

**HelmholtzZentrum münchen**

German Research Center for Environmental Health

# Lecture 3

Statistics for Human Genetics

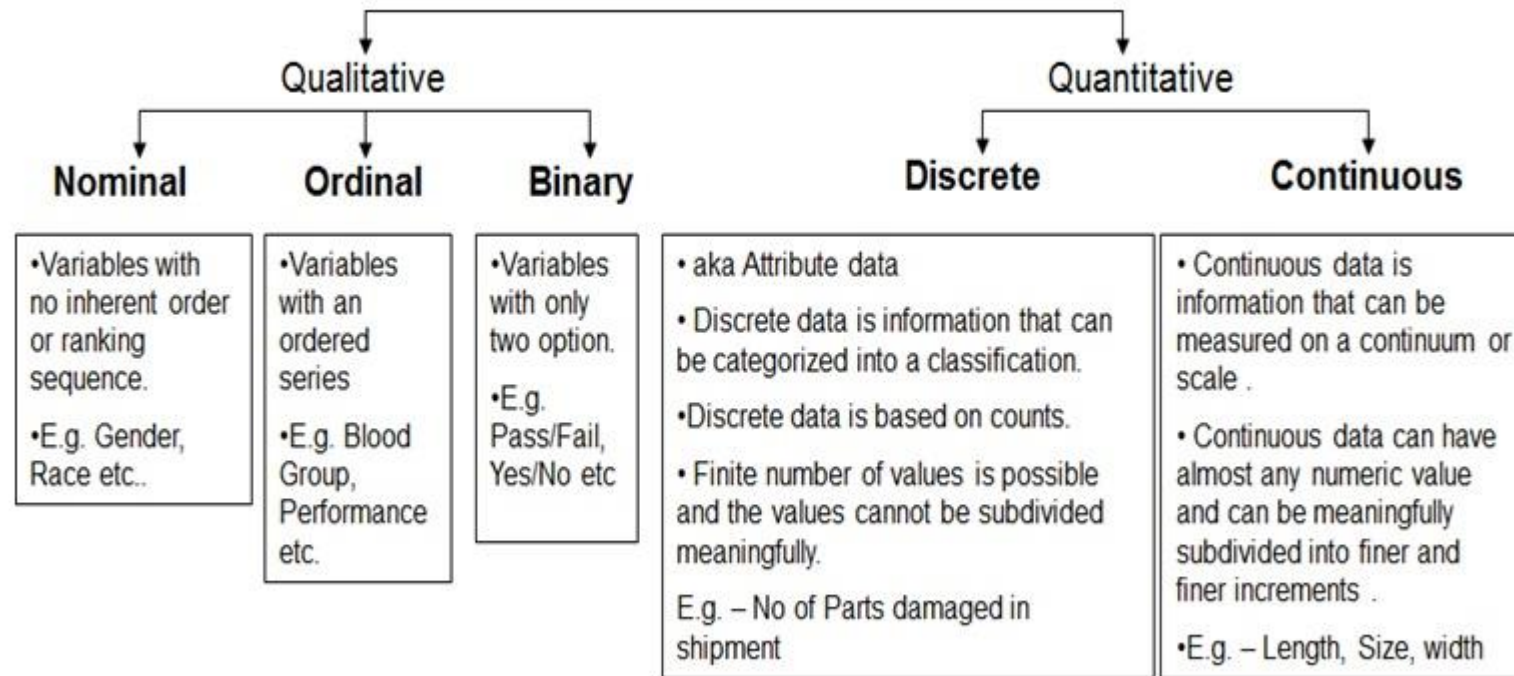
**HELMHOLTZ** RESEARCH FOR  
GRAND CHALLENGES

# What can we do with statistics?

- Estimation
- Hypothesis testing
- Modelling
- Predicting

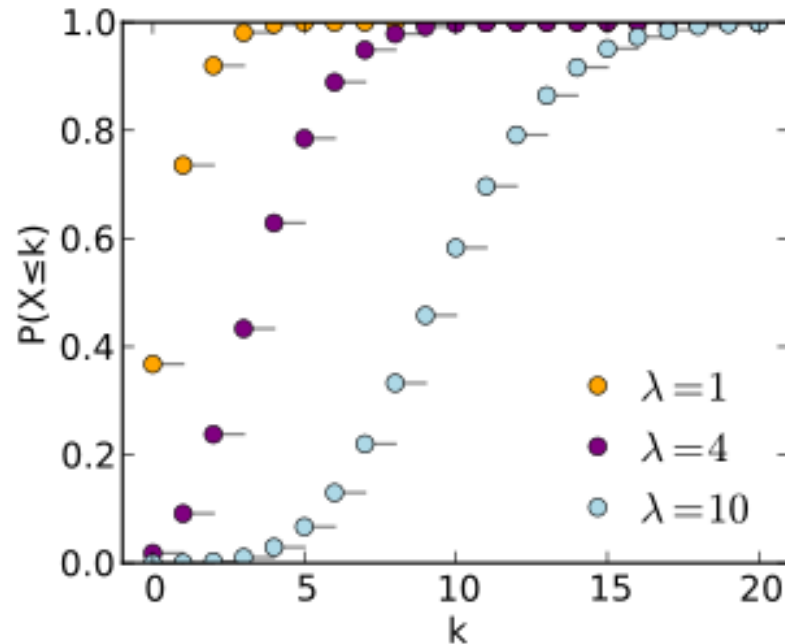
# Random Variables

- In statistics, we measure realizations of random variables
- Often, random variables follow a distribution
- They can be qualitative or quantitative, continuous or discrete



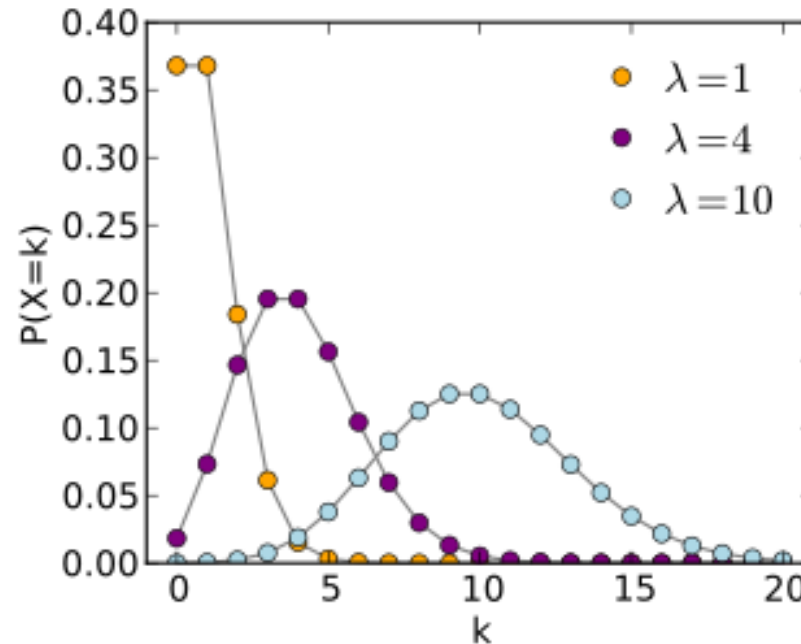
# Distributions

- Two ways to represent them :



## ***Cumulative distribution function (CDF)***

- $y = p(X \leq x)$
- Always growing
- Ideal way to represent, but hard to read
- All distributions look the same

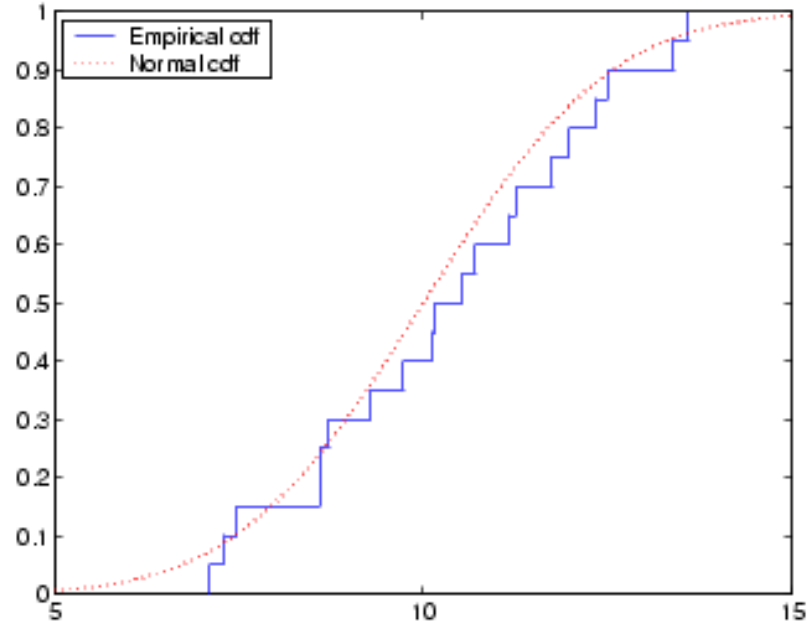


## ***Probability density function (PDF)***

- $y = p(X = x)$  for discrete
- Shows how values are distributed
- Nice visually, but mathematically hard to deal with

