cov}(\text{SNP}_x^F, \text{SNP}_x^M) = \mathbb{E}(\text{SNP}_x^F \text{SNP}_x^M) - \mathbb{E}(\text{SNP}_x^F) \mathbb{E}(\text{SNP}_x^M) \quad (26)$$

$$= \mathbb{E}(\text{SNP}_x^F \text{SNP}_x^M | \text{IBD}) \mathbb{P}(\text{IBD}) \quad (27)$$

$$+ \mathbb{E}(\text{SNP}_x^F \text{SNP}_x^M | \text{no IBD}) \mathbb{P}(\text{no IBD}) - (2p)^2 \quad (28)$$

$$= (4p^2 + 2(2p(1-p)))\Phi_X + 4p^2(1 - \Phi_X) - 4p^2 \quad (29)$$

$$= 4p(1-p)\Phi_X \quad (30)$$

We can now see that using the X chromosome kinship values from Table 5, the variance for a female and male SNP is indeed as calculated in Equations 10 and 13 after incorporating the self-kinship values. Thus, the variance for a given individual i for an X chromosome SNP is

$$\text{var}(y_i) = \beta_1^2 4p(1-p)\Phi_X + \sigma_X^2 + \sigma_A^2 + \sigma_\epsilon^2 \quad (31)$$

The parameter h_{snp}^2 indicates the heritability of the X chromosome SNP. It can be calculated from the equation

$$h_{\text{snp}}^2 = \frac{\beta_1^2 4p(1-p)\Phi_X}{\beta_1^2 4p(1-p)\Phi_X + \sigma_\epsilon^2 + \sigma_A^2 + \sigma_X^2} \quad (32)$$

where p is the allele frequency of the causal SNP. On the other hand, we can calculate the heritability of all SNPs on the X chromosome, which is

$$h_x^2 = \frac{\beta_1^2 4p(1-p)\Phi_X + \sigma_X^2}{\beta_1^2 4p(1-p)\Phi_X + \sigma_\epsilon^2 + \sigma_A^2 + \sigma_X^2} \quad (33)$$

Relatedness Estimation Using X Chromosome SNPs

For the proposed model to work, we first must convince ourselves that we can accurately estimate relatedness using genetic material from the X chromosome. Table 5 displays the autosomal and X chromosome kinship coefficients (KC) for a given pair of relatives. The autosomal KC is defined as the probability of sampling two alleles IBD from a given pair of individuals. The X chromosome KC, on the other hand, is defined as the probability of sampling one allele IBD from an X chromosome in a given pair of individuals. Thus, the X chromosome KC will differ from the usual autosomal KC. Furthermore, when calculating the theoretical X chromosome KC, we must take into account the sex of the individuals, whether the pair is related maternally or paternally, and the type of relationship.

X chromosome relatedness can be estimated between two individuals using SNP genotypes. Let j, k, l, m be four individuals such that j, k are male and l, m are female. The genotypes for each of the independent X chromosome SNPs in individual j are denoted by X_{ij} for SNPs $i \in \{1, \dots, N\}$ and are coded as 0, 1, 2 in females and 0, 2 in males. We can estimate the genetic relatedness (GR) using SNP genotypes, which is

twice the X chromosome KC in female-female pairs, $\sqrt{2}$ the X chromosome KC in female-male pairs and equal to the X chromosome KC in male-male pairs. The X chromosome genetic relatedness between two individuals is

$$\text{GR}_{FF} = \frac{1}{N} \sum_{i=1}^N \frac{(X_{il} - 2p_i)(X_{im} - 2p_i)}{2p_i(1 - p_i)} \quad (34)$$

$$\text{GR}_{MM} = \frac{1}{N} \sum_{i=1}^N \frac{(X_{ij} - p_i)(X_{ik} - p_i)}{p_i(1 - p_i)} \quad (35)$$

$$\text{GR}_{MF} = \frac{1}{N} \sum_{i=1}^N \frac{(X_{ij} - p_i)(X_{il} - 2p_i)}{\sqrt{2}p_i(1 - p_i)} \quad (36)$$

To investigate how accurate we can estimate the X chromosome KC, I simulated varying numbers of X chromosome SNPs for one iteration of the 16-sample pedigree shown in Figure 1. The allele frequency of the SNPs was set at 0.4. Figure 2 shows the difference between the estimated and theoretical X chromosome KC for all sample pairs for increasing numbers of SNPs. Figures 3, 4 and 5 show the results broken up by the composition of sex in the related pair. Figure 6 shows a histogram of the estimated X chromosome KC for each relationship type. The true value is shown with a red dotted line and the number of relationships at a given KC value is displayed in the plot title. Many of the relationships have a small sample size. Some relationships are underestimated but all are generally centered around the truth.

