

Association Testing with X Chromosome Data

An Application To HCHS/SOL

Caitlin McHugh, with Tim Thornton

Department of Biostatistics
University of Washington

5 Feb 2015

Outline

Intro

Step 1: Estimate Φ_X

Step 2: Determine Covariates

Step 3: Fit the Model

- ▶ For each genotyped *autosomal* SNP, we can fit a model of the form

$$Y = \beta_0 + \beta_1 \text{SNP} + g_A + \text{covariates} + \epsilon$$

where

$$g_A \sim \text{MVN}(0, \sigma_A^2 \Phi_A)$$

$$\epsilon \sim \text{MVN}(0, \sigma_\epsilon^2 \mathbb{I})$$

- ▶ When testing association on X chromosome SNPs, we propose to fit the model

$$Y = \beta_0 + \beta_1 \text{SNP}_X + g_A + g_X + \text{covariates} + \epsilon$$

where further

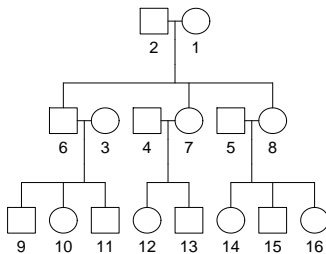
$$g_X \sim \text{MVN}(0, \sigma_X^2 \Phi_X)$$

- ▶ The X chromosome kinship coefficient between individuals i and j , Φ_{ij}^X , is defined as the probability of sampling one allele IBD at random from individual i and individual j on the X chromosome.
- ▶ *Note* for males there is no randomness in sampling, as there is only one allele at each location on the X chromosome.

Step 1: Estimate Φ_X

Step 2: Determine Covariates

Step 3: Fit the Model



		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	$\frac{1}{8}$	$\frac{6}{16}$
Paternal	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	$\frac{1}{8}$	$\frac{1}{2}$
	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

- We can estimate Φ_X using the following GRM equations:

$$GR_{FF} = \frac{1}{N} \frac{\sum_{i=1}^N (X_{il} - 2p_i)(X_{im} - 2p_i)}{\sum_{i=1}^N 2p_i(1 - p_i)}$$

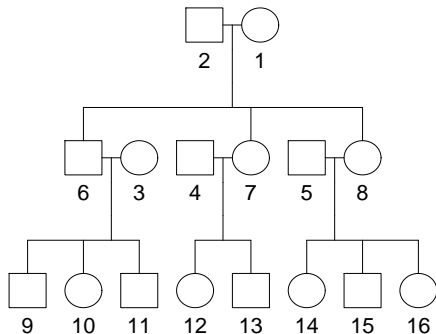
$$GR_{MM} = \frac{1}{N} \frac{\sum_{i=1}^N (X_{ij} - p_i)(X_{ik} - p_i)}{\sum_{i=1}^N p_i(1 - p_i)}$$

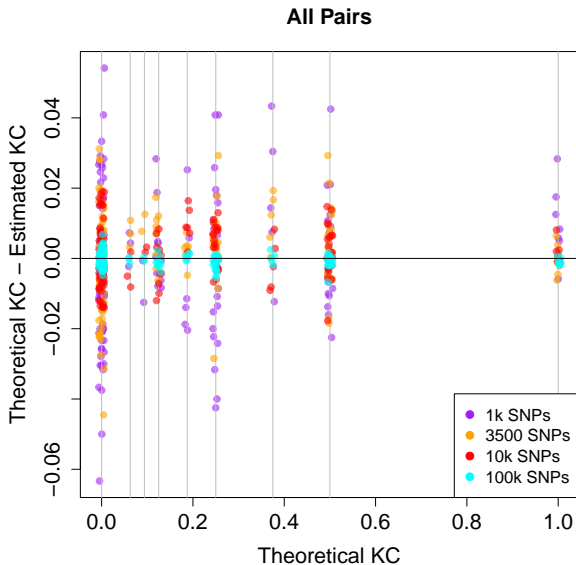
$$GR_{MF} = \frac{1}{N} \frac{\sum_{i=1}^N (X_{ij} - p_i)(X_{il} - 2p_i)}{\sum_{i=1}^N \sqrt{2}p_i(1 - p_i)}$$

