

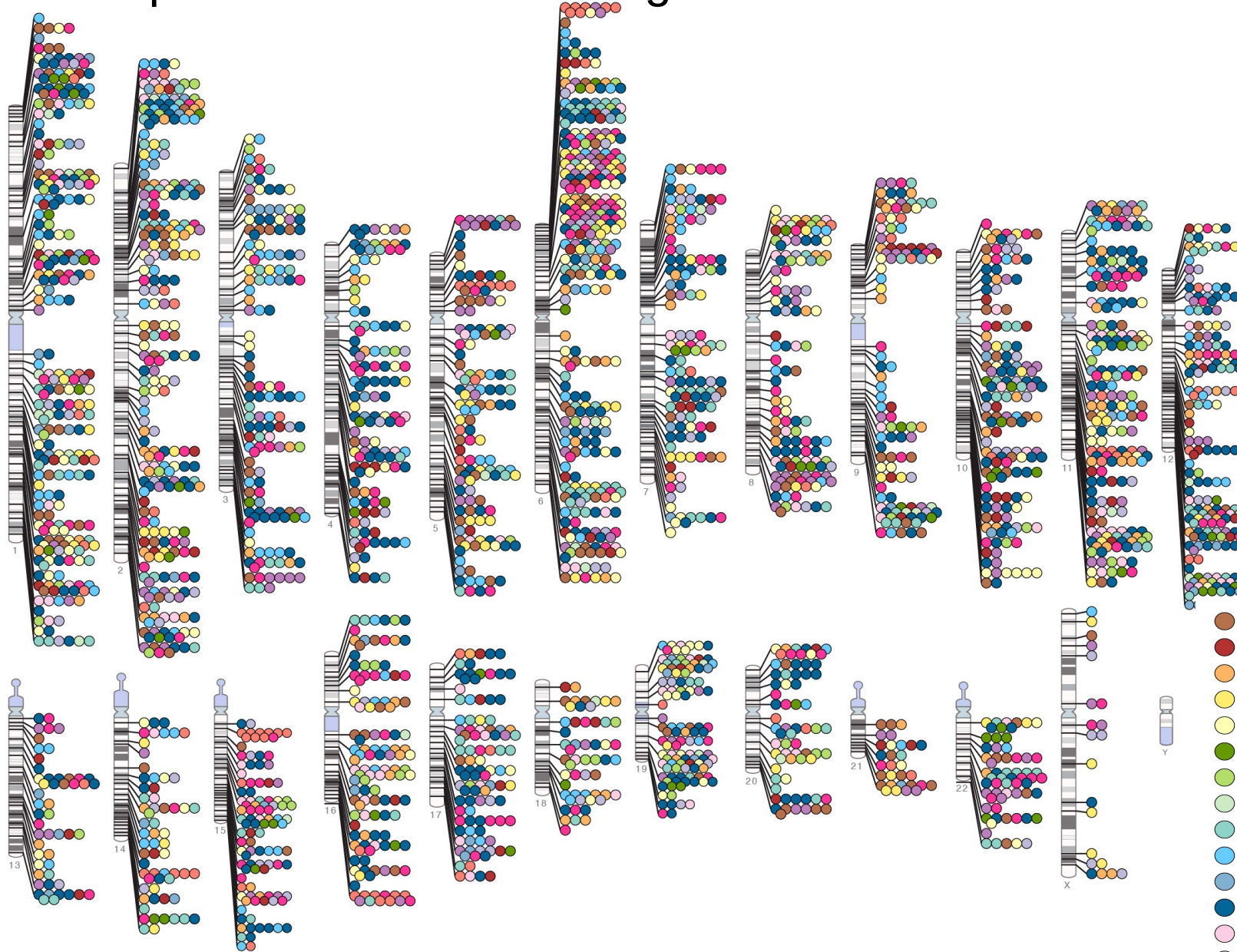
A minicourse on
Genomewide association studies
(GWAS)
Part I

Matti Pirinen
FIMM, University of Helsinki

December 4th 2012

Published Genome-wide associations 07/2012

$p < 5e-8$ for 18 trait categories



Contents:
1. Concepts
And Rationale

2. Technology
and Data

3. Statistics
and Analysis

4. Results and
Hot topics

- Digestive system disorder
- Cardiovascular disorder
- Metabolic disorder
- Immune system disorder
- Neurological disorder
- Liver enzyme measurement
- Lipid or lipoprotein measurement
- Inflammatory marker measurement
- Hematological measurement
- Body measurement
- Cardiovascular measurement
- Other measurement
- Chemical compound
- Biological process
- Cancer
- Other disease
- Other trait
- Trait mapping in progress

Concepts

- Individuals (humans, mice, flies, grass)
- Phenotypes, traits
 - Quantitative (height, cholesterol levels)
 - Discrete (case-control, Parkinson's disease)
- EXPLAIN
 - Understanding, mechanisms, therapeutics
 - Example: SORT1 (Musunuru 2010)
- PREDICT
 - Early intervention, agriculture

$$Y = \mu + G + E + (G \times E)$$

- Y phenotype, μ population mean
- E nvironment
 - Chemicals, temperature, food, physical activity...
- G enetics
 - DNA, A,C,G,T, $3e+9$ bases, 22 chrs + X + Y, diploid, meiosis
 - Genes \rightarrow proteins, $2e+4$, 1-2 % of DNA
 - Single nucleotide polymorphisms (SNPs)
- GxE not in this course

Linkage disequilibrium

- Non-independence of alleles at (two) SNPs in population
- r^2 and D' , (Wray et al. 2011, p4 BOX 2)
- Look at 1000 Genomes browser: rs4988235
- LD \rightarrow SNPs tag each other in a population sample
- Basis for gene mapping with a SNP panel

Heritability

- $Y=G+E$
- Simplest model without measured variants or environment
 - $G \sim N(0, s^2)$, $\text{cor}(\text{full-sibs}) > \text{cor}(\text{half-sibs})$
 - $E \sim N(0, t^2)$, $\text{cor}(\text{household}) > 0$
 - $H = \text{var}(G) / \text{var}(Y) = s^2 / (s^2 + t^2)$, heritability
 - $\text{var}(Y)$ can be decomposed to s^2 and t^2 with variance component models, but environment is a potential confounder

Example: ACE twin model

- $Y=A+C+E$
 - A, additive genetic component
 - C, common environment for twins
 - E, non-shared environment for twins
- From (many pairs of) monozygotic (identical) twins AND dizygotic twins, $\text{var}(A)$, $\text{var}(C)$ and $\text{var}(E)$ can be estimated → estimate of heritability as $2(\text{cor}(\text{MZ})-\text{cor}(\text{DZ}))$
 - But with strong assumptions!

Linkage analysis

- Families
 - Does trait correlate with genetic sharing at some parts of the genome?
 - Parametric Linkage analysis in pedigrees
 - Affected sib-pairs (non-parametric)
 - Need for families restricts sample size and thus size of genetic effects that can be found
 - Localisation is coarse since large blocks of genome are linked in close relatives

Association analysis

- If a large panel of SNPs can be genotyped, association between each SNP and the trait can be tested
 - Role of LD (HapMap project)
- Risch 2000, Fig 4
- Only samples from (homogeneous) population (not families) needed
- Genome coverage and localisation depend on #SNPs and LD in the population

Genotyping

- Chips
 - YOUTUBE: “Microarray method for genetic testing”
 - YOUTUBE: “DNA chips and microarrays”
 - Currently 50-200 euros per sample
- Intensities → genotype calls (Vukcevic p.36)
- Problems (Vukcevic p.38)

Genotype probabilities

```
rs4345758-128_B_F_1516310071 rs4345758 28663 A G 1 0 0 1 0 0 1 0 0 1
rs10399793-128_B_F_1501305891 rs10399793 39161 A G 0 0.9985 0.0015 0
rs2462492-128_T_R_1551017801 rs2462492 44539 A G 0 1 0 1 0 0 0 1 0 0
rs3107975-128_B_F_1551017858 rs3107975 45189 A G 0 1 0 0 0.0001 0.99
rs4420028-128_T_R_1551213529 rs4420028 63065 A G 1 0 0 1 0 0 1 0 0 1
rs2462495-128_B_F_1551083146 rs2462495 68896 A G 0 0 1 0 0 1 0 0 1 0
rs3878915-128_B_F_1551017866 rs3878915 70249 A C 0 0 1 0 0 1 0 0 1 0
rs4477212-128_B_R_1513978262 rs4477212 72017 A G 1 0 0 1 0 0 1 0 0 1
rs1856862-128_T_F_1516310110 rs1856862 110905 A T 0 0 1 0 0 1 0 0 1
rs4096703-128_T_R_1516308906 rs4096703 160593 C G 1 0 0 1 0 0 1 0 0
rs6670732-128_T_R_1551148139 rs6670732 228130 A G 0 0 1 0 0 1 0 0 1
rs3956613-128_T_R_1551213503 rs3956613 344638 A G 0 0 1 0 0 1 0 0 0
rs2808353-128_B_F_1551017760 rs2808353 515561 A G 0 0 1 0 0 1 0 0 1
rs28780398-128_T_R_1551213411 rs28780398 516424 A C 0 0.9890 0.0111
rs28863004-128_B_F_1550922896 rs28863004 516599 C G 0 0 0 0 0.9973 0
rs6680723-128_T_R_1551148123 rs6680723 524055 A G 1 0 0 1 0 0 1 0 0
rs6683466-128_B_F_1551017542 rs6683466 524446 C G 0 0 1 0 0 1 0 0 1
rs12025928-128_B_R_1501322495 rs12025928 536560 A G 0 1 0 0 1 0 0 1
rs9633435-128_B_R_1551087070 rs9633435 554454 A C 1 0 0 1 0 0 1 0 0
```