As expected, with a larger number of SNPs, we are able to more accurately estimate the true X chromosome KC. We note from all Figures that the estimated KC is at most 0.06 away from the true value. The OLGA genotyping set resulted in approximately 3,500 SNPs on the X chromosome after pruning. Perhaps this number of SNPs, shown in orange in the Figures, should be considered most realistic. We conclude that we are sufficiently able to estimate the X chromosome KC from 3,500 independent, genotyped X chromosome SNPs.

Variance Components Estimation

I estimated the variance components for the autosomes, the X chromosome, and the remaining effects, using the true kinship matrix in both cases. I then fit the mixed model for a quantitative trait on the X chromosome, testing the genotypes simulated on the X chromosome.

Initially, I investigated whether the estimated variance components were converging to what I expected. I simulated 10,000 independent pedigrees as shown in Figure 1 for a total of 16,000 individuals, of whom 5,000 are unrelated (5 founders per pedigree), and performed this simulation 500 times. The relatedness matrices used were the true values, not the estimated ones. I estimated variance components for three models (where the

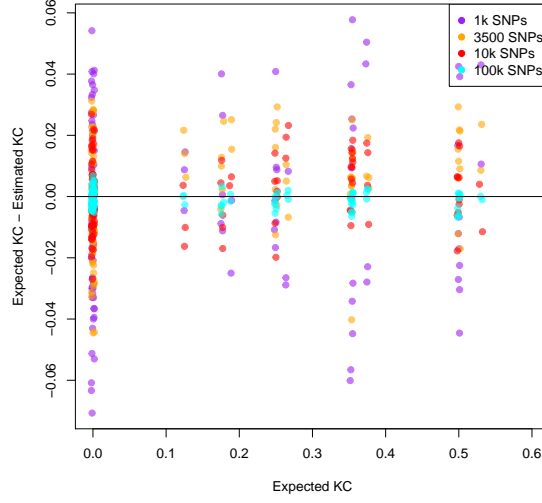


Figure 2: The expected X chromosome KC versus the difference between the expected and estimated X chromosome KC, for all sample pairs in the 16-person pedigree shown in Figure 1. The colors indicate the number of simulated X chromosome SNPs used to estimate the KC.

phenotype was simulated from the given model):

$$y = \beta_1 \text{SNP}_x + \epsilon \quad (37)$$

$$y = \beta_1 \text{SNP}_x + g_X + \epsilon \quad (38)$$

$$y = \beta_1 \text{SNP}_x + g_X + g_A + \epsilon \quad (39)$$

with the specifications of

$$g_A \sim \text{MVN}(0, 0.3\Phi_A)$$

$$g_X \sim \text{MVN}(0, 0.8\Phi_X)$$

$$\epsilon \sim N(0, 1)$$

$$\beta_1 = 0.8$$

$$p = 0.2$$

From computations as shown above, in each of the three models we expect the estimates for σ_X^2 , σ_A^2 and σ_ϵ^2 to be as displayed in Table 1 where $\sigma_{XT}^2 = \beta_1^2 4p(1-p) + \sigma_X^2$. Table 2 shows the results from the described simulation study. The standard deviations shown there are the mean lower and upper bounds as provided from Matt's confidence intervals, i.e. the confidence intervals calculated from each of the 500 iterations. We note that the simulations yield a mean value that is equal to the expected value. We conclude the variance components are being estimated accurately.

