

# Day 2

## Optional GWAS Practical

Luke Evans & Wei Zhou

Practical originally developed by Katrina Grasby and Lucía Colodro Conde

**Questions:** Please ask! You can also use the Q&A box under “Common & rare association”

**Learning Goals:**

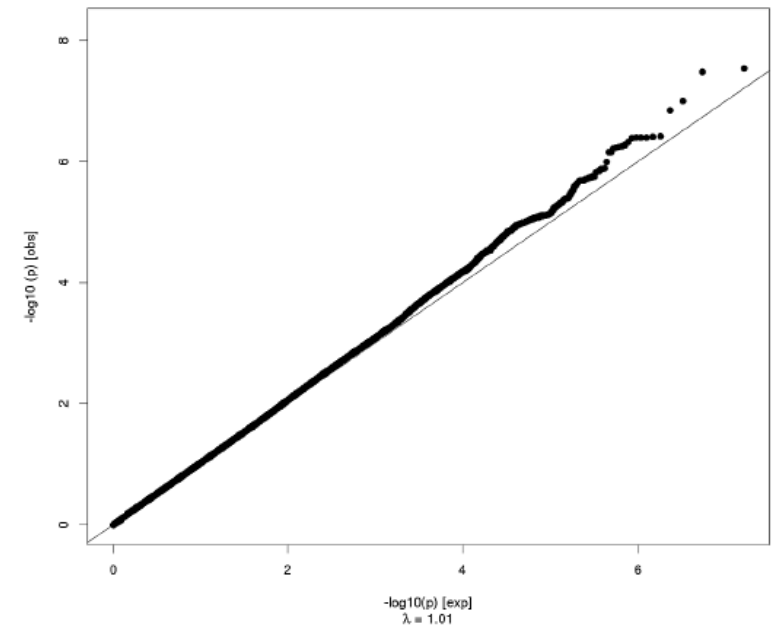
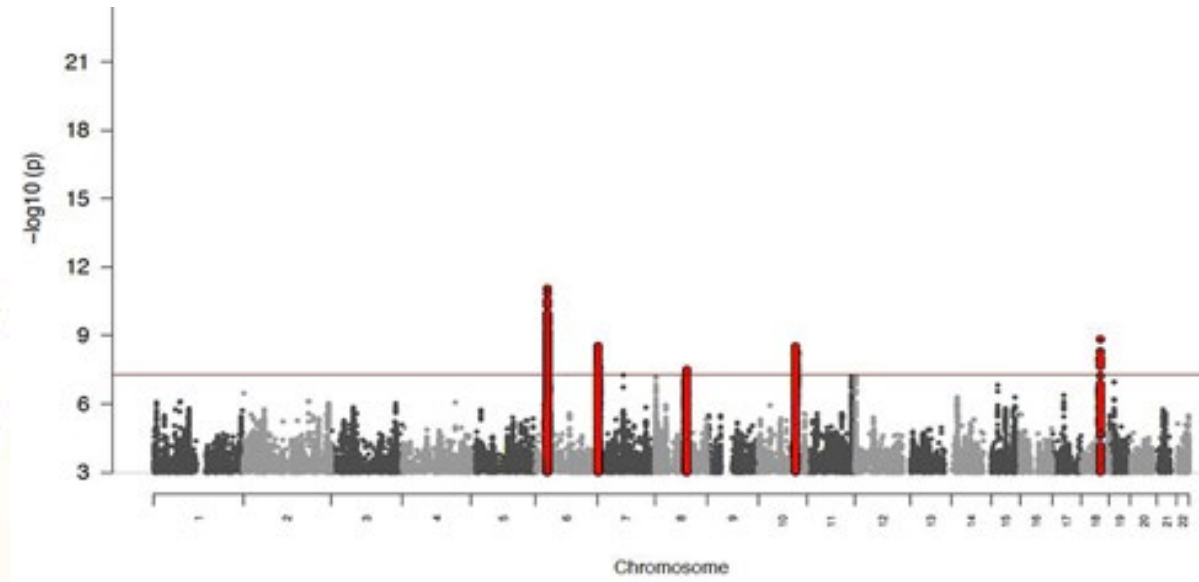
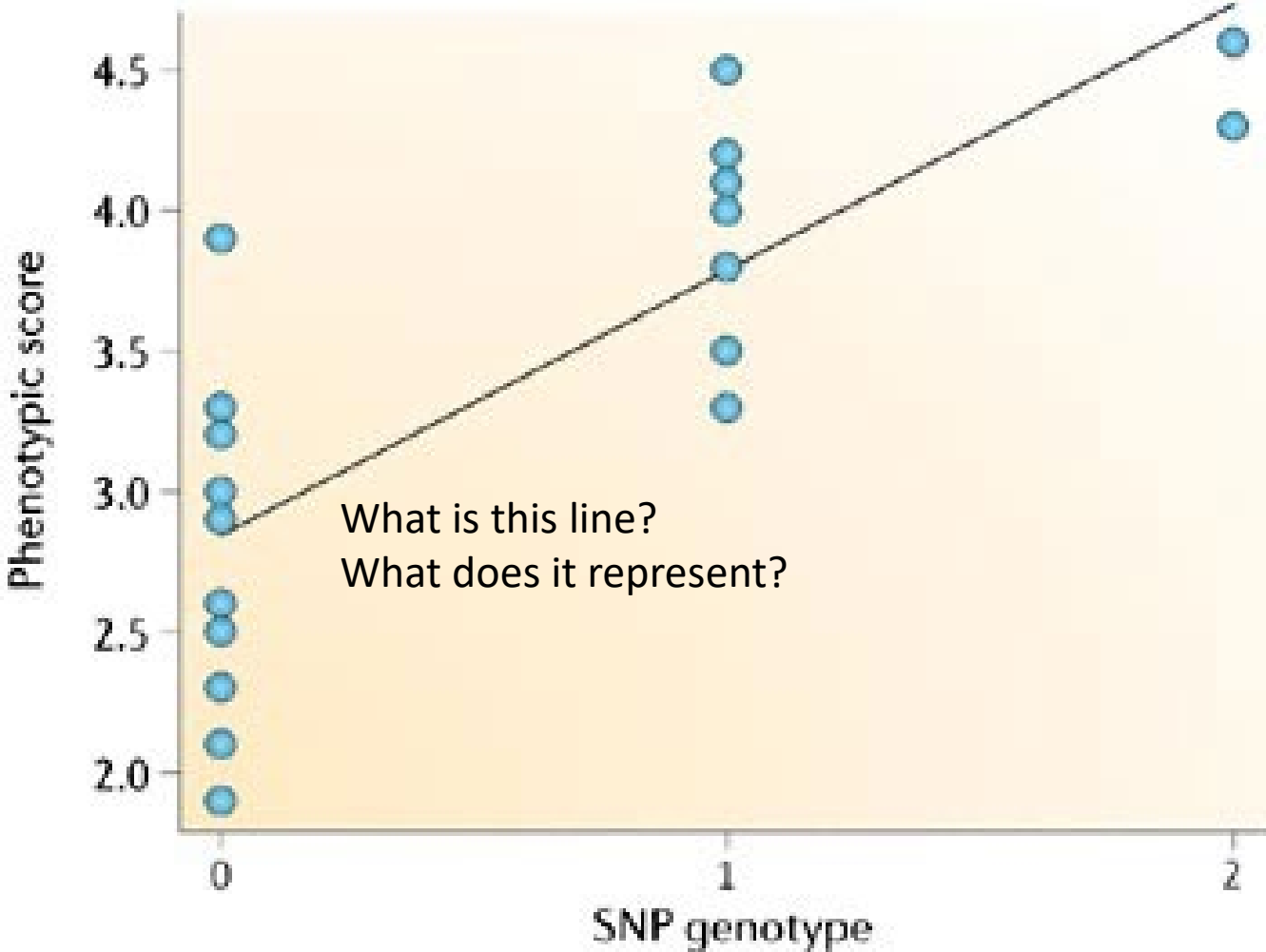
1. Run a GWAS using real data
2. Familiarize yourself with the standard types of output from a GWAS
3. Practice using the command line
4. Improve your familiarity with `plink2`, a common tool in genetics (beyond GWAS, too)
5. Interpret why peaks form in a Manhattan plot, and how to explore the genome around associated loci

