

Analyses of data from stratified populations: a methods comparison

LUPA Workshop in Statistical Methods for GWAS studies

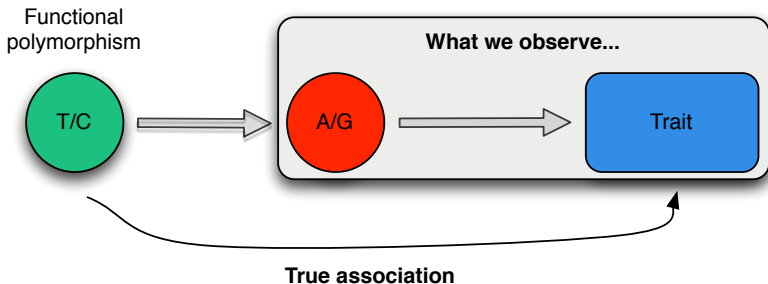
Marcin Kierczak*

*Computational Genetics Group
Faculty of Veterinary and Animal Breeding
Swedish University of Agricultural Sciences
Uppsala, SWEDEN

04-05 May 2011, Uppsala

Observed associations

Most of the observed associations are due to confounders!



Definition 1 – population (after: ABEL tutorial)

Individuals l_1 and l_2 belong to the same population P_1 and l_3 belongs to some other population iff:

$$\text{kin}(l_1, l_2) > 0 \ \&$$

$$\text{kin}(l_1, l_2) \ll \text{kin}(l_1, l_3) \ \&$$

$$\text{kin}(l_1, l_2) \ll \text{kin}(l_2, l_3)$$

Introduction

What population stratification is?

Definition 2 – population stratification (after: Price et al., 2006)

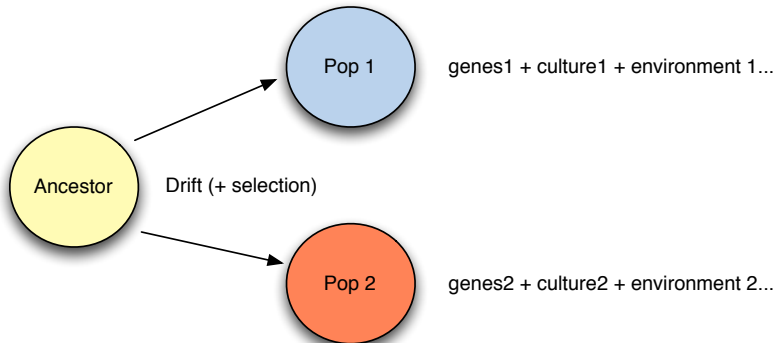
Population stratification – differences in allele frequency in cases and controls due to systematic ancestry differences.

Consequence

Possible spurious associations in GWAS!

