

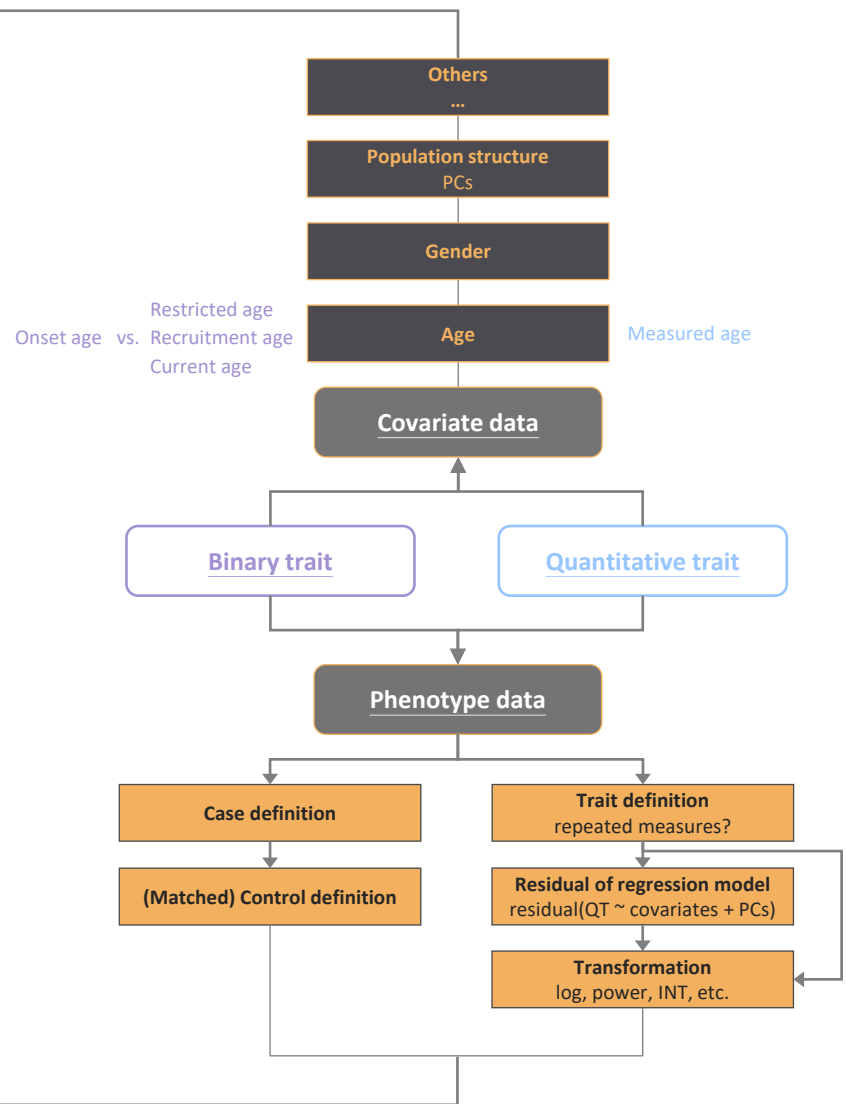
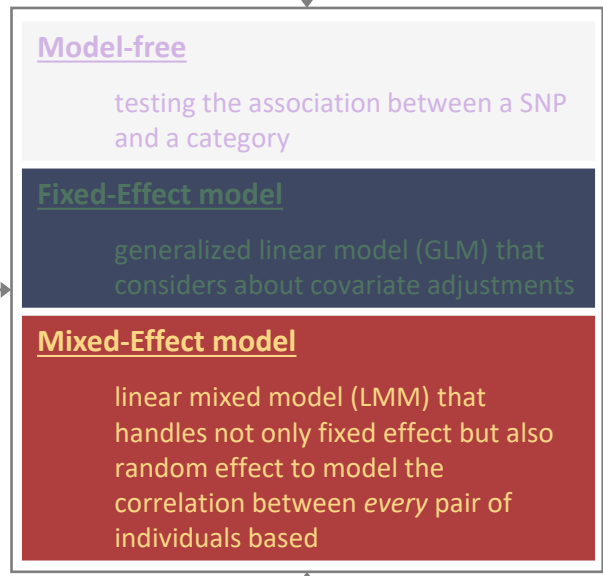
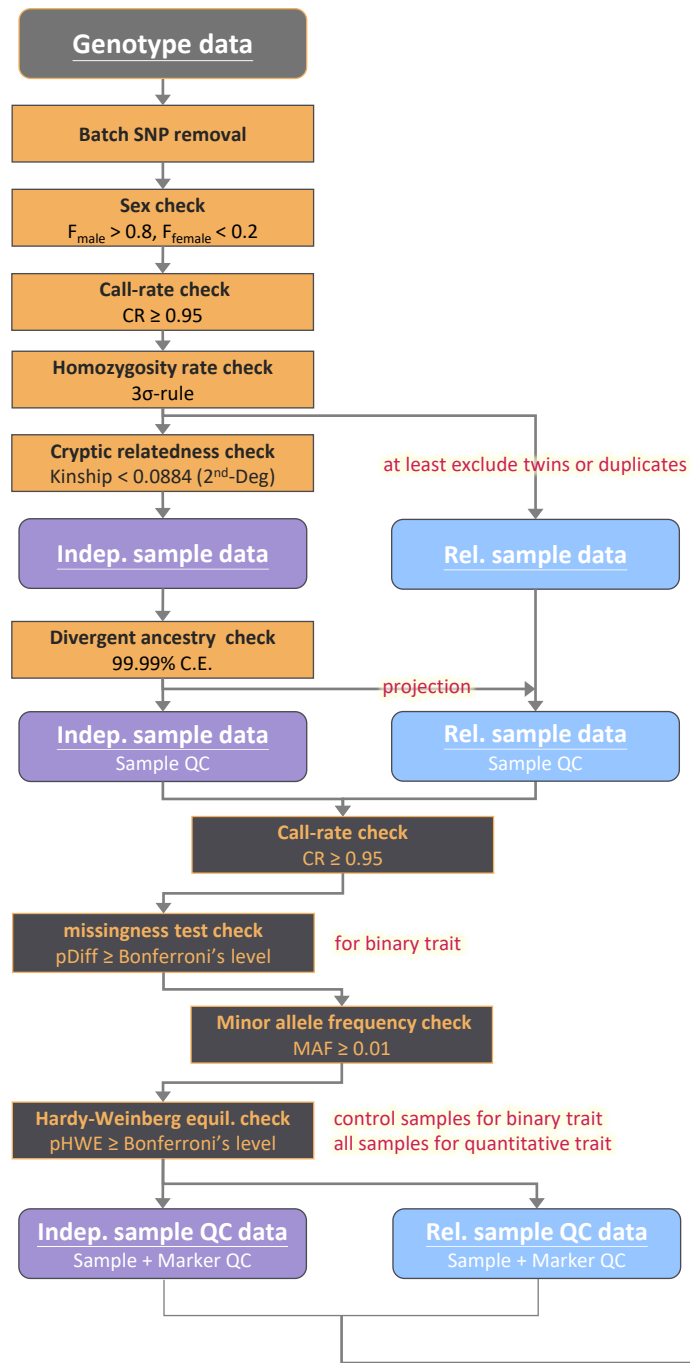
# Re:ZERO

Re: Analysis in a genomic world from zero

- Starting Analysis in Genomic World -

陳佳煒 (Jia-Wei Chen)

2026/04/07

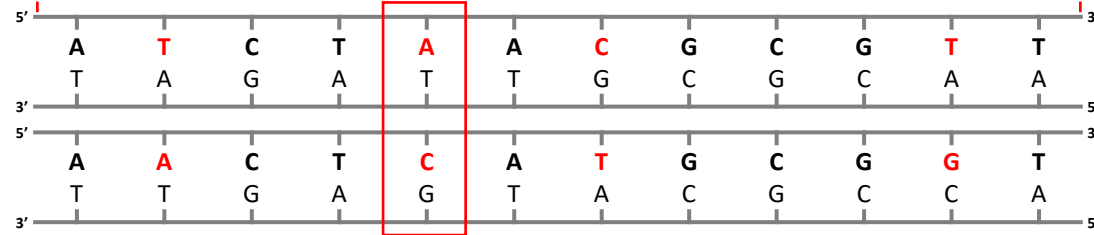
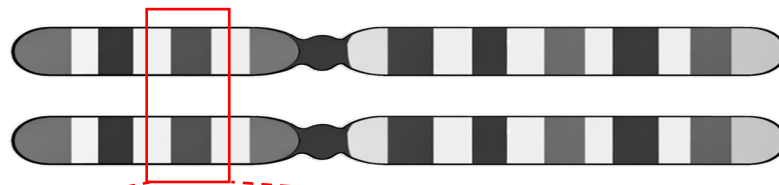


**:Terminology**



## Locus

The specific location (a base pair or a genomic region) on a chromosome



## Diploid

Having two complete sets of chromosomes, one from each biological parent

## Double strand

DNA molecule with two complementary polynucleotide chains held by base pairing

## Marker

A DNA position used to track variation (could be SNV, indel, etc.)

