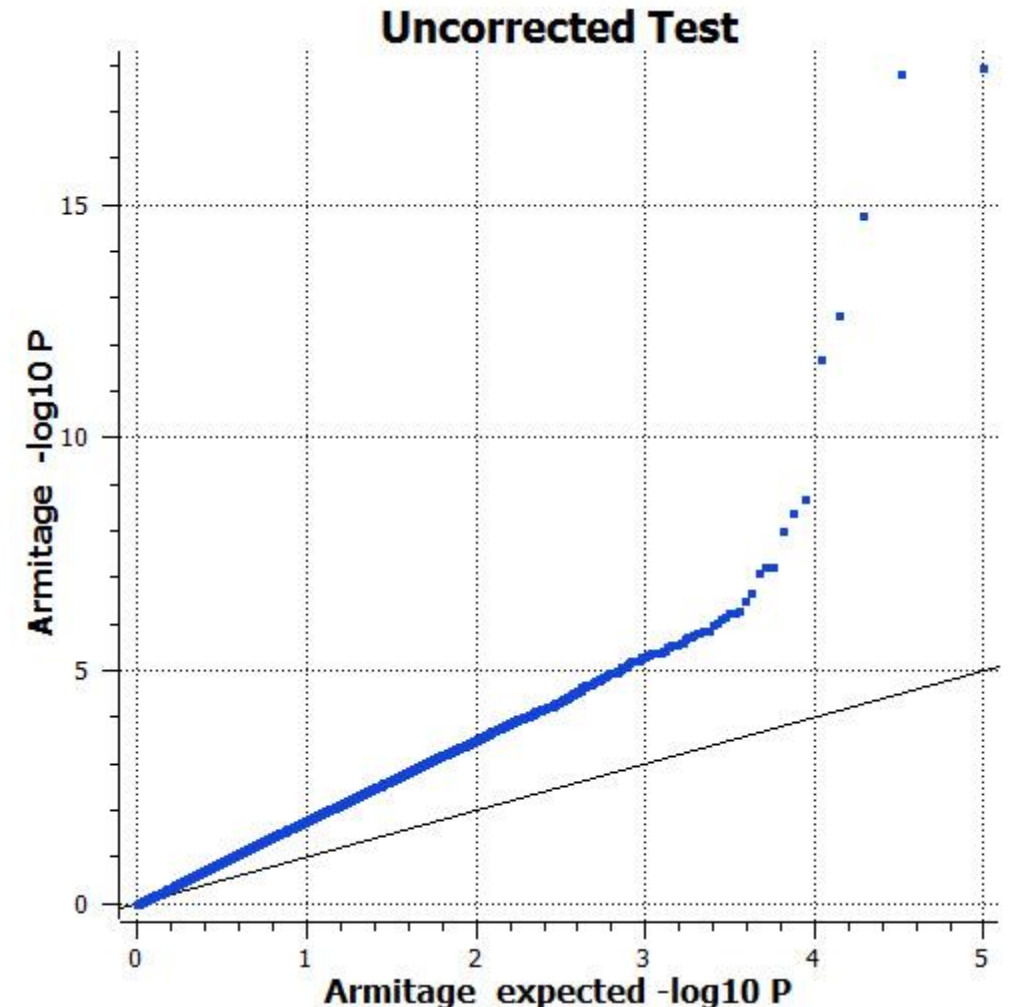


Population stratification

Why it is bad when we are doing association studies and what measures we can take to make sure it does not affect our results in ways we don't want it to affect them.