There are 3 exercises in the practical.

If you don't finish, that's OK – these resources are there for you to practice & use as you're working through your own data.

There are questions throughout the instructions – take some time to either look for the answer (e.g., in the output) or consider what the answers/results are telling you (e.g., why does a q-q plot look a certain way?)

# Recap of GWAS from yesterday



# Analyzing X Chromosome

- Often overlooked, but important to analyze, too
- Can analyze it, you just have to do it a little differently than the autosomes
- Imputation can be done and servers understand the different chromosomes and regions on them
- Dosage differences between sexes, dosage compensation and X inactivation are all important features
- X inactivation varies for different tissues



*Briefings in Bioinformatics*, 2022, 23(5), 1–9

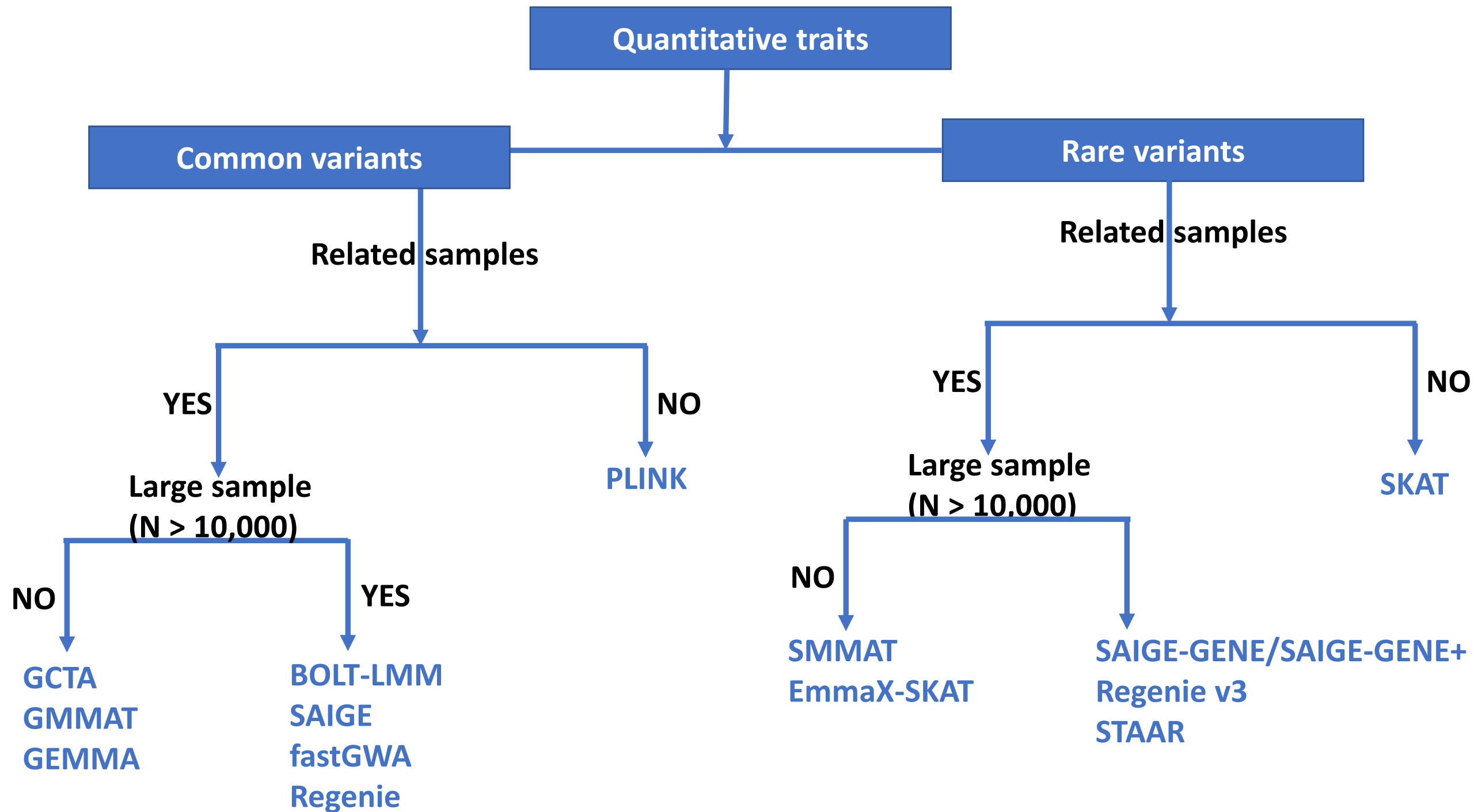
<https://doi.org/10.1093/bib/bbac287>

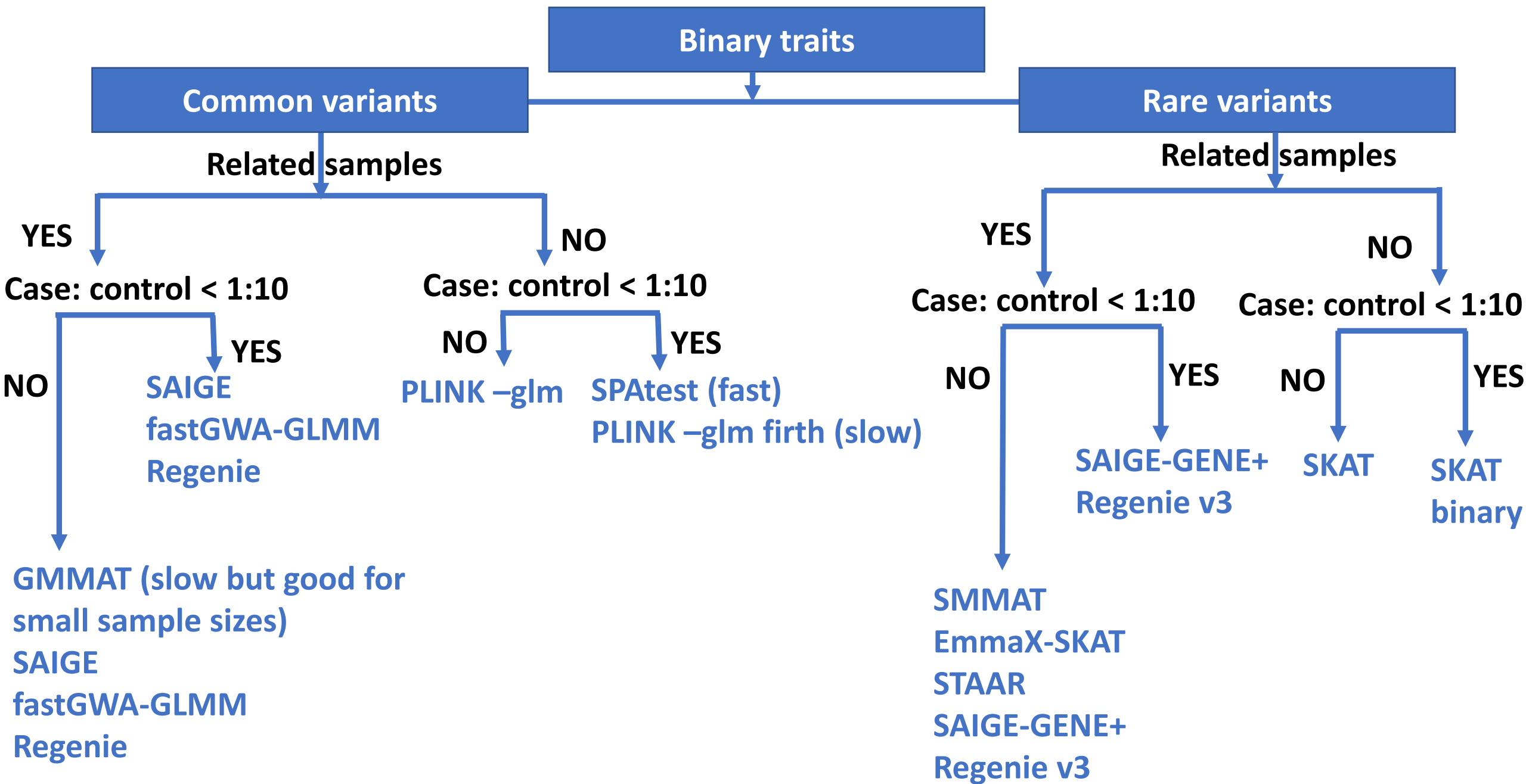