The data comes with 5 types of “header” information

We report the probability for each of the three genotypes

$P(AA) = 0$ $P(AG) = 0.9985$ $P(GG) = 0.0015$

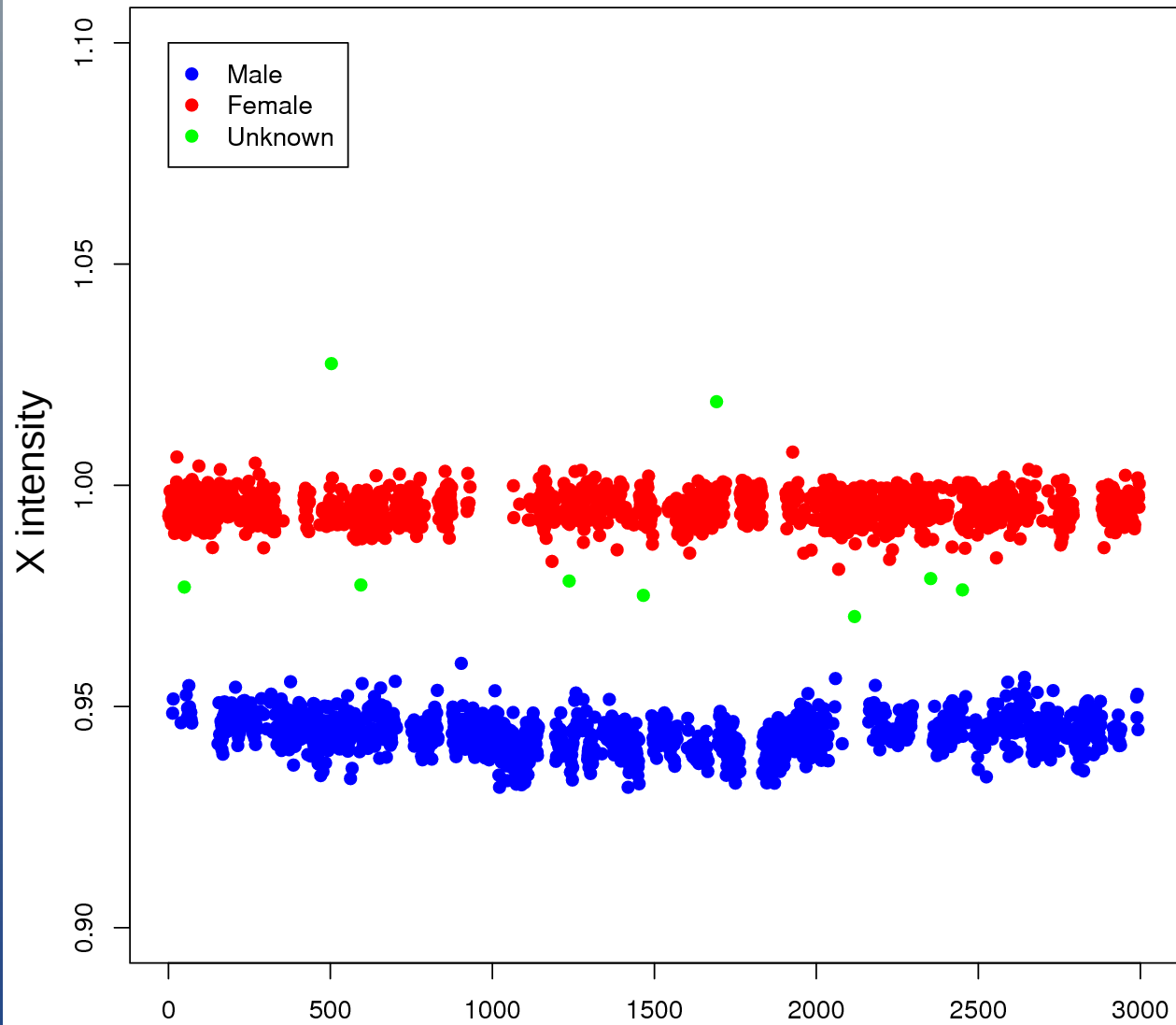
A null call is represent by three zeros

Quality control

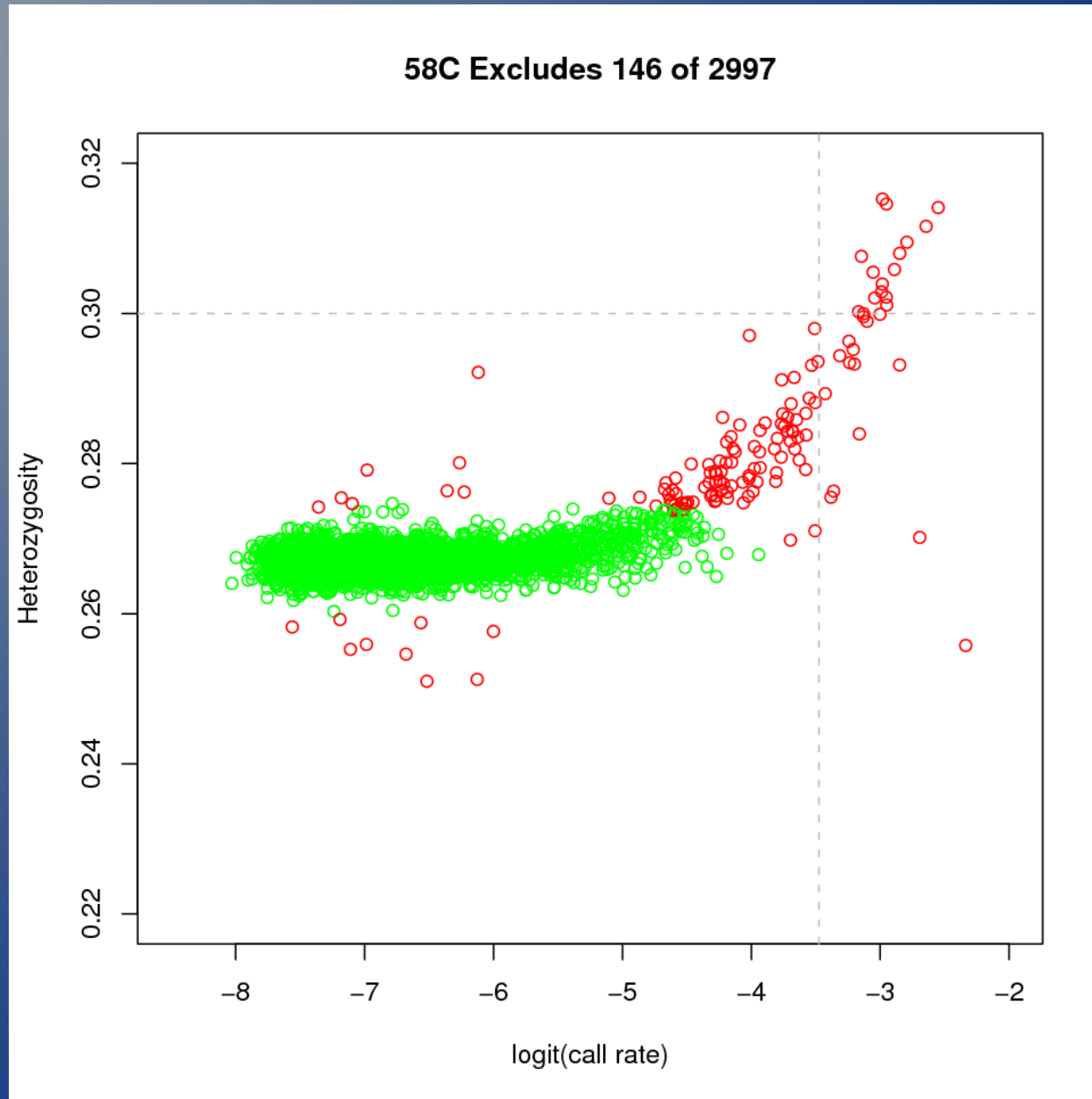
- Remove individuals and SNPs that show low quality
- Anderson et al. 2010
- Individuals: Sex discrepancies, missingness, heterozygosity, relatedness, ancestry
- SNPs: missingness, deviation from Hardy-Weinberg equilibrium, low minor allele frequency

Gender

58C Excludes 9 of 2997



Individual QC



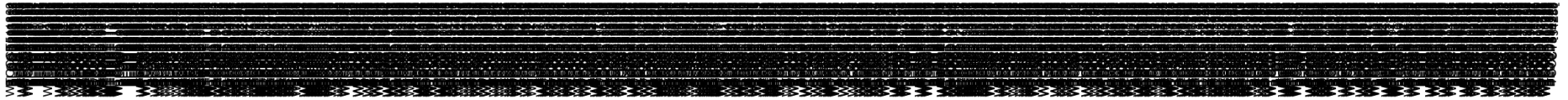
Relatedness

We look for relatedness in two stages:

- Calculate genome-wide allele sharing between all pairs
- For each individual estimate the portion of their genome IBD 0,1 and 2 with their closest relative

We exclude individuals until there is no pair of individuals with more than 5% IBD