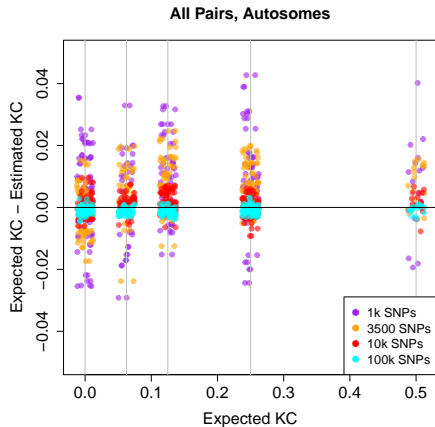
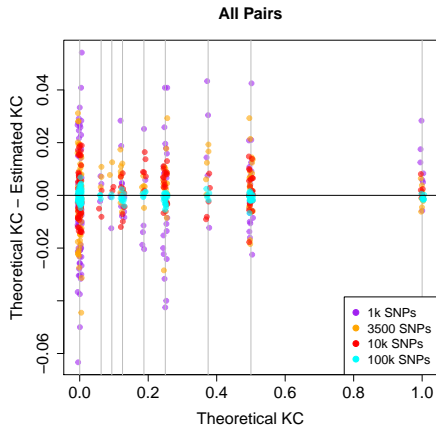
- ▶ When X chromosome genotypes are coded 0, 1, 2 for females and 0, 2 for males, the GRM equations for the X chromosome are the same for the autosomes.
- ▶ With this genotype coding, the covariance between two individuals i and j on the X chromosome is $4p(1 - p)\Phi_{ij}^X$, regardless of the sex of the individuals. {with proof, if needed}

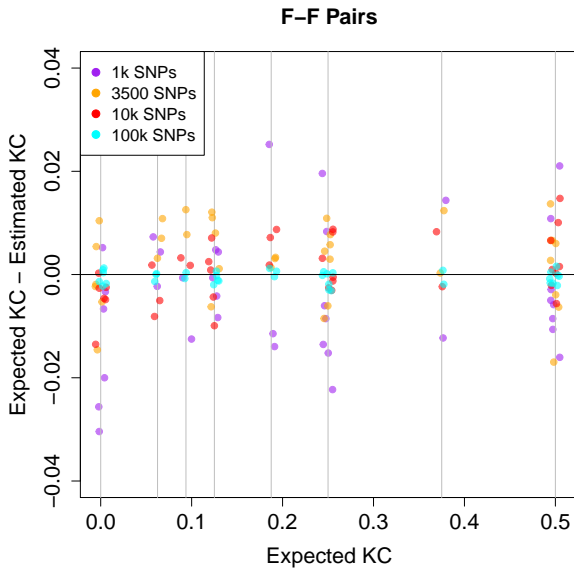
- ▶ Can we accurately estimate relatedness using X chromosome SNPs?
- ▶ After pruning, there are approximately 3,500 X chromosome SNPs in the SoL data.

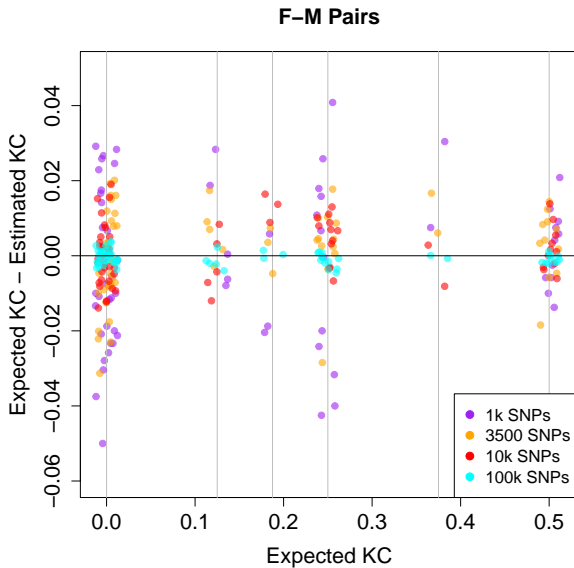
- ▶ Simulate increasing number of X chromosome SNPs for 16 individuals in this pedigree structure.
- ▶ *Note* this is assuming a homogeneous population.
- ▶ Individs 2 and 6 have $\Phi_{2,6}^X = 0$ but $\Phi_{2,6}^A = 1/4$.