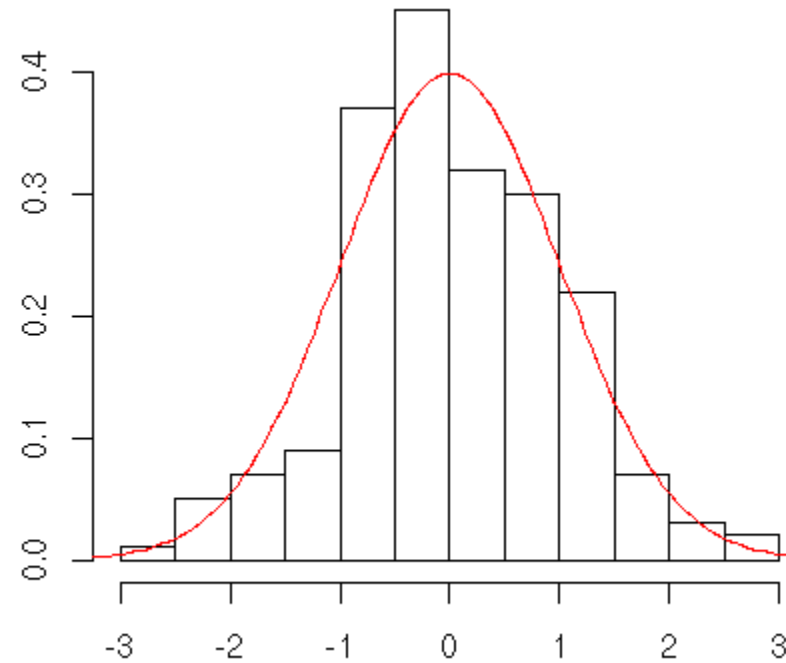
# Distributions

- How to estimate them:



## ***Empirical CDF***

- Rarely used
- Except when you want to compute empirical quantile functions



## ***Barplot (discrete)***

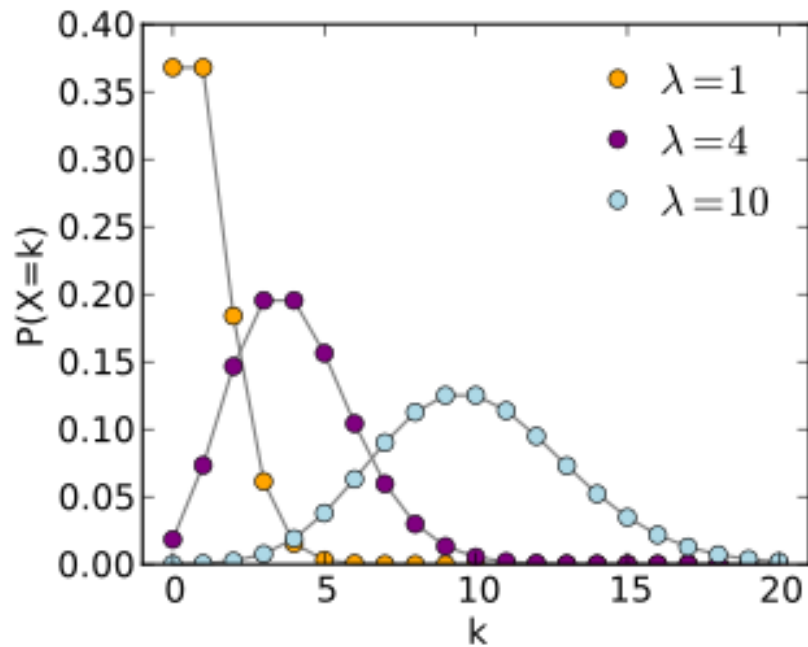
- For every value, count occurrences

## ***Histogram (continuous)***

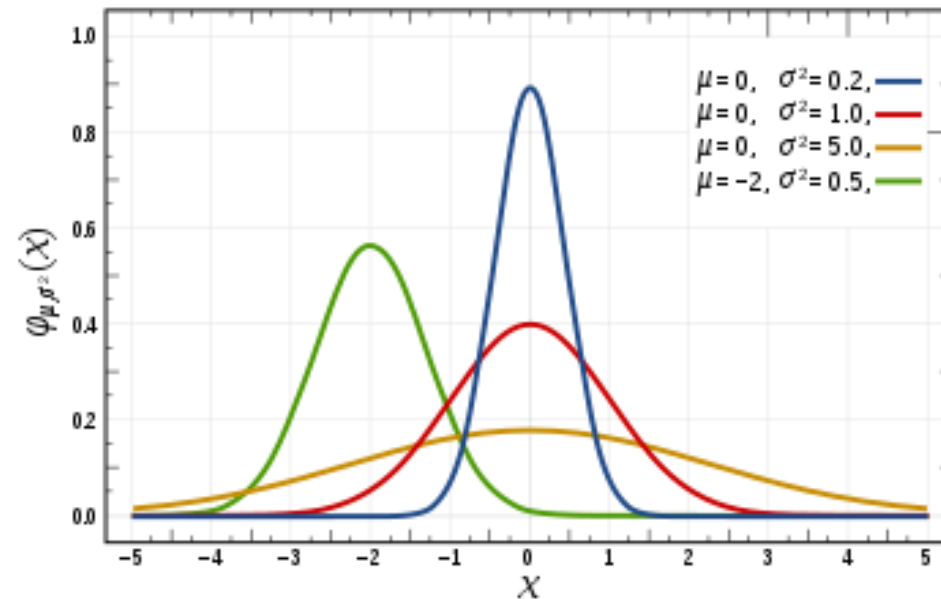
- Cut the interval into bins, count observations within bin

# Distributions

- Two broad types:
  - those followed by random variables (real world data)



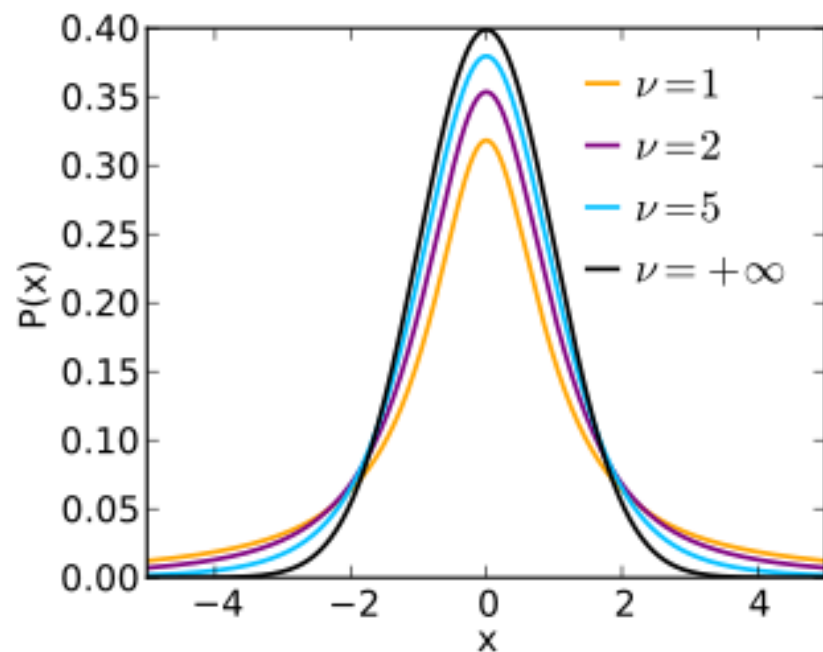
$$X \sim \text{Poisson}(\lambda)$$



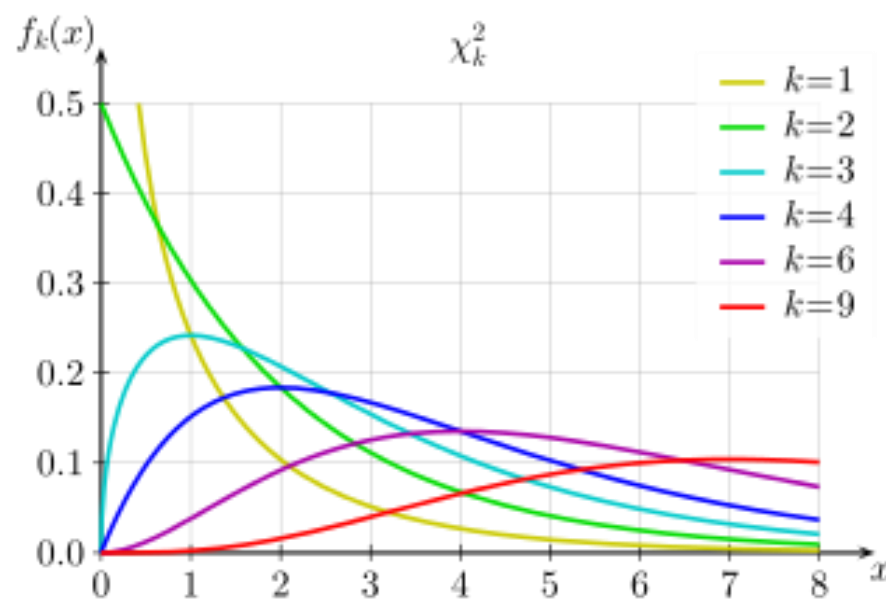
$$X \sim \mathcal{N}(\mu, \sigma^2)$$

# Distributions

- Two broad types:
  - those followed by test statistics



$$X \sim T(\nu)$$



$$X \sim \chi^2(k)$$

$\lambda, \mu, \sigma, \nu$  and  $k$  are ideal parameters. How to estimate them?