Review

## A systematic review of analytical methods used in genetic association analysis of the X-chromosome

Nick Keur, Isis Ricaño-Ponce, Vinod Kumar and Vasiliki Matzaraki

Corresponding author. Vasiliki Matzaraki, E-mail: [Vasiliki.Matzaraki@radboudumc.nl](mailto:Vasiliki.Matzaraki@radboudumc.nl)



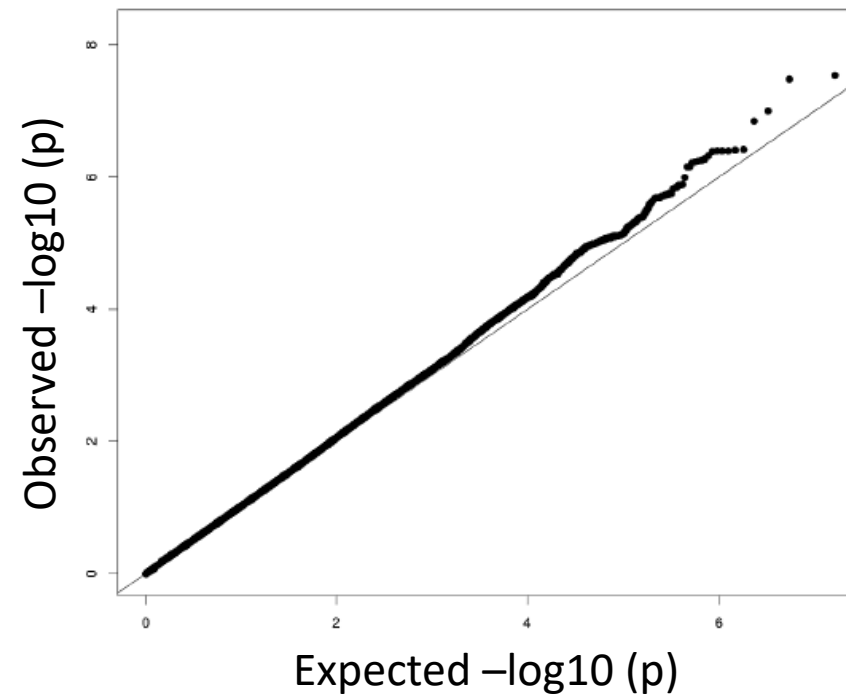


# There are many tools for GWAS

- The most appropriate method depends on the structure of your data: stratification? Relatedness? Sample size? Etc...
- Mixed linear models: SAIGE, BOLT-LMM, GCTA, FaST-LMM...
- We're going to use plink2 for today's practical: fast, simple command line program that is a general workhorse software for managing data and running analyses.
  - Original plink ped/map format
  - Binary plink format (plink1.9): bed/bim/fam
    - Fast, very useful
  - Plink2 format: pgen/psam/pvar
    - Very fast
    - Saves space compared to others
    - Can include dosage, phase, INFO (similar to VCF format)

