

A Hands-on Introduction to GWAS

Danièle Filiault

Popgen Intro Course

30 Sept 2020

The goals for today









- Give a practical, non-technical introduction of how to do a “standard” linear mixed model GWAS
- Address some FAQs
- Identify possible sticking points
- Understand the process enough to begin exploring GWAS in your organism of choice

Build a foundation for you to start learning more about GWAS on your own!

What is a GWAS?

Use genome-wide SNPs/variants
Associate phenotype and genotype

	...ATCA A TGGGGC G TAGAC CAT AC...
	...ATCA A TGGGGC G TAGAC CAT AC...
	...ATCA A TGGGGC C TAGAC CAT AC...
	...ATCA C TGGGGC C TAGAC AC...
	...ATCA A TGGGGC G TAGAC CAT AC...
	...ATCA C TGGGGC C TAGAC AC...
	...ATCA C TGGGGC G TAGAC AC...
	...ATCA A TGGGGC C TAGAC CAT AC...
	...ATCA A TGGGGC G TAGAC CAT AC...
	...ATCA C TGGGGC G TAGAC AC...

LD is our frenemy!



Linkage disequilibrium (LD) = tendency for SNP variants to be inherited together

Why? proximity, population structure, selection

Friend

We don't need to test every SNP (one SNP can tag many variants).

Enemy

LD can confound our results!

Complicated confounding patterns (especially if many alleles or loci underlie the trait)

****Population structure causes confounding genome-wide****

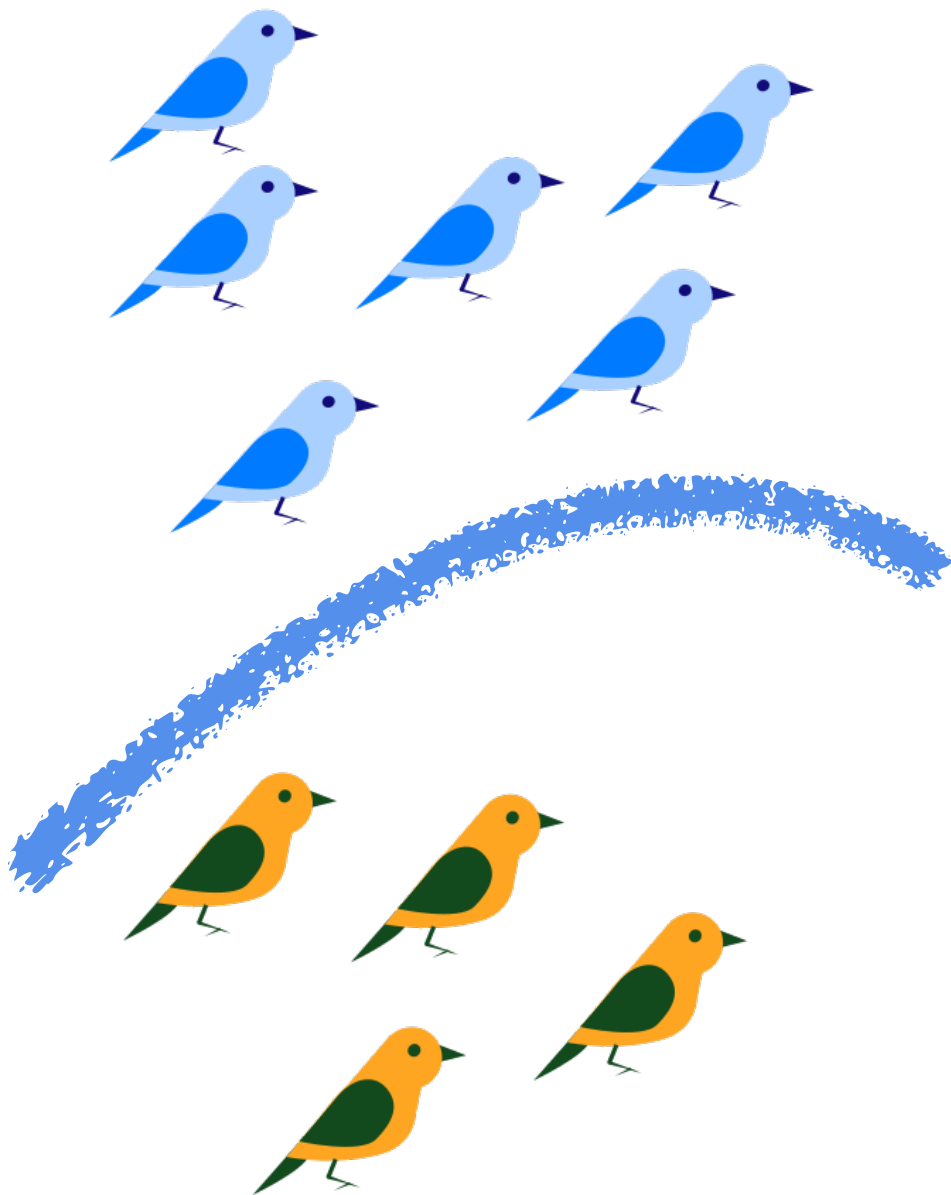
If we don't try to correct for population structure, we will end up with too many false positives.