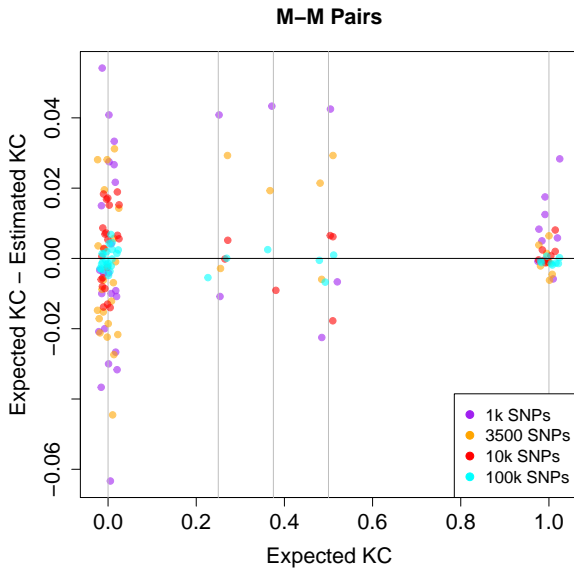






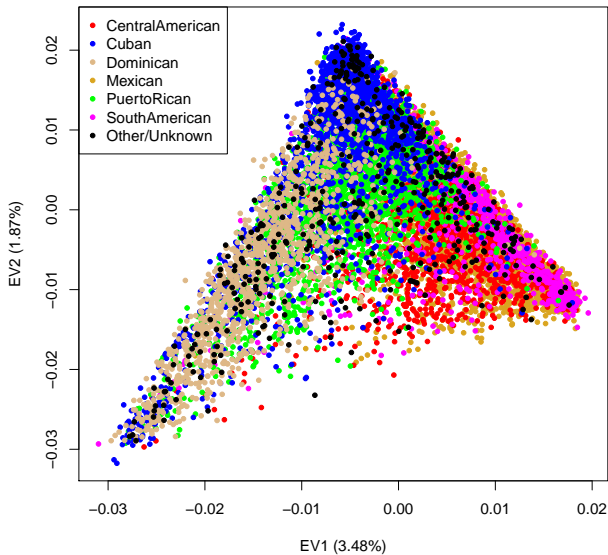


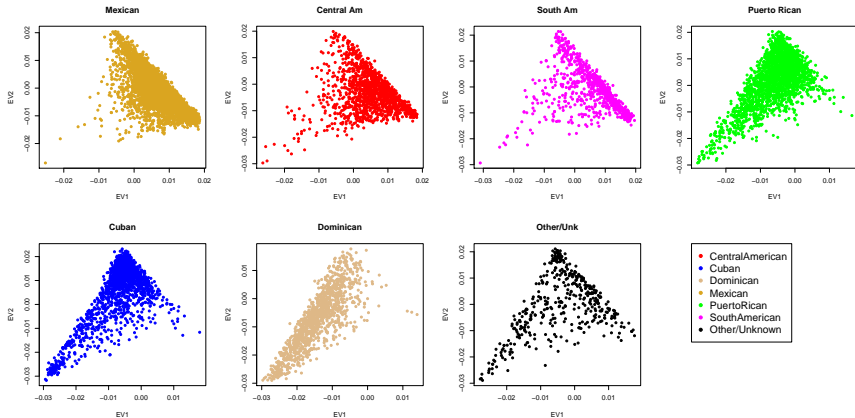


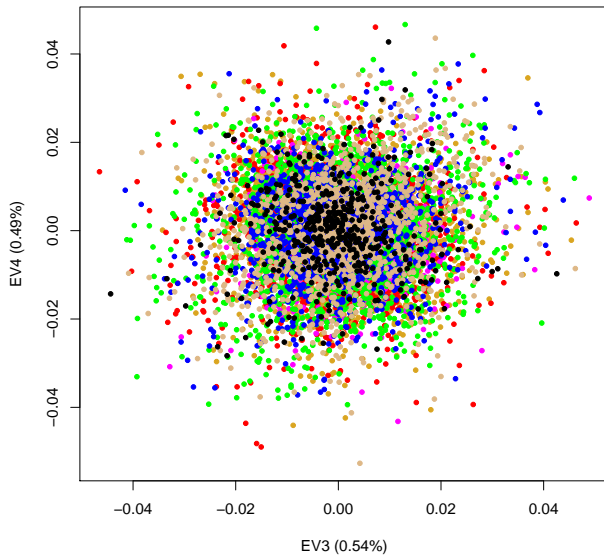


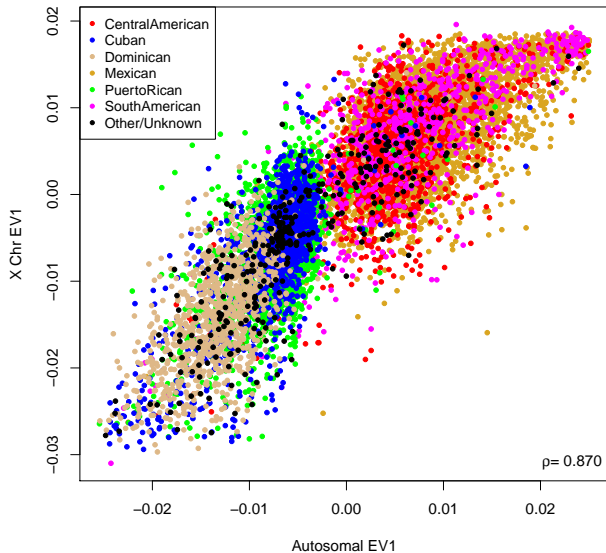
- ▶ What about samples from an admixed population?
- ▶ We need to determine if we can find 'ancestry adjusted' relatedness estimates on the X chromosome.
 - ▶ Using global ancestry, such as PCs estimated on the X.
 - ▶ Using local ancestry, such as average X chromosome local ancestry. {we need X chr local ancestry estimates}