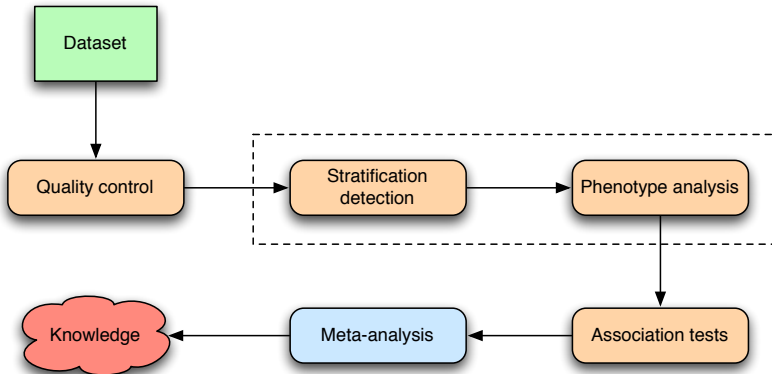
Introduction

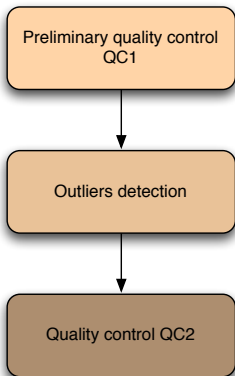
Population stratification



Workflow

Basic workflow in GWAS

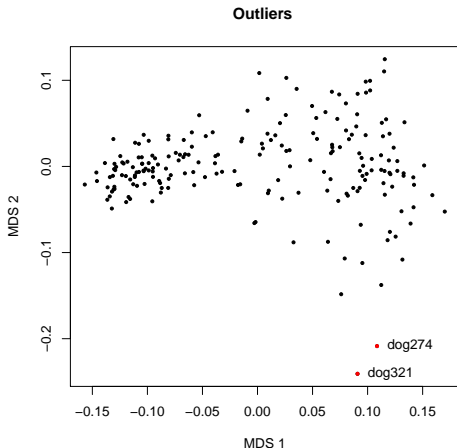




Three phases of QC

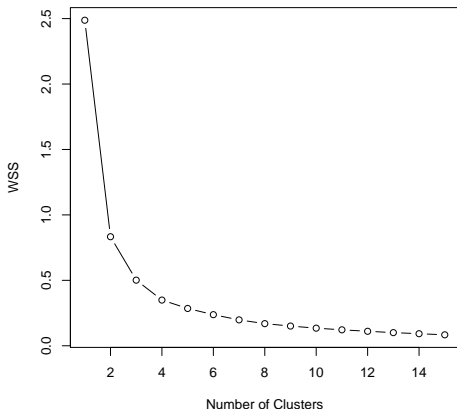
- **Preliminary (QC1)** – call rate, too close kinship, only strong departures from HWE.
- **Outliers detection** – removal of outliers using genomic kinship-based distances.
- **Final (QC2)** – departures from HWE, esp. in controls.

Qc1 and QC2 are iterative!



Outliers detection

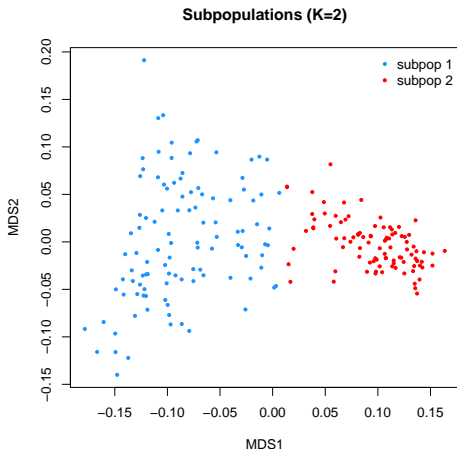
- Compute IBS matrix using autosomal markers (subset).
- Transform IBS matrix to dist. matrix.
- Run MDS (\Rightarrow 2D)
- Visualize and identify outliers.
- Remove outliers (avg. dist. from cluster centre?).



How many subpops?

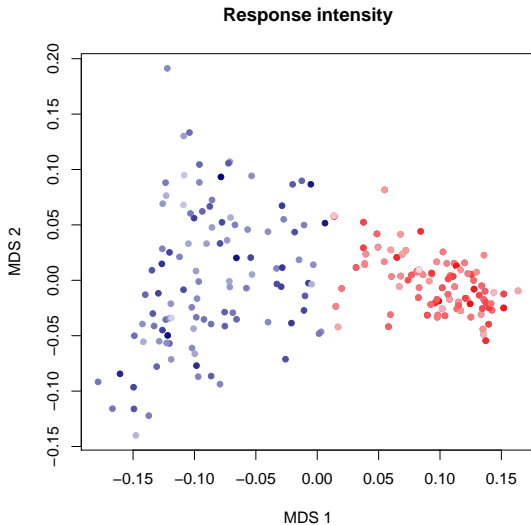
- Visual assessment.
- Clustering e.g., K-means.

In K-means clustering, plot WCSS for different K and look for a bend. Similar to a scree test.



