

# Introduction to **GWAS**

## Exploratory Data Analysis & Data Pre-Processing

Christian Werner

*(Quantitative geneticist and biostatistician)* **EiB**, **CIMMYT**, Texcoco (Mexico)

Filippo Biscarini

*(Biostatistician, bioinformatician and quantitative geneticist)* **CNR-IBBA**, Milan (Italy)



HerrFaloppio

Oscar González-Recio

*(Computational biologist and quantitative geneticist)* **INIA-UPM**, Madrid (Spain)



OscarGenomics



# Some basic data handling – **plink** (run in the shell)

Basic **plink** command structure:

***./plink --function specification***

Call program from path

Prefix for input files

dogs.ped and dogs.map are the basic input files

***./plink --dog --file dogs --recode vcf --out dogs***

Specify a nonhuman  
chromosome set:  
**--dog** = --chr-set 38

Recode .ped  
and .map file  
to .vcf file

Prefix for  
output files

# Some basic data handling – **plink** (run in the shell)

Basic **plink** command structure: ***./plink --function specification***

Problem with this command: `./plink --dog --file dogs --recode vcf --out dogs`

**POP1** **SAMPLE1** 0 0 0 -9 C C

**POP1** **SAMPLE2** 0 0 0 -9 C C

Vs.

#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NAME1 NAME2 NAME3....

`./plink --dog --file dogs --recode vcf-fid --out dogs`

`./plink --dog --file dogs --recode vcf-iid --out dogs`

'**vcf-fid**' and '**vcf-iid**' cause family IDs and within-family IDs respectively to be used for the sample IDs  
'**vcf**' merges both IDs and puts an underscore between them

# Some basic data handling – **plink** (run in the shell)

Basic **plink** command structure: ***./plink --function specification***

***When ped and map have different names:***

```
plink --dog --ped dogs.ped --map dogs.map --recode vcf --out dogs
```

***plink reads vcf too!***

```
./plink --vcf dogs.vcf --recode --out dogs
```

# Some basic data handling – **vcftools** (run in the shell)

Basic **plink** command structure: ***./vcftools --function specification***

```
./vcftools --vcf <path to vcf file> --plink --out <path to out file>
```

```
vcftools --vcf dogs.vcf --plink --out dogs_plink
```

(only biallelic markers will be in the output)

# **EDA: Exploratory Data Analysis**

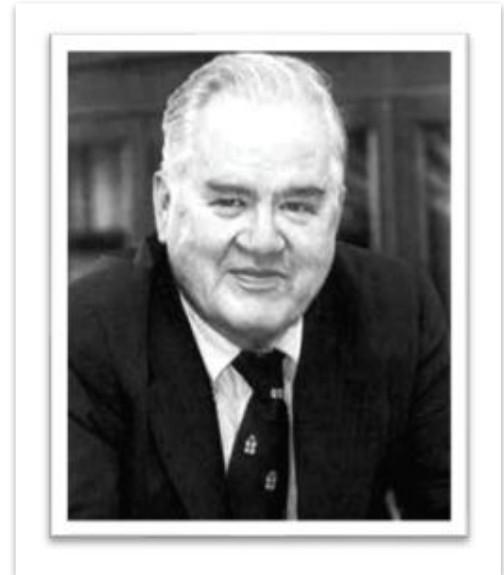
**(crap in, crap out...)**

# Exploratory data analysis – **EDA**

Exploratory data analysis was promoted by John Tukey to encourage statisticians to **explore the data**, and possibly formulate hypotheses that could lead to new data collection and experiments

**EDA is a fundamental step in all statistical and data analysis problems**

- approach to analyzing data sets to summarize their **main characteristics** often with visual methods.
- EDA is for seeing **what the data can tell us beyond the formal modeling or hypothesis testing task.**



# Exploratory data analysis – EDA

Before we conduct a GWAS, we have **two types of data** to explore

- **Genotypic data**
- **Phenotypic data**



# Exploratory data analysis – EDA

## Genotypic data

# Exploratory data – **genotypes**

Various metrics – statistics are performed either across SNP or across samples

Key concept: detect SNP and samples that should be removed prior to GWAS

Use of some metrics and thresholds is species- and population-specific

Still some level of subjectivity in the thresholds

# Exploratory data – **genotypes**

**Also referred to as Quality Control (QC)**

Some parameters to look at...