# Statistics

- A statistic is a meaningful quantity derived from the data
- Often, estimators are realization of distribution parameters
- Examples? Mean, proportion
- For simple distributions/parameters, there is a formula
- For more complex ones, we have to use other techniques (Monte-Carlo, Permutations...)

$$\hat{p} = \frac{x}{n}$$

$$(\hat{\mu} =) \bar{x} = \frac{1}{n} \sum_{k=0}^n x_k$$

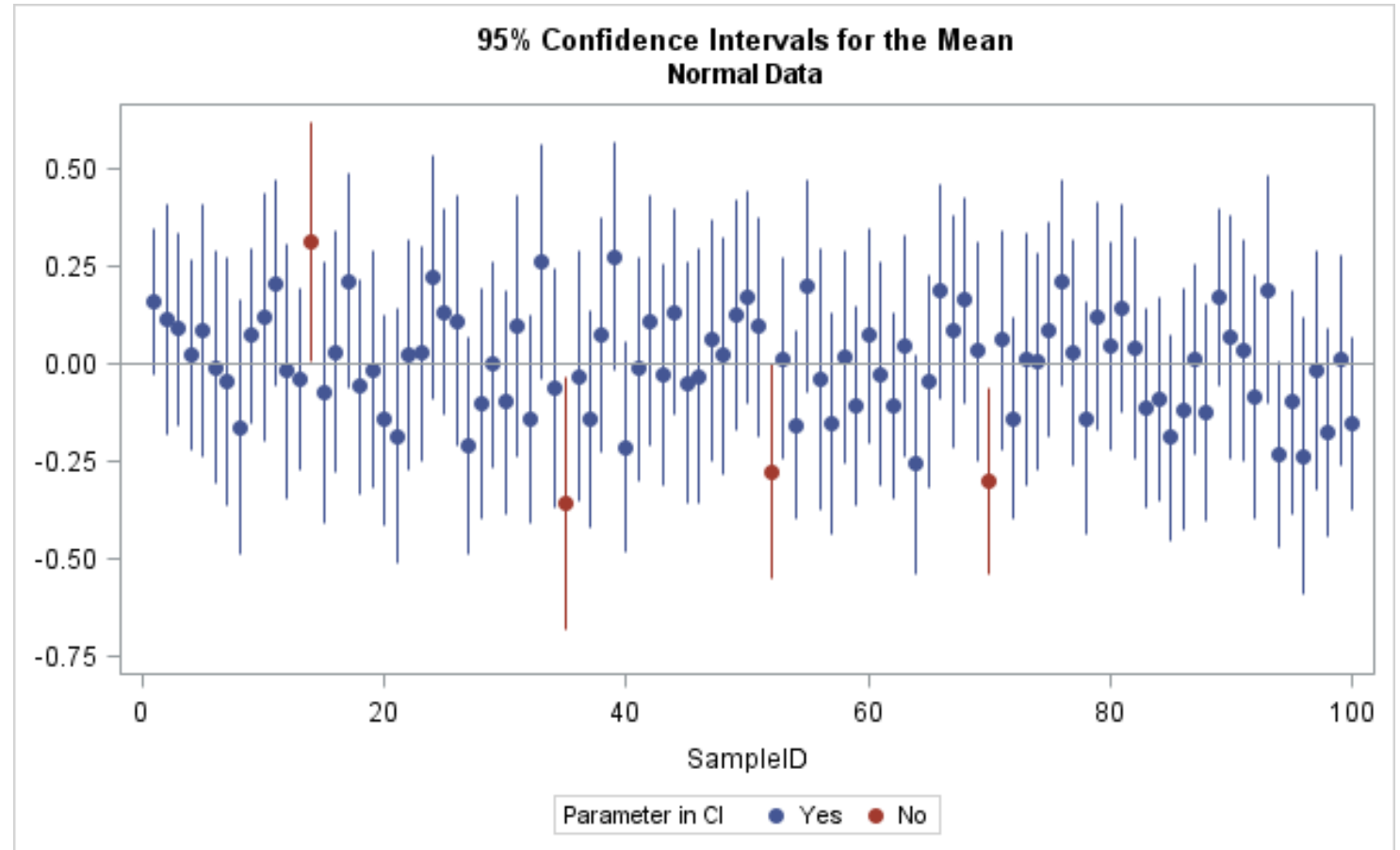
$$w = \frac{(\hat{\theta} - \theta_0)^2}{se(\hat{\theta})} \sim \mathcal{N}(0,1)$$

$$(\widehat{\sigma^2} =) s^2 = \frac{1}{N-1}$$



# One particular statistic: Confidence intervals

- $x\%$  confidence interval ( $x\%C.I.$ ) :  $x\%$  of the time when this interval is calculated, it will contain the true value of the parameter



# Hypothesis testing

- We want to measure whether the data gives sufficient evidence to reject a hypothesis
- Null/Alternative hypothesis ( $\mathcal{H}_0/\mathcal{H}_A$ )
- We prove that we can produce a statistic that follows a certain distribution if the null hypothesis is true = name of the test
- We calculate the statistic based on our data
- Because we know the distribution, we can calculate the CDF  $p(X \leq x)$
- = how unlikely it is that our measurement comes from the null : p-value
- Example: proportion test, t-test, chi-squared test...

**Summary statistic**  
(helps distinguish  $H_0$  and  $H_A$ )



**Test statistic**  
(standard distribution with no unknown parameters under  $H_0$ )

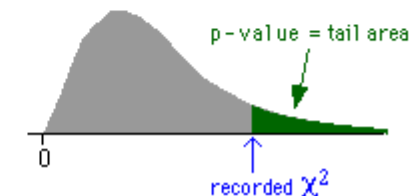


**P-value**  
(probability of more 'extreme' test statistic)

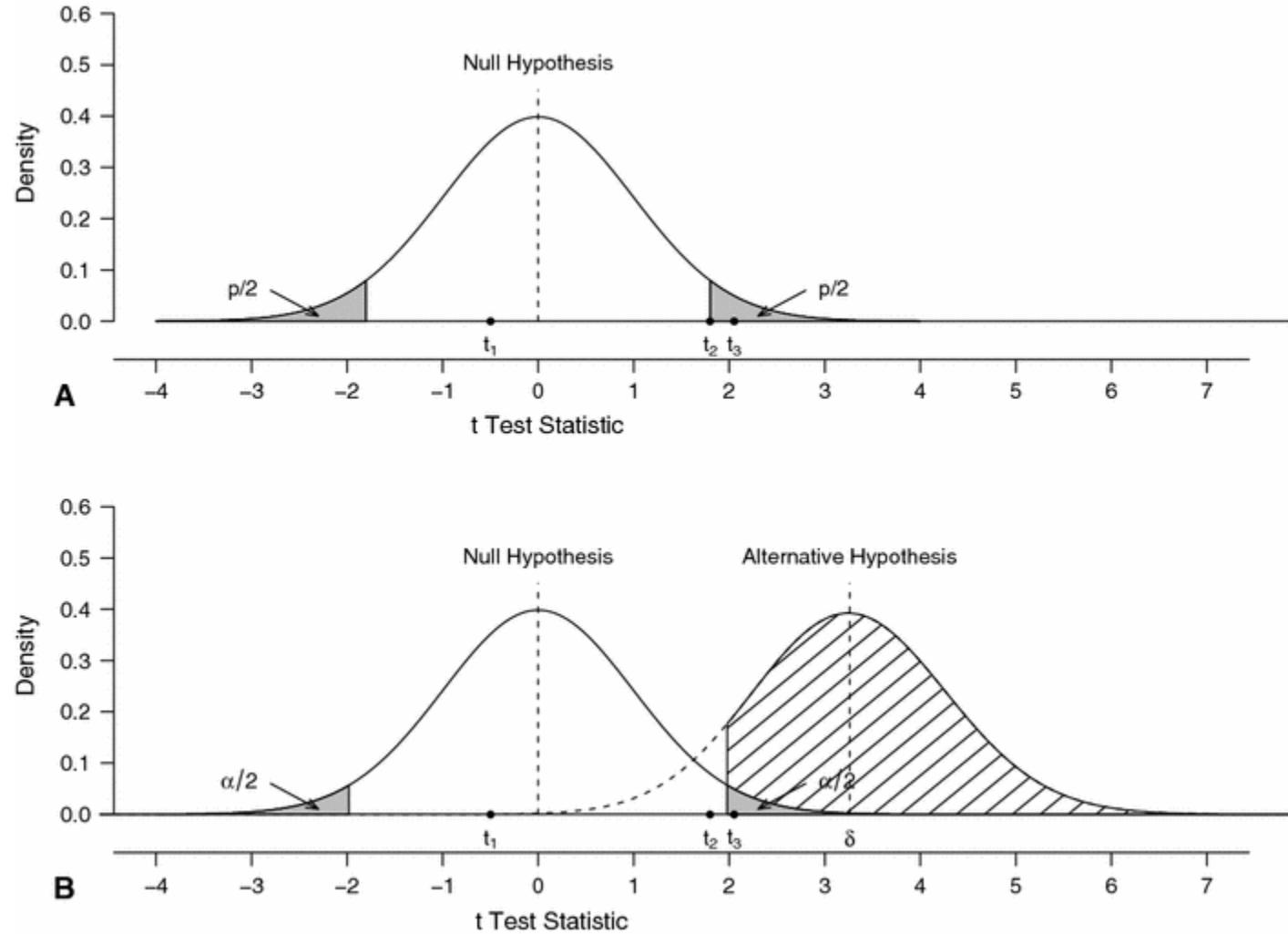
$$\chi^2 = \sum \frac{(n_{xy} - e_{xy})^2}{e_{xy}}$$