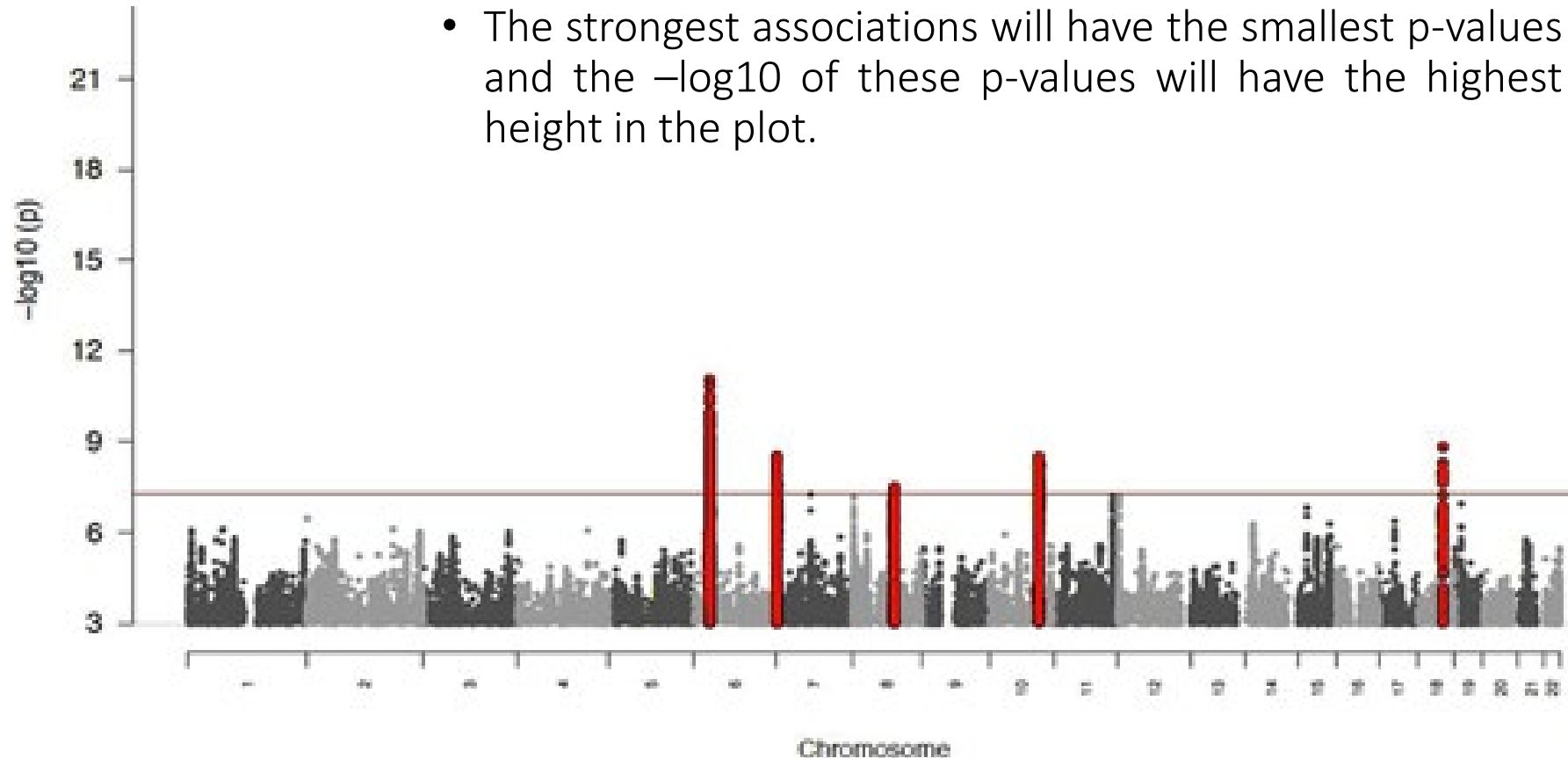
# QQ (quantile-quantile) plot

- Checks the overall distribution of test statistics or  $-\log_{10}$  p-values of our results with the expectation under the null hypothesis of no association (the diagonal line shows where the points should fall under the null).
- Evaluates systematic bias and inflation (undetected sample duplications, unknown familial relationships, gross population stratification, problems in QC...).