 - ▶ Or, could we use IBD estimates on the X chromosome from BEAGLE?

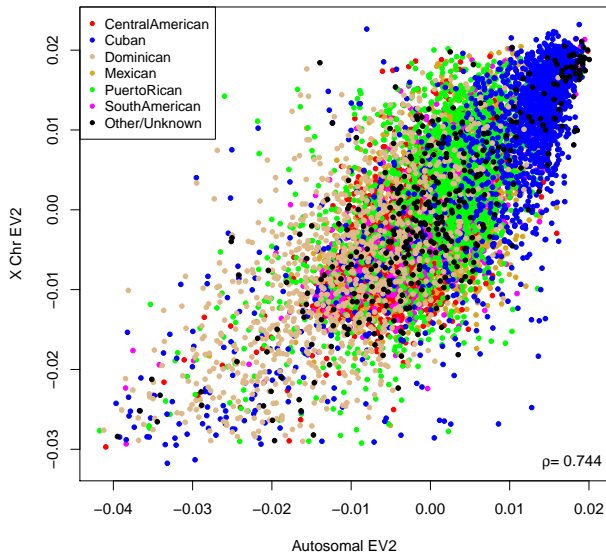
- ▶ We estimated PCs in the SOL subjects using 3,600 LD-pruned X chromosome SNPs and PC-AiR.
- ▶ The unrelated set `unrelated.pcair.deg4` of 10,272 samples as defined from the autosomes was set, and only study samples (`subj.plink & geno.cnt1==0`) excluding `gengrp6.outliers` were projected for a total of 12,747 samples.



