# Mixed Models: How to Effectively Account for Inbreeding and Population Structure in GWAS
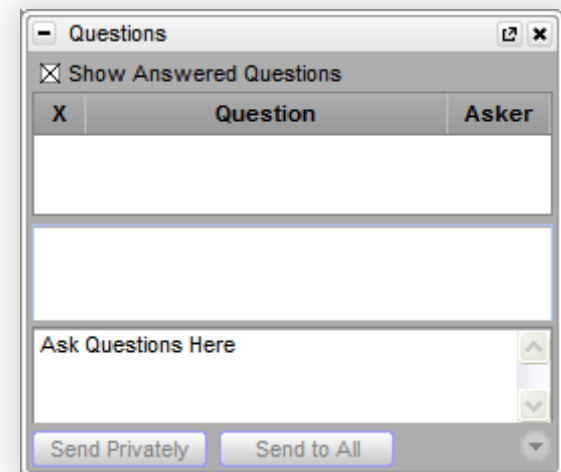
Greta Linse Peterson, Senior Statistician
June 5, 2013

GOLDEN HELIX
Accelerating the Quest for Significance™

# Questions During the Presentation

Use the Questions pane in your GoToWebinar window

# Agenda

**1** Background of GWAS Approaches

**2** Review of Select Mixed Model Methods

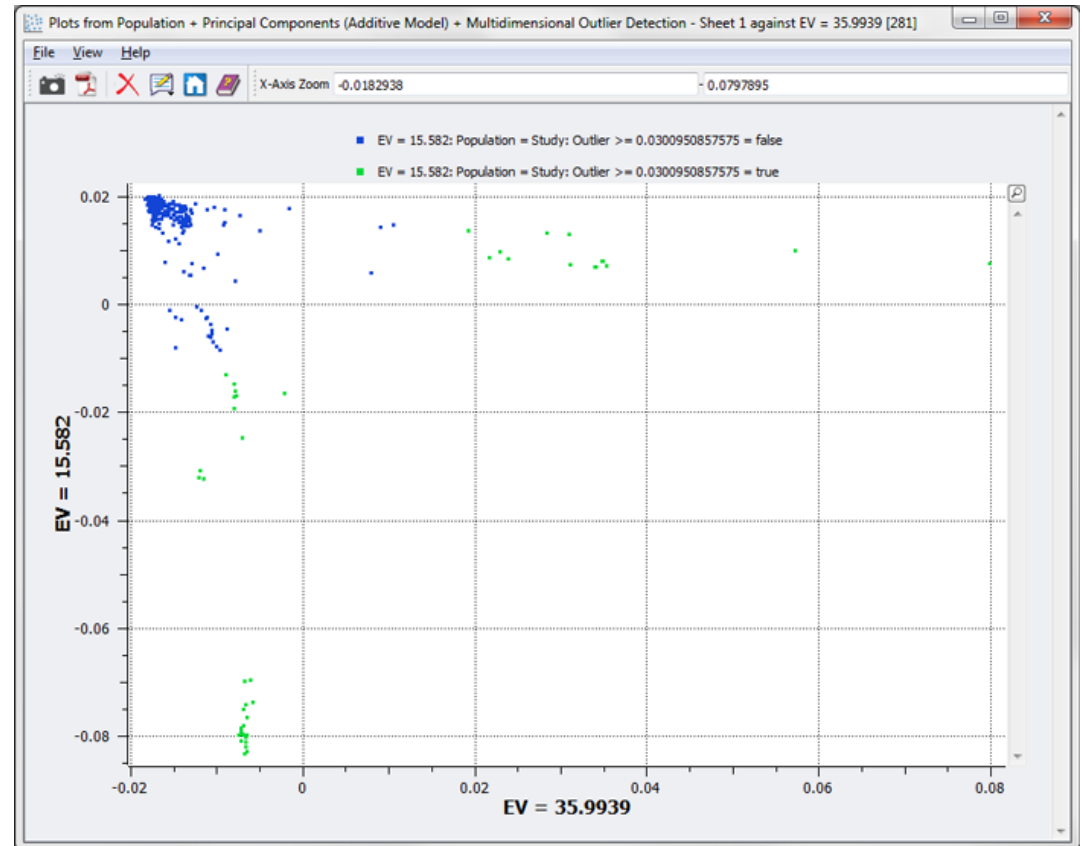**3** Mixed Models in SVS

**4** Demo

**5** Compare Results

[Poll: What category of species are you studying?]
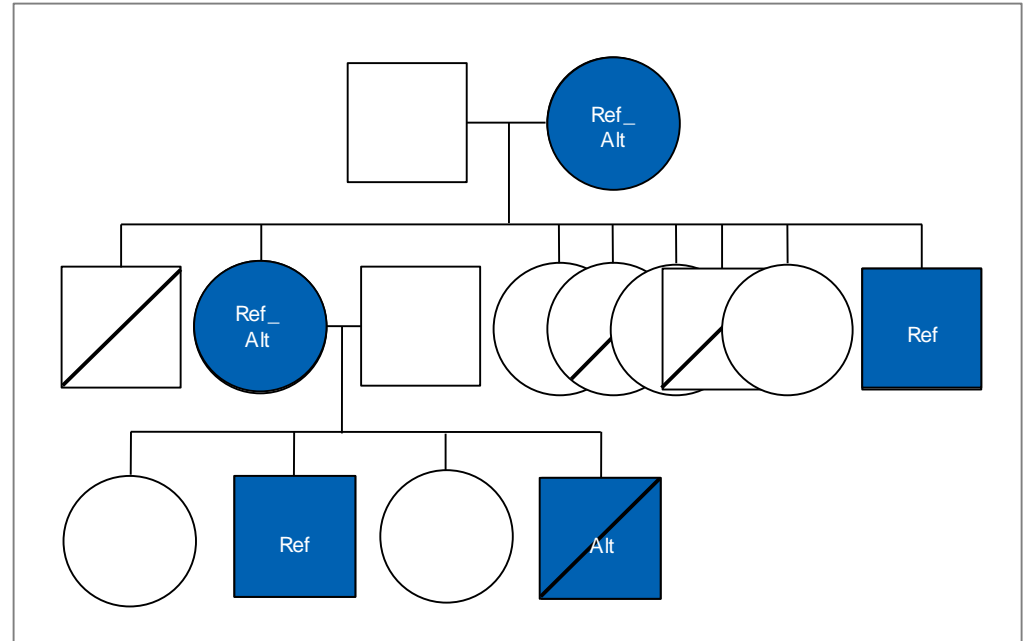
# A brief background of GWAS

- First the naïve approaches: Correlation/Trend Test, Linear/Logistic Regression

- Batch Effects, Population Structure and sharing of controls violated assumptions of the naïve approaches.

# Goal of better GWAS approaches

- Minimize false positives, obtain cleaner results, don't over correct the data to miss out on interesting results

- Handle population, family-based or mixed study designs.

- **Mixed Linear Model:**
  - $Y = X\beta + Z\mu + \epsilon$ where $\mu \sim N(0, \sigma_G^2 K^*)$, $\epsilon \sim N(0, \sigma_e^2 I)$ and $Cov(\mu, \epsilon) = 0$

- **Fixed Factors:**
  - Sex, age, known loci

- **Random Effects:**
  - Family or Population Structure, batch effects

- **Kinship Matrix:**
  - Any N x N matrix that describes the pairwise relationships between N samples

- **Null Hypothesis (generally):** $\sigma_G^2 = 0$

# Approximate Timeline

| Naïve GWAS | GWAS + Correcting for Population Stratification | Mixed Model Approaches |
|---|---|---|
| Corr/Trend Test | Genomic Control | EMMA (Kang 2008) |
| Regression Analysis | Structured Association (STRUCTURE) | BLUP/GBLUP Approaches for GWAS (Zhang 2008) |
| | PCA Correction (Eigenstrat Price 2006) | EMMAX (Kang 2010) |
| | | MLMM (Segura 2012) |