## Allele

A variant of the sequence of nucleotides at a particular location, or locus

- \* Allele 1 (A1) = A, Allele 2 (A2) = C
- \* Major / Minor allele
- \* Reference (Ref) / Alternative (ALT) allele
- \* Effect allele (EA)

## Genotype

Combination of alleles at a locus

- \* Homozygous: AA, CC
- \* Heterozygous: AC

## SNV

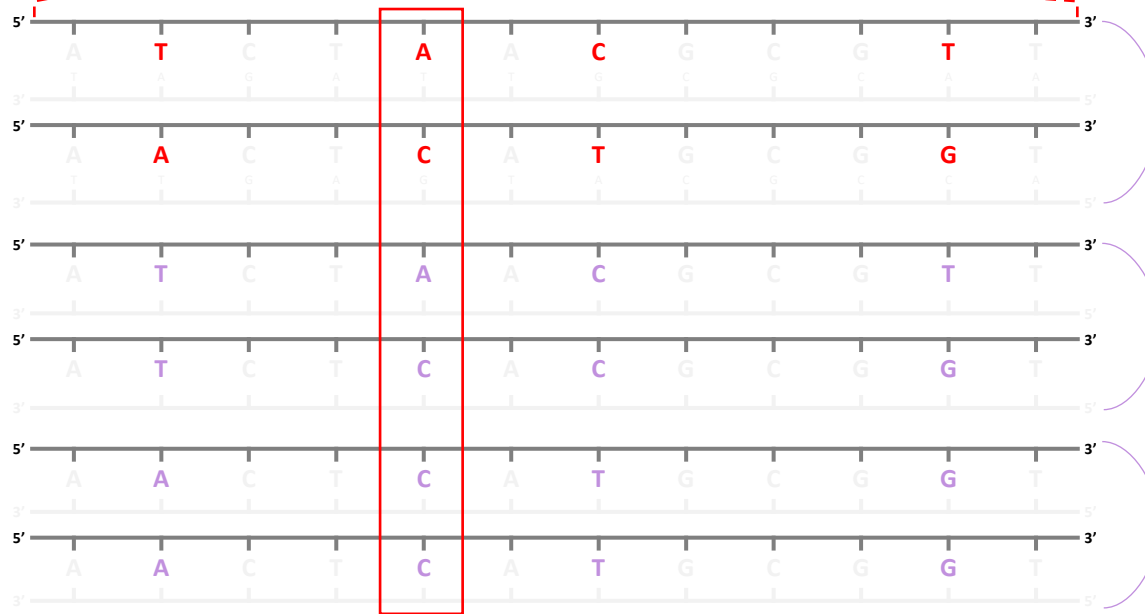
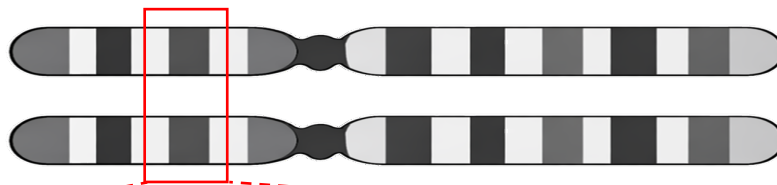
Single Nucleotide Variant: any single base change in DNA sequence

Common	MAF $\geq$ 0.05
SNP	MAF $\geq$ 0.01
Low-frequency	$0.01 \leq$ MAF $<$ 0.05
Rare	$0.001 \leq$ MAF $<$ 0.01
Ultra-rare	MAF $<$ 0.001

## Haplotype

Combination of alleles on a strand

A-T-C-T-A-A-C-G-C-G-T-T  
 A-A-C-T-C-A-T-G-C-G-G-T



### CR

Call-Rate: proportion of non-missing genotypes for a variant or a sample

### HR

Homozygous/Heterozygous Rate: proportion of homozygous/heterozygous genotypes for a variant or a sample

Homozygous rate (variant level)  
 $(0 + 0 + 1) / 3 = 0.33$   
Homozygous rate (genome-wide)  
 $(0 + 0 + 0 + 0) / 4 = 0$   
 $(1 + 0 + 1 + 0) / 4 = 0.5$   
 $(1 + 1 + 0 + 1) / 4 = 0.75$

### IBD

Identity By Descent: DNA segments inherited from a recent common ancestor without recombination

### IBS

Identity By State: DNA segments that match in sequence, regardless of shared ancestry (includes IBD + chance/population matches)

### GRM

Genetic Relationship Matrix: pairwise genetic relatedness estimates from genome-wide SNP data

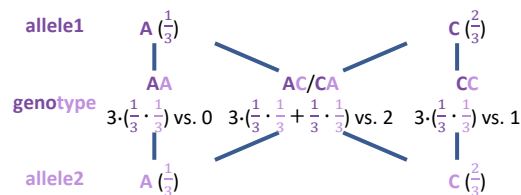
### MAF

Minor Allele Frequency: frequency of the less common allele at a variant site

$AF(A) = (1 + 1 + 0) / (3 \times 2) = 0.33$   
 $AF(C) = (1 + 1 + 2) / (3 \times 2) = 0.67$   
 $MAF = AF(A) = 0.33$

