### Genome-Wide Association Studies

#### Caitlin Collins, Thibaut Jombart

MRC Centre for Outbreak Analysis and Modelling Imperial College London

Genetic data analysis using 30-10-2014



### Outline

- Introduction to GWAS
- Study design
  - GWAS design
  - Issues and considerations in GWAS
- Testing for association
  - Univariate methods
  - Multivariate methods
    - Penalized regression methods
    - Factorial methods

# Genomics & GWAS

 $\bullet$   $\bullet$   $\bullet$ 

# The genomics revolution

#### Sequencing technology

- o 1977 Sanger
- o 1995 1<sup>st</sup> bacterial genomes
  - < 10,000 bases per day per machine
- o 2003 1<sup>st</sup> human genome
  - > 10,000,000,000,000
     bases per day per machine

#### GWAS publications

2005 – 1<sup>st</sup> GWAS

 Age-related macular degeneration

 2014 – 1,991 publications

 14,342 associations





# A few GWAS discoveries...



## So what is GWAS?

- Genome Wide Association Study

   Looking for SNPs...
   associated with a phenotype.
- Purpose:
  - o **Explain** 
    - Understanding
    - Mechanisms
    - Therapeutics
  - Predict

Genomics & GWAS

- Intervention
- Prevention
- Understanding not required





## Association

#### Definition

 Any relationship between two measured quantities that renders them statistically dependent.

Heritability

 The proportion of variance explained by genetics

 $\circ \mathsf{P} = \mathsf{G} + \mathsf{E} + \mathsf{G}^*\mathsf{E}$ 

Heritability > 0





### The case of the missing heritability

Genomics & GWAS

6

# Why?

- Environment, Gene-Environment interactions
- Complex traits, small effects, rare variants
- Gene expression levels
- GWAS methodology?



#### The case of the missing heritability

Study Design

# **GWAS** design

#### Case-Control

- Well-defined "case"
- o Known heritability

#### Variations

- Quantitative phenotypic data
  - Eg. Height, biomarker concentrations
- Explicit models
  - Eg. Dominant or recessive

### Issues & Considerations

- Data quality
   0 1% rule
- Controlling for confounding

   Sex, age, health profile
   Correlation with other variables
- Population stratification\*
- Linkage disequilibrium<sup>\*</sup>



# **Population stratification**

#### Definition

- "Population stratification" = population structure
- Systematic difference in allele frequencies btw. subpopulations...
  - ... possibly due to different ancestry



#### Problem

- Violates assumed population homogeneity, independent observations
  - → Confounding, spurious associations
- Case population more likely to be related than Control population
  - → Over-estimation of significance of associations

## Population stratification II

- Solutions
   O Visualise
  - Phylogenetics
  - PCA

#### Correct

- Genomic Control
- Regression on Principal Components of PCA



# Linkage disequilibrium (LD)

### Definition

- Alleles at separate loci are NOT independent of each other
- Problem?
  - Too much LD is a problem
    - → noise >> signal
  - Some (predictable) LD can be beneficial
    - → enables use of "marker" SNPs





# Testing for Association

 $\bullet$   $\bullet$   $\bullet$ 

## Methods for association testing

- Standard GWAS

   Univariate methods
- Incorporating interactions

   Multivariate methods
  - Penalized regression methods (LASSO)
  - Factorial methods (DAPC-based FS)



# Univariate methods

- Approach
  - Individual test statistics
  - Correction for multiple testing



- Variations
  - o **Testing** 
    - Fisher's exact test, Cochran-Armitage trend test, Chisquared test, ANOVA
    - Gold Standard—Fischer's exact test
  - Correcting
    - Bonferroni
    - Gold Standard—FDR

Testing for Association

### Univariate – Strengths & weaknesses

### Strengths

- Straightforward
- Computationally fast
- Conservative
- Easy to interpret

#### Weaknesses

- Multivariate system, univariate framework
- Effect size of individual SNPs may be too small
- Marginal effects of individual SNPs ≠ combined effects



### Interactions

- White White Epistasis White White o "Deviation from li general linear n × AAbb aaBB  $Y_i = w_0 + w_1 A_i + w_2 B_{i} + w_3 A_i B_{i}$ Purple F1 • With p predictors, t All AaBb Purple •  $\binom{p}{k} = \frac{p^k}{k!}$  k-way interactions
  - p = 10,000,000 → 5 x 10<sup>11</sup>
     That's **500 BILLION** possible pair-wise interactions!

# Need some way to limit the number of pairwise interactions considered...

• Testing for Association

•

## Multivariate methods

#### **Penalized Regression**

LASSO penalized regression Ridge regression

#### **Bayesian Approaches**

Bayesian partitioning Bayesian Logistic Bayesian Epistasis Regression with Association Mapping Stochastic Search Variable Selection

#### **Factorial Methods**

Sparse-PCA Supervised-PCA DAPC-based FS (snpzip) Odds-ratiobased MDR

**Neural Networks** Genetic programming optimized neural networks **Logic Trees** Logic feature selection Monte Carlo Logic regression Logic Regression Modified Logic **Regression-Gene Expression Programming** Genetic Programming for Set association Association Studies approach **Non-parametric Methods** Random forests Restricted partitioning method Combinatorial partitioning method

## Multivariate methods (ii)

Penalized regression methods
 LASSO penalized regression

 Factorial methods
 DAPC-based feature selection

•23

# Penalized regression methods

#### Approach

- Regression models multivariate association
- Shrinkage estimation → feature selection

#### Variations

- o LASSO, Ridge, Elastic net, Logic regression
  - Gold Standard—LASSO penalized regression

# LASSO penalized regression

o Generalized linear model ("glm")

# Penalization L1 norm

•

 $\circ \text{ Coefficients } \rightarrow 0$ 

• Feature selection!



Testing for Association

## LASSO – Strengths & weaknesses

### Strengths

- Stability
- Interpretability
- Likely to accurately select the most influential predictors
  - Sparsity

#### Weaknesses

- Multicollinearity
- Not designed for high-p
- Computationally intensive
- Calibration of penalty parameters

   User-defined → variability
- Sparsity
- NO p-values!

### Factorial methods

#### Approach

- Place all variables (SNPs) in a multivariate space
- Identify discriminant axis → best separation
- Select variables with the highest contributions to that axis

#### Variations

- Supervised-PCA, Sparse-PCA, DA, DAPC-based FS
- Our focus—DAPC with feature selection (snpzip)

# DAPC-based feature selection



## DAPC-based feature selection

#### Where should we draw the line?

 $\circ \rightarrow$  Hierarchical clustering



### Hierarchical clustering (FS)





\$FS\$`Number of selected vs. unselected alleles`
[1] 1 4

\$FS\$`List of selected alleles`
[1] 1

(Hooray!

\$FS\$`Names of selected alleles`
[1] "a"

\$FS\$`Contributions of selected alleles to discriminant<sup>\*</sup>axis
1
0.39

## DAPC – Strengths & weaknesses

### Strengths

- More likely to catch all relevant SNPs (signal)
- Computationally quick
- Good exploratory tool

### Weaknesses

- Sensitive to n.pca
- N.snps.selected varies
- No "p-value"
- Redundancy > sparsity
- Redundancy > sparsity

### Conclusions

- Study design
  - GWAS design
  - Issues and considerations in GWAS
- Testing for association
  - Univariate methods
  - Multivariate methods
    - Penalized regression methods
    - Factorial methods

# Thanks for listening!

• • •

Questions?

#### • • •