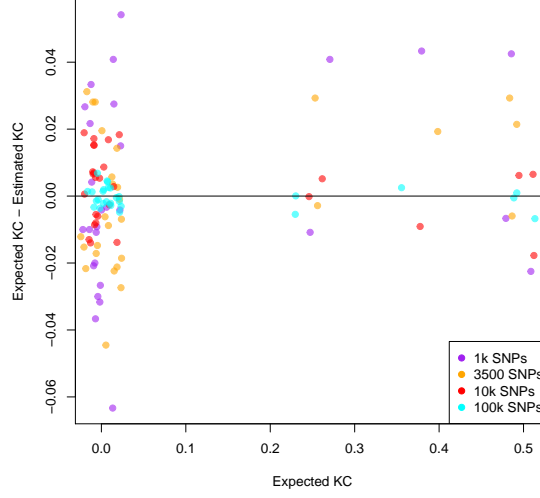


Figure 3: The expected X chromosome KC versus the difference between the expected and estimated X chromosome KC, for all male-male sample pairs in the 16-person pedigree shown in Figure 1. The colors indicate the number of simulated X chromosome SNPs used to estimate the KC.

Association Testing on the X Chromosome

I estimated the variance components for the autosomes, the X chromosome, and the remaining effects, using the known, theoretical autosomal and X chromosome kinship matrices. I fit the mixed model for a quantitative trait on the X chromosome, testing the genotypes simulated on the X chromosome.

We can evaluate the type I error and power when the true model is

$$y = \beta_1 \text{SNP}_x + g_X + g_A + \epsilon \quad (40)$$

and when fitting the misspecified, usual model, $y = \beta_1 \text{SNP}_x + g_A + \epsilon$. I also fit the misspecified model $y = \beta_1 \text{SNP}_x + g_X + \epsilon$ for comparison. From 1,000 iterations using 8,000 samples (500 pedigrees as displayed in Figure 1) of which 5,500 are related and 2,500 are unrelated, we set the parameters as shown in Table 3. Because the usual model is not properly calibrated, we compare the false positive rate to the true positive rate. In this manner, we can identify how many true positives (the power) we are able to detect for a given false positive rate (type I error).

When the effect size reaches 0.18, all three models tested always found the true signal. Plots showing the The type I error, however, differed in each of the models. Table 4 shows the type I errors when fitting the two models.

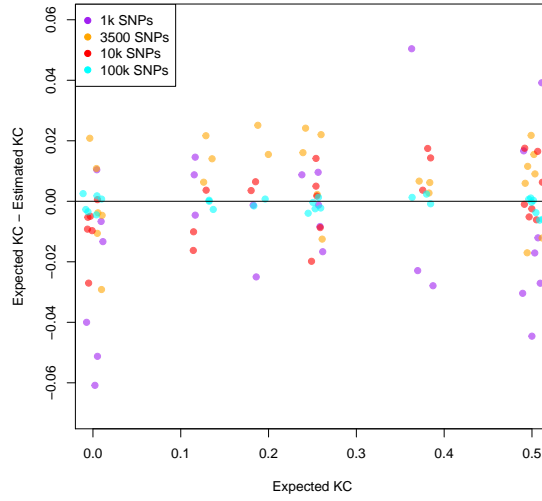


Figure 4: The expected X chromosome KC versus the difference between the expected and estimated X chromosome KC, for all female-female sample pairs in the 16-person pedigree shown in Figure 1. The colors indicate the number of simulated X chromosome SNPs used to estimate the KC.

Performing Mixed Model Association Testing on the X Chromosome

To perform mixed model association testing on X chromosome SNPs, there are a few steps to take.

1. Estimate relatedness on the X chromosome, call the results Φ_X . Use independent (pruned) SNPs, excluding the pseudoautosomal regions. In OLGA, the number of pruned X chromosome SNPs was approximately 3,500. Histograms of the results for theoretical vs estimated in 3,500 simulated X chromosome SNPs for the pedigree in Figure 1 are shown in the file 'hist_xchrKC_byRelType.pdf.'
2. Run Matt's MLM program, including as a random effect the relatedness matrix on the X chromosome, Φ_X . We are investigating the implications of including Φ_A as well, when testing a SNP on the X chromosome for association.

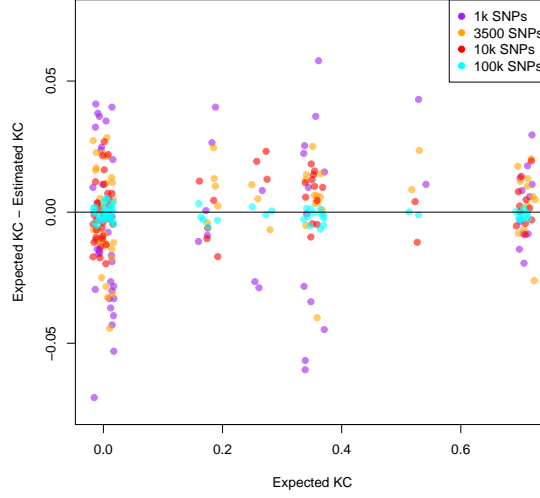


Figure 5: The expected X chromosome KC versus the difference between the expected and estimated X chromosome KC, for all female-male sample pairs in the 16-person pedigree shown in Figure 1. The colors indicate the number of simulated X chromosome SNPs used to estimate the KC.