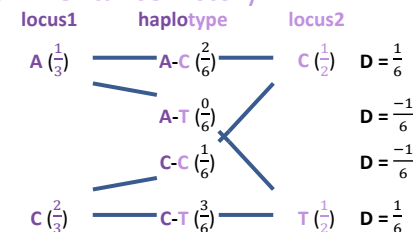
### HWE

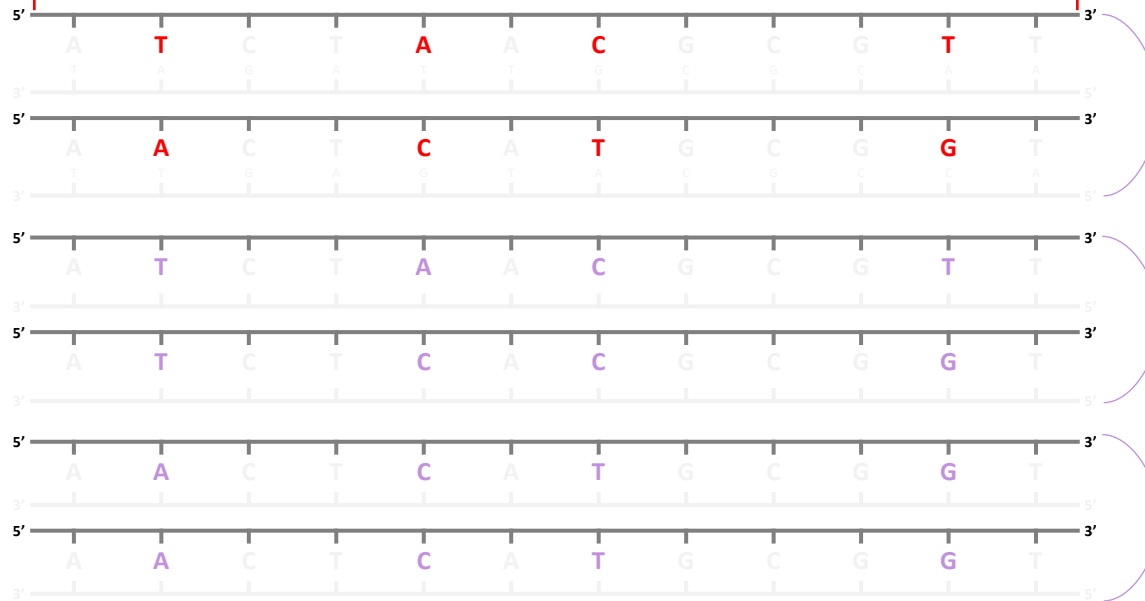
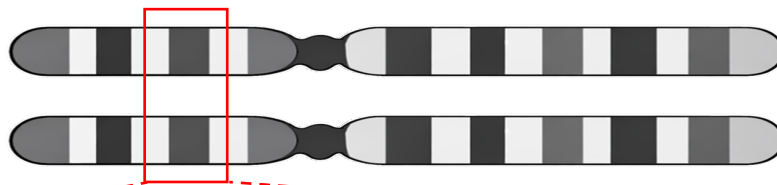
Hardy-Weinberg Equilibrium: allele/genotype frequencies stable across generations without evolutionary forces



### LD

Linkage Disequilibrium: non-random association of alleles at different loci due to shared inheritance history





	Binary trait	Quantitative trait
FID/IID	FID = 0 or FID = IID	
PID/MID	PID = MID = 0	
SEX	0 = unknown, 1 = male, 2 = female	
PHENO	-9/0 = missing 1 = unaffected 2 = affected	at least two different values
A1/A2	0 = missing 1 = A, 2 = C, 3 = G, 4 = T non-zero integers or any characters	

CHR	1-22, 23 = X, 24 = Y, 25 = XY, 26 = MT
ID	probe set ID (Ax-***): unique id for probe sequence affy SNP ID (Affy-***): unique id for CHR, POS, REF, and ALT dbSNP RS ID (rs***): unique id for a genome build
cM	genetic distance (centiMorgan)
POS	plain integers check genome build

### .ped (pedigree + genotype)

Family ID	Individual ID	Paternal ID	Maternal ID	SEX	Phenotype	variant1		variant2		variant3		variant4		...
						A1	A2	A1	A2	A1	A2	A1	A2	
FID	IID	PID	MID		PHENO	A1	A2	A1	A2	A1	A2	A1	A2	
FAM001	IND1	0	0	1	2	T	A	A	C	C	T	T	G	
FAM001	IND2	0	0	1	2	T	T	A	C	C	C	T	G	
FAM001	IND3	0	0	2	1	A	A	C	C	T	T	G	G	
	⋮													

### .map

Chromosome	Variant ID	Genetic distance	Physical position
CHR	ID	cM	POS
1	var1	0	565433
1	var2	0	752566
1	var3	0	753541
	⋮		

## bfile

2-bit genotype codes

00	Homozygous for first allele in .bim file
01	Missing genotype
10	Heterozygous
11	Homozygous for second allele in .bim file

### .bim (extended map)

A2	A	C	T	(major allele)
A1	T	A	C	(minor allele)
POS				
cM				
ID				
CHR				

FID	IID	PID	MID	SEX	PHENO

.fam (pedigree)

var1	var2	var3	var4	
10	10	10	10	...
00	10	00	10	
11	11	11	11	
:				

.bed (binary genotype)

## pfile

### .pvar (extended map)

ALT	A	C	T	(alternative allele)
REF	T	A	C	(reference allele)
POS				
cM				
ID				
CHR				

#FID	IID	PAT	MAT	SEX	PHENO

.psam (pedigree)

var1	var2	var3	var4

.pgen (binary genotype)

### .ped (pedigree + genotype)

FID	IID	PID	MID	SEX	PHENO	var1		var2		var3		var4		...
						A1	A2	A1	A2	A1	A2	A1	A2	
FAM001	IND1	0	0	1	2	T	A	A	C	C	T	T	G	...
FAM001	IND2	0	0	1	2	T	T	A	C	C	C	T	G	
FAM001	IND3	0	0	2	1	A	A	C	C	T	T	G	G	
:														

### .map

CHR	ID	cM	POS
1	var1	0	565433
1	var2	0	752566
1	var3	0	753541
:			

# bfile

2-bit genotype codes

00	Homozygous for first allele in .bim file
01	Missing genotype
10	Heterozygous
11	Homozygous for second allele in .bim file

## .bim (extended map)

A2	A	C	T	(major allele)
A1	T	A	C	(minor allele)
POS				
cM				
ID				
CHR				

FID IID PID MID SEX PHENO

--	--	--	--	--	--	--

.fam (pedigree)

var1 var2 var3 var4

10	10	10	10	...
00	10	00	10	
11	11	11	11	
:				

.bed (binary genotype)

PHENO	<p><b>Binary</b> 2 = Case, 1 = Control, 0/-9/NA = Unknown (default) 1 = Case, 0 = Control, NA = Unknown</p>
	<p><b>Quantitative</b> any numeric values</p>
COVAR	<p><b>Demographic</b> age, gender</p> <p><b>Subpopulation structure</b> PC1, PC2, ...</p> <p><b>Others</b> BMI, abc-covariates, etc.</p>

## .phe (phenotype)

FID	IID	PHENO1	PHENO2	...
FAM001	IND1			
FAM001	IND2			
FAM001	IND3			
:				

## .cov (covariate)

FID	IID	COVAR1	COVAR2	...
FAM001	IND1			
FAM001	IND2			
FAM001	IND3			
:				

## .phecov (phenotype + covariate)

FID	IID	PHENO1	PHENO2	...	COVAR1	COVAR2	...
FAM001	IND1						
FAM001	IND2						
FAM001	IND3						
:							

**:Toolset**

# LINUX

## Terminal (终端機)

```
$ uname -m
x86_64 # 64-bit
i686 i386 # 32-bit
```

## Terminal (终端機)

```
$ lscpu | grep avx2
... avx2 ... # supports AVX2
# not support AVX2
```

## Terminal (终端機)

```
# execute from the specific folder
$ cd path/of/plink
$ plink

# execute everywhere
$ copy path/of/plink /usr/local/bin/
$ plink
```

# PLINK

a free, open-source whole genome association analysis toolset, designed to perform a range of basic, large-scale analyses in a computationally efficient manner.

## TECHNICAL NOTE

Open Access

## Second-generation PLINK: rising to the challenge of larger and richer datasets

Christopher C Chang<sup>1,2\*</sup>, Carson C Chow<sup>3</sup>, Laurent CAM Tellier<sup>2,4</sup>, Shashaank Vattikuti<sup>3</sup>, Shaun M Purcell<sup>3,5,6,7,8</sup> and Jamin **PLINK 1.9**

<https://www.cog-genomics.org/plink/1.9/>

The screenshot shows the PLINK website with two main sections: PLINK 1.9.0 beta and PLINK 2.0.0 alpha. The PLINK 1.9.0 beta section includes a table of binary downloads for Linux (64-bit, 32-bit), macOS (64-bit), and Windows (64-bit, 32-bit). The PLINK 2.0.0 alpha section includes a table of binary downloads for Linux AVX2 AMD, Linux AVX2 Intel, Linux 64-bit Intel, Linux 32-bit, macOS M1, macOS AVX2, macOS 64-bit, Windows AVX2, Windows 64-bit, and Windows 32-bit.