- Genotype calling and signal intensities (not covered here...)
- Marker allele frequencies
- Missing rate per marker and per individual
- Hardy-Weinberg equilibrium
- Heterozygosity

# Exploratory data – **genotypes**

## **Marker allele frequencies**

- Allele counts & genotype counts
- Minor allele frequency (MAF)
  - Some SNPs will be monomorphic
  - One of the alleles may be at very low frequency
  - Might be due to genotyping errors
  - Power to detect the association is very low

Common MAF thresholds are between 1 – 5%

In samples with known group structure, MAF should be checked within groups

# Exploratory data – **genotypes**

## **Missing rate per marker & per individual**

**SNPs might be of poor quality if their genotyping failed in many individuals**

- Should be investigated separately for all study groups (if known)
- Common thresholds are 2 – 5% (based on sample size & SNP number)

**Sample DNA might be of poor quality if there are many missing SNPs in an individual**

- Too many missing SNPs per individual can be an indication of poor DNA quality
- (or true deletions...)
- Common call rate thresholds are between 2 – 5% (based on sample size & SNP number)
- Includes monomorphic SNP !!

# Exploratory data – genotypes

Marker	Individuals												
	1	2	3	4	5	6	7	8	9	10	11	12	13
2	2	2	2	2	2	2	2	2	2	2	2	2	2
2	2	2	2	2	2	-1	2	2	2	2	2	2	2
2	-1	2	2	2	2	2	1	0	2	1	2	2	2
2	2	2	2	2	2	2	2	2	2	2	2	2	2
0	1	0	0	0	0	-1	0	0	0	0	0	0	0
0	0	0	0	0	0	-1	0	0	0	0	0	0	0
2	1	1	0	0	0	-1	0	0	1	0	0	0	1
0	2	1	1	2	2	2	1	1	2	0	1	1	2
1	0	1	0	0	0	0	0	0	0	0	0	0	0
1	2	1	1	2	2	2	1	2	2	1	2	1	1

# Exploratory data – genotypes

Monomorphic marker

Marker	Individuals												
	1	2	3	4	5	6	7	8	9	10	11	12	13
1	2	2	2	2	2	2	2	2	2	2	2	2	2
2	2	2	2	2	2	-1	2	2	2	2	2	2	2
3	2	-1	2	2	2	2	1	0	2	1	2	2	2
4	2	2	2	2	2	2	2	2	2	2	2	2	2
5	0	1	0	0	0	-1	0	0	0	0	0	0	0
6	0	0	0	0	0	-1	0	0	0	0	0	0	0
7	2	1	1	0	0	-1	0	0	1	0	0	0	1
8	0	2	1	1	2	2	1	1	2	0	1	1	2
9	1	0	1	0	0	0	0	0	0	0	0	0	0
10	1	2	1	1	2	2	1	2	2	1	2	1	1

# Exploratory data – genotypes

Marker	Individuals												
	1	2	3	4	5	6	7	8	9	10	11	12	13
	2	2	2	2	2	2	2	2	2	2	2	2	2
	2	2	2	2	2	-1	2	2	2	2	2	2	2
	2	-1	2	2	2	2	1	0	2	1	2	2	2
	2	2	2	2	2	2	2	2	2	2	2	2	2
	0	1	0	0	0	-1	0	0	0	0	0	0	0
	0	0	0	0	0	-1	0	0	0	0	0	0	0
	2	1	1	0	0	-1	0	0	1	0	0	0	1
	0	2	1	1	2	2	1	1	2	0	1	1	2
	1	0	1	0	0	0	0	0	0	0	0	0	0
	1	2	1	1	2	2	1	2	2	1	2	1	1

Monomorphic marker

Low MAF



# Exploratory data – genotypes

Monomorphic marker

Individuals

	2	2	2	2	2	2	2	2	2	2	2	2	2	2	
	2	2	2	2	2	-1	2	2	2	2	2	2	2	2	
	2	-1	2	2	2	2	1	0	2	1	2	2	2	2	
	2	2	2	2	2	2	2	2	2	2	2	2	2	2	
	0	1	0	0	0	-1	0	0	0	0	0	0	0	0	
	0	0	0	0	0	-1	0	0	0	0	0	0	0	0	
	2	1	1	0	0	-1	0	0	1	0	0	0	0	1	
	0	2	1	1	2	2	1	1	2	0	1	1	1	2	
	1	0	1	0	0	0	0	0	0	0	0	0	0	0	
	1	2	1	1	2	2	1	2	2	1	2	1	1	1	