Relatedness

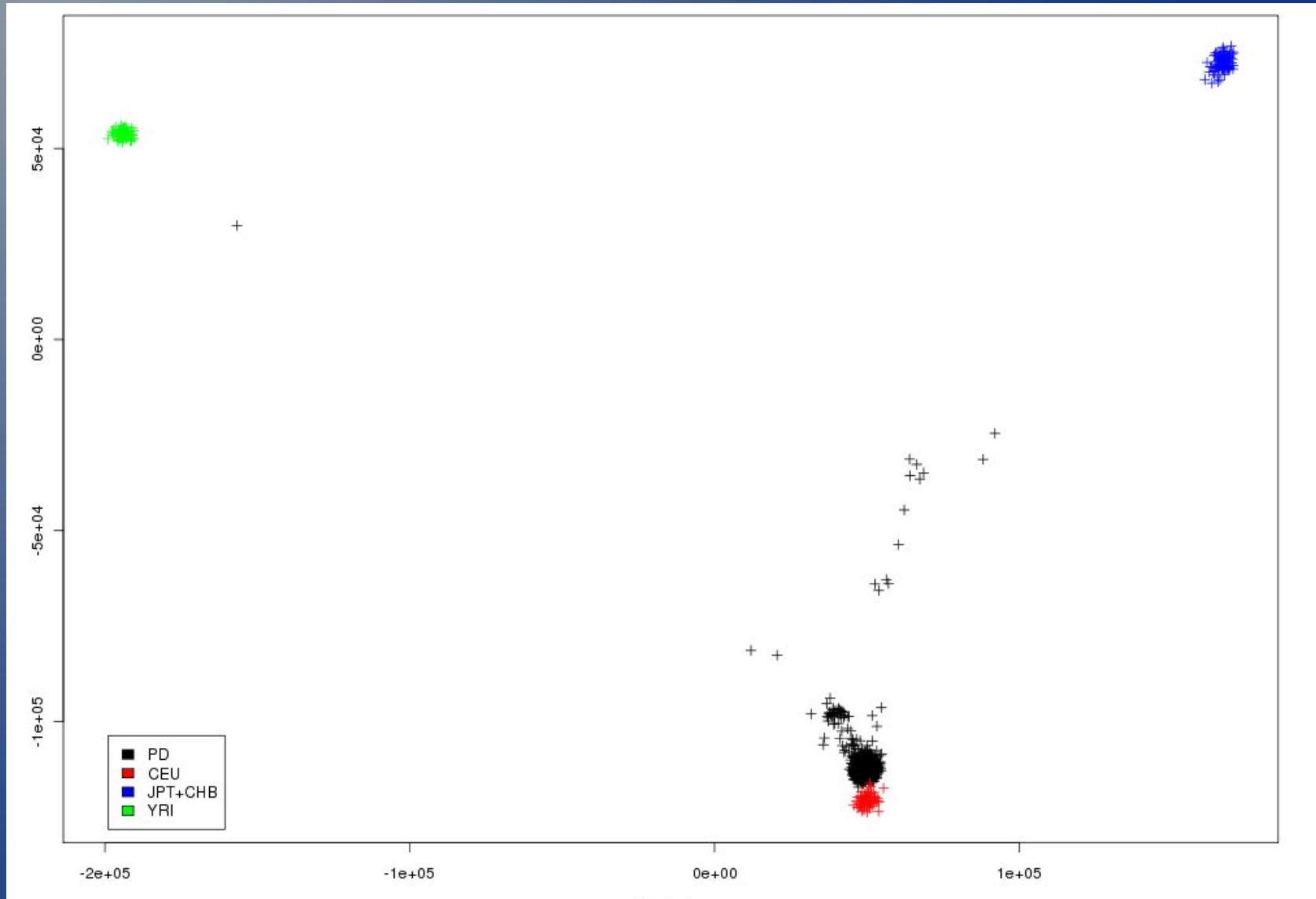


Population structure

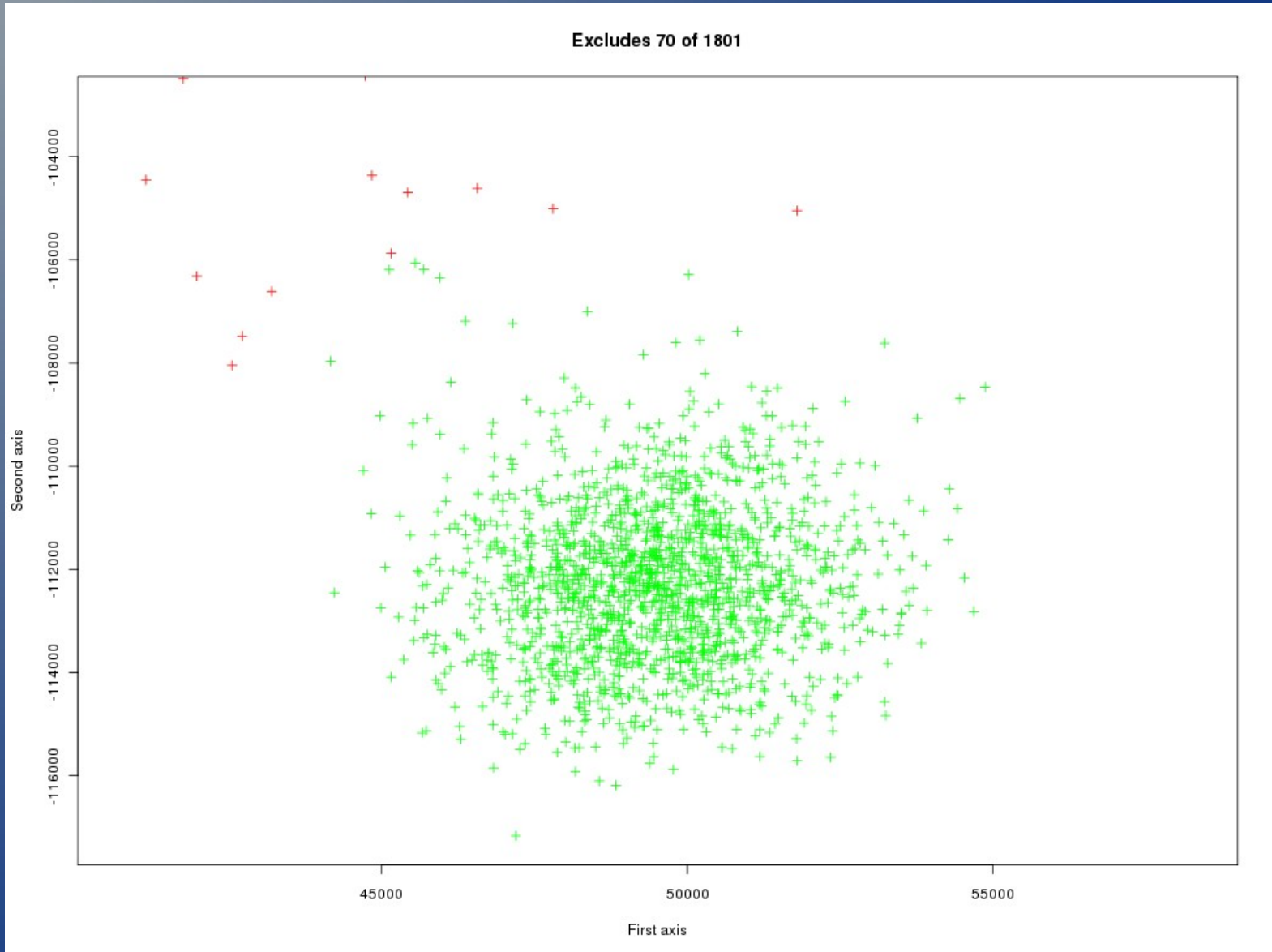
For all collections we project individuals onto axis of variation which are chosen to explain the diversity in HapMap

We cluster individuals and exclude those who's ancestry outlier with respect to the rest of the sample

Population structure



Population structure



QC genotype calls

The optimal approach is either to model the calling errors or to look at all cluster plots