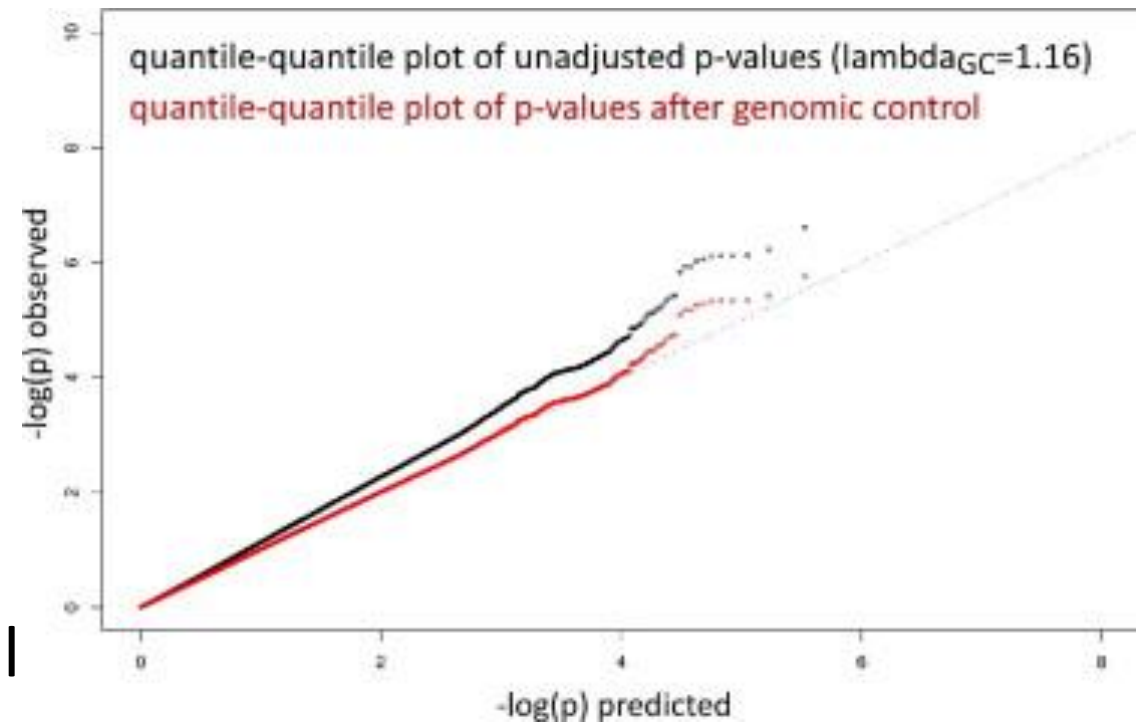
What is population stratification and why is it bad?

- Differences between cases and controls that happen not-quite-by chance but are picked up by association studies.
- Will affect overall distribution of test statistics.
- Association models often assume HWE or something similar:
 - Constant allele frequencies
 - Random mating
 - Infinite population size
 - No mutations
 - No migration/flow
- Of these, random mating is most often violated and introduces population structure/stratification.



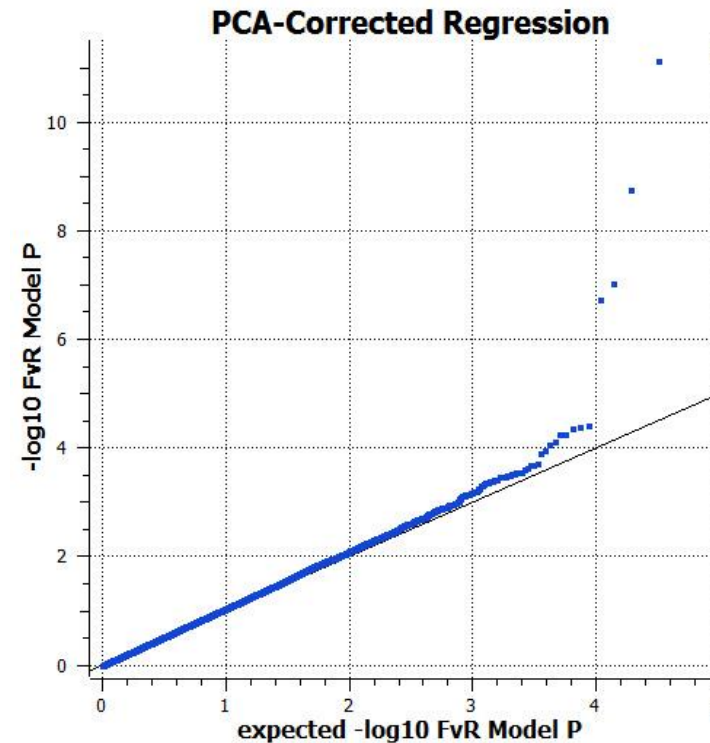
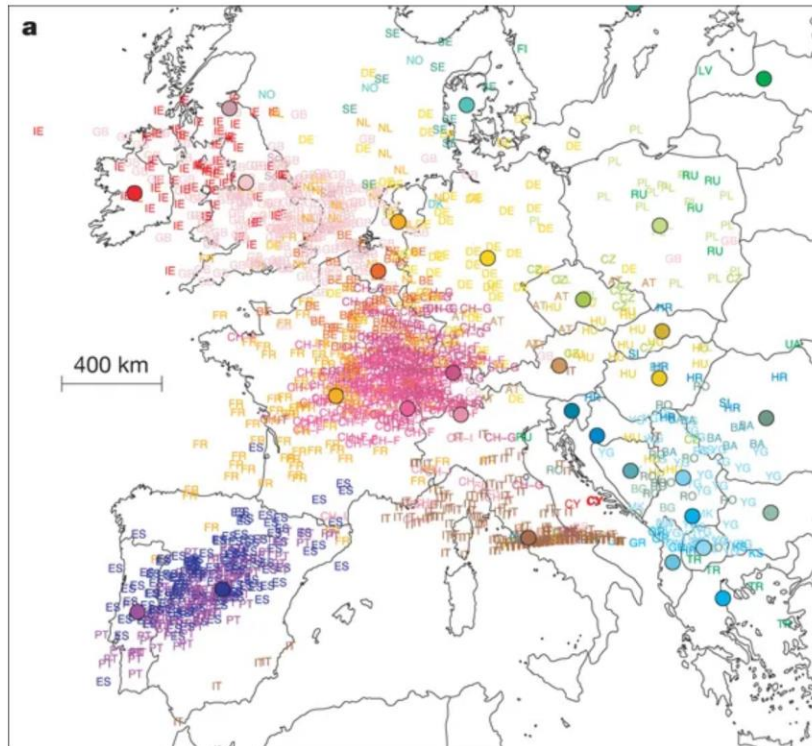
Genomic control