Low MAF

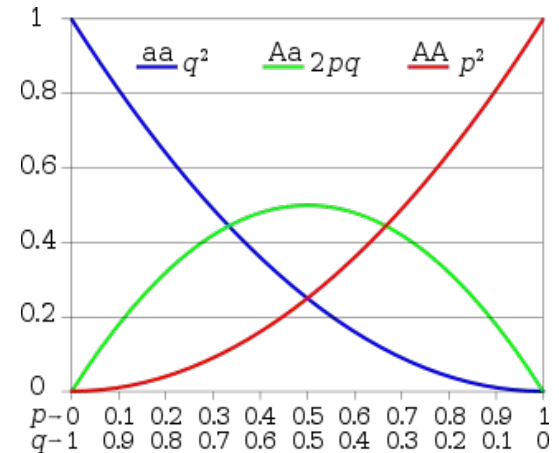
Many missing markers

# Exploratory data – genotypes

Relationship between allele and genotype frequencies

$$(p + q)^2 = p^2 + 2pq + q^2$$

- diploid genomes
- autosomal loci
- large population
- random mating
- equal frequencies in both sexes
- no selection
- no migration
- no mutations



# Exploratory data – **genotypes**

## Deviations from HW equilibrium

- systematic genotyping errors
- violation of assumptions

## Test for HW deviation

- chi-squared ( $\chi^2$ ) test
- Fisher's (exact) test
- many, many more ...

**But also selection, assortative mating, population structure and inbreeding cause deviations from HWE!**

**HWE is, in most cases, NOT a reasonable assumption...**

# Exploratory data – **genotypes**

## **Heterozygosity**

- Proportion of heterozygotes
- Heterozygosity can be checked per locus & per marker

### **Very high sample heterozygosity can be an indication of DNA contamination**

- But also could be that a small proportion of samples are truly very different from the rest...
- Removal of samples that depart  $\pm 3$  SD from the mean

### **Very high heterozygosity per marker could also indicate poor DNA quality, but also be due to...**

- the breeding scheme (e.g. hybrid breeding in plants, or very low heterozygosity in lines)
- Genome duplications

# Exploratory data analysis – EDA

## Phenotypic data

# Exploratory data – phenotypes

## Data type

- Continuous (e.g. height)
- Binary (e.g. case/control)
- Categorical (e.g. scores (ordered), eye colour (ordered))

**Measure of centrality:** mean, mode, median

**Measures of dispersion:** range, variance, standard deviation

## Distribution of the data

- Distribution of values as expected? Outliers?
- representative sample of the population?
- Other explanatory covariables

# Exploratory data – phenotypes

## Covariables

Are there any variables which may have a relationship with the phenotype?

- E.g. sex, breed, age, treatments, year effects, ... (population structure)

**The data needs to be corrected for these effects.** Otherwise they can be confounded with allelic variants with an effect on the phenotype which we try to identify.

Covariables with significant effects on the phenotype can be identified using ANOVA (requires balanced datasets - ANOVA is outdated) or **linear mixed models**. However, a comprehensive preparation of phenotypic data including model comparison is not covered here...

# Exploratory data – **phenotypes**

## **Assumptions for continuous variables** (different for binary traits...)

- Normally distributed **residuals** (prerequisite of GWAS model assumptions)
- Homogeneity of variance (differences in variance might indicate a factor that has not been included in the phenotype processing)



# Exploratory data – phenotypes

## Outliers

- Apparently rare phenotypes are often a result of errors or poor models rather than true outliers
- However, the values might be real - outliers should be investigated thoroughly rather than relying on statistical tests

## Data transformation

- **Positively skewed distributions of the residual** with the long tail into the positive direction can be corrected with a logarithmic or square root transformation.
- **Negatively skewed distributions of the residual** that have a long tail in the negative direction can be corrected with cubing or squaring
- Transformations only if really necessary...