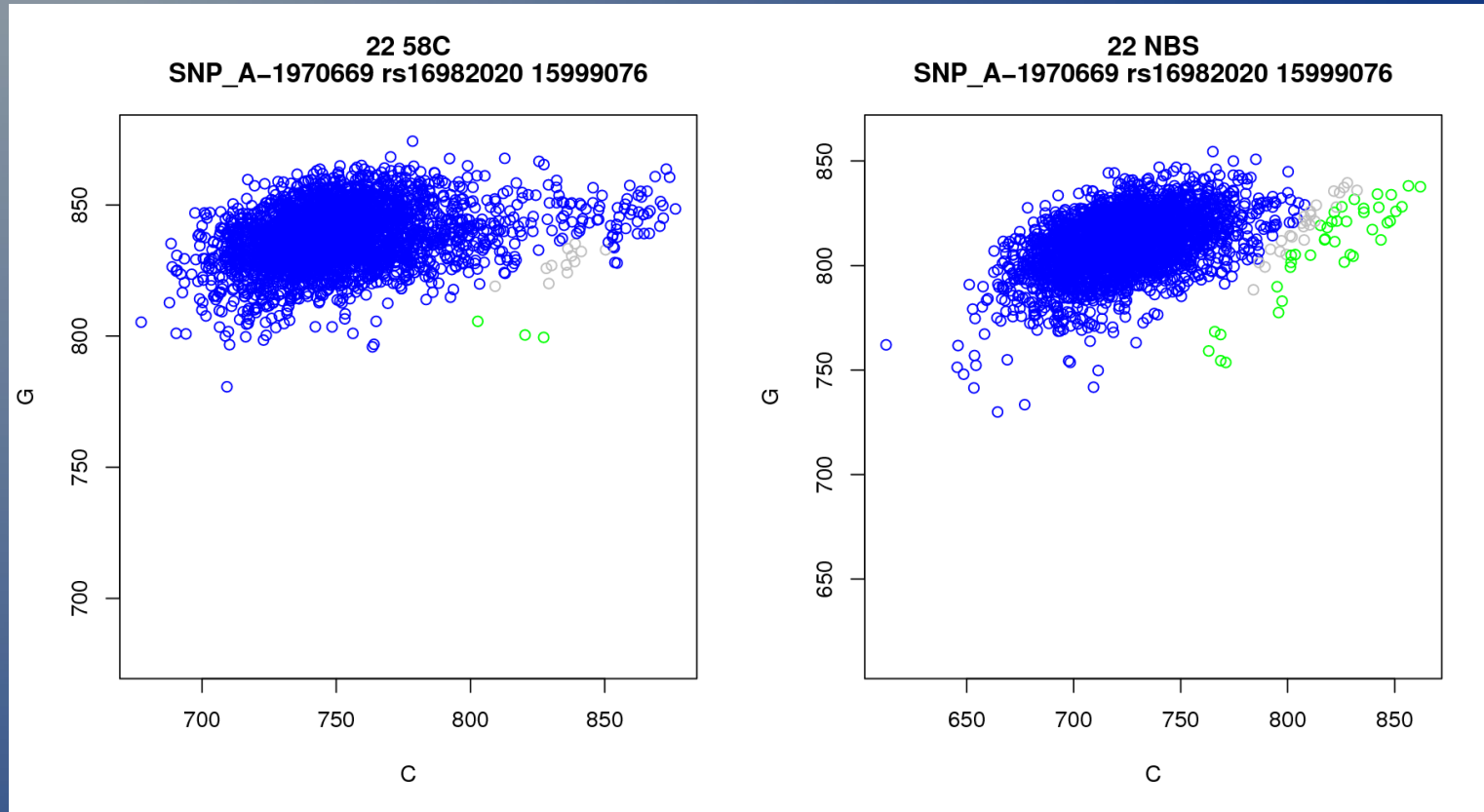
SNP filtering is then a short cut

The level of SNP filtering is therefore a trade-off

Software for looking at cluster plots using binary data
<http://sourceforge.net/projects/evoker/>

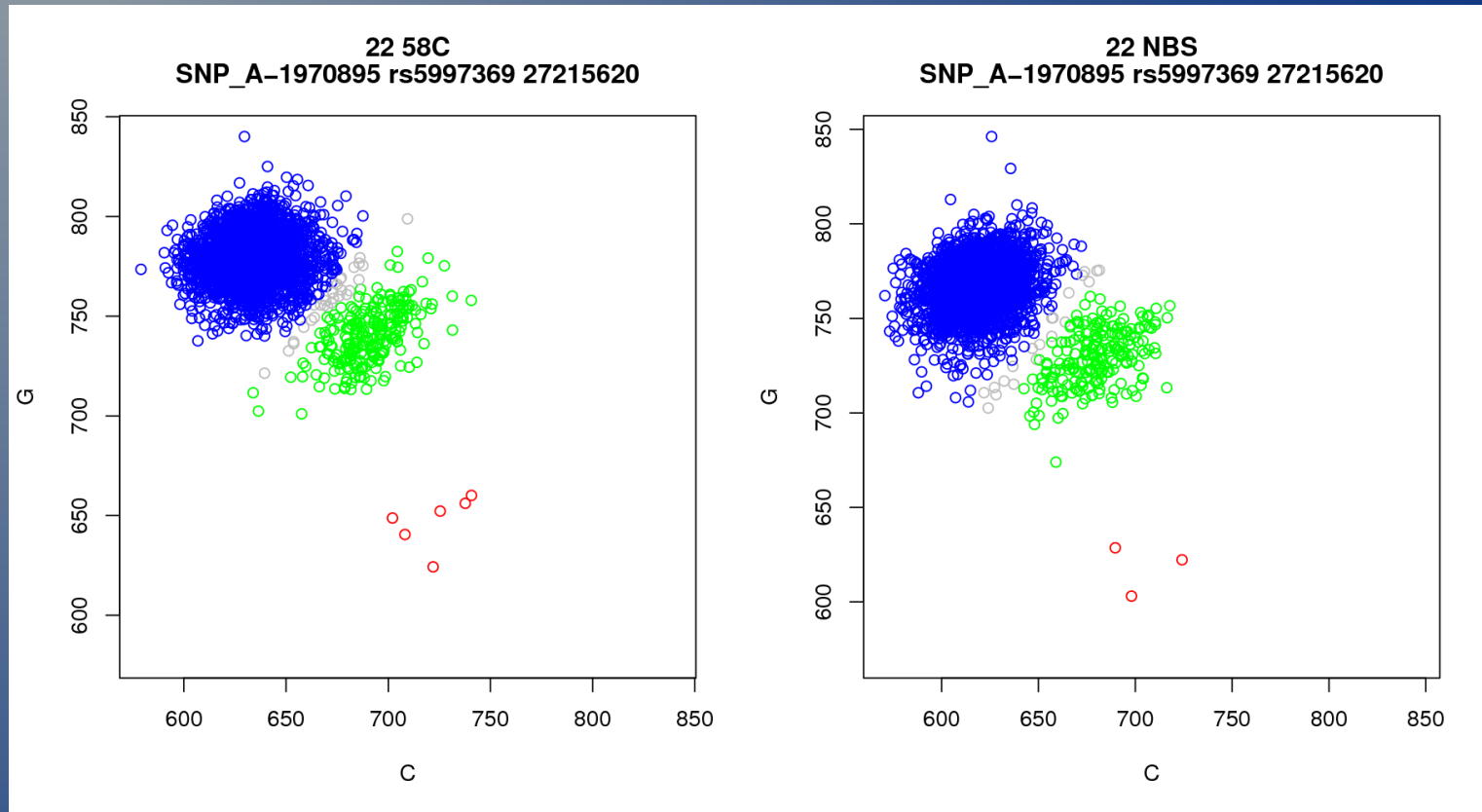
(Jeff Barrett)

Minor allele frequency



Genotype calling algorithms like to model the extra clusters

Low Information



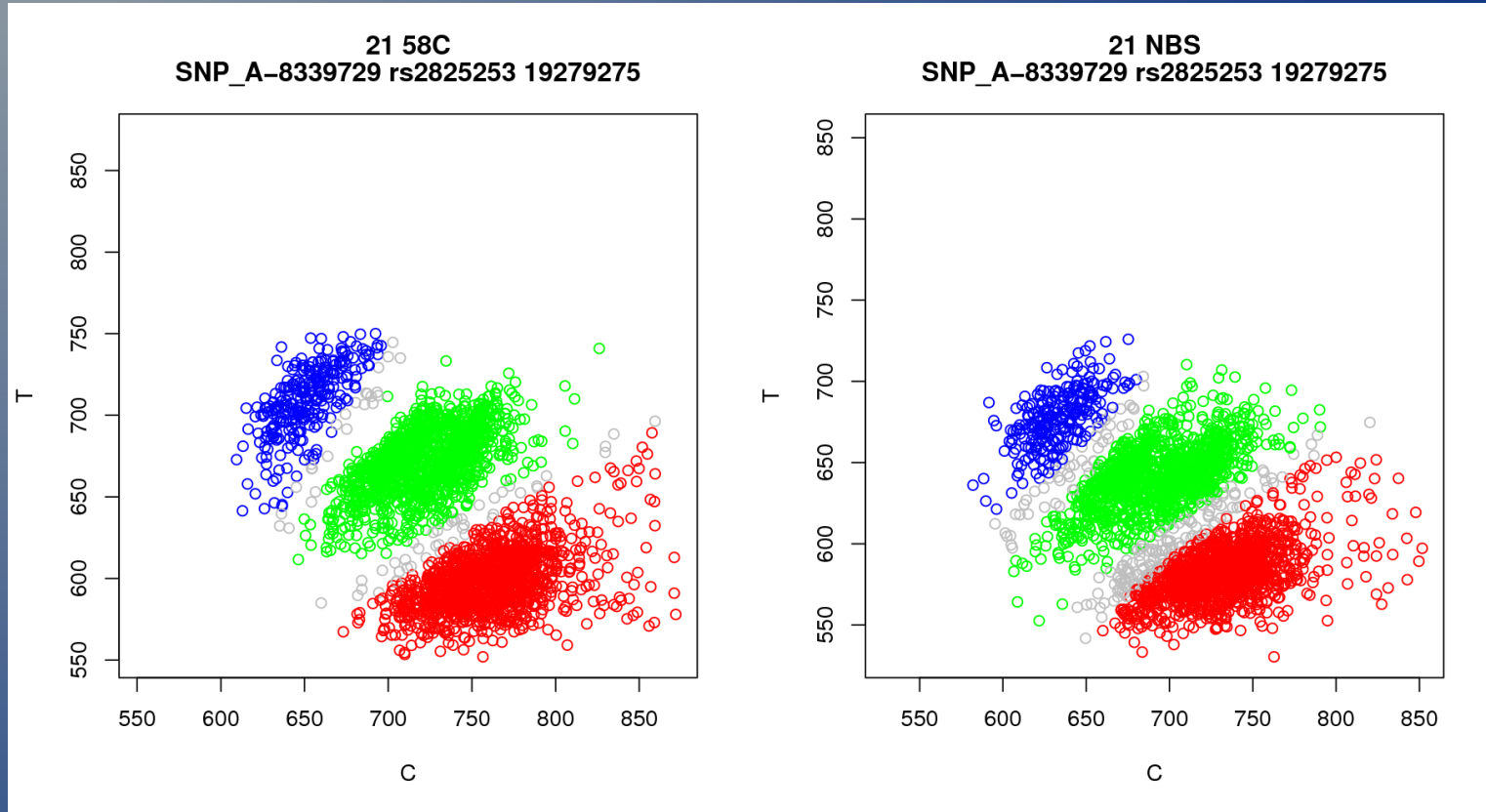
Information is good a identify SNPs with insufficient resolution of genotypes