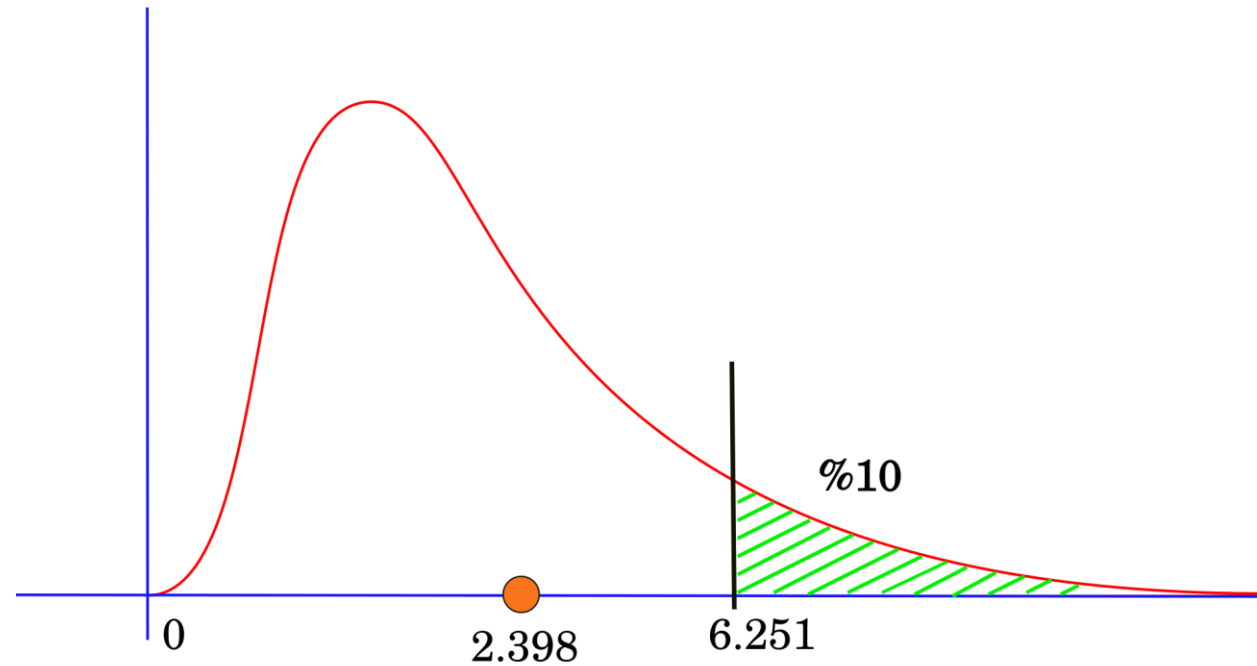
$\chi^2 \sim \text{chi-squared } ((r-1)(c-1) \text{ df})$



# Hypothesis testing



# Two tails or one



# Exercise I : Proportion test

- In a population, we observe 41,009 potentially damaging variants among 14,281,180 variants
- What is the proportion? In a very large reference population, we observe a proportion of  $1.52 \times 10^{-3}$ . Is it significantly different? (prop.test, binom.test)

# Multiple testing

## Family-wise (FWER)

- Bonferroni correction
- Simple to implement, harder to interpret

$$p_{critical} = \frac{0.05}{m}$$

- *“If all tests are under the null, probability that **one or more** of them is a false positive.”*

## False-discovery based (FDR)

- Benjamini-Hochberg procedure
- Harder to implement, easy to understand

$$p_{critical} = \operatorname{argmax}(p < \frac{i}{m} Q)$$

- $i$ =rank,  $Q$ =FDR.
- *“Proportion of significant tests that are false positives.”*

# Multiple testing

## Family-wise (FWER)

- Bonferroni correction
- Simple to implement, harder to interpret

$$p_{critical} = \frac{0.05}{m}$$

- *“If all tests are under the null, probability that **one or more** of them is a false positive.”*

## False-discovery based (FDR)

- Benjamini-Hochberg procedure
- Harder to implement, easy to understand

$$p_{critical} = \operatorname{argmax}(p < \frac{i}{m} Q)$$

- $i$ =rank,  $Q$ =FDR.
- *“Proportion of significant tests that are false positives.”*

### When to use which depends on

- 1) best practices
- 2) relative cost of a false negative/positive

# Checking results

- Statistical significance:
  - One test:  $p < 0.05$
  - Genome-wide: one test per variant and per phenotype
  - But all variants are not independent, in reality, we account for LD
  - $5 \times 10^{-8}$  for GWAS,  $10^{-9}$  for sequencing-based



## Exercise II

- If the adjusted genome-wide significance threshold is  $5 \times 10^{-8}$  for GWAS, how many “effective” variants are there in a genotyped human genome?
- You are writing an article about a GWAS for 16 different traits. What will be your threshold for declaring significance?

# Modelling and predicting

- If we estimate the effect of one variable on another variable, we do modelling
- When we apply this effect to new observations of the variable, we do prediction
- Process is called machine learning, predictive modelling or predictive analysis
- In human genetics, main task is to model effect of genotypes on phenotypes

$$\textit{phenotype} \sim \beta \times \textit{genotype} + \epsilon$$

$$\begin{bmatrix} \textit{pheno}_0 \\ \vdots \\ \textit{pheno}_n \end{bmatrix}$$

$$\begin{bmatrix} A/T \\ \vdots \\ T/T \end{bmatrix}$$

# Modelling and predicting

- If we estimate the effect of one variable on another variable, we do modelling
- When we apply this effect to new observations of the variable, we do prediction
- Process is called machine learning, predictive modelling or predictive analysis
- In human genetics, main task is to model effect of genotypes on phenotypes

$$\textit{phenotype} \sim \beta \times \textit{genotype} + \epsilon$$

$$\begin{bmatrix} \textit{pheno}_0 \\ \vdots \\ \textit{pheno}_n \end{bmatrix} \quad \begin{bmatrix} A/T \\ \vdots \\ T/T \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 2 \end{bmatrix}$$

$= \{0,1\}$  (case-control)  
 $\in \mathbb{R}$  (quantitative)  $\sim \mathcal{N}(0,1)$

$= \{0,1,2\}$  (genotype, directly typed)  
 $\in [0,2]$  (dosage, imputed)

$$\begin{bmatrix} 0.965 \\ \vdots \\ 1.816 \end{bmatrix}$$

# Modelling and predicting

- If we estimate the effect of one variable on another variable, we do modelling
- When we apply this effect to new observations of the variable, we do prediction
- Process is called machine learning, predictive modelling or predictive analysis
- In human genetics, main task is to model effect of genotypes on phenotypes

Usually, we do not predict (except PRS)

$$\textit{phenotype} \sim \beta \times \textit{genotype} + \epsilon$$

$$\begin{bmatrix} \textit{pheno}_0 \\ \vdots \\ \textit{pheno}_n \end{bmatrix} \quad \begin{bmatrix} A/T \\ \vdots \\ T/T \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 2 \end{bmatrix}$$

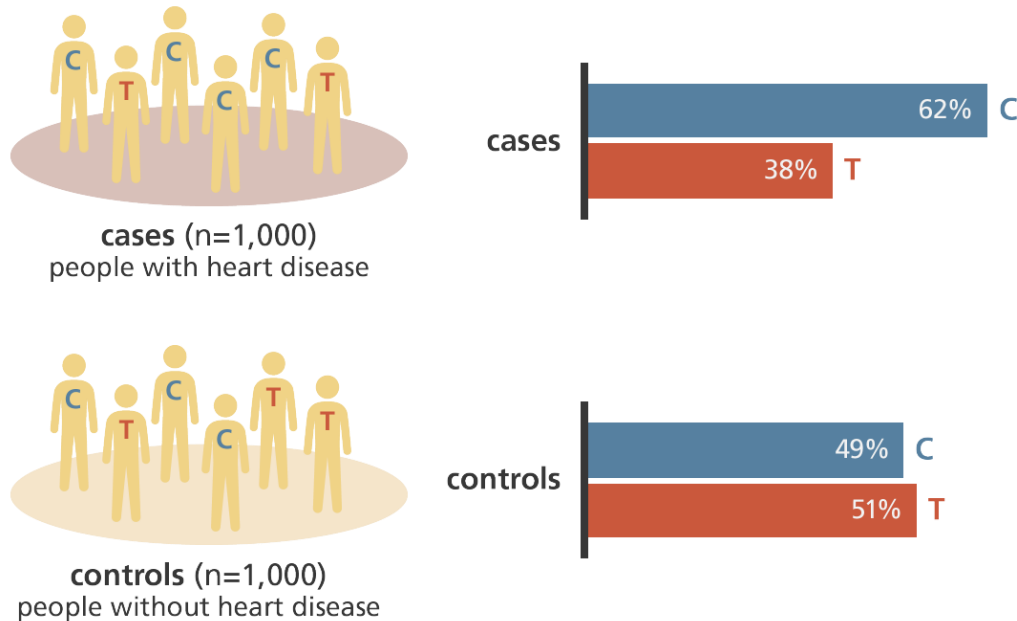
$= \{0,1\}$  (case-control)  
 $\in \mathbb{R}$  (quantitative)  $\sim \mathcal{N}(0,1)$

$= \{0,1,2\}$  (genotype, directly typed)  
 $\in [0,2]$  (dosage, imputed)

$$\begin{bmatrix} 0.965 \\ \vdots \\ 1.816 \end{bmatrix}$$

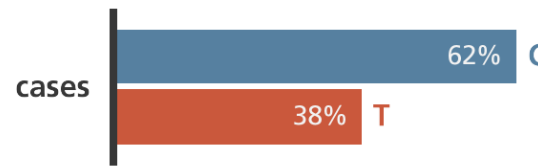
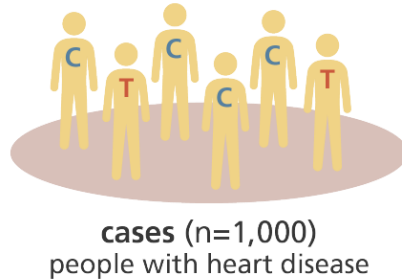
# Case/control

- Estimated effect: odds ratio (OR)  
“how much more likely are you to be a case if you carry the risk allele?”  
per genotype, calculate the odds  $O = \frac{p}{1-p}$