# Methods for MLMs in GWAS

# Agenda

**1**  Background of GWAS Approaches

**2**  **Review of Select Mixed Model Methods**

**3**  Mixed Models in SVS

**4**  Demo

**5**  Compare Results

# Methods Overview

- **Regression with PCA Correction**
  - Accounts for the relationship between samples with Principal Components
  - Need to know how many components to correct for

- **EMMAX**
  - Adjusts for the relationship between samples using a kinship matrix
  - Approximates the variance components and uses the same variance for all probes
  - Tests a single loci at a time

- **MLMM**
  - Adjusts for the relationship between samples using a kinship matrix
  - Approximates the variance components and uses the same variance for all probes, but re-computes at every step
  - Stepwise EMMAX, assumes multiple loci are associated with the phenotype

- **GBLUP**
  - Adjusts for the relationship between samples using a kinship matrix
  - Computes allele substitution effects to determine best genomic predictors of the phenotpye

# Method Comparison

| | Population Structure as Fixed Effect | Multiple Loci | EMMA used | Uses Kinship as Random Effect | Output Random Effect Component | Compute Allele Substitution Effects | Compute P-Value |
|---|---|---|---|---|---|---|---|
| Regression with PCA | Yes | No | No | No | No | No | Yes |
| EMMAX | Yes | No | Yes | Yes | No | No | Yes |
| MLMM | Yes | Yes | Yes | Yes | No | No | Yes |
| GBLUP | No | No | Yes* | Yes | Yes | Yes | No |

\* Uses EMMA for REML estimates

# Regression with PCA method overview

- **First compute the principal components**
  - Assumes the first few components are associated with the largest batch effects including population structure, plate effects, etc.

- **Decide how many components to correct for**

- **Either run regression on PCA corrected data or on genotype data including top principal components as fixed factors**

# EMMAX method overview

- Published in 2010 by the authors of EMMA

- Assumes a complex disease and that all SNP loci have a small effect on the phenotypic trait by themselves

- Instead of re-computing the variance components for every SNP (under the Alternative Hypothesis) computes it once under the Null Hypothesis

- Null Hypothesis: $\sigma_G^2 = 0$ ;

## Variance component model to account for sample structure in genome-wide association studies

Hyun Min Kang[1,2,8], Jae Hoon Sul[3,8], Susan K Service[4], Noah A Zaitlen[5], Sit-yee Kong[4], Nelson B Freimer[4], Chiara Sabatti[6], and Eleazar Eskin[3,7]

[1]Center for Statistical Genetics, Department of Biostatistics, University of Michigan, Ann Arbor, Michigan, USA

[2]Center for Computational Medicine and Bioinformatics, The University of Michigan Medical School, Ann Arbor, Michigan, USA

[3]Computer Science Department, University of California, Los Angeles, California, USA

[4]Center for Neurobehavioral Genetics, University of California, Los Angeles, California, USA

[5]Department of Epidemiology and Biostatistics, Harvard School of Public Health, Boston, Massachusetts, USA

[6]Department of Health Research and Policy, Stanford University School of Medicine, Stanford, California, USA

[7]Department of Human Genetics, University of California, Los Angeles, California, USA

## Abstract

Although genome-wide association studies (GWASs) have identified numerous loci associated with complex traits, imprecise modeling of the genetic relatedness within study samples may cause substantial inflation of test statistics and possibly spurious associations. Variance component approaches, such as efficient mixed-model association (EMMA), can correct for a wide range of sample structures by explicitly accounting for pairwise relatedness between individuals, using high-density markers to model the phenotype distribution; but such approaches are computationally impractical. We report here a variance component approach implemented in publicly available software, EMMA eXpedited (EMMAX), that reduces the computational time for analyzing large GWAS data sets from years to hours. We apply this method to two human GWAS data sets, performing association analysis for ten quantitative traits from the Northern Finland Birth Cohort and seven common diseases from the Wellcome Trust Case Control Consortium. We find that EMMAX outperforms both principal component analysis and genomic control in correcting for sample structure.