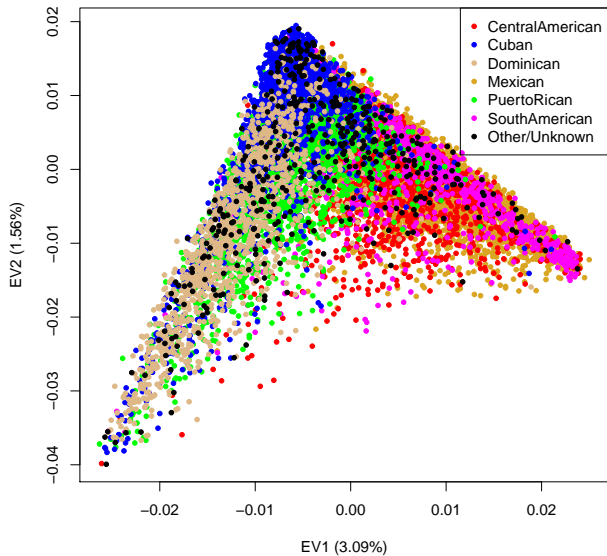


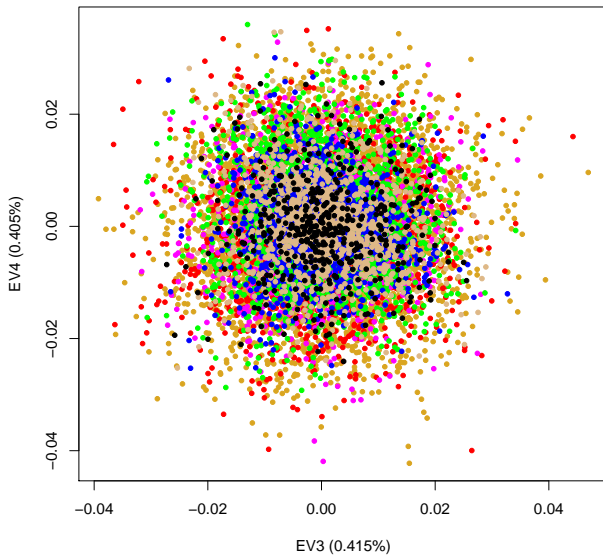


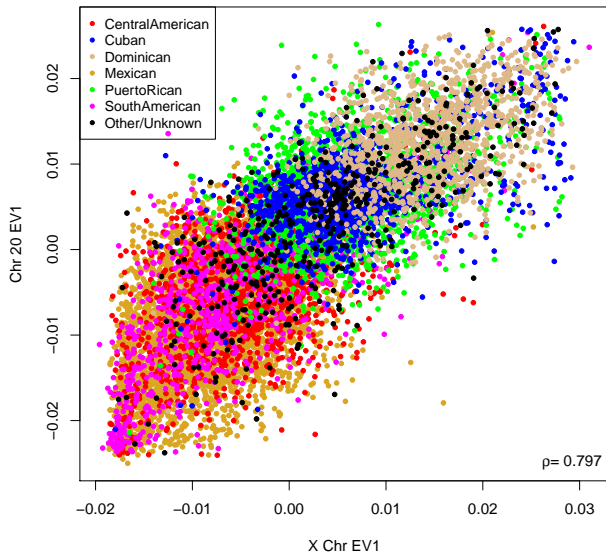


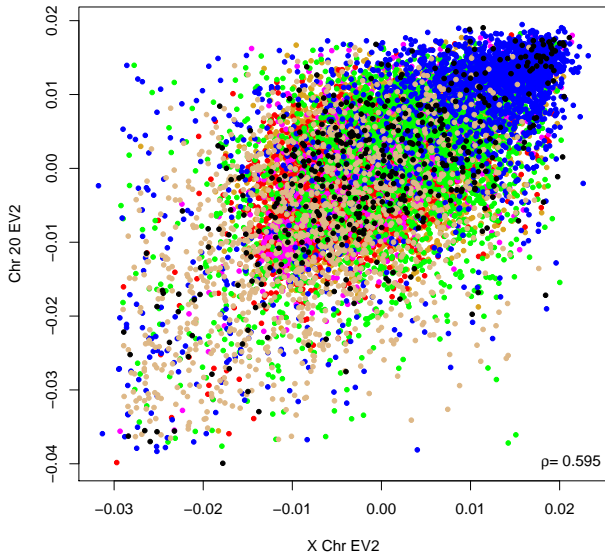


- ▶ We compare these results which use only 3,600 X chromosome SNPs to a pruned set of 4,413 chromosome 20 SNPs.