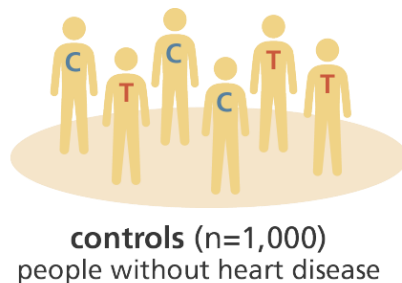


# Case/control

- Estimated effect: odds ratio (OR)  
“how much more likely are you to be a case if you carry the risk allele?”  
per genotype  $g$  and for a disease  $Y$ , calculate the odds  $O = \frac{p_{Y=1|g}}{1-p_{Y=1|g}}$



	case s	controls
T	380	490
C	620	510



$$O_T = \frac{380/n_T}{490/n_T} \quad O_C = \frac{620/n_C}{510/n_C}$$

$$OR_{C/T} = \frac{620 \times 490}{510 \times 380} = 1.56$$

# Case/control

## Dominant

Marker allele	Affected	Unaffected
DD+Dd	$n_{2A} + n_{1A}$	$n_{2U} + n_{1U}$
dd	$n_{0A}$	$n_{0U}$

$$OR = \frac{n_{affected\ carriers} \times n_{healthy\ non-carriers}}{n_{healthy\ carriers} \times n_{affected\ non-carriers}}$$

## Recessive

Marker allele	Affected	Unaffected
DD	$n_{2A}$	$n_{2U}$
Dd+dd	$n_{1A} + n_{0A}$	$n_{1U} + n_{0U}$

## Additive

Marker genotype	Affected	Unaffected
DD	$n_{2A}$	$n_{2U}$
Dd	$n_{1A}$	$n_{1U}$
dd	$n_{0A}$	$n_{0U}$

$$OR = \frac{(2 \times n_{2A} + n_{1A}) \times (2 \times n_{0U} + n_{1U})}{(2 \times n_{0A} + n_{1A}) \times (2 \times n_{2U} + n_{1U})}$$

Allelic odds-ratio

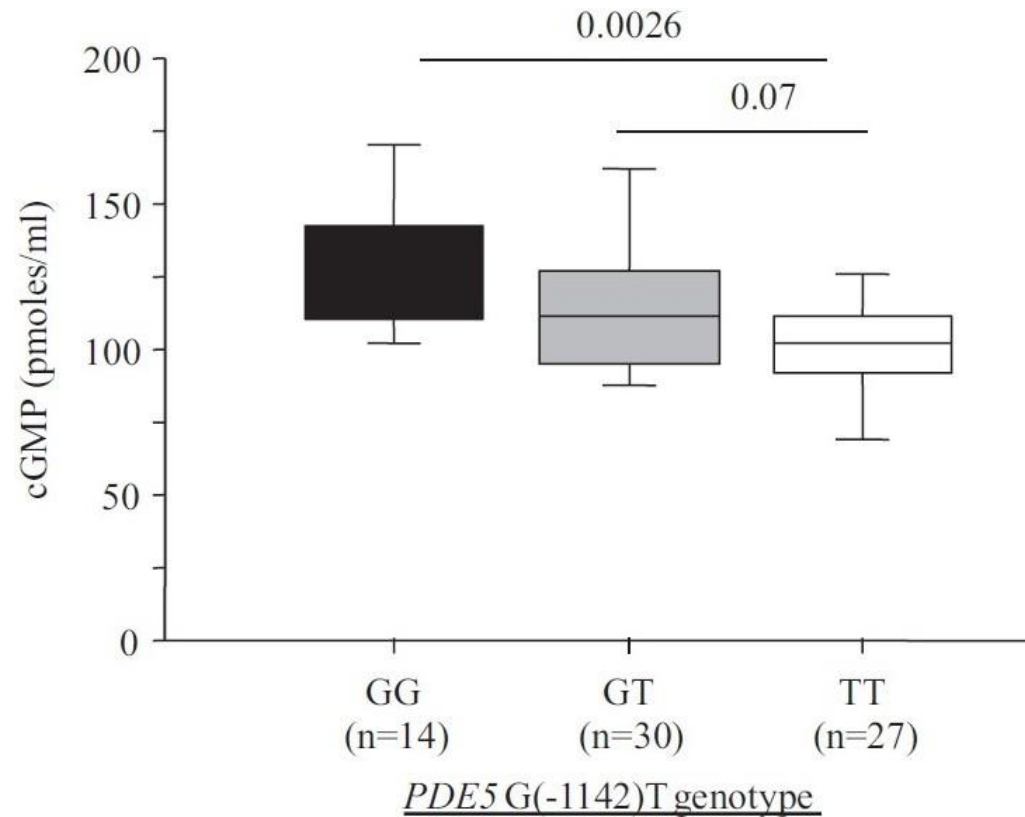
# Case/control

- Output: OR and 95% confidence interval of the OR
- Test: is it significantly different from 1?
- Tests: Fisher's exact test or Chi-squared
- In case of dosages or covariates: logistic regression

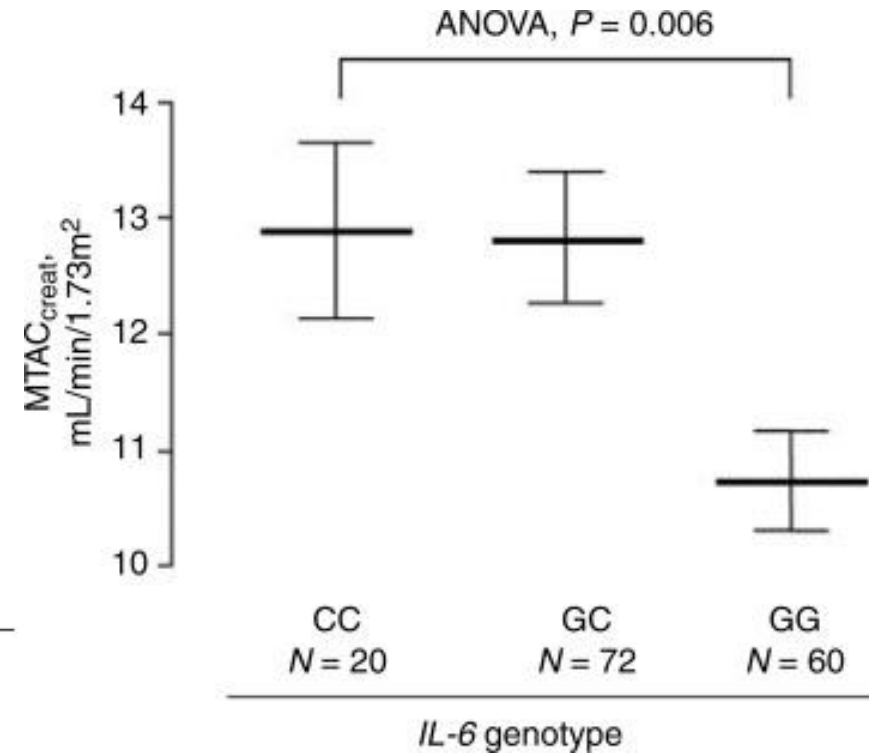


# Continuous trait

- For directly typed (0,1,2): ANOVA



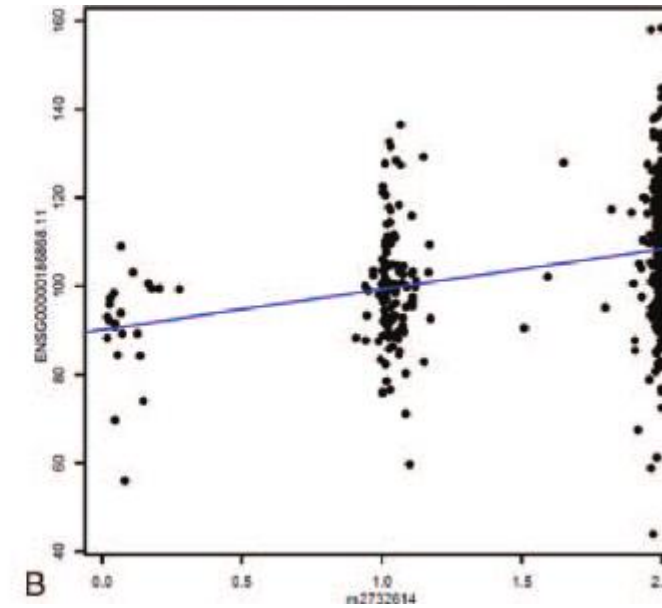
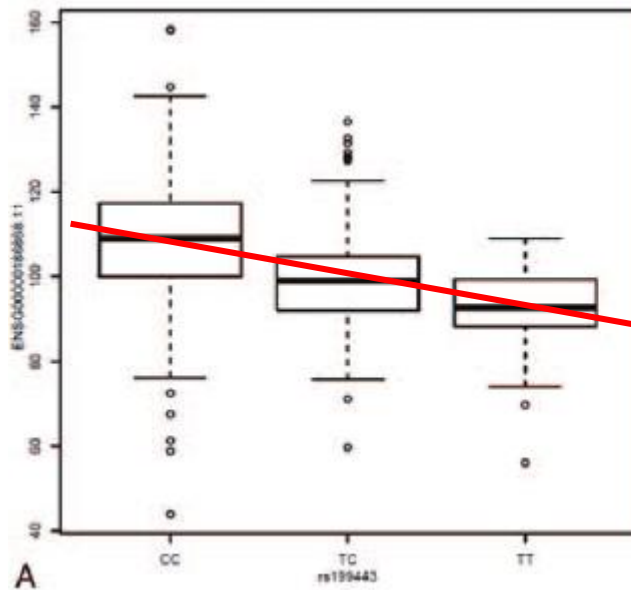
additive



recessive

# Continuous trait

- For dosages (imputed quantity of minor allele  $d \in [0,1]$ ) : linear regression
- In general: generalized linear model



# Continuous trait

A linear regression model is defined as

$$y = x\beta_1 + \beta_0 + \varepsilon$$

Data:

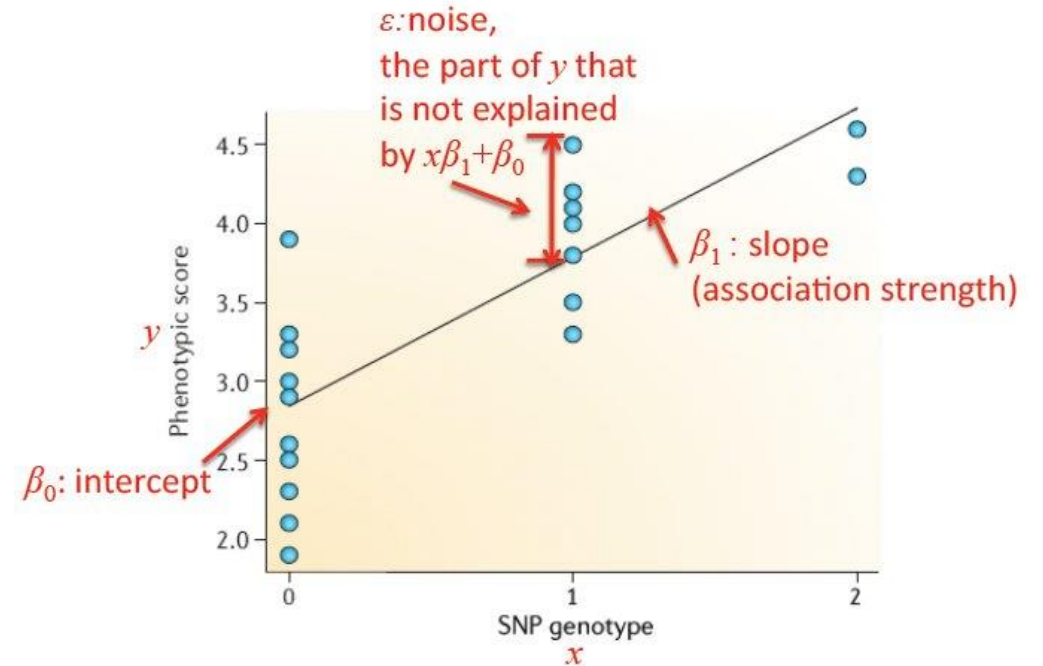
- $y$ : a continuous trait
- $x$ : SNP genotype at a given locus

Parameters:

- $\beta_1$ : regression coefficient, represents the strength of association between  $x$  and  $y$
- $\beta_0$ : intercept term (is 0 or ignored)
- $\varepsilon$ : noise or the part of  $y$  that is not explained by  $x$  (e.g., environmental effect)

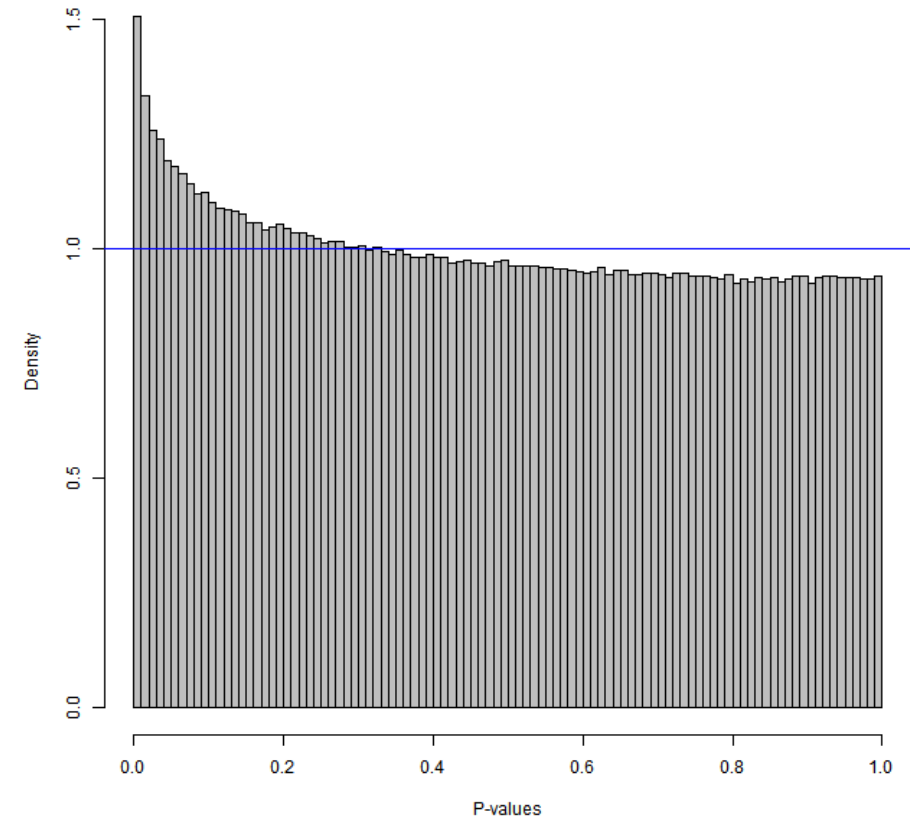
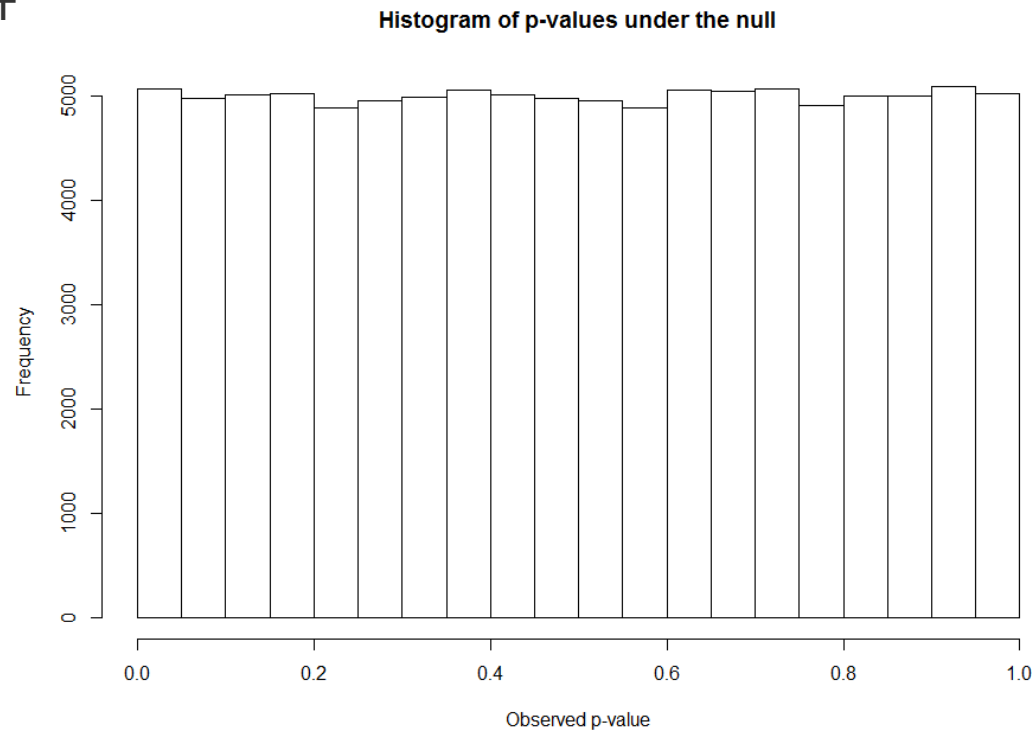
Assumptions:

- The individuals in the study are not related
- The phenotype  $y$  has a normal distribution



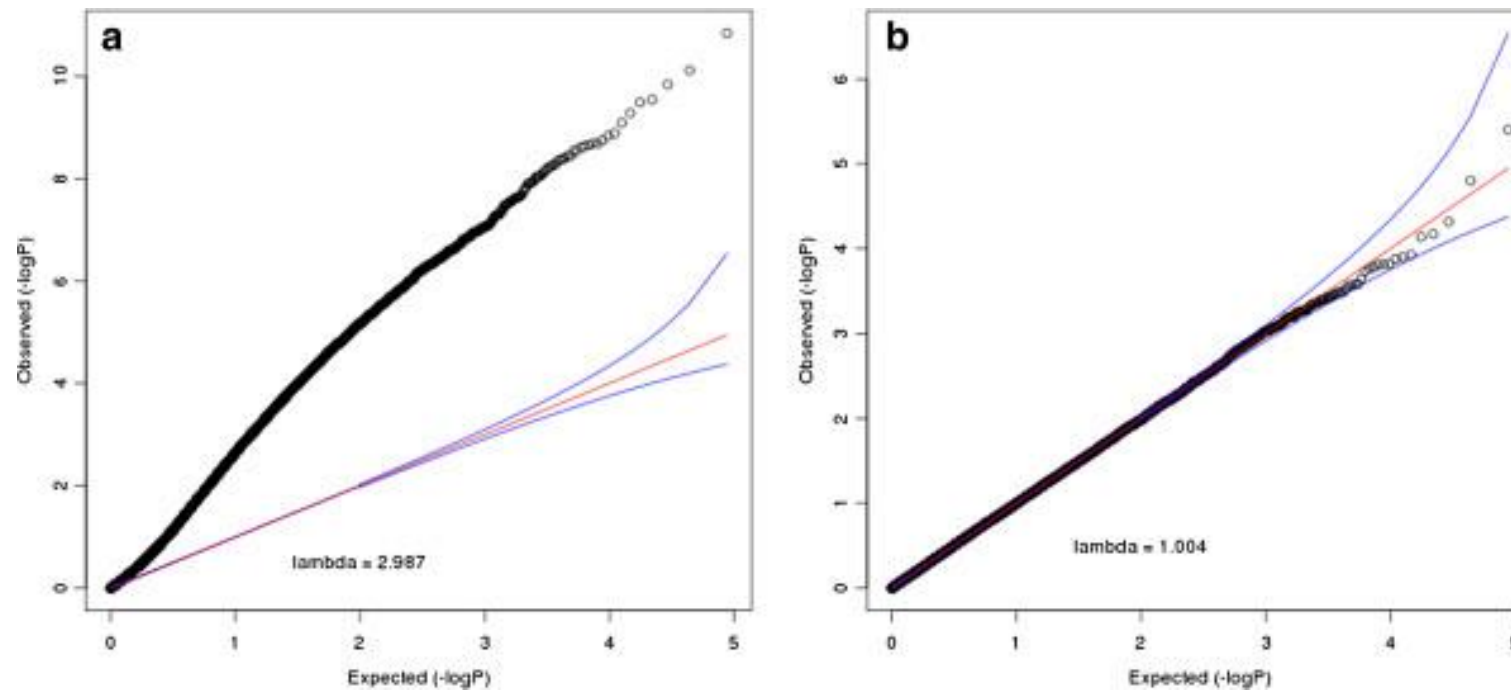
# Checking results

- QQ-plot
  - Distribution of p-values is uniform  $[0,1]$  under the null
  - If we have much signal, more around 0
  - Compare quantiles with expected ones : QQ-plot
  - In R: `qqunif`