**PLINK 2.0**  
<https://www.cog-genomics.org/plink/2.0/>

# WINDOWS

## Command Prompt (命令提示字元)

```
> wmic os get osarchitecture
64-bit # 64-bit
32-bit # 32-bit
```

## Command Prompt (命令提示字元)

```
> wmic cpu get caption
Family 6 Model 60 (Intel) # supports AVX2
```

## Command Prompt (命令提示字元)

```
# execute from the specific folder
$ cd path/of/plink
$ plink
```

**--keep / --remove:** Filters individuals in/out based on a list of Family and Individual IDs

**--extract / --exclude:** Filters variants (SNPs) in/out based on a list of IDs

**--het:** Calculates the heterozygosity of individuals to identify sample contamination or inbreeding

**--geno:** Filters out SNPs with a missing call rate higher than a specified threshold (e.g., 0.05)

**--hwe:** Filters out SNPs that deviate significantly from Hardy-Weinberg Equilibrium

**--indep-pairwise:** Performs Linkage Equilibrium (LD) pruning, identifying SNPs that are highly correlated to reduce redundancy

**--genomes:** Calculates Identity-by-Descent (IBD) to identify cryptic relatedness or sample duplicates

**--pheno:** Points to an external file containing your phenotype data

**--covar:** Loads a file with covariates to include in the regression model to prevent confounding

**--zero-cluster:** Sets all genotypes to missing for members of a specific cluster (used for specific cleaning tasks)

```
Terminal (終端機)
# basic execution
$ plink --bfile path/of/files --flags [arg.]
```

**--king-cutoff:** Uses the KING-robust kinship estimator to remove closely related individuals, ensuring sample independence

**--make-bed / --make-pgen:** Creates a PLINK 1.9 binary triplet (.bed, .bim, .fam) / PLINK 2.0 binary triplet (.pgen, .pvar, .psam)

**--glm:** Performs association testing using Generalized Linear Models (linear, logistic, or Firth regression)

**--threads:** Sets the number of CPU cores for the process. Vital for speeding up calculations on large datasets like the 90k samples you're analyzing

**--within:** Loads family or cluster assignments from a file, ensuring analysis is performed within these specified groups

**--missing:** Generates statistics on missing data per individual and per SNP

**--check-sex:** Compares the sex assigned in the dataset to the actual heterozygosity on the X chromosome to identify mislabeled samples

**--maf:** Filters out SNPs with a Minor Allele Frequency below a specific threshold (e.g., 0.01)

**--score:** Calculates Polygenic Risk Scores (PRS) by summing weighted alleles across the genome

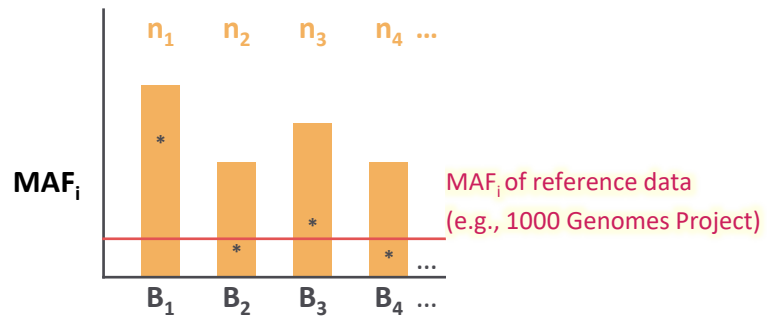
**--pca:** Conducts Principal Component Analysis to visualize and adjust for population stratification

**:Quality Control**

## Batch SNP removal

### Why batch SNP removal?

- False positive associations (e.g., cases and controls are in different batches)



### How to do?

- Compare AFs of SNPs to ones in public database (e.g., 1000 Genomes Project) to identify the SNPs with weird Afs
- Instead of excluding the SNP, marking the genotype calls of samples from problematic batches as missing for this specific SNP

## Input

### dat.clust

FID	IID	Batch
FAM001	IND1	1
FAM001	IND2	1
FAM001	IND3	2
⋮		

The batch information can be requested from typing center

### batchSNPs.zero

SNP	Batch
var1	1
var2	3
var3	3
⋮	

It may be a challenge for non-programmer to find out batch SNPs

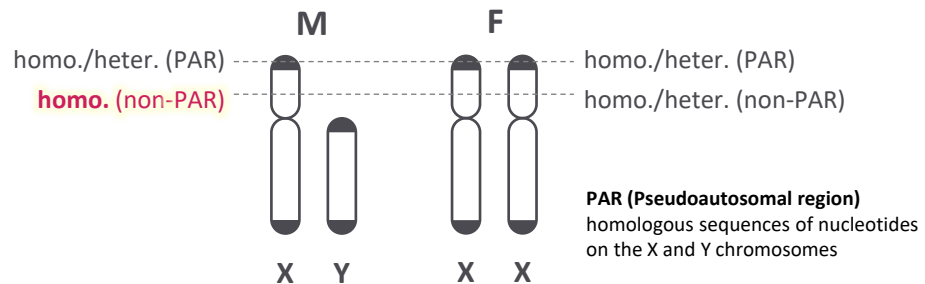
### # batch SNP removal (by batches)

```
$ plink --bfile path/of/your/file --zero-cluster batchSNPs.zero --within dat.clst --make-bed --out dat_wg
$ plink --bfile dat_wg --chr 1-22 --make-bed --out dat_as
```

## Sex check

### Why sex check?

- Chromosome aberration
- Covariate identification
- Sample misidentification



### How to do?

- Using either  
**homozygosity rate** : M ( $\geq 0.9$ ) & F ( $< 0.9$ ) or  
**inbreeding coefficient** : M ( $> 0.8$ ) & F ( $< 0.2$ )  
of X chromosome to check for the gender
- If EHR data is accessed, directly comparing EHR  
gender to genetic gender

### # sex check

```
$ plink --bfile dat_wg --check-sex --out dat_wg  
$ grep "PROBLEM" dat_wg.sexcheck > rmlnd_sex.txt  
$ plink2 --bfile dat_as --remove rmlnd_sex.txt --write-samples --out dat_as_1
```



### Output

dat\_wg.sexcheck

FID	IID	PESEX	SNPSEX	STATUS	F
FAM001	IND1	1	1	OK	0.9588
FAM001	IND2	1	1	OK	0.9616
FAM001	IND3	1	1	OK	0.9588
⋮					
FAM001	IND17	2	1	PROBLEM	0.9539
⋮					
FAM001	IND28	1	2	PROBLEM	-0.05586
⋮					

rmlnd\_sex.txt

FAM001 ind17  
FAM001 ind28

Remove samples with inconsistent genders (**STATUS = PROBLEM**)

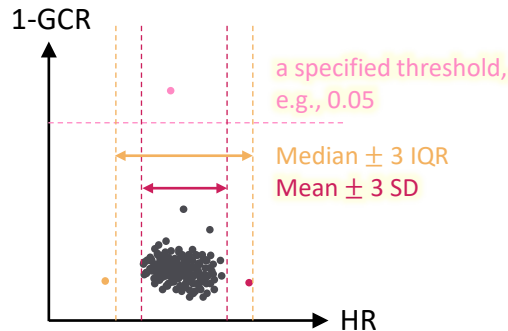
## Call-rate & Homozygosity rate check

### Why genotype call-rate check?

- Low DNA quality or concentration

### Why homozygosity rate check?

- Inbreeding (excess Homozygosity)
- DNA contamination (excess Heterozygosity)



$N$  = # of markers

$N_i$  = # of nonmissing markers for individual  $i$

$O_i$  = # of homozygous markers for individual  $i$

$$GCR_i = \frac{N_i}{N}$$

$$HR_i = \frac{N_i - O_i}{N_i}$$

### How to do?

- Calculate genotype call-rate (GCR) and heterozygosity rate (HR) for **autosomes**
- A specific cutoff (95% or even 99%) for CR and use 3- $\sigma$  rule or IQR rule for HR

## # call-rate & homozygosity rate check

```
$ plink --bfile dat_as --keep dat_as_1.id --missing --het --out dat_as_1
$ plink --bfile dat_as --keep dat_as_1.id --remove rmInd_missing_het.txt --write-samples --out dat_as_2
```