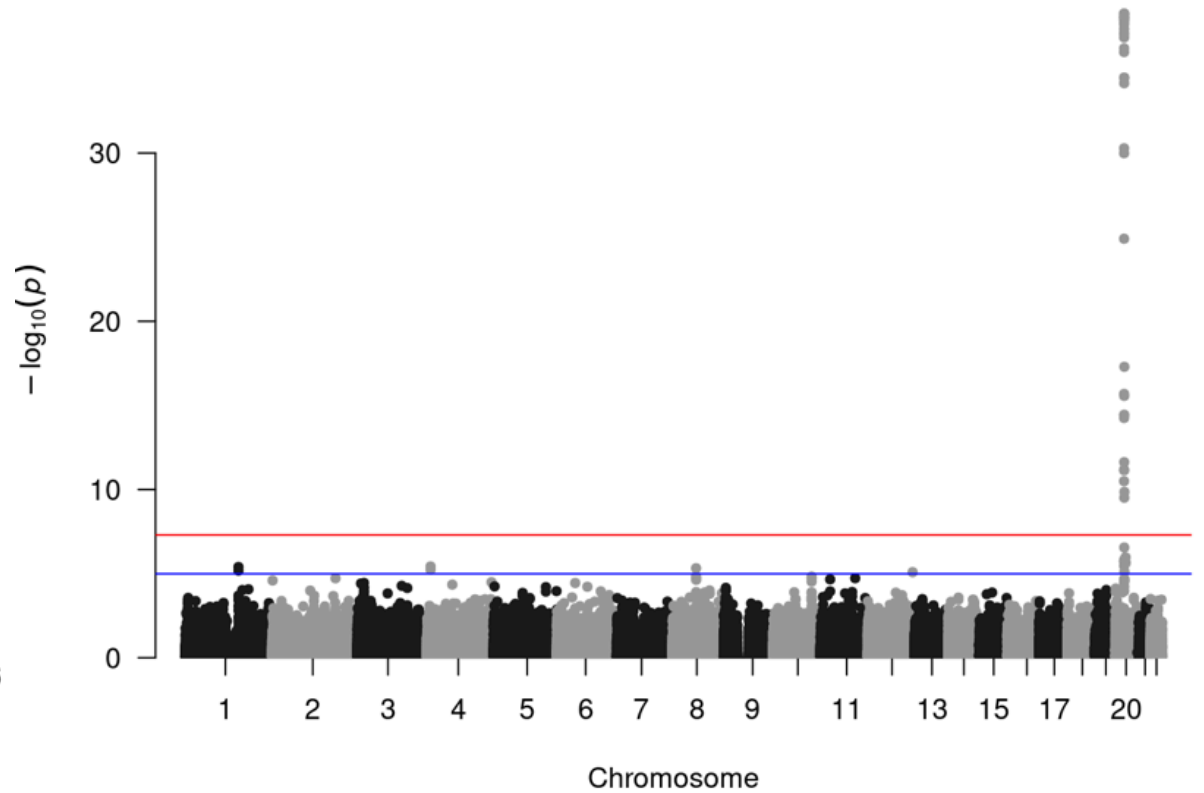
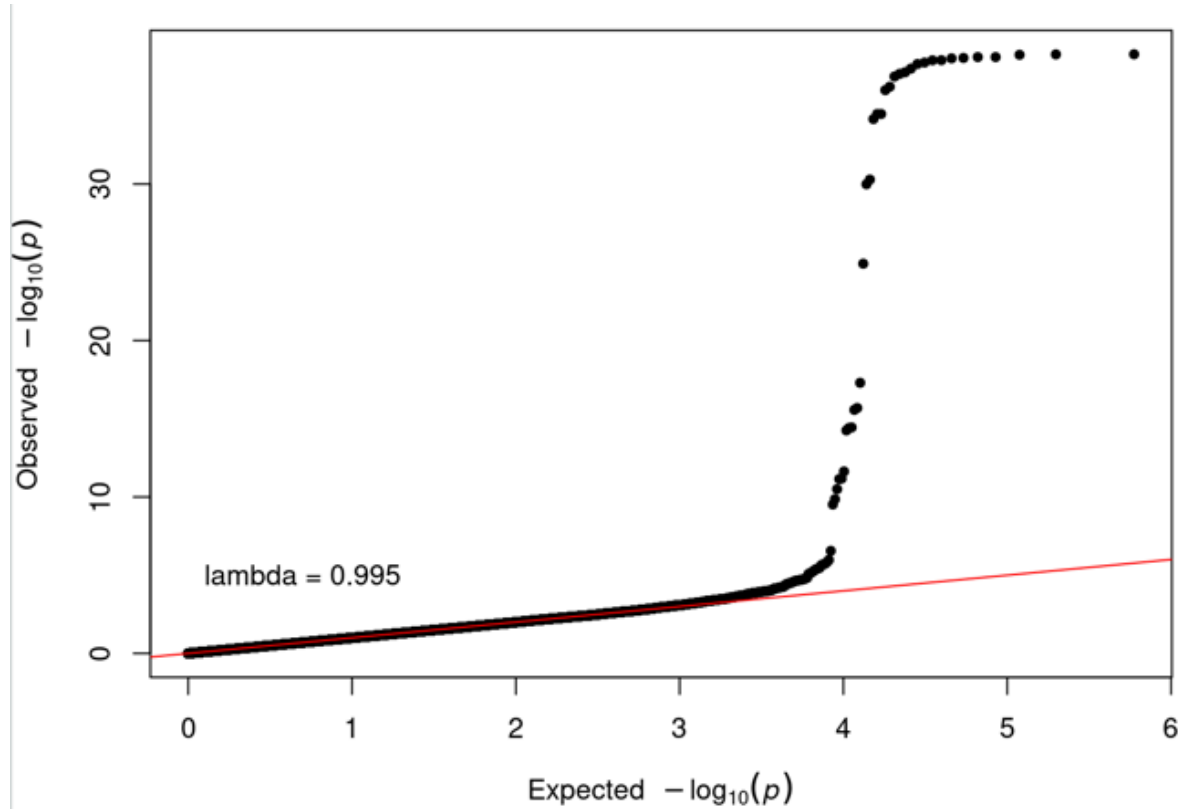
# Manhattan plot