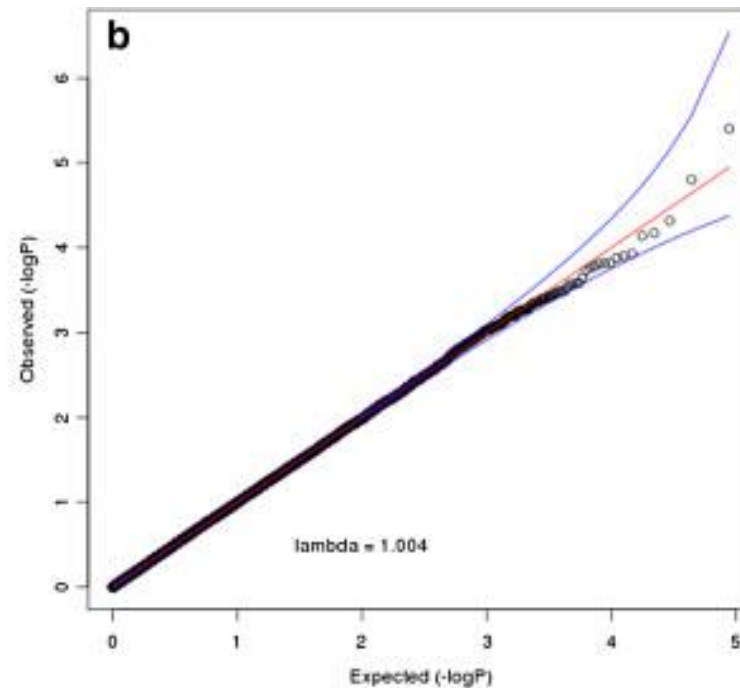
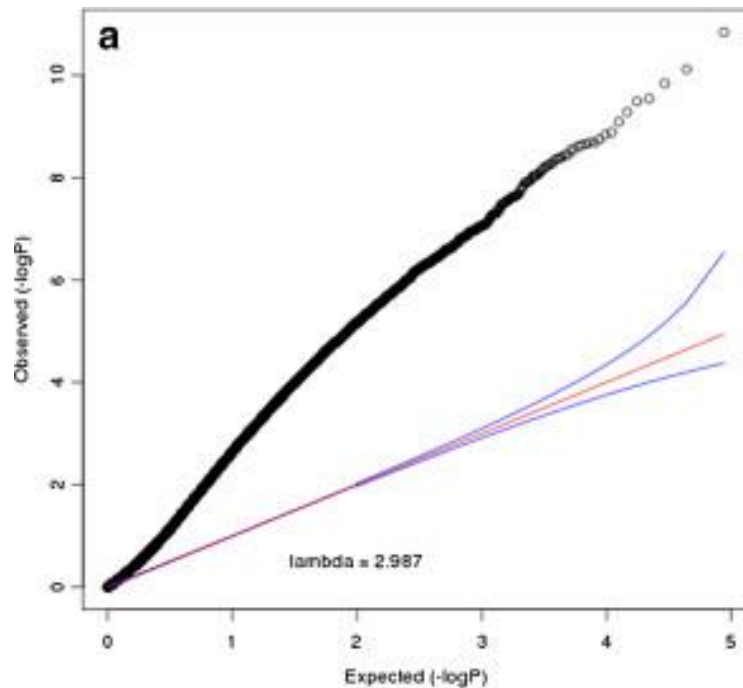
# Checking results

- QQ-plot
  - Distribution of p-values is uniform  $[0,1]$  under the null
  - If we have much signal, more around 0
  - Compare quantiles with expected ones : QQ-plot
  - In R: `qqunif`



# Checking results

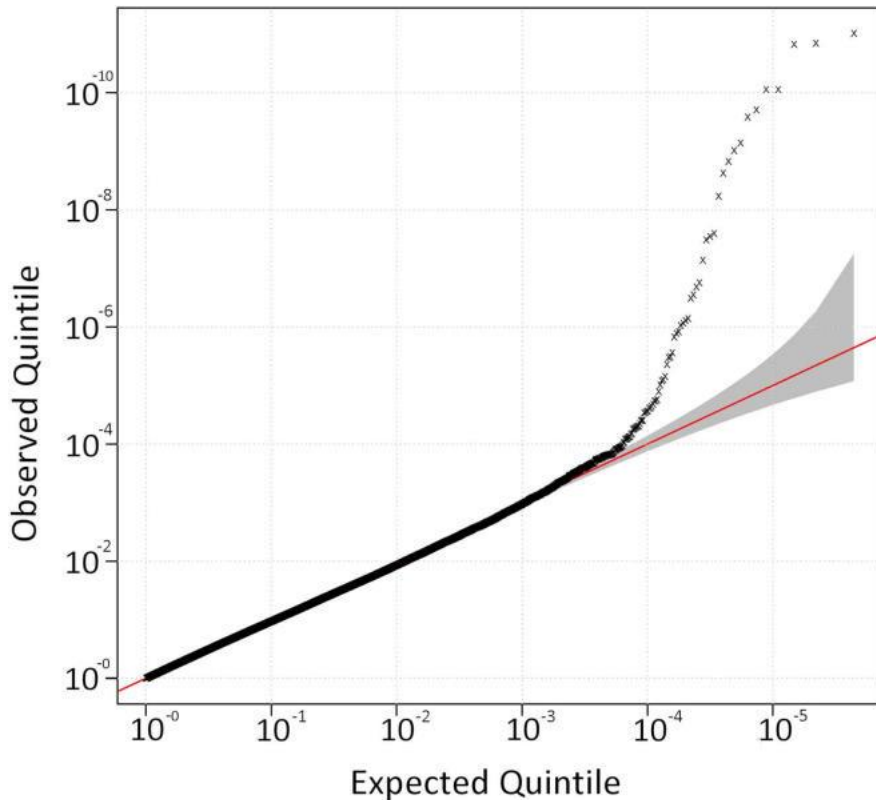
- QQ-plot
  - Distribution of p-values is uniform [0,1] under the null
  - If we have much signal, more around 0
  - Compare quantiles with expected ones : QQ-plot
  - In R: qqunif



- Inflation: too much signal
- Measured visually, but also lambda

# Checking results

- QQ-plot
  - Distribution of p-values is uniform [0,1] under the null
  - If we have much signal, more around 0
  - Compare quantiles with expected ones : QQ-plot
  - In R: qqunif



Appearances can be deceiving:

- A QQ-plot can look inflated when it isn't (just a lot of signal)
- And conversely
- We calculate the genomic inflation factor

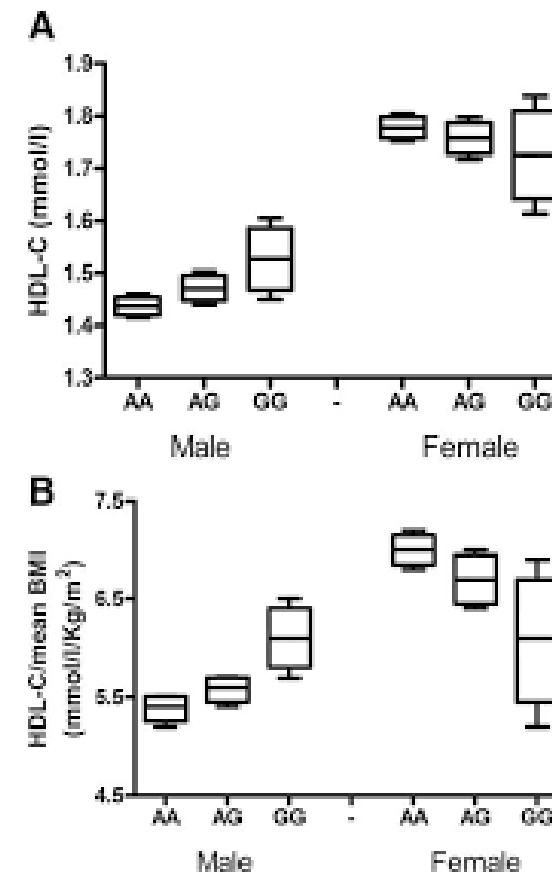
$$\lambda = \frac{\text{median}(Q_{\chi^2}(p))}{0.45}$$

(median of  $\chi^2$  test statistics divided by median of  $\chi^2_1$ )

- Ideally, want to correct in the model
- Can also adjust: GC correction (divide by lambda)