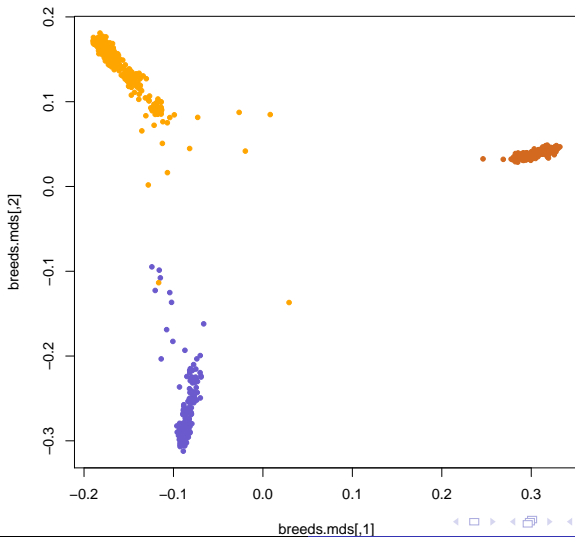
What else can be done?

- Check which individuals fall into the clusters.
- Compare sex distribution in the clusters.
- Compare distribution of cases/controls in the clusters.
- More strata?



Workflow

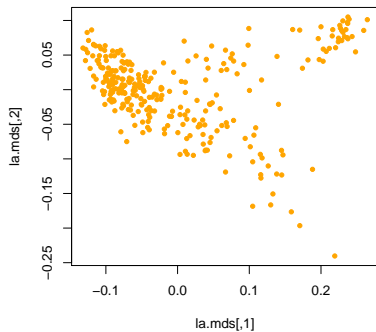
Population structure – jumping over the fence?



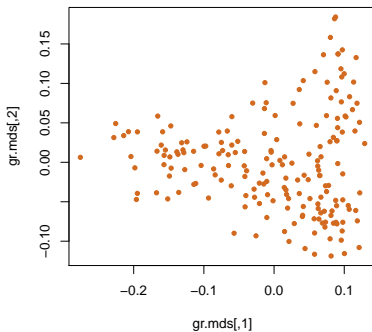
Workflow

Population structure – separately for each breed

Labradors



Golden retrievers



It is always worth to include all available information and to make comparisons between clusters.

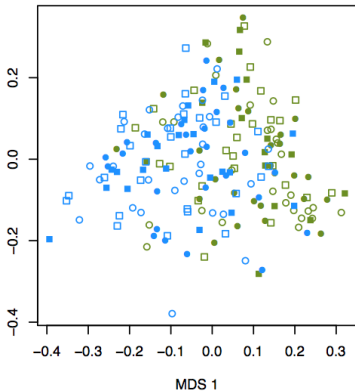
Consider

- Geographic location.
- Kennel club.
- Age (generation).
- Type of dog, e.g.: working, show.
- Import due to, e.g.: fashion.

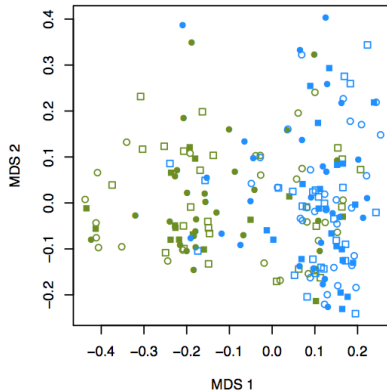
Workflow

Population structure – separately for each chromosome

Chromosome 33



Chromosome 34



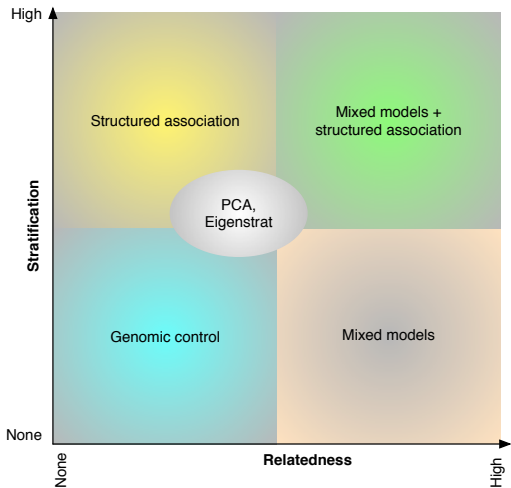
Different approaches exist to correct for stratification:

- Simple association tests.
- Genomic control.
- PCA-based correction – Eigenstrat.
- Mixed models.
- Structured association.
- Combinations of the above!
- Use of co-variate, e.g.: breed, address, etc.