Model	σ_{XT}^2	σ_A^2	σ_ϵ^2
1	0.4096	-	1
2	1.2096	-	1
3	1.2096	0.3	1

Table 1: Values of simulated variance components.

Model	σ_{XT}^2	σ_A^2	σ_ϵ^2
1	0.4098 (0.3667, 0.4529)	-	1.001 (0.9685, 1.033)
2	1.210 (1.136, 1.284)	-	0.9995 (0.9623, 1.037)
3	1.211 (1.103, 1.319)	0.3035 (0.1303, 0.4768)	1.000 (0.9525, 1.048)

Table 2: Mean (mean CI bounds) from simulation results for three models, where the variance components were simulated as shown in Table 1. The simulation included 16,000 samples, of which 5,000 were unrelated.

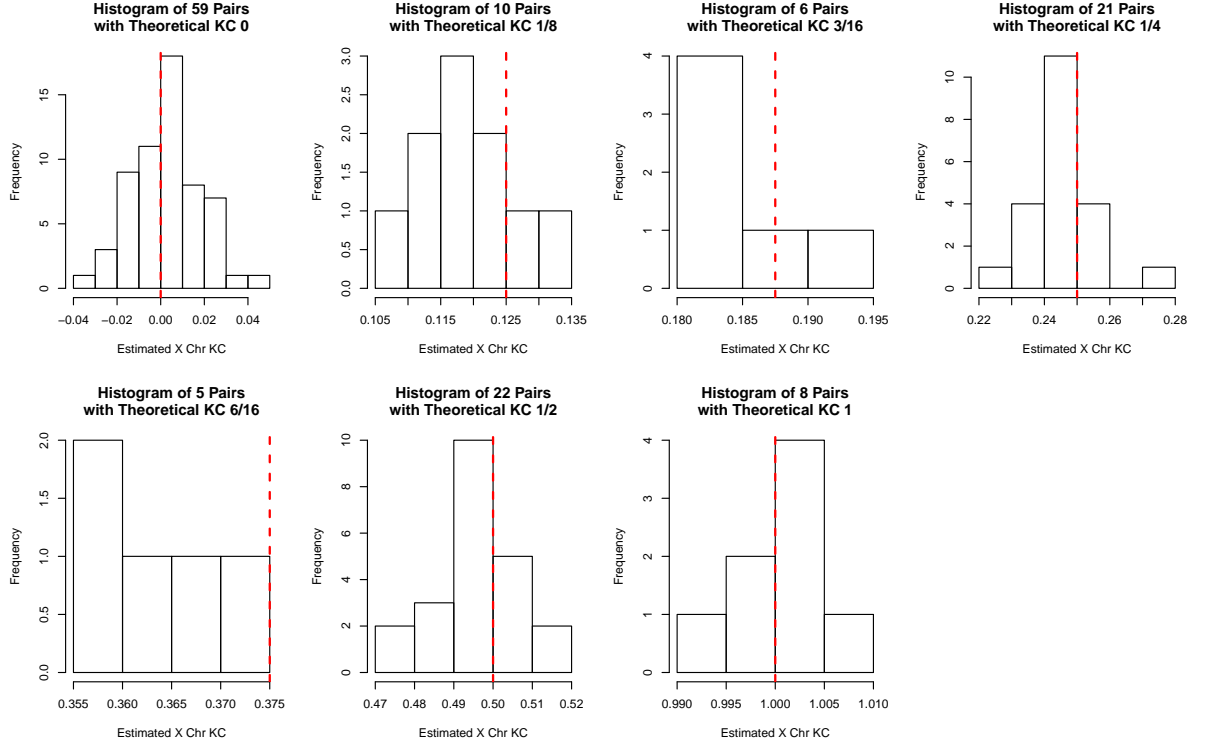


Figure 6: Histograms of the estimated X chromosome KC, by relationship type, for all pairs of relatives shown in the pedigree in Figure 1. The true value is shown with a red dotted line. The estimates were calculated from 3,500 simulated SNPs.