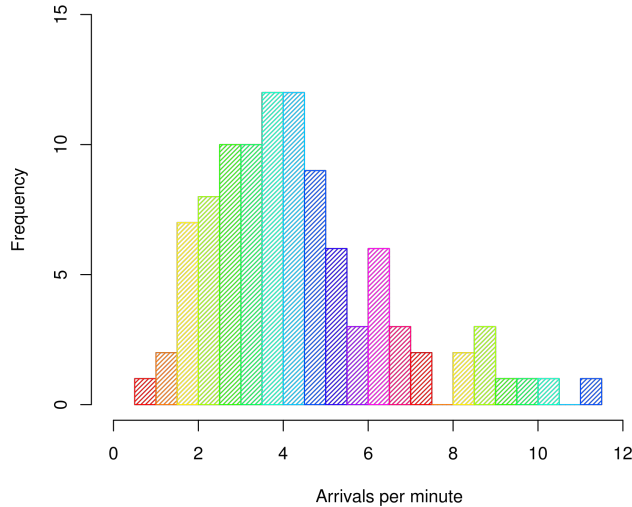
# Exploratory data – phenotypes

Histograms,

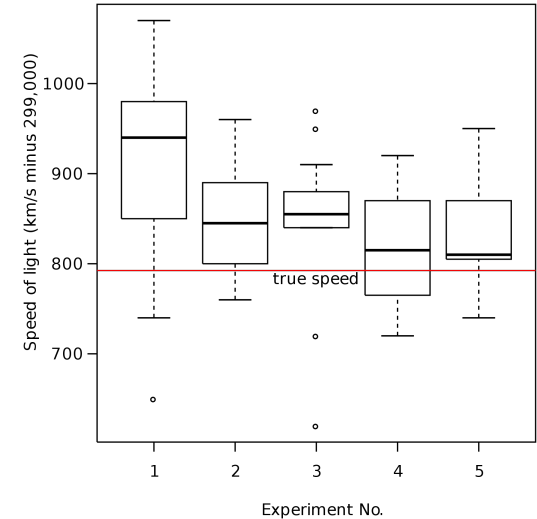
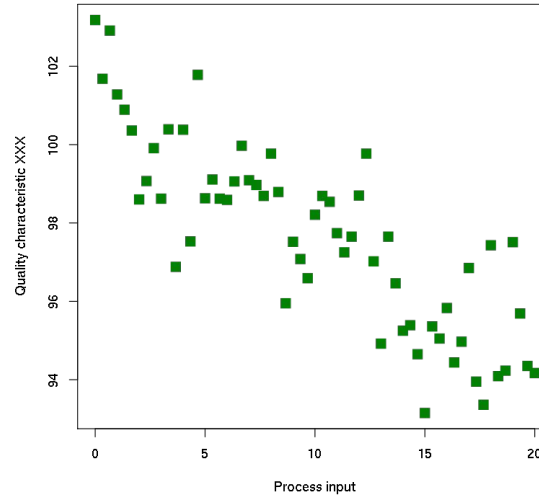
Scatter plots,

Box plots,

Histogram of arrivals



Scatterplot for quality characteristic XXX



# Exploratory data – phenotypes

## Normality of residuals

- Look at the data (Histogram, QQ-Plot)
- Don't rely on Shapiro-Wilk test

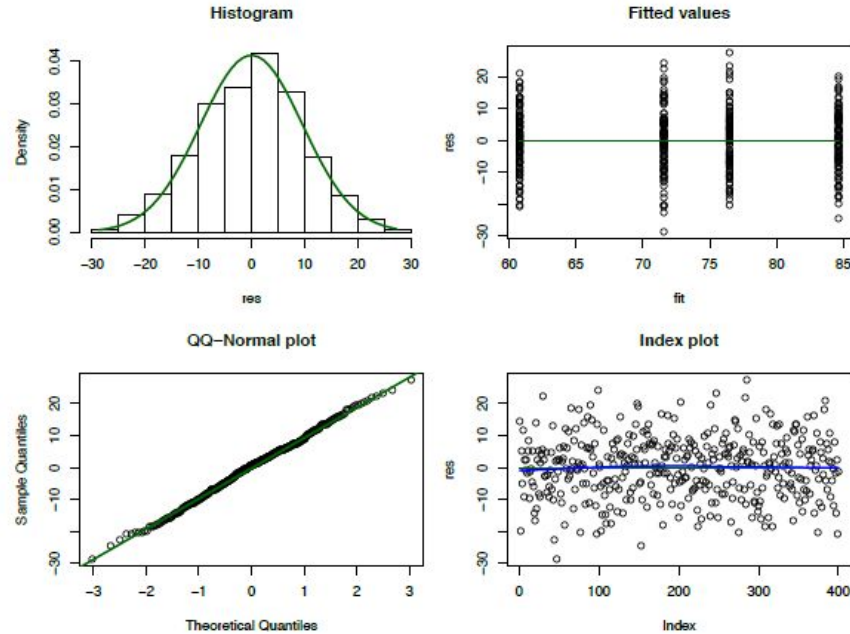
## Variance homogeneity

- Look at the data (Scatterplot)
- Don't rely on Leven's test

Tests might be very conservative and might indicate a violation of the assumptions of normality and variance homogeneity in a suitable “real-world” dataset.