Hardy-Weinberg



Clustering errors, or putative CNP often distorts
HW equilibrium

Missingness



High missingness is also indicative of clustering failure or poor cluster resolution

SNP exclusions

The application of SNP exclusions is aimed at capturing three classes:

- Poor clustering
- Low signal to noise
- Non-SNP like variation (known CNVs)

Guide to looking at cluster plots in the Supplementary material of WTCCC 2007 paper

Information metric

A natural way to assess the utility of the calls at a SNP is to assess the information about the allele frequency. (info column in SNPTTEST output)

Compare two measures of information

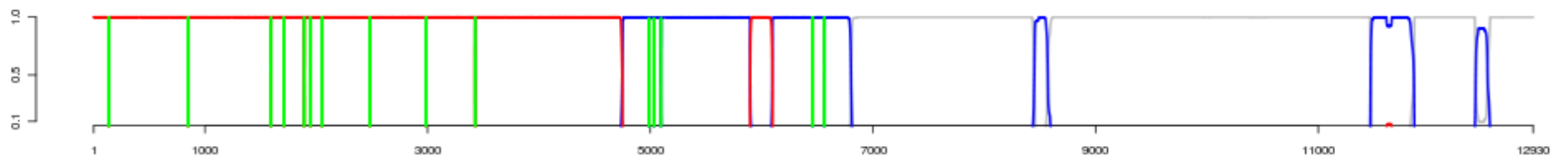
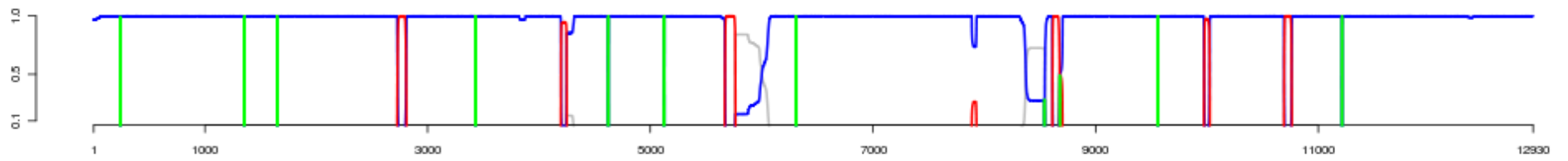
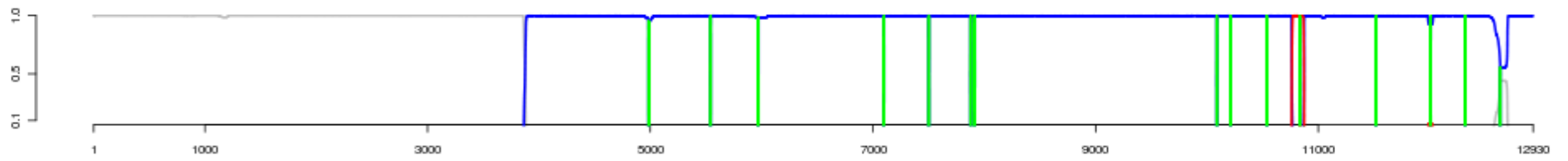
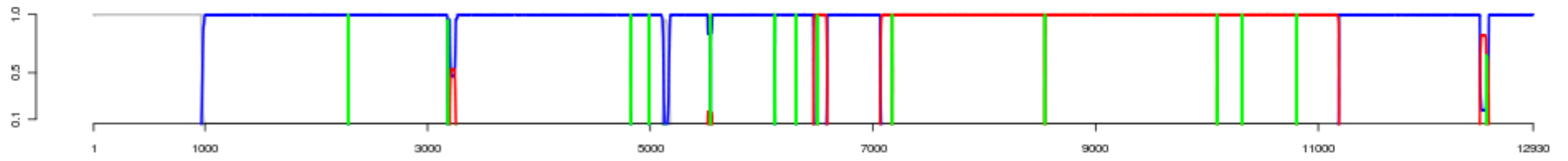
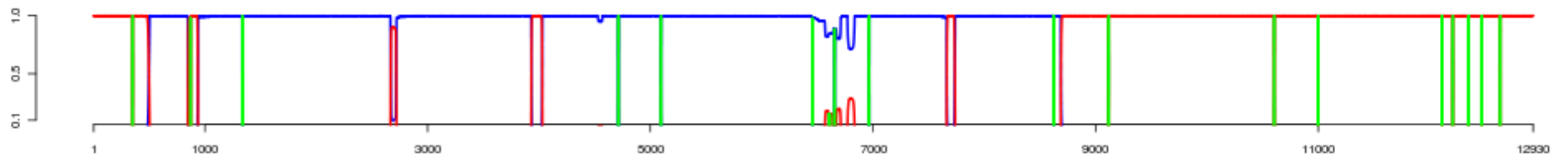
- Using expectation of genotype calls
- Accounting for probabilistic genotypes

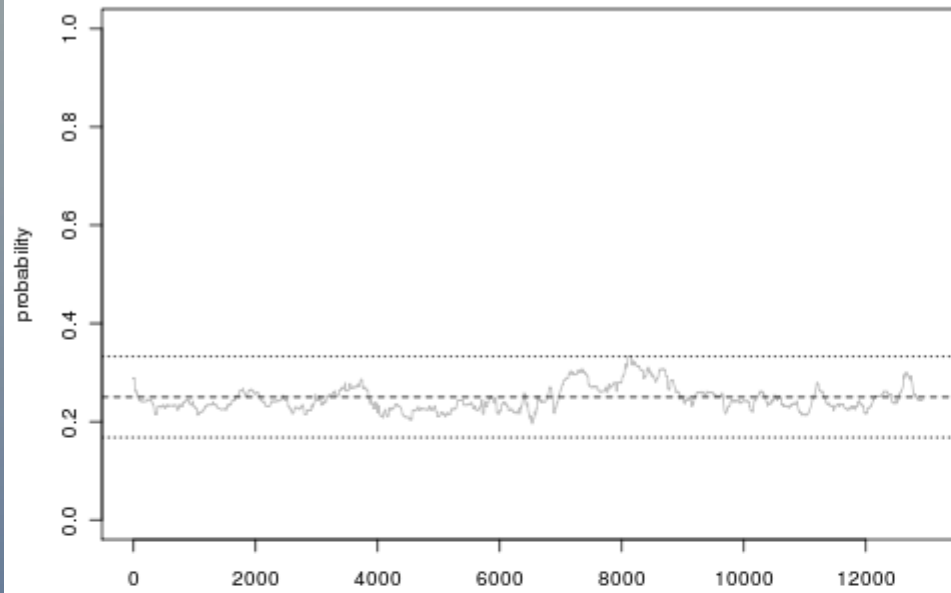
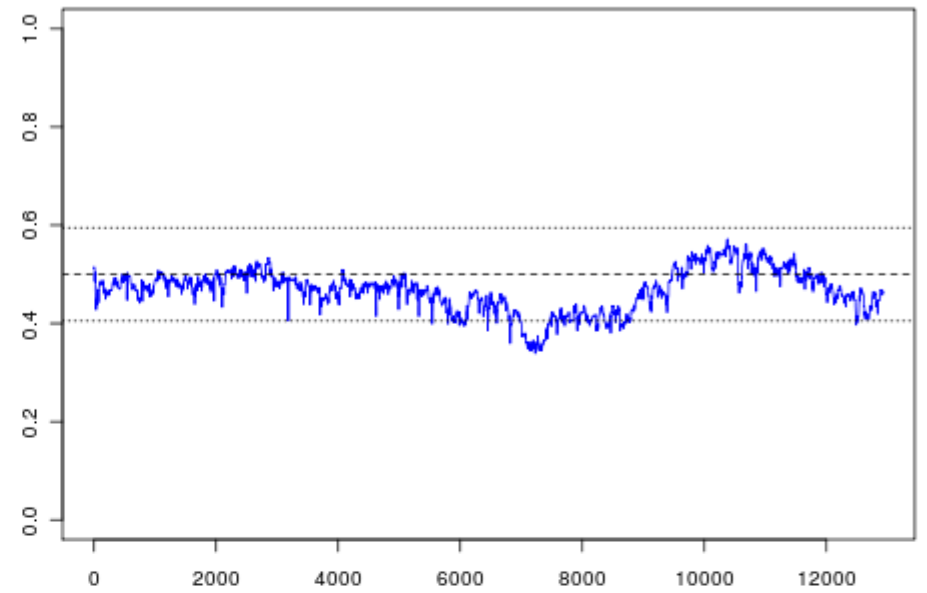
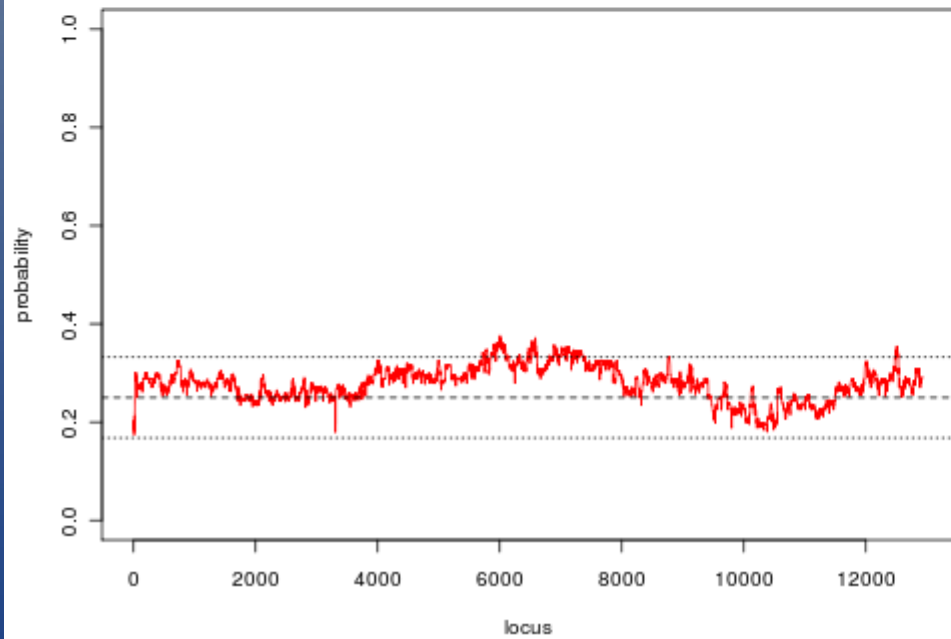
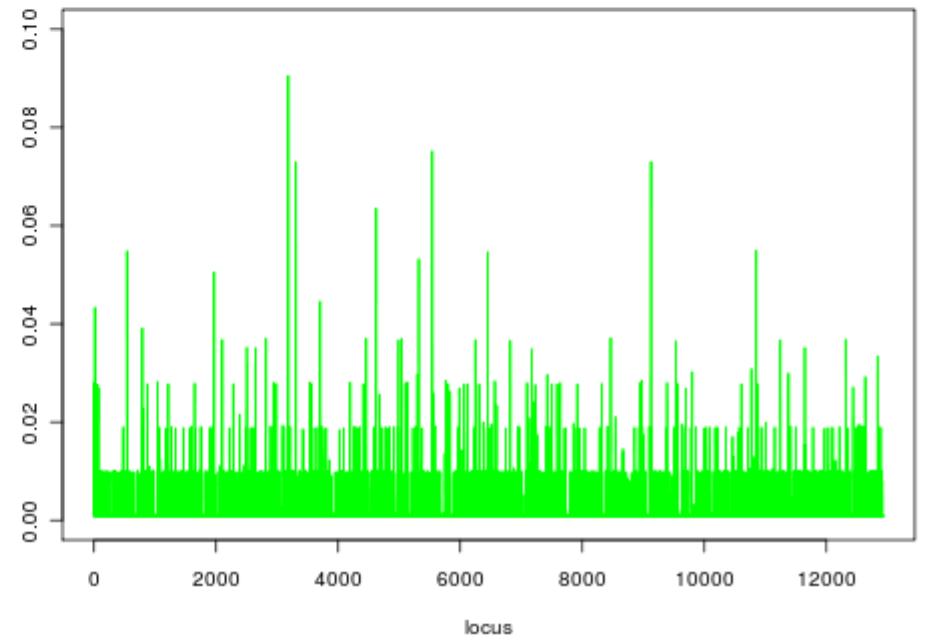
SNP QC