- The classical approach (Devlin and Roeder 1999)
- Define a **genomic inflation factor**
 $\lambda = \text{median}(\chi^2)/0.456$
and adjust the test statistics:
 $\chi^2_{\text{adj}} = \chi^2/\lambda$
- Has some assumptions that may fail
 - "Few" SNPs are truly causal
 - All SNPs have the same inflation factor



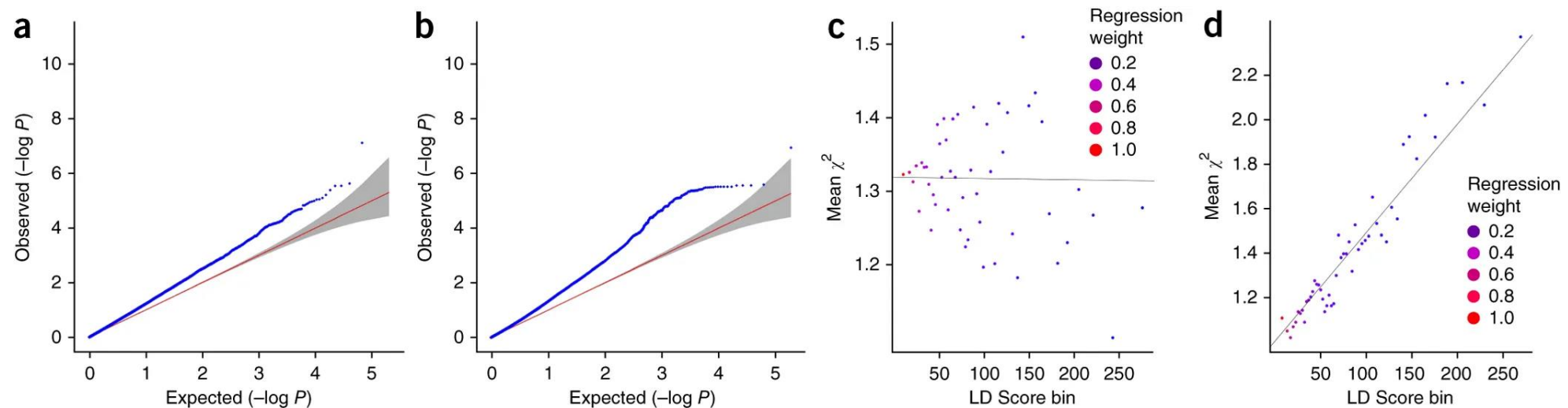
Principal components

- Principal components from the population are added as covariates to the model.
- Principal components usually capture large parts of population stratification and cryptic relatedness.