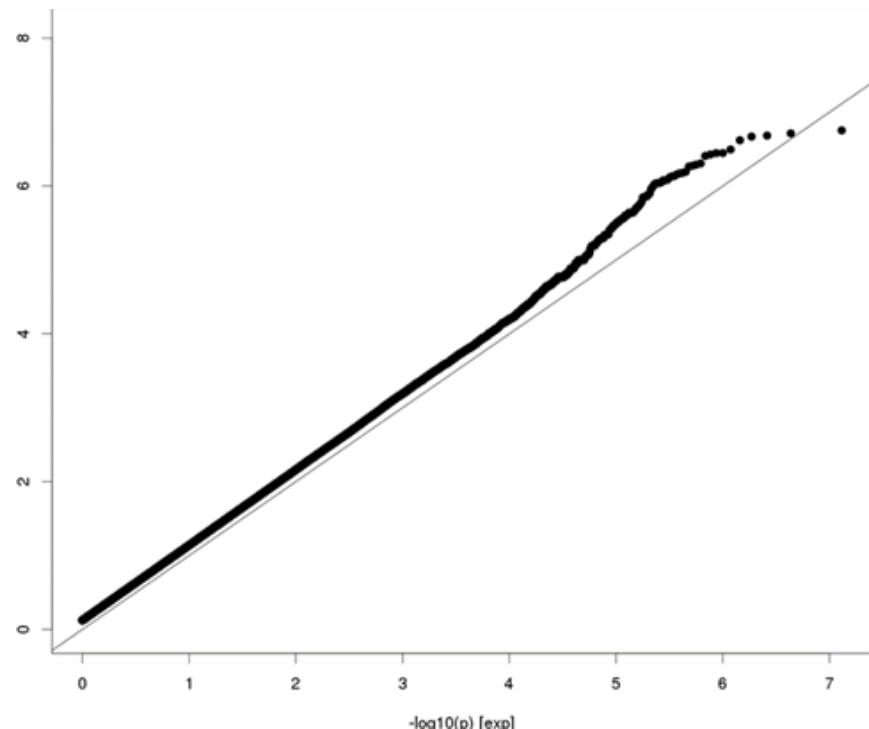
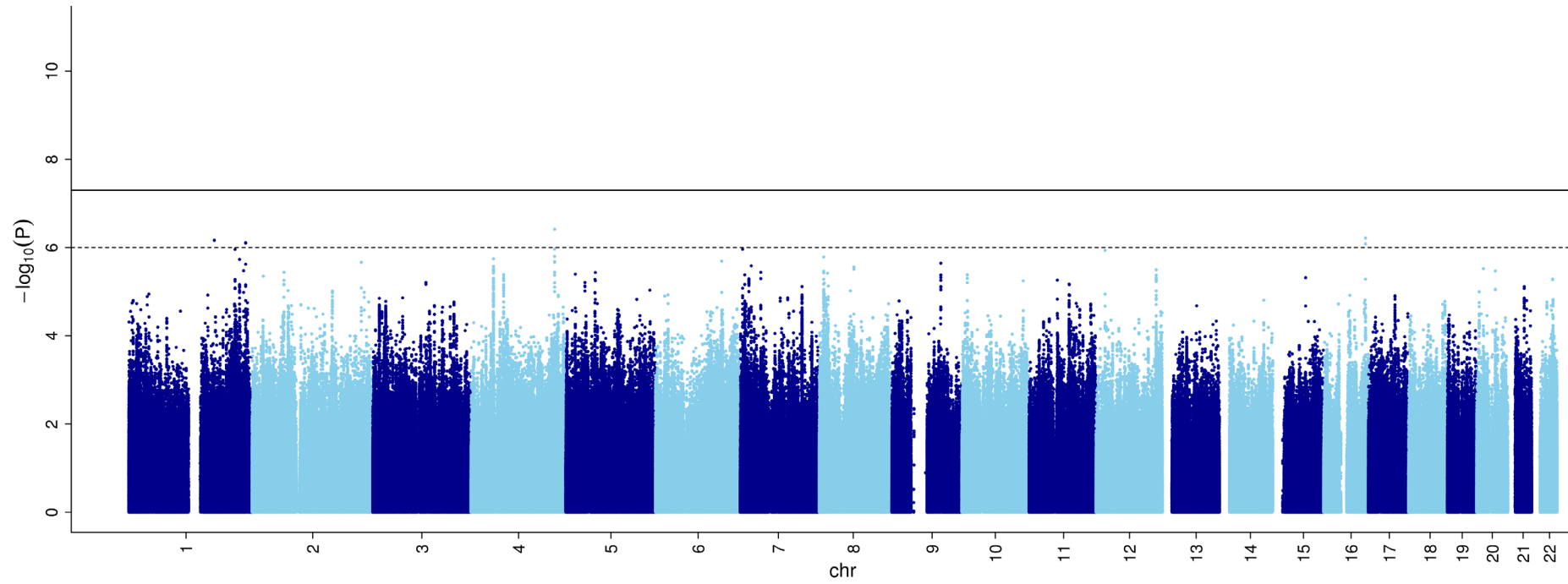
- Plots the  $-\log_{10}$  of the association p-value for each SNP against the genomic coordinates.
- The strongest associations will have the smallest p-values and the  $-\log_{10}$  of these p-values will have the highest height in the plot.