### Output

dat\_as\_1.imiss

FID	IID	MISS_PHENO	N_MISS	N_GENO	F_MISS
FAM001	IND1	Y			0.00158
FAM001	IND2	Y			0.00181
FAM001	IND3	Y			0.1375
:					

rmInd\_missing\_het.txt

FAM001	IND3	CR
FAM001	IND2	HET

Remove samples with CR = 1 - F\_MISS smaller than a threshold, say 0.95

dat\_as\_1.het

FID	IID	O(HOM)	E(HOM)	N(NM)	F
FAM001	IND1	426964		624027	
FAM001	IND2	588153		623882	
FAM001	IND3	425156		624154	
:					

Remove samples with HR =  $\frac{O(\text{HOM})}{N(\text{NM})}$  out of mean(HR)  $\pm$  3 sd(HR)

## Cryptic relatedness check

### Why cryptic relatedness check?

- Biostatistical Independence (e.g., assumption of GLM, bias AF estimation)
- Pedigree structure (e.g., family-based study)

	Proportion IBD	Kinship coeff.
duplicate/MZ twin	1	> 0.354
1 <sup>st</sup> degree	0.5	[0.177, 0.354]
2 <sup>nd</sup> degree	0.25	[0.0884, 0.177]
halfway of 2 <sup>nd</sup> & 3 <sup>rd</sup> degrees	> 0.1875	
3 <sup>rd</sup> degree	0.125	[0.0442, 0.0884]

### How to do?

- using **independent** markers (pair correlation  $r^2 < 0.2$ ) to calculate the relatedness of individuals by either **proportion IBD** or **kinship coefficient** to check for the relatedness

### # cryptic relatedness check (preferred)

```
$ plink --bfile dat_as --keep dat_as_2.id --indep-pairwise 50 5 0.2 --out dat_as_2
$ plink --bfile dat_as --keep dat_as_2.id --extract dat_as_2.prune.in --king-cutoff 0.0442 --out dat_as_2
$ plink --bfile dat_as --keep dat_as_2.id --remove dat_as_2.king.cutoff.out.id --write-samples --out dat_as_3
```

### # cryptic relatedness check

```
$ plink --bfile dat_as --keep dat_as_2.id --indep-pairwise 50 5 0.2 --out dat_as_2
$ plink --bfile dat_as --keep dat_as_2.id --extract dat_as_2.prune.in --genome --out dat_as_2
$ plink --bfile dat_as --keep dat_as_2.id --remove rmlnd_rel.txt --write-samples --out dat_as_3
```



### Output

dat\_auto\_2.genome

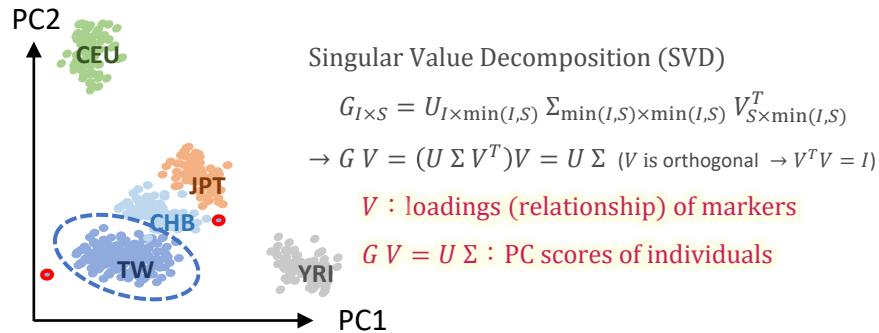
FID1	IID1	FID2	IID2	RT	EZ	Z0	Z1	Z2	PI_HAT	PHE	DST	PPC	RATIO
:	:	:	:	:	:	:	:	:					

Remove samples with lower CR in a pair of samples if **PI\_HAT > 0.1875**  
It may be a challenge for a non-programmer to find the optimal sample set

## Divergent ancestry check

### Why divergent ancestry check?

- Population stratification (e.g., [Hamer, D. , & Sirota, L. \(2000\)](#))
- Pedigree structure (e.g., family-based study)



### How to do?

- merge study and reference data
  - exclude ambiguous markers → prune correlated markers
  - PCA → exclude outliers
- PCA on reference data → extract loadings (marker weights)
  - project study data onto reference data → exclude outliers

## # divergent ancestry check

```
$ https://www.cog-genomics.org/plink/2.0/resources#phase3\_1kg # 1000 genomes phase 3 data

$ awk '{print $1,$4,$4,$2}' dat_as.bim > lst_var.txt
$ plink2 --zst-decompress all_hg38.pgen.zst > all_hg38.pgen
$ plink2 --pfile all_hg38 vzs --allow-extra-chr --extract bed1 lst_var.txt --snps-only --max-alleles 2 --set-all-var-ids @:# --rm-dup exclude-all --make-bed --out all_hg38_extract
$ sed 's/#IID/IID/g;s/Population/population/g' all_hg38.psam > relationships_w_pops.txt

$ plink2 --bfile all_hg38_extract --maf 0.01 --freq counts --pca biallelic-var-wts --out all_hg38_extract
$ plink2 --bfile dat_as --keep dat_as_3.id --read-freq all_hg38_extract.accounts --score all_hg38_extract.eigenvec.var 2 4 header read variance-standardize no-mean-imputation --score-col-nums 5-14 --out dat_as_3
$ plink --bfile dat_as --keep dat_as_3.id --remove rmlnd_divAncestry.txt --make-bed --out dat_as_sampQC
```

### Output

all\_hg38\_extract.account

#CHROM	ID	REF	ALT	ALT_CTS	OBS_CT
	var1				
	var2				
	var3				
	:				

Allele counts of variants

all\_hg38\_extract.eigenvec.var

#CHROM	ID	MAJ	NONMAJ	PC1	PC2	...	PC10
	var1						
	var2						
	var3						
	:						