Parameter	Sim1	Sim2	Sim3
β_1	0.8	0.3	0.3
σ_A^2	0.3	0.3	0.3
σ_X^2	0.8	0.8	0.3
σ_ϵ^2	1	1	1
σ_{XT}^2	1.2096	0.8576	0.3576

Table 3: Parameters used in simulations.

α	0.05	0.01	0.005	0.001	0.0001
X + Auto adj	0.04604	0.00834	0.00450	0.00092	0
Auto adj only	0.06823	0.01960	0.01034	0.00133	0
X chr adj only	0.04446	0.00817	0.00509	0.00100	0

Table 4: Type I error from 12,000 iterations of an 8,000 sample simulation. The true model is as described in Equation 39 and the results are shown from fitting the true model, the model without fitting X chromosome effects, $y = \beta_1 \text{SNP}_x + g_A + \epsilon$, and the model fitting only X chromosome effects $y = \beta_1 \text{SNP}_x + g_X + \epsilon$.

α	Auto + X	X	Auto
0.05 (0.0496, 0.0504)	0.04604 (0.04214, 0.04994)	0.04446 (0.04056, 0.04836)	0.06823 (0.06433, 0.07213)
0.01 (0.00960, 0.0104)	0.00834 (0.00656, 0.01012)	0.00817 (0.00639, 0.00996)	0.01960 (0.01782, 0.02138)
0.005 (0.00460, 0.0054)	0.00450 (0.00324, 0.00577)	0.00509 (0.00383, 0.00635)	0.01034 (0.00908, 0.01161)
0.001 (6.00e-04, 0.0014)	0.00092 (0.00035, 0.00148)	0.00100 (0.00044, 0.00157)	0.00133 (0.00077, 0.00190)
5e-04 (9.98e-05, 0.0009)	0 (-0.00040, 0.00040)	0.00017 (-0.00023, 0.00057)	0 (-0.00040, 0.00040)

		Autosomes	X Chromosome
Maternal	Self, Female	$\frac{1}{2}$	$\frac{1}{2}$
	Self, Male	$\frac{1}{2}$	1
	Mother-Daughter	$\frac{1}{4}$	$\frac{1}{4}$
	Mother-Son, Father-Daughter	$\frac{1}{4}$	$\frac{1}{2}$
	Father-Son	$\frac{1}{4}$	0
	Full sisters	$\frac{1}{4}$	$\frac{6}{16}$
	Full brothers	$\frac{1}{4}$	$\frac{1}{2}$
	Sister-Brother	$\frac{1}{4}$	$\frac{1}{4}$
	Aunt-Niece	$\frac{1}{8}$	$\frac{3}{16}$
	Aunt-Nephew	$\frac{1}{8}$	$\frac{6}{16}$
	Uncle-Niece	$\frac{1}{8}$	$\frac{1}{8}$
	Uncle-Nephew	$\frac{1}{8}$	$\frac{1}{4}$
	Grandma-Granddaughter	$\frac{1}{8}$	$\frac{1}{8}$
	Grandma-Grandson	$\frac{1}{8}$	$\frac{1}{4}$
	Grandpa-Granddaughter	$\frac{1}{8}$	$\frac{1}{4}$
	Grandpa-Grandson	$\frac{1}{8}$	$\frac{1}{2}$
Paternal	Aunt-Niece	$\frac{1}{8}$	$\frac{1}{8}$
	Aunt-Nephew	$\frac{1}{8}$	0
	Uncle-Niece	$\frac{1}{8}$	0
	Uncle-Nephew	$\frac{1}{8}$	0
	Grandma-Granddaughter	$\frac{1}{8}$	$\frac{1}{4}$
	Grandma-Grandson	$\frac{1}{8}$	0
	Grandpa-Granddaughter	$\frac{1}{8}$	0
	Grandpa-Grandson	$\frac{1}{8}$	0
Maternal-Maternal	First cousins, Girl-Girl	$\frac{1}{16}$	$\frac{3}{32}$
	First cousins, Girl-Boy	$\frac{1}{16}$	$\frac{3}{16}$
	First cousins, Boy-Boy	$\frac{1}{16}$	$\frac{6}{16}$
Paternal-Paternal	First cousins, Girl-Girl	$\frac{1}{16}$	$\frac{1}{32}$
	First cousins, Girl-Boy	$\frac{1}{16}$	0
	First cousins, Boy-Boy	$\frac{1}{16}$	0
Paternal-Maternal	First cousins, Girl-Girl	$\frac{1}{16}$	$\frac{1}{16}$
	First cousins, Girl-Boy	$\frac{1}{16}$	0
	First cousins, Boy-Boy	$\frac{1}{16}$	0