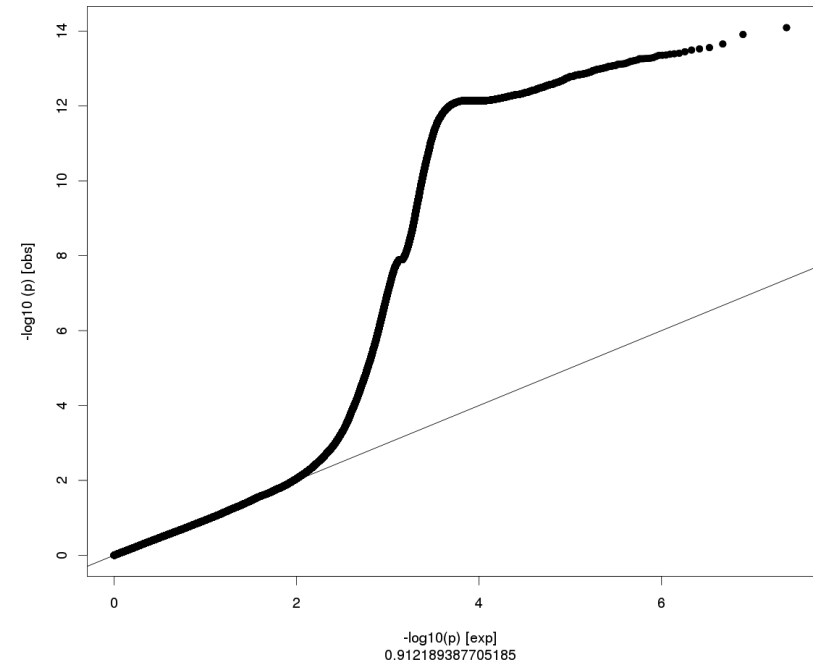
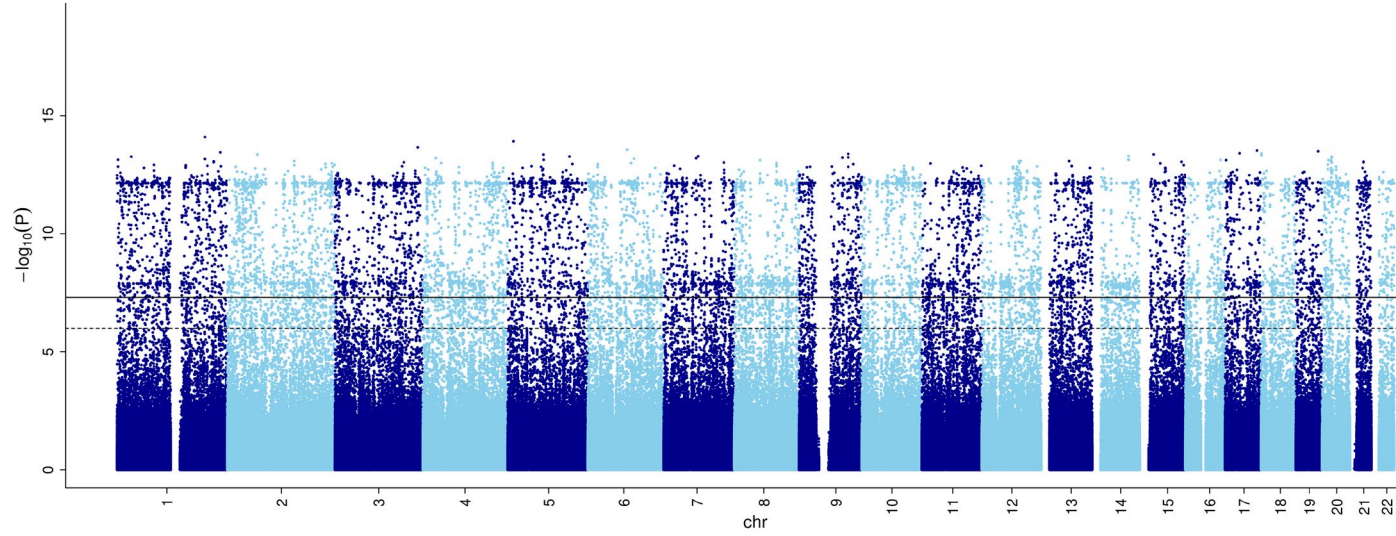


# What are these telling us?





# And these? What are these telling us?



# Regional plot



NATIONAL CANCER INSTITUTE  
Division of Cancer Epidemiology & Genetics

[Home](#) [LDassoc](#) [LDhap](#) [LDmatrix](#) [LDpair](#) [LDproxy](#) [SNPchip](#) [SNPclip](#) [API Access](#) [Help](#)

## Welcome to LDlink!

LDlink is a suite of web-based applications designed to easily and efficiently interrogate linkage disequilibrium in population groups. All population genotype data originates from Phase 3 (Version 5) of the 1000 Genomes Project and variant RS numbers are indexed based on dbSNP 151. Where coordinates are specified, GRCh37/hg19 is used. Only bi-allelic variants are permitted as input. LDlink includes the following modules:

**LDassoc:** Interactively visualize association p-value results and linkage disequilibrium patterns for a genomic region of interest. Input is a tab or space delimited association output file and a population group.

**LDhap:** Calculate population specific haplotype frequencies of all haplotypes observed for a list of query variants. Input is a list of variant RS numbers (one per line) and a population group.

**LDmatrix:** Create an interactive heatmap matrix of pairwise linkage disequilibrium statistics. Input is a list of variant RS numbers (one per line) and a population group.

**LDpair:** Investigate correlated alleles for a pair of variants in high LD. Input is two RS numbers and a population group.

**LDproxy:** Interactively explore proxy and putatively functional variants for a query variant. Input is an RS number and a population group.

**SNPchip:** Find commercial genotyping platforms for variants. Input is a list of variant RS numbers (one per line) and desired arrays.

**SNPclip:** Prune a list of variants by linkage disequilibrium. Input is a list of variant RS numbers (one per line) and a population group.