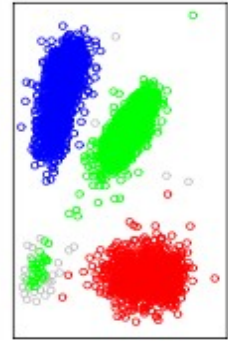
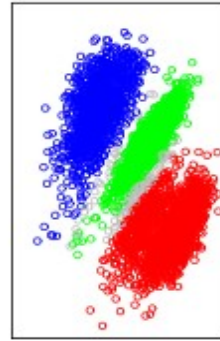
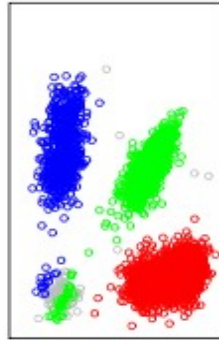
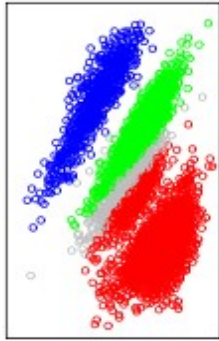
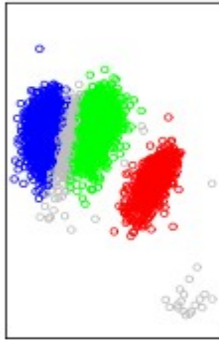
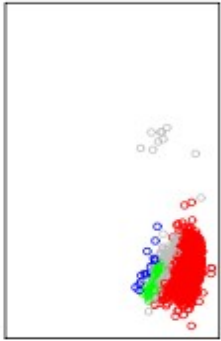
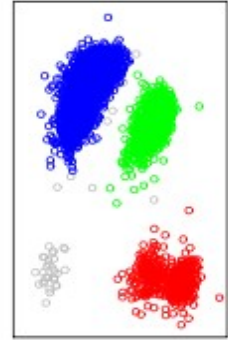
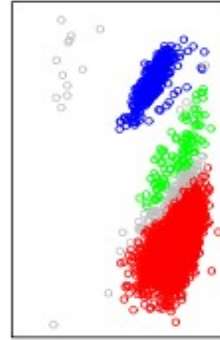
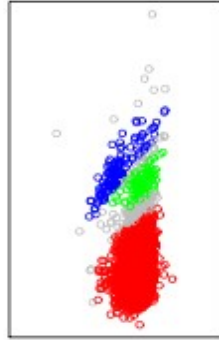
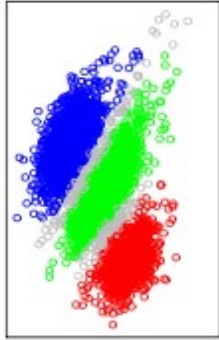
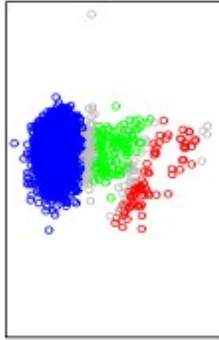
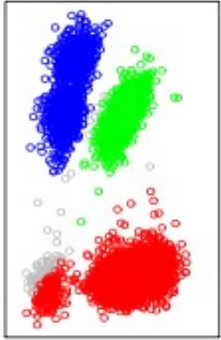
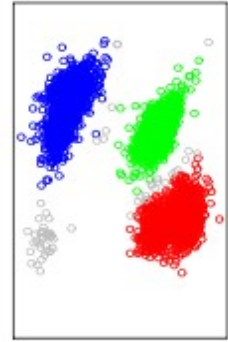
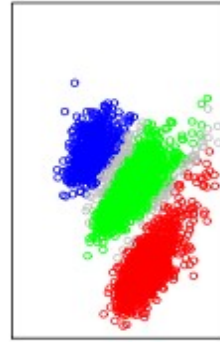
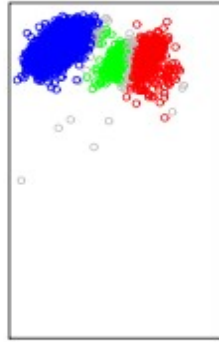
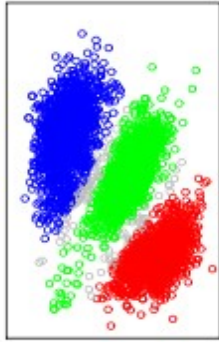
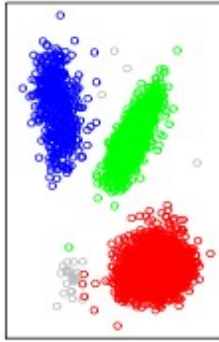
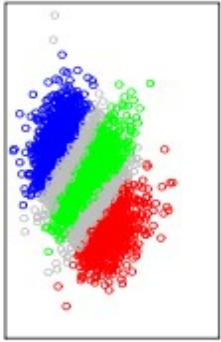
Most poor cluster plots can be identified by these simple filters

We have chosen a set of SNP filters on the basis of the control-control comparison

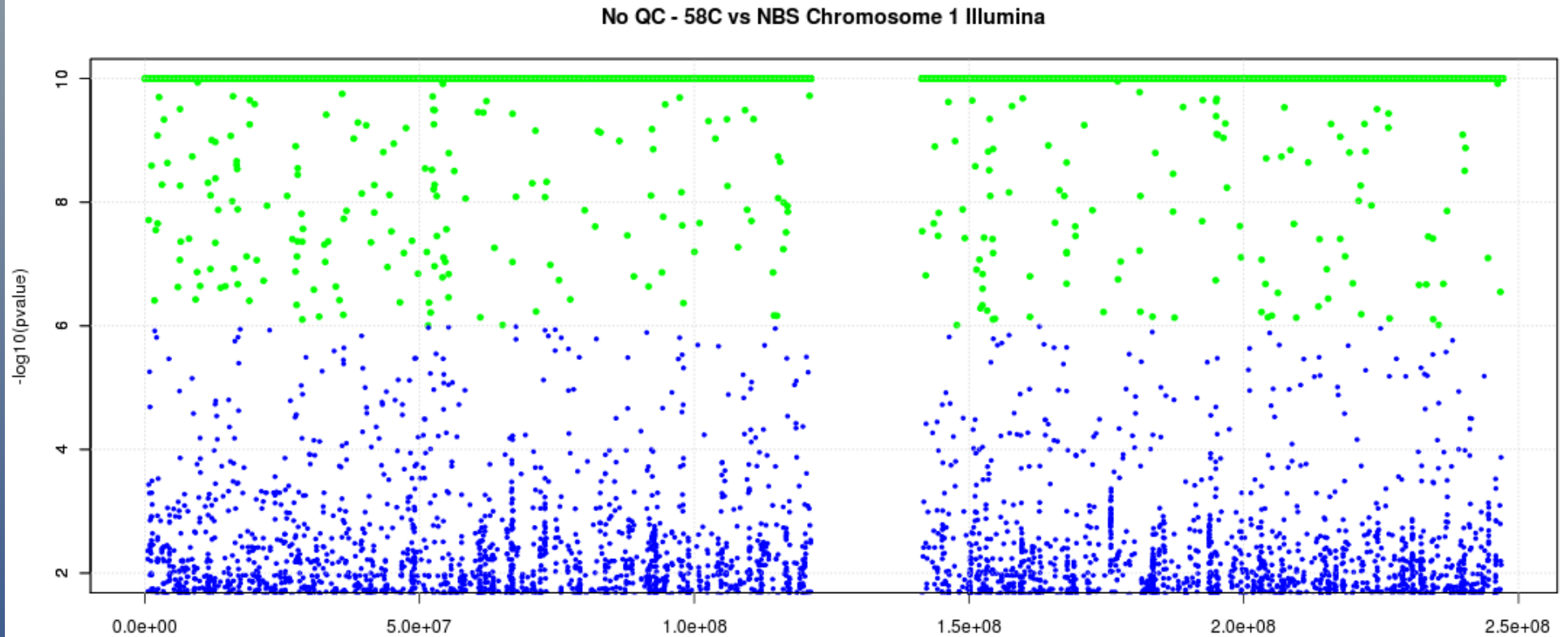
Filters depend on data quality and requires an iterative approach for each data set...