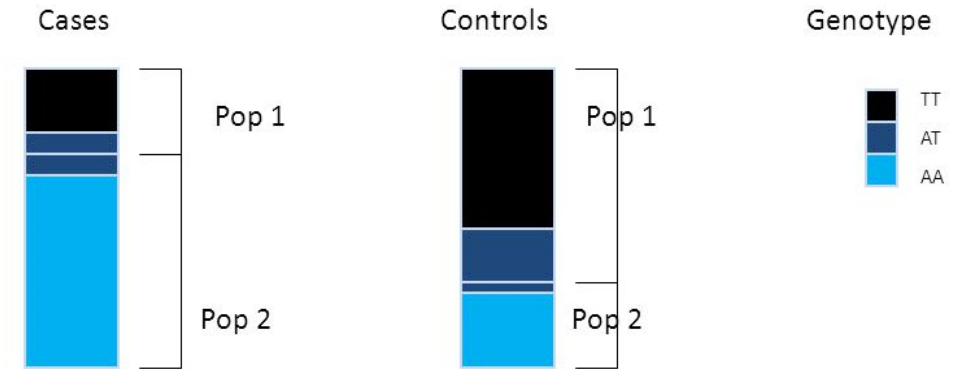
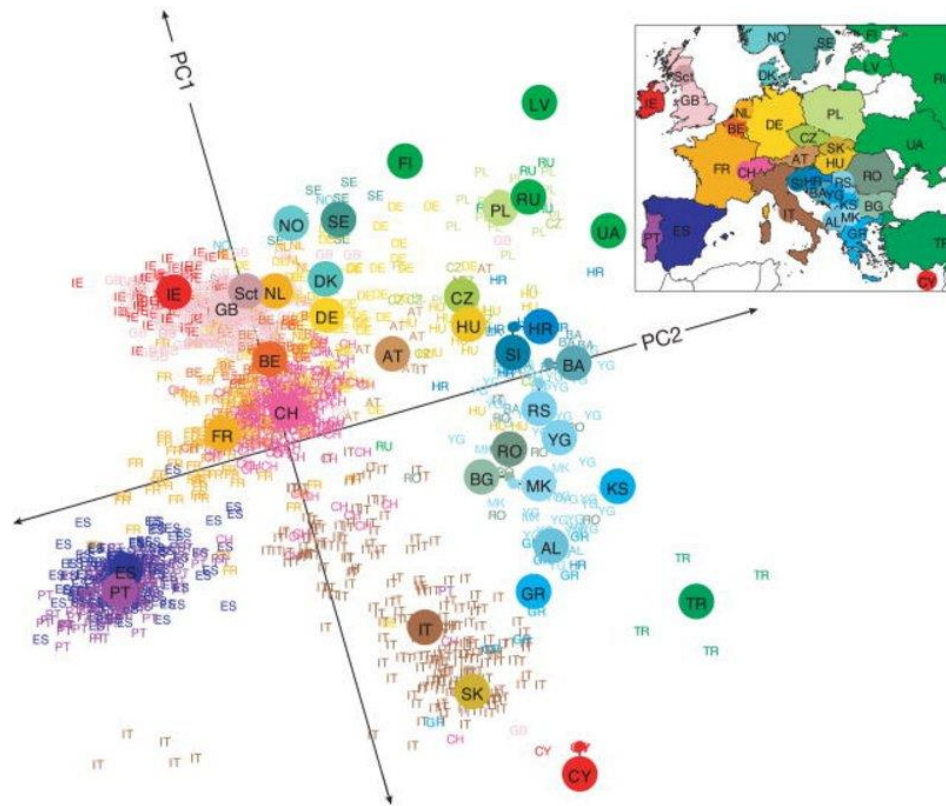
# Problems

- Covariates: sex, batch effects, chip effects





# Problems

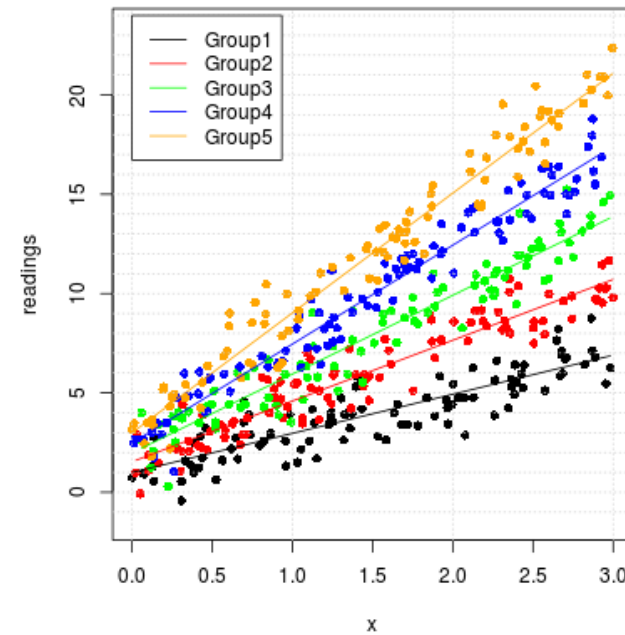
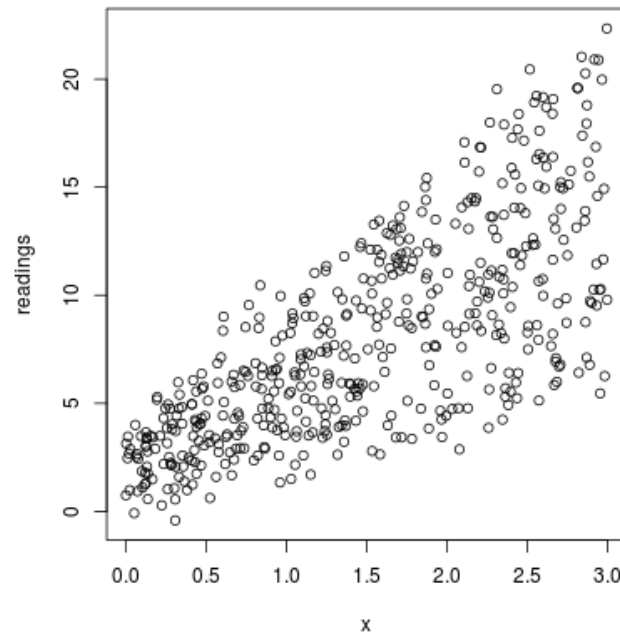


- Covariates: sex, batch effects, chip effects
- Structure: villages or subpopulations: linear mixed models can model the intra-group effect

# Problems

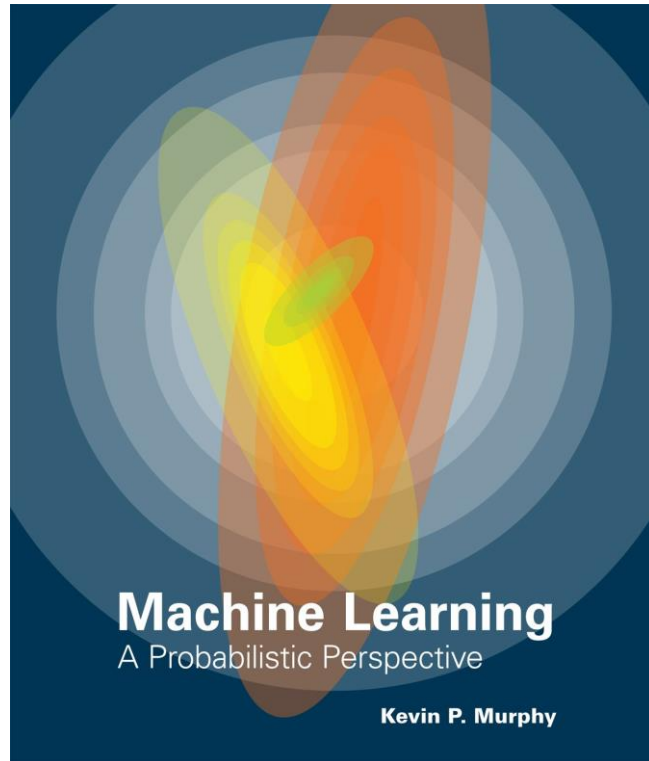
- Covariates: sex, batch effects, chip effects
- Structure: villages or subpopulations: linear mixed models can model the intra-group effect

$$\text{phenotype} \sim \beta \times \text{genotype} + \beta_1 \times \text{covariates} + \beta_2 \times \text{structure} + \epsilon$$

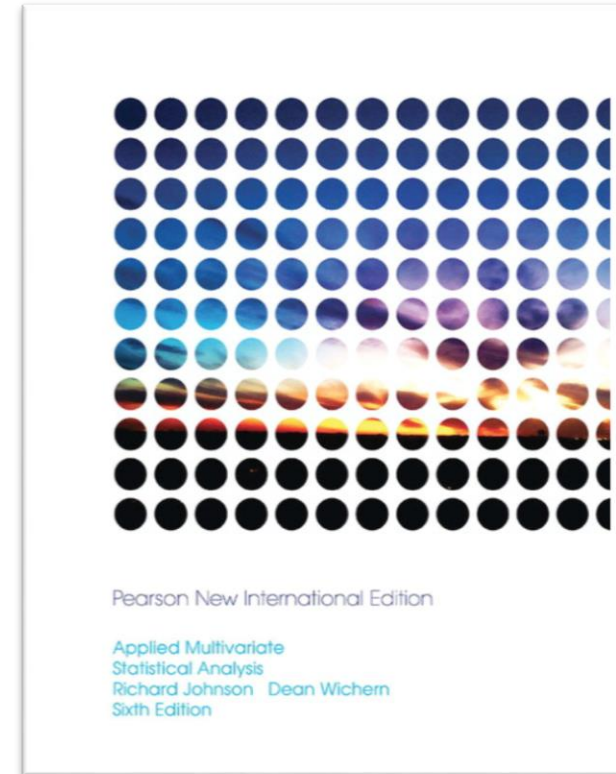


# References

*These references are for probability and statistical theory only  
They are not specific to statistical genetics*



- Murphy K, **Machine learning: A probabilistic perspective.** (MIT Press)



- Johnson and Wichern, **Applied multivariate statistical analysis.** (Pearson)