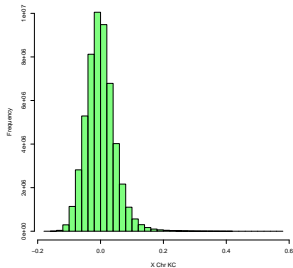
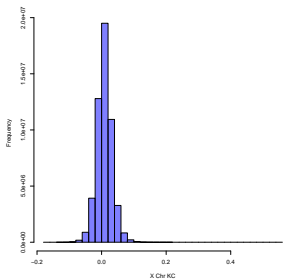
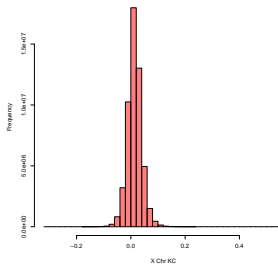
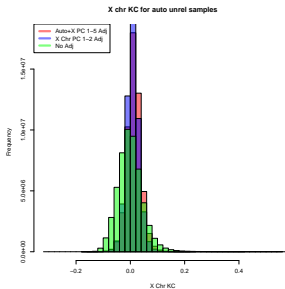


We estimated Φ_X for all OLGA samples using the following scenarios:

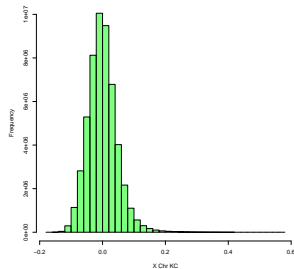
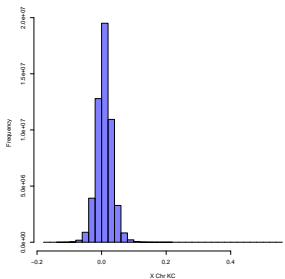
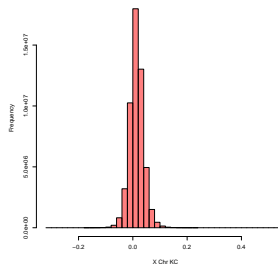
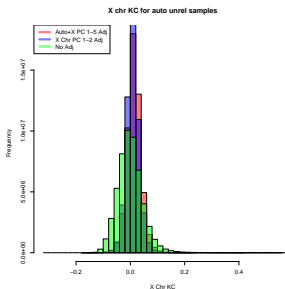
	SNP set	PCA run	EVs used
1	all X chr	-	-
2	all X chr	autosomes + X chr	1-5
3	all X chr	autosomes	1-5
4	all X chr	X chr	1-2
5	pruned X chr	-	-
6	pruned X chr	autosomes + X chr	1-5
7	pruned X chr	autosomes	1-5
8	pruned X chr	X chr	1-2
9	pruned autosomal	autosomes	1-5

All settings used the autosomal unrelated set of 10,272 samples and estimated Φ_X for 12,734 study samples posted to dbGaP with no X chromosome anomalies.

Estimate of Φ_X in 10,272 autosomal-unrelated samples models 1-4.



Estimate of Φ_X in 10,272 autosomal-unrelated samples models 5-8.



- ▶ What covariates should we include in the model?
- ▶ Should we include PCs calculated across the autosomes?
- ▶ Should we include PCs calculated on the X chromosome?
ADMIXTURE estimates from the X chromosome? Local
ancestry estimates averaged across the X chromosome?

- What happens if we simply fit our usual autosomal model when testing X chromosome markers?

- ▶ Simulate X chromosome genotypes for 8,000 samples = 500 iterations of the 16-person pedigree; 2,500 unrelated + 5,500 relatives.
- ▶ *Note* these samples are from a homogeneous population.
- ▶ Simulate phenotypes that follow the model

$$Y = \beta_0 + \beta_1 \text{SNP}_X + g_A + g_X + \epsilon$$

$$g_A \sim \text{MVN}(0, \sigma_A^2 \Phi_A)$$

$$g_X \sim \text{MVN}(0, \sigma_X^2 \Phi_X)$$

$$\epsilon \sim \text{MVN}(0, \mathbb{I})$$

where $\sigma_A^2 = 0.3$, $\sigma_X^2 = 0.8$ and for various β_1 values.