Which method to use?

Workflow

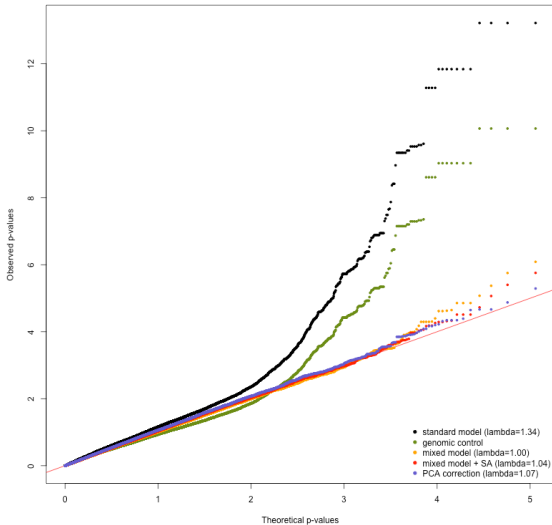
Applicability of different methods



after: Aulchenko et. al., 2010 (ABEL Tutorial)

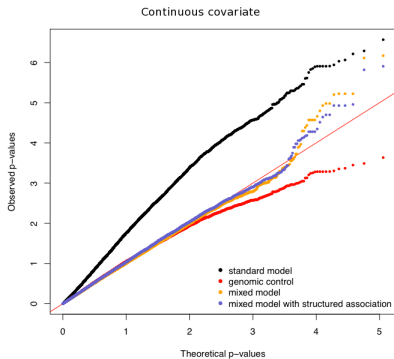
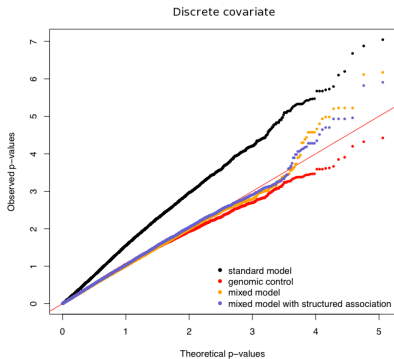
Workflow

Comparing methods...



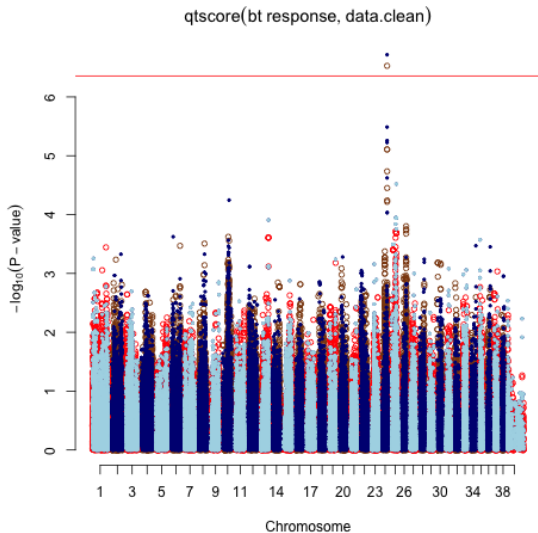
Workflow

Comparing methods – using co-variate



Workflow

Comparing methods – Manhattan plots



Conclusion

There is no single universal approach. Each dataset is specific and has to be examined in detail! Look at your data and experiment with different approaches!

Occam's razor

The simplest approach that gives good correction should always be chosen! BUT: do not fall into a trap of choosing a method that corrects "by chance". Do not use GC when you know your data is stratified!

Thank You! and:

- Leif Andersson
- Yurii Aulchenko
- Örjan Carlborg
- Hille Fieten
- Dirk Jan de Koonig
- Kerstin Lindblad-Toh
- Xia Shen
- Katarina Tengvall