GWASs may utilize either case-control cohorts to test for associations with diseases or population cohorts to identify associations with quantitative traits. In both cases, it is

# MLMM method overview

- "Multiple-Loci Mixed Models"; stepwise EMMAX

- Assumes complex diseases where multiple loci are associated with the phenotype

- Cofactors are selected in a stepwise fashion by choosing the probe with the smallest p-value

- Since EMMAX is used, genetic and error are computed once for each step.

- Genetic and error variances are then re-estimated at for every step

# GBLUP method overview

- Best Linear Unbiased Predictor (BLUP) provides residual errors
  - Residual Breeding Values for plant/animal studies

- Estimates of allele substitution effects

- Pseudo-heritability estimate can be used to compare the method with other methods

- Uses a genomic relationship matrix which computes faster than IBS

- Have a dataset with inbreeding or some population structure

- Dataset is filtered down to samples and SNPs with:
  - "Good" Call Rate
  - SNP MAF > 0.05 (common variants)

- Whole Genome Sequencing data is fine if looking for common variants

- NOT for RARE VARIANT ANALYSIS!!!!

# Which Model to Use?

| | |
|---|---|
| **Regression with PCA** | • Homogeneous datasets or datasets with simple structure |
| **EMMAX** | •Complex diseases on a structured population, assuming all loci have a small effect on the trait |
| **MLMM** | • Complex diseases on a structured population, assuming there are several loci that have a large effect on the trait and the rest have small effects on the trait |
| **GBLUP** | • Obtain estimated breeding values, rank allele substitution effects to find QTL or find genomic relationship matrix in structured populations |

# Agenda

**1** Background of GWAS Approaches

**2** Review of Select Mixed Model Methods

**3** **Mixed Models in SVS**

**4** Demo

**5** Compare Results

# Mixed Models in SVS

- Derived from the mixmogam python package

- By B. Vilhjalmsson, coauthor of MLMM paper*

- Note, GBLUP also uses utilities from mixmogam



* V. Segura et al. "An efficient multi-locus mixed model approach for genome-wide association studies in structured populations" (Nat Genetics, 2012)

# SVS Implementation

- Provides user friendly interface for:
  - GBLUP
  - Mixed Linear Models Analysis

- Runs directly from a spreadsheet and has an options dialog where you can select your fixed factors and other parameters

- Visualization of results in SVS' Genome Browser is quick and easy

- Optimized so that analyses run fast

# Agenda

**1**    Background of GWAS Approaches

**2**    Review of Select Mixed Model Methods

**3**    Mixed Models in SVS

**4**    **Demo**

**5**    Compare Results

# Why Sheep? What about Humans or…?

- The Sheep HapMap dataset was chosen because of
  - the large number of samples and
  - the large number of breeds

- The dataset was available for public use on request from the ISGC

- The dataset was sufficiently structured enough to demonstrate all of the four methods

- Other datasets used by Mixed Model papers include:
  - WTCCC (all diseases including RA and T1D)
  - NFBC66
  - Arabidopsis thaliana dataset
  - Zea maize dataset
  - Various cattle datasets

- Mixed models used on datasets not expected to have family structure (WTCCC and NFBC66)

# First a little about the dataset…

- Sheep HapMap SNP50_Breedv1 dataset

- Illumina 50k SNP array

- 49,034 markers were left after filtering by the consortium

- 110 unmapped markers

- Only 1 marker in Chr Y

# Sample Statistics/Filtering

- **Removed samples from Boreray & Soay breeds**
  - 72 Breeds & Cross-Breeds left

- **Imputed gender from heterozygosity rates in the X chromosome**
  - Males: 1611
  - Females: 1081

# IBS and PCA on Marker Subset

- Filtered down to $MAF \geq 0.05$

- LD pruned
  - $R^2 \geq 0.5$ (CHM method)
  - Window of 50 markers
  - Step size of 5 markers

- Left 45,117 total markers (44,057 autosomal markers)