- ▶ Fit three models:

$$Y = \beta_0 + \beta_1 \text{SNP}_X + g_A + g_X + \epsilon$$

$$Y = \beta_0 + \beta_1 \text{SNP}_X + g_X + \epsilon$$

$$Y = \beta_0 + \beta_1 \text{SNP}_X + g_A + \epsilon$$

α	Adj for X + Auto	Adj for X	Adj for Auto
0.05	0.04983	0.04867	0.07325
0.01	0.00931	0.00913	0.02080
0.005	0.00494	0.00534	0.01095
0.001	0.00094	0.00102	0.00191

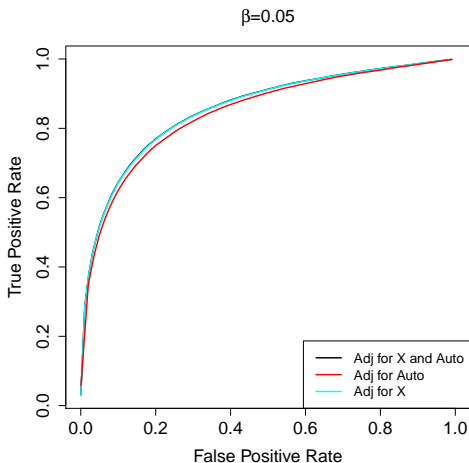
Table: Type I error rate as calculated from 22,000 simulation iterations where $\beta_1=0$.

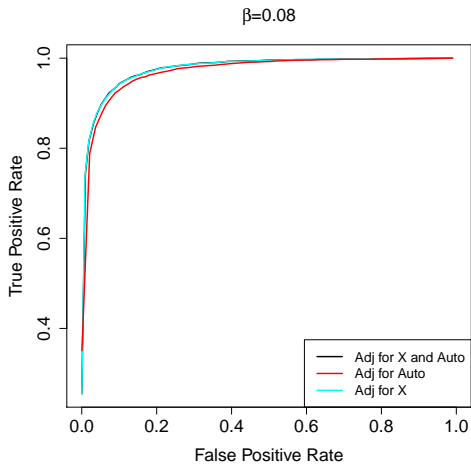
- ▶ When only adjusting for autosomal effects, no type I error CIs include the α under consideration.

α	Estimate	95% CI
0.05	0.07325	(0.07040, 0.07610)
0.01	0.02080	(0.01949, 0.02210)
0.005	0.01095	(0.01003, 0.01188)
0.001	0.00191	(0.00150, 0.00233)

Table: Estimate and 95% CI for type I error rate as calculated from 22,000 simulations where $\beta_1=0$ and fitting the model $Y = \beta_0 + \beta_1 \text{SNP}_X + g_A + \epsilon$.

- For each α value ranging from $1e-04$ to 1, and 7,500 simulation iterations, we calculate the true and false positive rate for each model.





Supplementary Slides

$$\begin{aligned} \text{cov}(F, M) &= \mathbb{E}(F, M) - \mathbb{E}(F)\mathbb{E}(M) \\ &= \mathbb{E}(FM|\text{IBD})\Phi_X + \mathbb{E}(FM|\text{not IBD})(1 - \Phi_X) - (2p)^2 \\ &= [4p^2 + 4p(1 - p)]\Phi_X + [4p^3 + 4p^2(1 - p)](1 - \Phi_X) - 4p^2 \\ &= 4p(1 - p)\Phi_X \end{aligned}$$

Type I Error CIs

α	Estimate	95% CI
0.05	0.04867	(0.04582, 0.05152)
0.01	0.00913	(0.00783, 0.01043)
0.005	0.00534	(0.00442, 0.00627)
0.001	0.00102	(0.00061, 0.00144)

Table: Estimate and 95% CI for type I error rate as calculated from 22,000 simulations where $\beta_1=0$ and fitting the model

$$Y = \beta_0 + \beta_1 \text{SNP}_X + g_X + \epsilon.$$

Type I Error CIs

α	Estimate	95% CI
0.05	0.04983	(0.04698, 0.05268)
0.01	0.00931	(0.00801, 0.01061)
0.005	0.00494	(0.00402, 0.00587)
0.001	0.00094	(0.00052, 0.00135)

Table: Estimate and 95% CI for type I error rate as calculated from 22,000 simulations where $\beta_1=0$ and fitting the model

$$Y = \beta_0 + \beta_1 \text{SNP}_X + g_A + g_X + \epsilon.$$