IBD 0**IBD 1****IBD 2****ERROR**

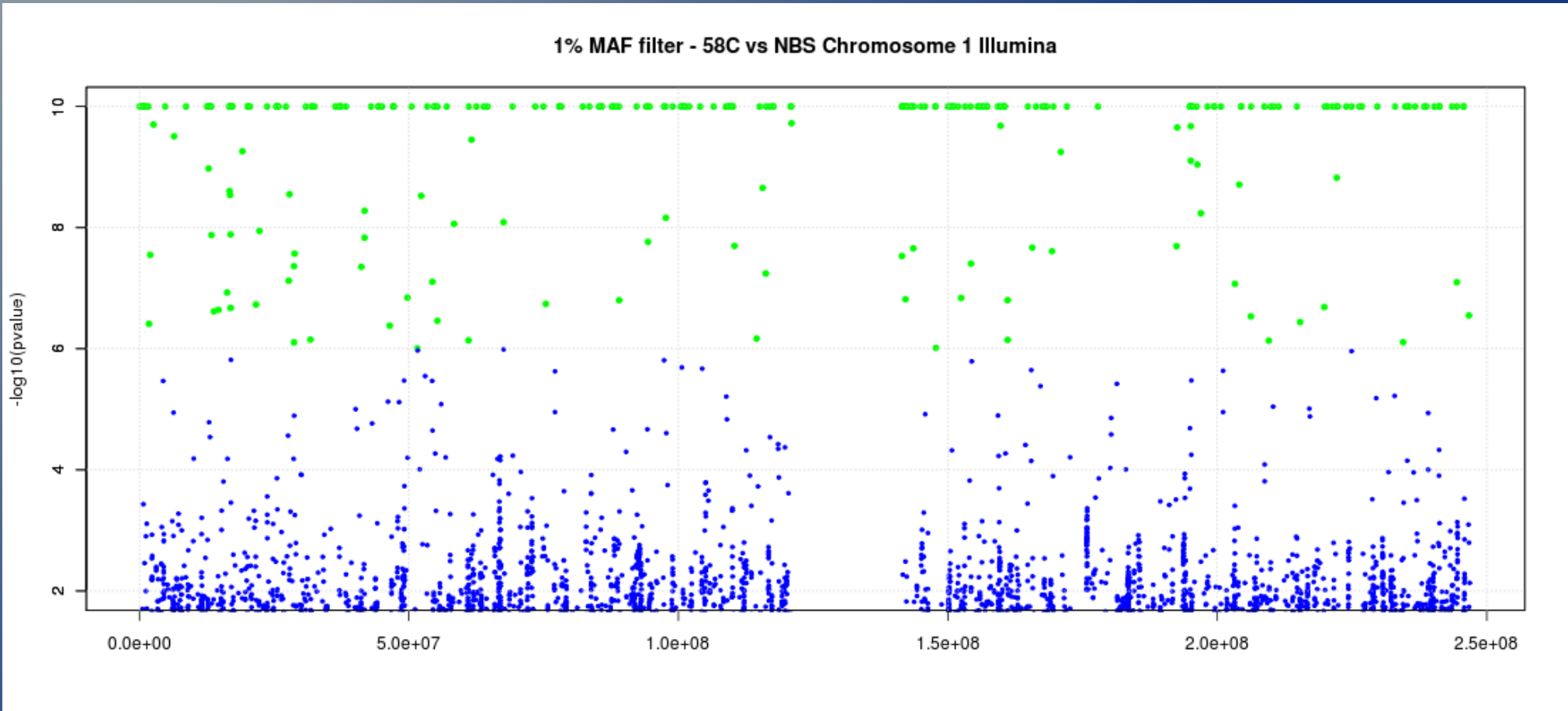


Testing for association



100% of SNPS

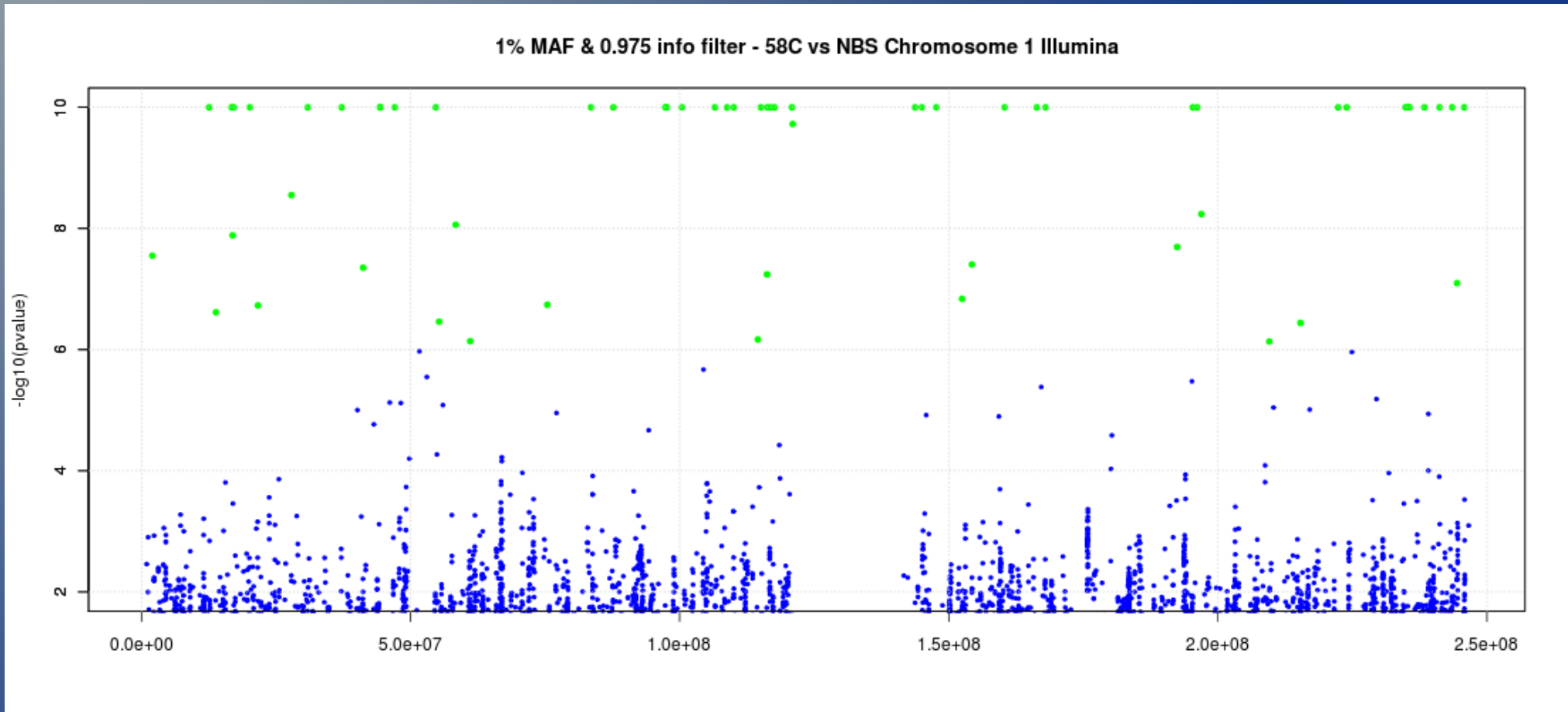
Testing for association



80.69% of SNPS

1% MAF

Testing for association

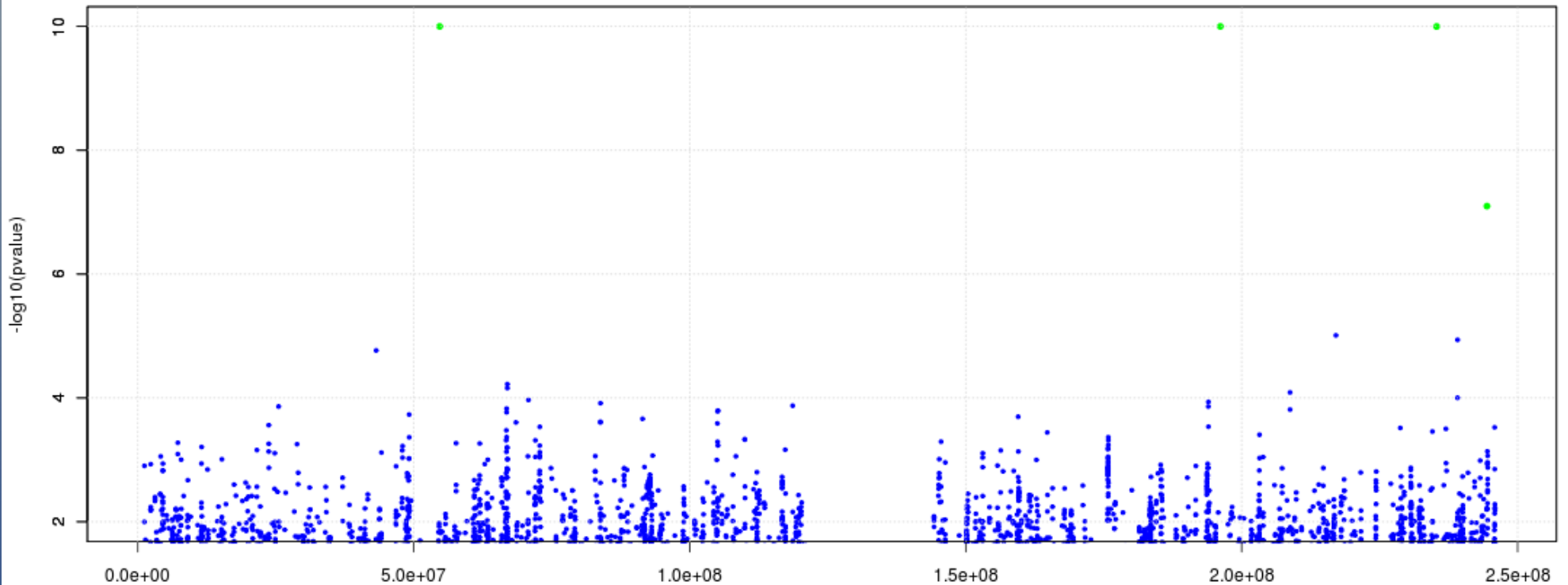


78.36% of SNPS

MAF > 1% & info > 0.975

Testing for association

1% MAF & 0.975 info & 1e-20 HWE filter - 58C vs NBS Chromosome 1 Illumina

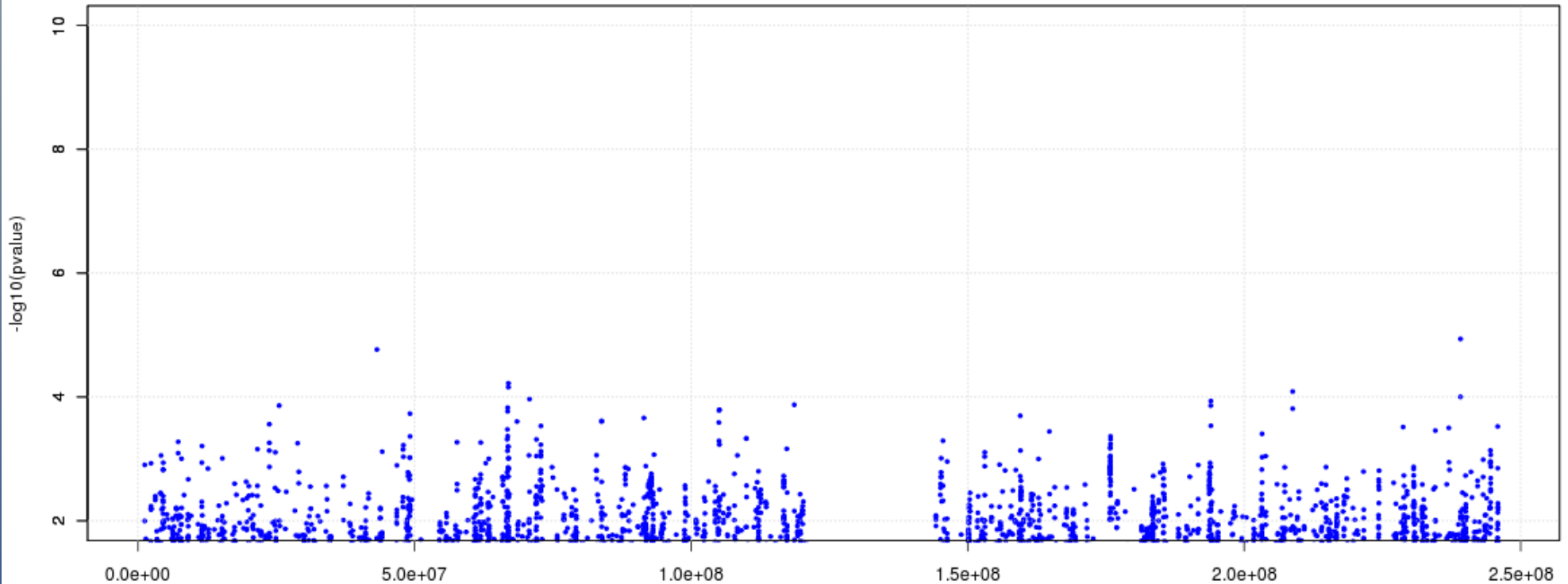


78.36% of SNPS

MAF > 1% & info > 0.975 & HW < 1e-20

Testing for association

1% MAF & 0.975 info & 1e-20 HWE & 2% missing filter - 58C vs NBS Chromosome 1 Illumina



77.92% of SNPS

MAF > 1% & info > 0.975 & HW < 1e-20 miss < 2%

A “clean” data set

Applying these filters results in a typically only 10 to 100 signals of association (e.g. $p < 1e-7$)