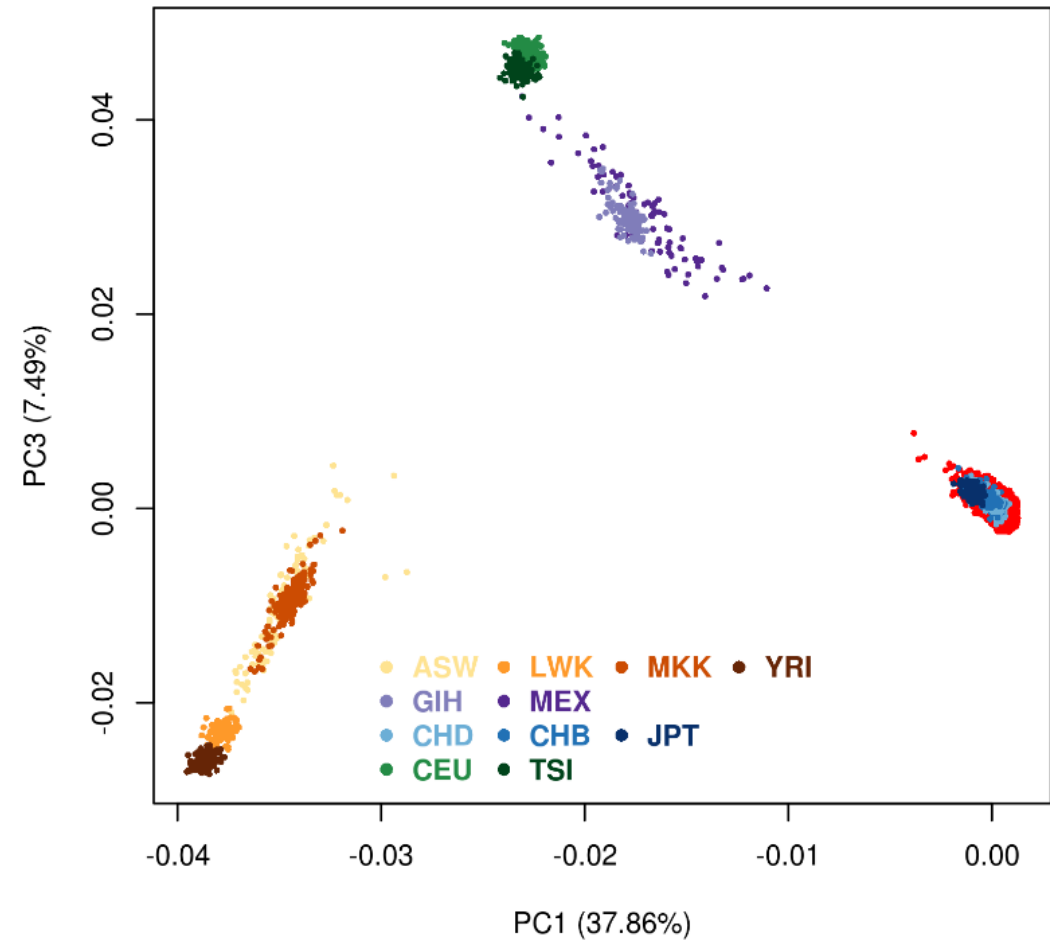
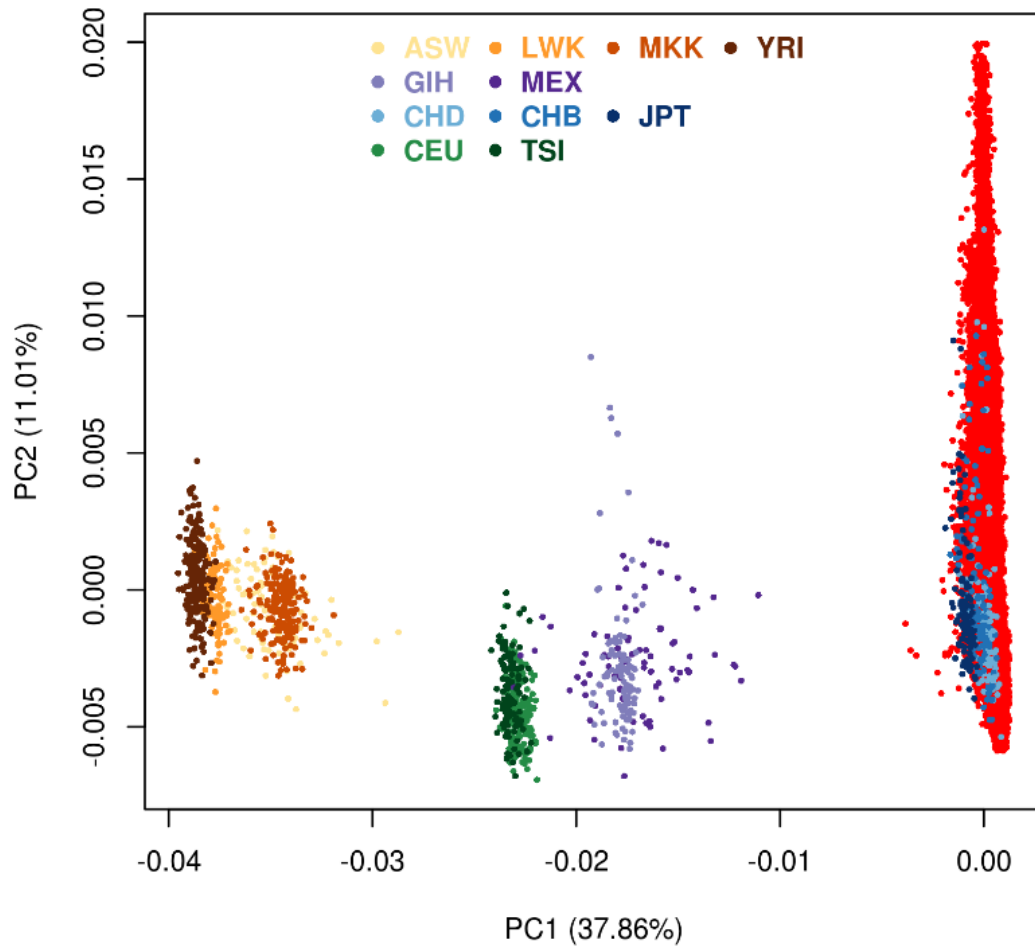
PC loadings of variants

dat\_as\_3.score

#FID	IID	ALLELE_CT	NAMED_ALLELE_DOSAGE_SUM	SCORE1_AVG	...	SCORE10_AVG
FAM001	IND1					
FAM001	IND2					
FAM001	IND3					
	:					

PC scores of individuals

## Divergent ancestry check



When sample size of target data is much larger than that of ancestry data, one of the first two PCs may be dominant by the variation in the target data

## Marker QC

### How to do?

- A specific cutoff (95% or even 99%) for CR
- A specific cutoff (0.01 or 0.05) for MAF
- A specific cutoff or Bonferroni's level for HWE; HWE test is applied to controls for case/control study and to all for quantitative-trait study

$$\begin{array}{ccccccc} n_{AA} & + & n_{AB} & + & n_{BB} & + & n_{Miss} & = & N \\ \hline & & n_A & + & n_B & & & = & 2(N - n_{Miss}) \\ \hline & & \text{CR} & & \text{MAF} & & \text{HWE} & & \end{array}$$

### Why marker QC?

- Technical artifacts (lower CR reflects assay issues since a failure across many samples)
- Statistical power (low MAF yields unstable variance and spurious associations)
- Biological Plausibility (HWE violation suggests miscalling or hidden population structures)

### # call-rate check

```
$ plink --bfile dat_as_sampQC --geno 0.05 --write-snplist --out dat_as_5
```

### # case/control nonrandom missingness check

```
$ plink --bfile dat_as_sampQC --keep dat_as_5.snplist --test-missing --out dat_as_5
$ n=$(wc -l < dat_as_5.missing)
$ awk -v n=$n '$5 < 0.05/(n-1) {print $2}' dat_as_5.missing > dat_as_5.snplist
```

### # MAF check

```
$ plink --bfile dat_as_sampQC --keep dat_as_5.snplist --maf 0.01 --write-snplist --out dat_as_6
```

### # HWE check

#### # quantitative trait

```
$ bonf=$(bc -l <<< "0.05/$(wc -l < dat_as_6.snplist)")
$ plink --bfile dat_as_sampQC --keep dat_as_6.snplist --hwe ${bonf} --make-bed --out dat_as_sampQC_varQC
```

#### # binary trait

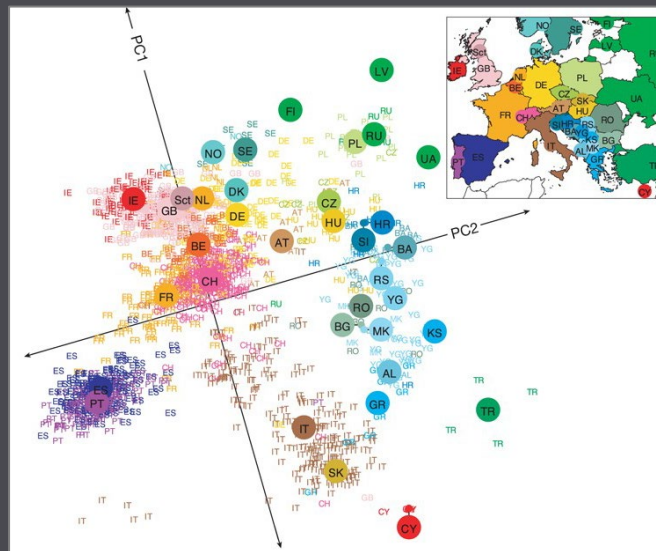
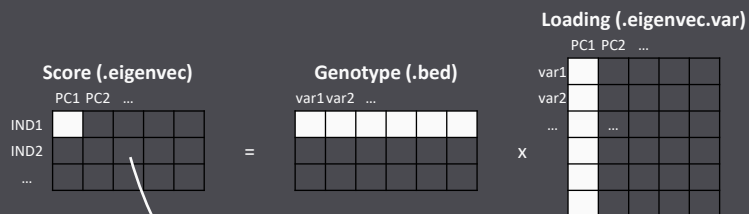
```
$ bonf=$(bc -l <<< "0.05/$(wc -l < dat_as_6.snplist)")
$ plink2 --bfile dat_as_sampQC --pheno dat.phecov --1 --keep-if "PHENO==0" dat_as_6.snplist --hwe ${bonf} --write-snplist --out dat_as_7
$ plink2 --bfile dat_as_sampQC --keep dat_as_7.snplist --make-bed --out dat_as_sampQC_varQC
```