Table 5: The theoretical kinship coefficients (KC) stratified by X chromosome and autosomes. The autosomal KC value is $\frac{1}{2}\kappa_2 + \frac{1}{4}\kappa_1$, where κ_1 and κ_2 are the probabilities of sampling one and two alleles IBD, respectively. The X chromosome KC value is the probability of sampling one allele IBD on the X chromosome in a given pair of individuals.

		Autosomes	X Chromosome
	Self, Male	$\frac{1}{2}$	1
	Self, Female	$\frac{1}{2}$	$\frac{1}{2}$
	Mother-Son, Father-Daughter	$\frac{1}{4}$	$\frac{1}{2}$
	Full brothers	$\frac{1}{4}$	$\frac{1}{2}$
Maternal	Grandpa-Grandson	$\frac{1}{8}$	$\frac{1}{2}$
	Full sisters	$\frac{1}{4}$	$\frac{6}{16}$
Maternal	Aunt-Nephew	$\frac{1}{8}$	$\frac{6}{16}$
	Mother-Daughter	$\frac{1}{4}$	$\frac{1}{4}$
	Sister-Brother	$\frac{1}{4}$	$\frac{1}{4}$
Maternal	Uncle-Nephew	$\frac{1}{8}$	$\frac{1}{4}$
Maternal	Grandma-Grandson	$\frac{1}{8}$	$\frac{1}{4}$
Maternal	Grandpa-Granddaughter	$\frac{1}{8}$	$\frac{1}{4}$
Paternal	Grandma-Granddaughter	$\frac{1}{8}$	$\frac{1}{4}$
Maternal	Aunt-Niece	$\frac{1}{8}$	$\frac{3}{16}$
Maternal-Maternal	First cousins, Girl-Boy	$\frac{1}{16}$	$\frac{3}{16}$
Maternal-Maternal	First cousins, Boy-Boy	$\frac{1}{16}$	$\frac{6}{16}$
Maternal	Uncle-Niece	$\frac{1}{8}$	$\frac{1}{8}$
Maternal	Grandma-Granddaughter	$\frac{1}{8}$	$\frac{1}{8}$
Paternal	Aunt-Niece	$\frac{1}{8}$	$\frac{1}{8}$
Maternal-Maternal	First cousins, Girl-Girl	$\frac{1}{16}$	$\frac{3}{32}$
Paternal-Maternal	First cousins, Girl-Girl	$\frac{1}{16}$	$\frac{1}{16}$
Paternal-Paternal	First cousins, Girl-Girl	$\frac{1}{16}$	$\frac{1}{32}$
	Father-Son	$\frac{1}{4}$	0
Paternal	Aunt-Nephew	$\frac{1}{8}$	0
Paternal	Uncle-Niece	$\frac{1}{8}$	0
Paternal	Uncle-Nephew	$\frac{1}{8}$	0
Paternal	Grandma-Grandson	$\frac{1}{8}$	0
Paternal	Grandpa-Granddaughter	$\frac{1}{8}$	0
Paternal	Grandpa-Grandson	$\frac{1}{8}$	0
Paternal-Paternal	First cousins, Girl-Boy	$\frac{1}{16}$	0
Paternal-Paternal	First cousins, Boy-Boy	$\frac{1}{16}$	0
Paternal-Maternal	First cousins, Girl-Boy	$\frac{1}{16}$	0
Paternal-Maternal	First cousins, Boy-Boy	$\frac{1}{16}$	0

Table 6: The theoretical kinship coefficients (KC) stratified by X chromosome and autosomes. The autosomal KC value is $\frac{1}{2}\kappa_2 + \frac{1}{4}\kappa_1$, where κ_1 and κ_2 are the probabilities of sampling one and two alleles IBD, respectively. The X chromosome KC value is the probability of sampling one allele IBD on the X chromosome in a given pair of individuals.