Comparing number of exclusions depends heavily on low MAF (which depends on population etc)

Important to apply filters to each set of collections or calls separately

Of those remaining...still lots of checking to do

Fingerprint
markers, gender
checks

Genotype QC, intensity
outlier

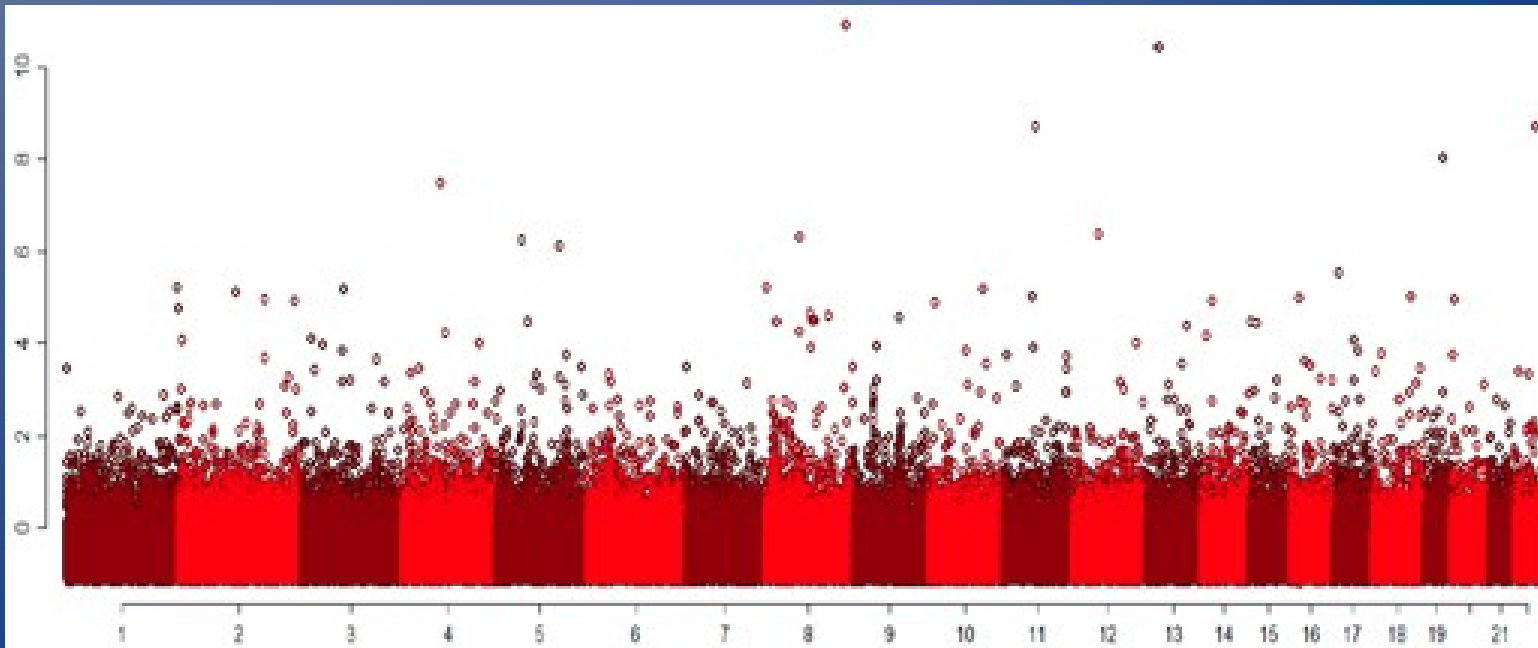
Given an individual's disease status, their genotype at a SNP, is independent of, and identically distributed to, other samples conditional on the underlying genotype frequencies in cases and controls

Duplicate checks,
relatedness

Population structure
analysis

A retracted paper

- Sebastiani et al. “Genetic signatures of exceptional longevity in humans” Science July 2010
 - Was retracted in July 2011 because QC had not been done properly!



Sebastiani
et al. 2010
Science