**:PopulationStructure**

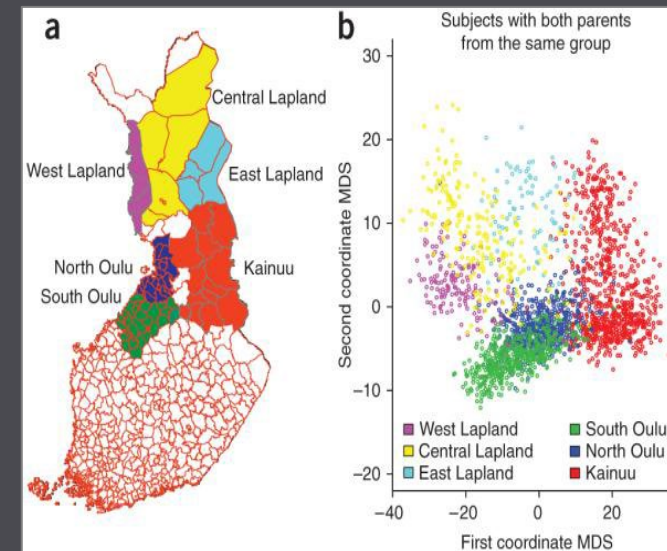
# Subpopulation structure

## Why subpopulation structure?

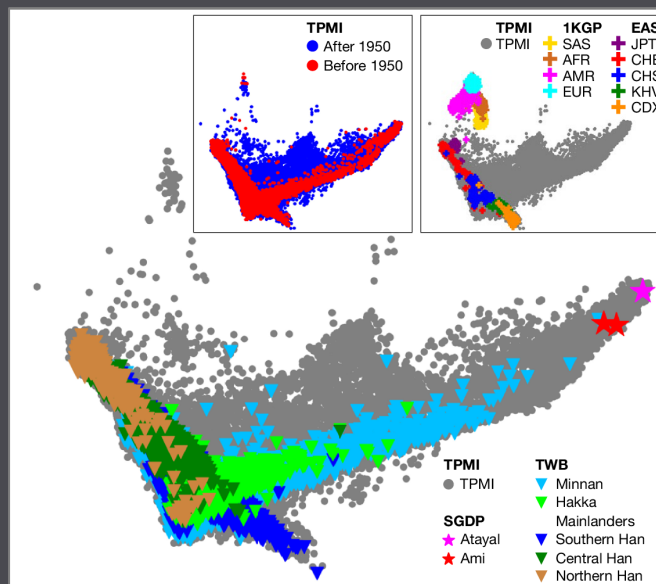
- A **confounder** effect that AF differences between subpopulations are correlated with trait differences can induce spurious associations



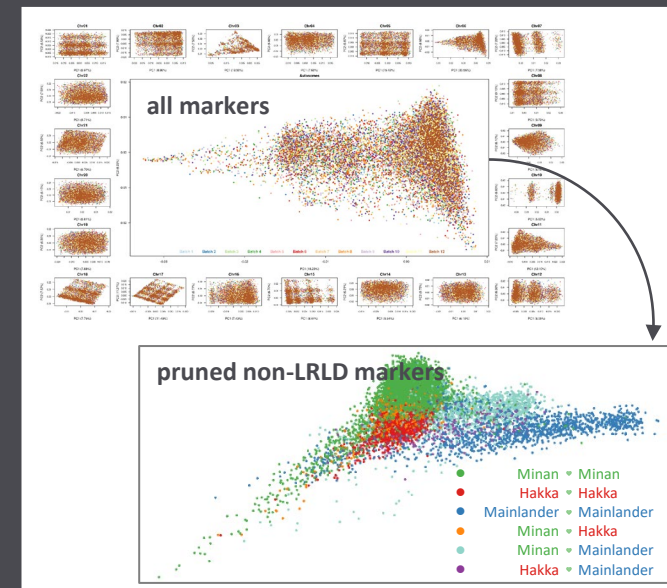
Novembre, J. et al. *Nat.* (2008) European



Sabatti, C. et al. *Nat. Genet.* (2009) Finland



Yang et al. *Nat.* (2025) Taiwan Precision Medicine Initiative (TPMI)



TaiWan Biobank (TWB1)



**:AssociationTest**



## Binary trait

A clear dichotomous outcome (or state)

- \* combine ICD codes with lab tests, medications, procedures to improve phenotype accuracy
- \* super control vs. standard control vs. matched control
- \* extreme case/control imbalance may require firth logistic regression

## Quantitative trait

A continuous (or near-continuous) values where distances are meaningful

- \* technically discrete values may treat as quantitative if it has many levels
- \* under model assumption, transformation is usually needed but may not be easy for interpretation

Inverse normal transformation (INT):

$$\text{INT}(x) = \Phi^{-1} \left( \frac{\text{rank}(x) - c}{n - 2c + 1} \right) \begin{cases} c = \frac{3}{8} & (\text{Blom}, 1958) \\ c = \frac{1}{3} & (\text{Tukey}, 1962) \\ c = \frac{1}{2} & (\text{Bliss}, 1967) \end{cases}$$

## Covariates

Variables are included in the model to control confounding or improve accuracy

- \* age: onset age vs. recruitment age vs. restricted age vs. current age
- \* gender
- \* subpopulation structure: PCs
- \* others

## Model free

Compare **case vs control frequencies directly**, without specifying a regression model

Total (alleles)	A	B		AA $x_0 = 0$	AB $x_1 = 1$	BB $x_2 = 2$	Total (genotypes)
2R	$2r_0 + r_1$	$r_1 + 2r_2$	<b>case</b>	$r_0$	$r_1$	$r_2$	R
2S	$2s_0 + s_1$	$s_1 + 2r_2$	<b>control</b>	$s_0$	$s_1$	$s_2$	S
2N	$2n_0 + n_1$	$n_1 + 2n_2$		$n_0$	$n_1$	$n_2$	N

### Allele-based test

$$\chi^2 = \frac{2N[(2r_0+r_1)(s_1+2r_2)-(r_1+2r_2)(2s_0+s_1)]}{(2R)(2S)(2n_0+n_1)(2n_2+n_1)} \sim \chi^2(1)$$

### Genotype-based test

$$\chi^2 = \frac{1}{RS} \sum_{j=0}^2 \frac{(r_jS - s_jR)^2}{n_j} \sim \chi^2(2)$$

### Cochran–Armitage trend test

$$T = \sum_{j=0}^2 x_j (r_jS - s_jR), \frac{T^2}{\text{var}(T)} \sim \chi^2(1)$$

## Fixed-effect model

Account for **covariates** like Age, Sex, or Principal Components (PCs), and assume all predictors' impact are constant across the entire study population (**fixed Effects**)

covariates target-SNP

$$y = X\alpha + g\beta + \epsilon$$

BT: Firth logistic regression

QT: Linear regression

## Mixed-effect model

Extend linear model by adding **random effects** to model **genetic similarity** (complex population structures and cryptic relatedness)