- Performed IBS & PCA analysis on remaining samples and markers

# Sheep HapMap PCA Plot

# Simulated Phenotype

- Filtered markers down to those within predicted gene regions

- Randomly selected 25 causal markers

- Set $h^2 = 0.4$

- Used a $\chi^2$ distribution for the effect sizes

- Added an error adjustment from a skewed normal distribution

# Analysis steps

| Marker Filtering | Sample Filtering | Compute Kinship Matrix | Compute Principal Components | Perform Mixed Model Analysis | Visualize Results |
|---|---|---|---|---|---|
| - Call Rate<br>- MAF<br>- LD Prune for Kinship, PCA | - Call Rate | - IBS, or<br>- IBD, or<br>- Gen Rel Matrix | - On filtered, LD pruned Dataset<br>- X Chr filtered out | - EMMAX<br>- MLMM<br>- GBLUP | - P-value plot<br>- Venn Diagram |

GOLDEN HELIX
*Accelerating the Quest for Significance™*

[Demo]

Golden Helix
Accelerating the Quest for Significance™

# Agenda

| 1 | Background of GWAS Approaches |
|---|---|

| 2 | Review of Select Mixed Model Methods |
|---|---|

| 3 | Mixed Models in SVS |
|---|---|

| 4 | Demo |
|---|---|

| 5 | **Compare Results** |
|---|---|

# Compare the methods

- **Top 1000 markers show some overlap in results**

# QQ Plots of methods

# Conclusion

- Mixed models can be a useful tool when applied appropriately.

- Use the method best suited for your data.

- Mixed models are not the "cure all" for bad data.

- Watch for a blog post to come out later this week on more mixed model methods!

# Acknowledgements

- Bjarni Vihjálmsson

- Christopher Seabury

- John McEwan of ISGC

# References

- Kang HM, et al (2008). 'Efficient control of population structure in model organism association mapping', Genetics, 178, 1709–1723.
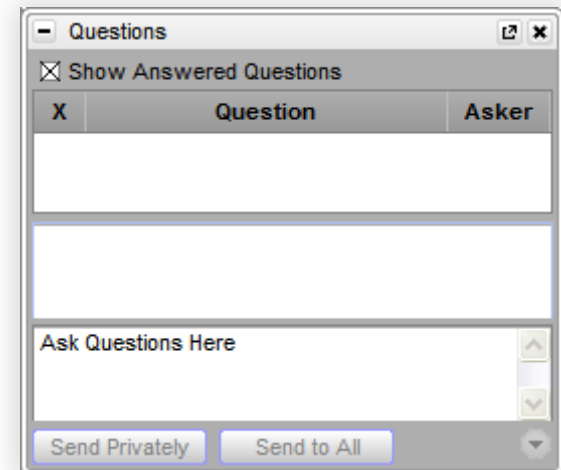
- Kang HM, et al (2010). 'Variance component model to account for sample structure in genome-wide association studies', Nature Genetics 42, 348–354.

- Patterson N, Price AL, Reich D (2006) Population Structure and Eigenanalysis PLoS Genet 2(12): e190. doi:10.1371/journal.pgen.0020190.

- Segura V, Vihjálmsson BJ, Platt A, Korte A, Seren Ü, et al. (2012) 'An efficient multi-locus mixed-model approach for genome-wide association studies in structured populations', Nature Genetics, 44, 825–830.

- Taylor, J.F. (2013) 'Implementation and accuracy of genomic selection', Aquaculture, http://dx.doi.org/10.1016/j.aquaculture.2013.02.017

- VanRaden, P.M. (2008) 'Efficient Methods to Compute Genomic Predictions', J. Dairy Sci, 91, pp. 4414–4423.

- Kijas JW, Lenstra JA, Hayes B, Boitard S, Porto Neto LR, San Cristobal M, Servin B, McCulloch R, Whan V, Gietzen K, Paiva S, Barendse W, Ciani E, Raadsma H, McEwan J, Dalrymple B; International Sheep Genomics Consortium Members. "Genome-wide analysis of the world's sheep breeds reveals high levels of historic mixture and strong recent selection." PLoS Biol. 2012 Feb;10(2):e1001258. doi: 10.1371/journal.pbio.1001258. Epub 2012 Feb 7

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# IBS vs Genomic Relationship Matrix