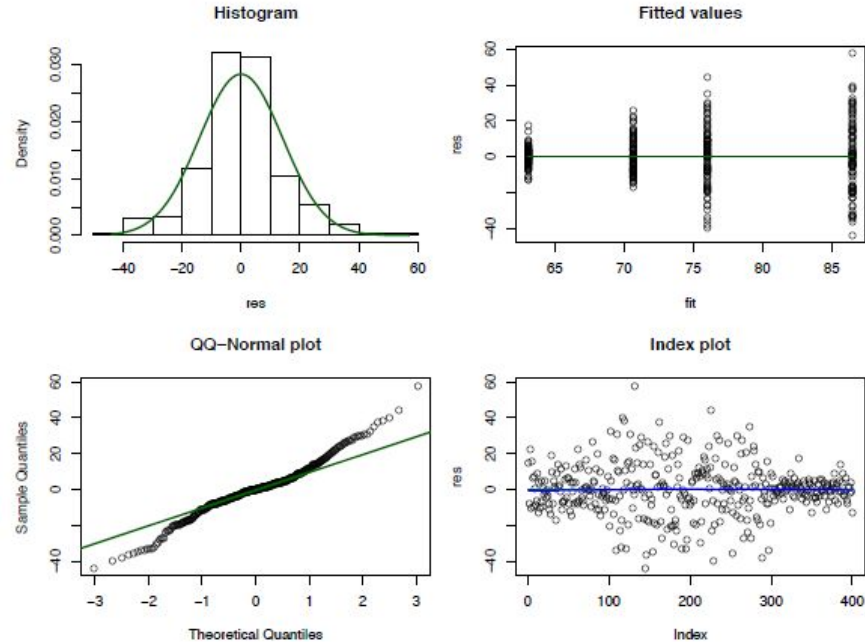
# Exploratory data – phenotypes

Almost perfectly distributed residuals



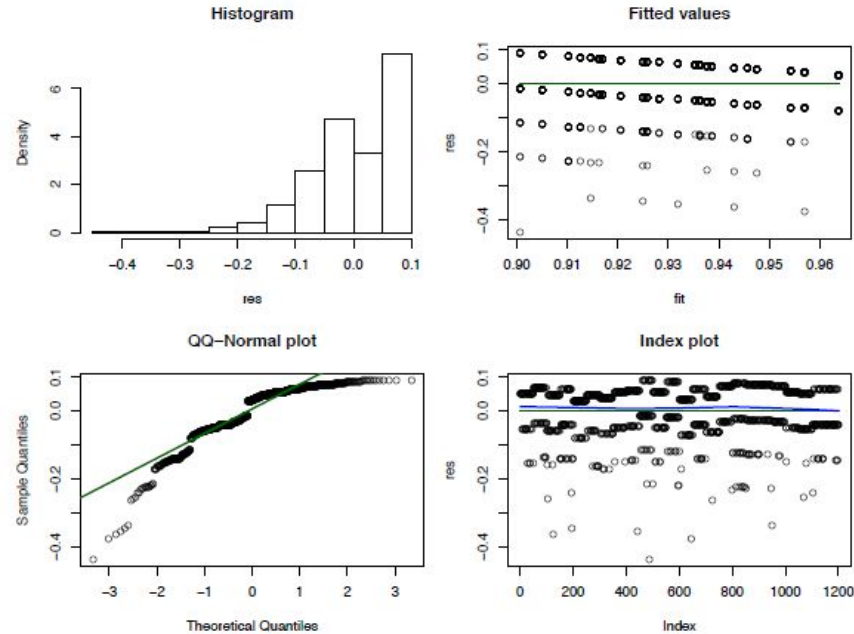
# Exploratory data – phenotypes

Still good...



# Exploratory data – phenotypes

There might be something wrong here....



# Summary of pre-processing for genotype data

# Pre-processing genotypic data – **Standard procedures**

## **missing rate**

- per-sample (e.g. max 10% missing SNP genotypes per sample)
- per-site (e.g. max 5% missing genotype per variant/SNP)
- stricter/looser thresholds depending on data/experiment (e.g. SNP array, GBS, quality of reference sequence, ...)

MAF

Hardy Weinberg equilibrium

others



# Pre-processing genotypic data – Standard procedures

missing rate

## MAF

- remove monomorphic variants → non-informative
- remove variants at low frequency (“rare”) → spurious associations
- threshold depends on sample size
- usually (re)done after imputation

Hardy Weinberg equilibrium

others

# Pre-processing genotypic data – **Standard procedures**

missing rate

MAF

## **Hardy Weinberg equilibrium**

- set **low threshold** for p-value (e.g.  $\exp(-10)$ )
- **questionable**: some of the forces driving out of HW equilibrium are what we are usually after (selection, migration, mutation, adaptation etc.)

others

# Pre-processing genotypic data – Standard procedures

missing rate

MAF

Hardy Weinberg equilibrium

**others**

- sex chromosomes (might need to be removed / analyzed separately)
- Mendelian errors
- quality scores (vcf files)
- relatedness (between samples - check for duplicates)