Infinitesimal (polygenic) all SNPs

fixed random

population stratification  
familial relatedness

$$y = X\alpha + g\beta + X_G\gamma + \epsilon$$

Proximal contamination

$$y = X\alpha + g\beta + X_{LOCO}\gamma' + \epsilon$$

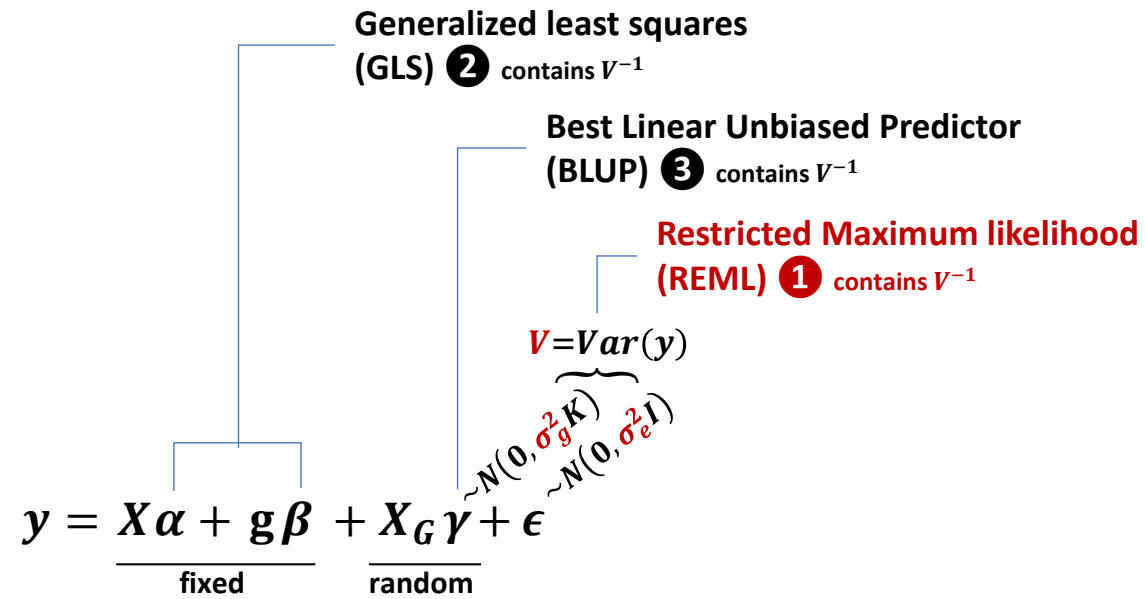
BT: SAIGE REGENIE

QT: BOLT-LMM SAIGE REGENIE

## Genomic Control

To reduce false positives inflated by population structure in GWAS (Devlin B & Roeder K. 1999)

$$\lambda_{GC} = \frac{\text{median}(\chi_{obs}^2)}{\text{median}(\chi_{null}^2)} \rightarrow \chi_{adj}^2 = \frac{\chi_{obs}^2}{\lambda_{GC}}$$



### BOLT-LMM

Loh P.R., et al. (2015) Nat Genet.

- \* a matrix-vector product + conjugate gradient
- \* REML (no explicit  $V$ )
- \* Quantitative trait (transformed by yourself)
- \* Binary trait (not recommended)
- \* Precompiled executable

### SAIGE

Zhou W., et al. (2018) Nat Genet.

- \* sparse GRM
- \* approximate REML
- \* Quantitative trait (INT provided)
- \* Imbalance case/control (SPA)
- \* Install from source codes

### REGENIE

Mbatchou, J., et al. (2021) Nat Genet.

- \* no GRM
- \* block ridge regression prediction
- \* Quantitative trait (INT provided)
- \* Imbalance case/control (Firth logistic regression)
- \* Precompiled executable

```
plink --bfile dat_as_sampQC_varQC --pheno
dat.phecov --1 --pheno-name PHENO1 --assoc --
out result_bt_allele
```

```
plink --bfile dat_as_sampQC_varQC --pheno
dat.phecov --1 --pheno-name PHENO1 --model --
out result_bt_model
```

```
regenie \
--step 1 \
--bed dat_as_sampQC_varQC \
--phenoFile dat.phecov \
--covarFile dat.phecov \
--covarCollist age,sex,PC1,PC2,PC3,PC4,PC5 \
--bt \
--bsize 1000
--lowmem --lowmem-prefix tmp_regenie \
--out regenie_step1
```

```
regenie \
--step 2 \
--bed dat_as_sampQC_varQC \
--phenoFile dat.phecov \
--covarFile dat.phecov \
--covarCollist age,sex,PC1,PC2,PC3,PC4,PC5 \
--bt --firth --approx --pThresh 0.01 \
--bsize 400 \
--pred regenie_step1_pred.list \
--out result_bt_regenie
```

```
plink2 \
--bfile dat_as_sampQC_varQC \
--pheno dat.phecov --1 \
--covar dat.phecov \
--covar-name age,sex,PC1,PC2,PC3,PC4,PC5 \
--covar-variance-standardize \
--glm no-x-sex hide-covar
cols=chrom,pos,ref,alt,a1freq,a1freqcc,gcountcc,n
obs,orbeta,se,tz,p,firth,err \
--out result_bt_plink2
```

Binary  
trait

Quantitative  
trait

```
bolt \
--bfile=dat_as_sampQC_varQC \
--phenoFile=dat.phecov --phenoCol=PHENO1 \
--covarFile=dat.phecov \
--covarCol=sex --qCovarCol=age,PC{1:5} \
--lmm \
--LDscoresFile=LDSCORE.1000G_EAS_m.tab \
--geneticMapFile=genetic_map_hg38_withX.txt
--numThreads=8 \
--statsFile=result_qt_bolt.txt \
```

```
Rscript step1_fitNULLGLMM.R \
--plinkFile=dat_as_sampQC_varQC \
--phenoFile=dat.phecov --phenoCol=PHENO1 \
--covarCollist=age,sex,PC1,PC2,PC3,PC4,PC5 \
--traitType=binary \
--outputPrefix=saige_null \
--nThreads=8
```

```
Rscript step2_SPAtests.R \
--bgenFile=dat_as_sampQC_varQC.bgen \
--sampleFile=dat_as_sampQC_varQC.sample \
--GMMATmodelFile=saige_null.rda \
--varianceRatioFile=saige_null.varianceRatio.txt \
--LOCO=TRUE \
--SAIGEOutputFile=result_bt_saige.txt
```

```
plink2 \
--bfile dat_as_sampQC_varQC \
--pheno dat.phecov \
--covar dat.phecov \
--covar-name age,sex,PC1,PC2,PC3,PC4,PC5 \
--covar-variance-standardize \
--glm no-x-sex hide-covar
cols=chrom,pos,ref,alt,a1freq,nobs,orbeta,se,tz,p \
--out result_qt_plink2
```

### .model

```
CHR SNP A1 A2
TEST (GENO, TREND, ALLELIC, DOM, REC)
AFF UNAFF
CHISQ DF P
```

### .regenie

```
CHROM GENPOS ID ALLELE0 ALLELE1
A1FREQ INFO
N TEST BETA SE CHISQ LOG10P
EXTRA
```

### .glm.logistic.hybrid

```
#CHROM POS ID REF ALT
CASE_NON_A1_CT CASE_HET_A1_CT
CASE_HOM_A1_CT CTRL_NON_A1_CT
CTRL_HET_A1_CT CTRL_HOM_A1_CT
A1_FREQ A1_CASE_FREQ A1_CTRL_FREQ
FIRTH? OBS_CT OR LOG(OR)_SE Z_STAT P
ERRCODE
```

Binary  
trait

Quantitative  
trait

```
SNP CHR BP GENPOS ALLELE1 ALLELE0
A1FREQ F_MISS
CHISQ_LIINREG P_LIINREG
BETA SE
CHISQ_BOLT_LMM_INF P_BOLT_LMM_INF
CHISQ_BOLT_LMM P_BOLT_LMM
```

### (binary trait)

```
CHR POS SNPID Allele1 Allele2
AC_Allele2 AF_Allele2 imputationInfo
N BETA SE Tstat Var
p.value (Saddlepoint approximation, SPA)
p.value.NA (Normal approximation)
ls.SPA.converge
varT varTstar AF.Cases AF.Controls
```

### (quantitative trait)

```
CHR POS SNPID Allele1 Allele2
AC_Allele2 AF_Allele2 imputationInfo
N BETA SE Tstat p.value varT varTstar
```

### .glm.linear

```
#CHROM POS ID REF ALT
A1_FREQ
OBS_CT BETA SE T_STAT P
```

**:Figure**





**Thanks for your attention!!**

**<(\_ \_)>**