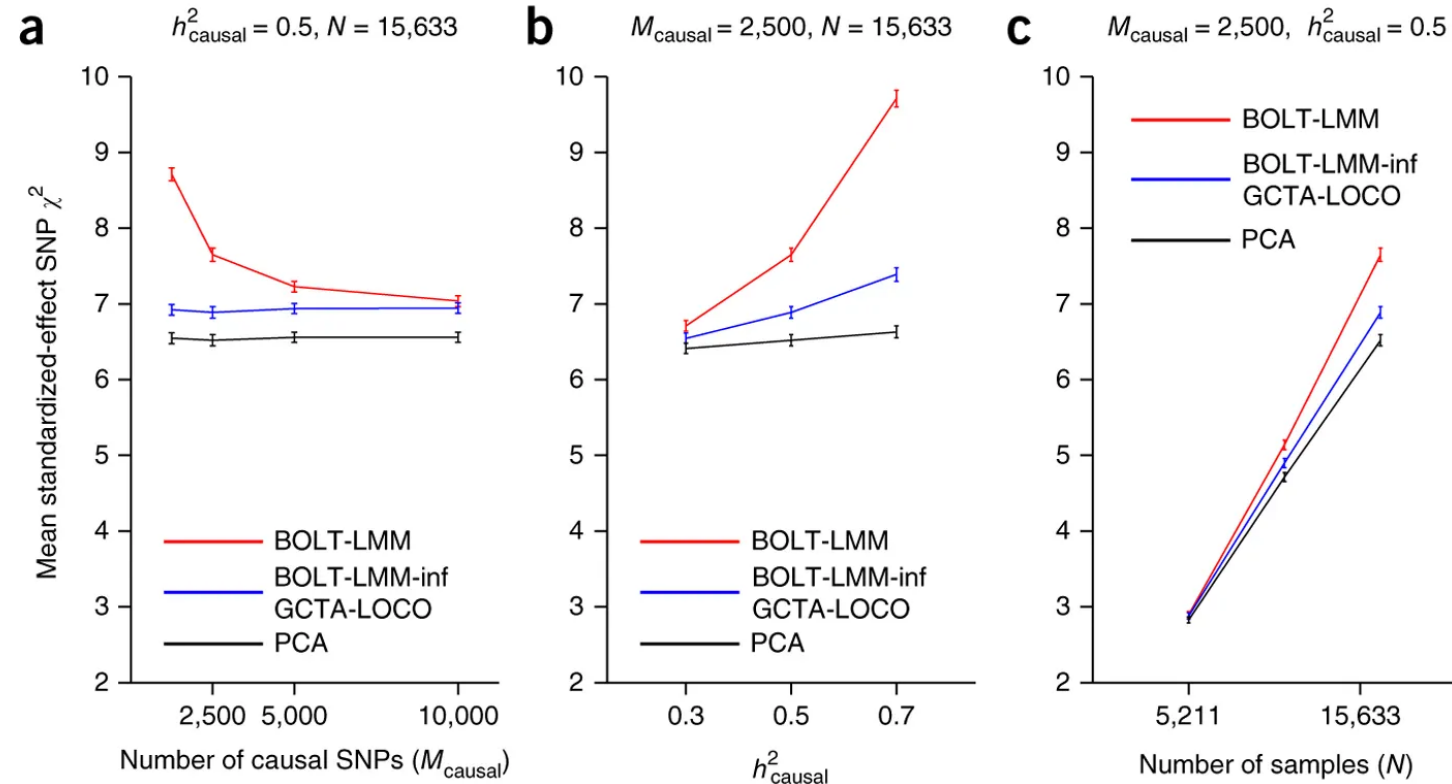
LDscore regression

- In the polygenic model, $E[\chi^2 | s_j] = Nh^2s_j/M + Na + 1$, where
 - N is the number of individuals;
 - M is the number of SNPs, such that h^2/M is the average heritability per SNP;
 - a is the contribution of confounding biases;
 - $s_j = \sum_k r_{jk}^2$ is the LD-score of variant j .
- Individual SNPs weighted by their LD-score
- Distinguishes between polygenicity and bias:



Mixed models

- Mixed models allow us to
 - model relatedness/population stratification,
 - jointly model all SNPs.
- We model phenotype as $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{W}\mathbf{u} + \boldsymbol{\varepsilon}$,
with $\text{Var}(\mathbf{y}) = \mathbf{W}\mathbf{W}^T\sigma_u^2 + \mathbf{I}\sigma_\varepsilon^2$.
- BOLT-LMM works around the computational difficulties.
- Assuming a non-infinitesimal model greatly increases power.
- SNPs assumed to have mixed-normal prior effect sizes



But what are the biases and limitations we're missing?

- We don't really know what all but the top few PCs are capturing
- Ascertainment bias in race/geography/socio-economic status
- Indirect genetic effects: We do not include information on families
- And probably like a million more