How to correct for population structure?

Individuals within a population are more related than those between populations.

Here, birds in a population share not only ***causative*** variants, but also ***non-causative*** variants that are more common in the population (genetic background).

We take this background into account in GWAS to try to reduce the significance of non-causative variants. The **K-matrix** represents this background relatedness.



What do we need to run a GWAS?

What the GWAS does:

Fit a linear mixed model for each SNP

$$Y = X\beta + u + \varepsilon$$

Phenotype = Genotype + Background relatedness (K) + error

(Genotype is a fixed effect, K is a random effect)

What we get as results:

P-value for each SNP (from comparing the full model to a model without the genotype term)

$$Y = u + \varepsilon$$

The effect of the SNP on the phenotype (β)

Phenotypes and Accessions



GWAS is testing an underlying hypothesis!

Genetic variation underlies the phenotypic variation we observe among these accessions.

Phenotype and accession choice shape this hypothesis (and interpretation of results).

Phenotypes need to be heritable traits.

Simple vs. complex phenotypes

Quantitative (can do binary or categorical with different models)

Remember, we are doing linear modeling, so there is the assumption that residuals will be normal.

Genotypes



Common genotyping methods

Whole-genome resequencing

RADseq/Genotyping by Sequencing

mRNA-Seq

Exome sequencing

SNP arrays

Understanding LD is key in selecting/interpreting results given a particular dataset.

How much of the genome are you actually assaying?

Coding versus non-coding variation?

SNPs aren't the only type of polymorphism.

K matrix



Genetic relationship between pairs of accessions

Pairwise relatedness matrix / IBS matrix

Used for population structure correction

Many programs for estimation: limix, plink, R packages, etc.

Today we will use a pre-calculated matrix - calculation is slow.

Sometimes population structure correction will take other forms:

- PCA (especially in humans)
- output from the Structure program

Options for running GWAS



Online options

EasyGWAS (<https://easygwas.ethz.ch/>)

GWA-Portal (<https://gwas.gmi.oeaw.ac.at/>)

- easy interface
- limited data availability
- limited types of analysis
- not flexible

Do-it-yourself

limix (<https://github.com/limix/limix>)

gemma (<https://github.com/genetics-statistics/GEMMA>)

and MANY others...

- requires coding skills
- maximum flexibility and options
- allows for more complicated GWAS analyses

The plan for hands-on work today



1. We will walk through GWAS together using a flowering time phenotype (3 jupyter notebooks).
2. You will do some GWAS by yourself with two additional phenotypes. This independent exercise is outlined in:

[GWAS_course_independent_work_instructions.pdf](#)

The GWAS jupyter notebooks



1_phenotype_exploration.ipynb

Explore the phenotype we will use

2_GWAS.ipynb

Prepare input variables, run GWAS, and output results

3_GWAS_interpretation.ipynb

Visualize and understand GWAS results

Jupyter commands:

Shift-enter or Command-enter to run a cell

enter -> edit mode

esc -> command mode

After esc, M (markdown - text) or Y (script) modes

I did a GWAS, now what?



- **Association is not causality!**
- **Most significant SNP is not necessarily causative** (confounding, incomplete knowledge of variants).
- **The “peakiness” of a peak is not necessarily a sign of its importance**

Interested in candidate genes?

- Use biology, annotation, expression, etc. Be creative!
- Sanger sequence candidate genes/regions
- Test causality by QTL and/or transgenics with different alleles in common background

Interested in evolutionary or ecological inferences?

- Consider the consequences of your choice of accessions and phenotypes
- Associations with selective pressures (environment/climate)?
- Other signs of selection?
- The reference allele is not necessarily the ancestral allele.

Moving beyond simple GWAS



1. Generalized linear models for different data distributions
2. Adding covariants
3. Incorporating genotype-by-environment interactions
4. Analyzing correlated traits together (increases power)
5. Bayesian approaches
6. Polygenic scores
7. Genomic prediction

There are endless packages/programs available to do GWAS, with more published all the time. Limix and gemma are great places to start.

Human GWAS is a bit of a different animal (no experiments, huge sample sizes, environmental confounding).

Final Questions/Comments?