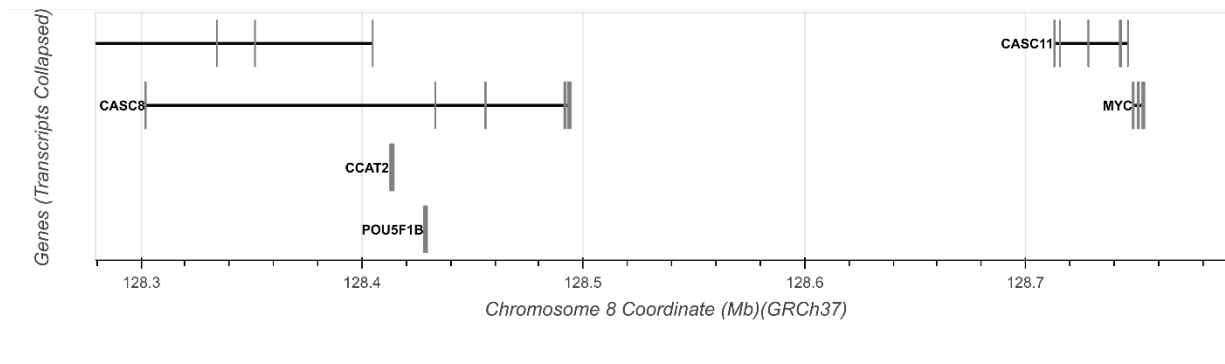
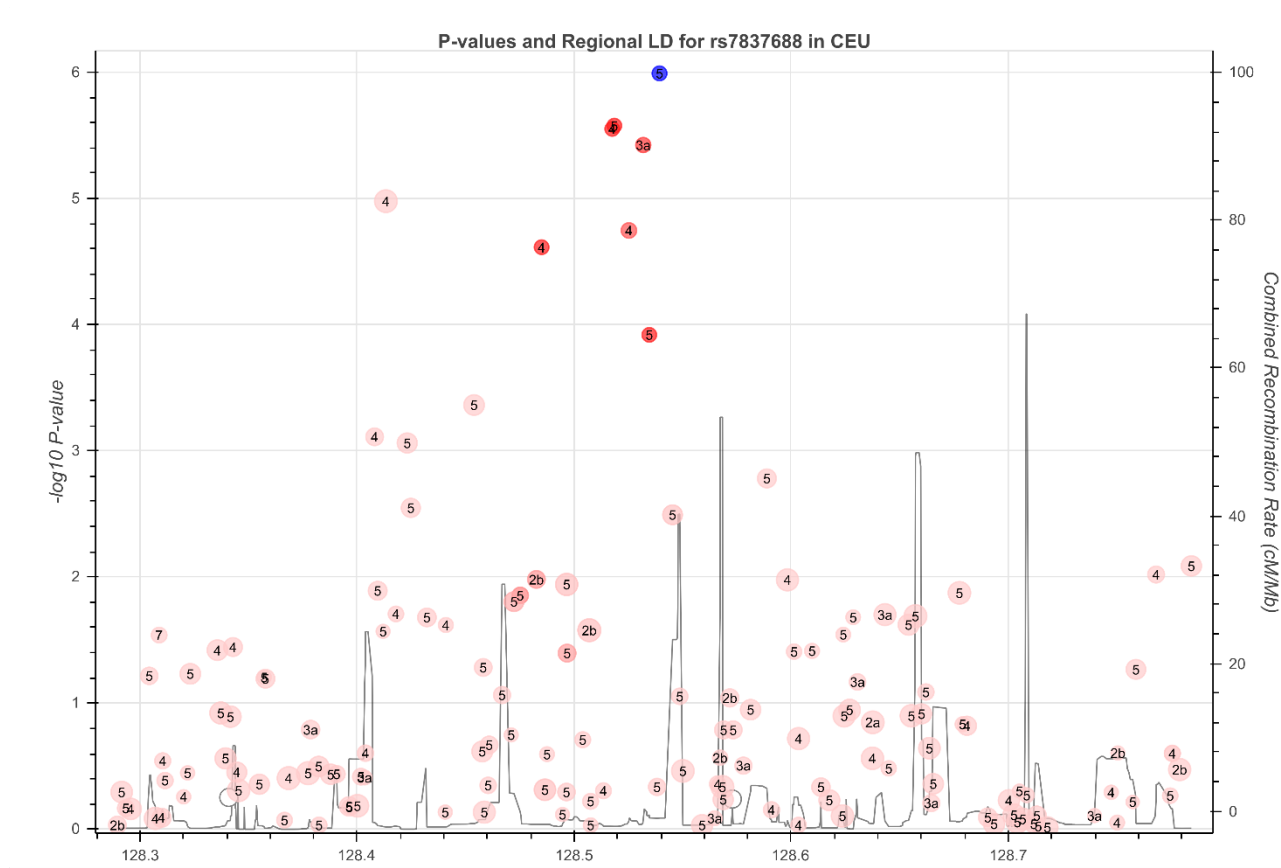
<https://ldlink.nci.nih.gov/>



## LDassoc Tool <sup>i</sup>

Interactively visualize association p-value results and linkage disequilibrium patterns for a genomic region of interest.

<input type="text" value="prostate_example.txt"/> <input type="button" value="Browse"/>	<input type="button" value="Variant ▾"/>	<input type="button" value="CEU ▾"/>	<input type="button" value="Calculate"/>
<input checked="" type="checkbox"/> Use example GWAS data	<input type="text" value="rs7837688"/>	LD Measure: <input type="button" value="R&lt;sup&gt;2&lt;/sup&gt;"/> <input type="button" value="D'"/>	
<input type="button" value="Chromosome: chr column"/>	± <input type="text" value="500000"/> base pair window	Collapse transcripts: <input type="button" value="Yes"/> <input type="button" value="No"/>	
<input type="button" value="Position: pos column"/>		RegulomeDB annotation: <input type="button" value="Yes"/> <input type="button" value="No"/>	
<input type="button" value="P-Value: p column"/>			



Today's practical data are freely available from:

**ARTICLE**

---

## Genetic Control of Human Brain Transcript Expression in Alzheimer Disease

Jennifer A. Webster,<sup>1,2,3,16</sup> J. Raphael Gibbs,<sup>4,5,16</sup> Jennifer Clarke,<sup>6</sup> Monika Ray,<sup>7</sup> Weixiong Zhang,<sup>7,8</sup>  
Peter Holmans,<sup>9</sup> Kristen Rohrer,<sup>4</sup> Alice Zhao,<sup>4</sup> Lauren Marlowe,<sup>4</sup> Mona Kaleem,<sup>4</sup>  
Donald S. McCorquodale III,<sup>10</sup> Cindy Cuello,<sup>10</sup> Doris Leung,<sup>4</sup> Leslie Bryden,<sup>4</sup> Priti Nath,<sup>4</sup>  
Victoria L. Zismann,<sup>1,2</sup> Keta Joshipura,<sup>1,2</sup> Matthew J. Huentelman,<sup>1,2</sup> Diane Hu-Lince,<sup>1,2</sup>  
Keith D. Coon,<sup>1,2,11</sup> David W. Craig,<sup>1,2</sup> John V. Pearson,<sup>1,2</sup> NACC-Neuropathology Group,<sup>12</sup>  
Christopher B. Heward,<sup>13,17</sup> Eric M. Reiman,<sup>1,2,14</sup> Dietrich Stephan,<sup>1,2,14</sup> John Hardy,<sup>4,5</sup>  
and Amanda J. Myers<sup>10,15,\*</sup>

Source:

<https://med.miami.edu/faculty/amanda-myers-phd>

<https://xzmxbgsv808roffneicreq.on.driv.tw/www.lfun/LFUN/LFUN/DATA.html>

# Set up – we'll go through the first few commands together

We will use the text editor, the terminal, R Studio, and the browser.

Copy all the files by typing in your working directory:

```
cp -r /faculty/luke/2023/gwas2/ ./
```

The practical is “Day2\_GWAS\_Practical\_INSTRUCTIONS.pdf”, and is available on the file sharing site.

A plain text version is also on the file sharing site

The rubric is “Day2\_GWAS\_Practical\_RUBRIC.pdf”, and is available on the file sharing site.

Some commands in the instructions are missing (indicated with XXX) – it's up to you & your group to determine what to include there.

Questions are in the instructions – some are about the specific results (what is N?) while others are asking you to think about what the results are telling you (QQ plot look OK?)