

# Population Stratification Practical

# Overview

- Use the SNP data to empirically investigate the ancestry of the samples
- Assign individuals to groups
- Capture Chinese/Japanese distinction?
- Generate a multi-dimensional scaling plot

# 1. Create Independent Set of SNPs

- clustering/MDS/PCA/etc... work best on sets of SNPs that are independent
- plink option
  - --indep-pairwise  $\langle N \rangle$   $\langle M \rangle$   $\langle r^2 \rangle$ 
    - N = windows size
    - M = window shift
    - $r^2$  = threshold for removal

## *Filter SNPs based on LD*

```
plink --bfile wgas3  
      --indep-pairwise 50 10 0.2  
      --out prune1
```

*Output: prune1.pruned.in*

```
rs4040617  
rs4075116  
rs9442385  
rs11260562  
rs6685064  
rs7519837  
rs3855951  
....
```

*Output: run1.pruned.out*

```
rs3094315  
rs3766180  
rs6603791  
rs3737628  
rs7511905  
rs3107157  
...  
....
```

## 2. Measure IBS between Individual

- IBS = Identity-by-State
  - Measure of total genotype similarity
  - plink commands
    - --genome
    - --extract prune1.prune.in      *or*
- exclude prune1.prune.out



# ibs1.genome columns

- RT = relationship from pedigree
- EZ = expected IBD sharing
- PI\_HAT = observed IBD sharing  
=  $P(\text{IBD}=2) + 0.5 * P(\text{IBD}=1)$
- PHE = pairwise phenotypic code  
(1,0,-1 = AA, AU and UU pairs)
- **RATIO = HETHET : IBS 0 SNPs (expected 2)**  
(HETHET = Number of IBS2 het/het SNP pairs)
- PPC = p-value of deviation of RATIO from 2

# Clustering

- Use the IBS values to perform clustering
- plink commands
  - `--read-genome <genome file>`
  - `--cluster`
  - `--mds-plot <N>`
- Add restrictions
  - PPC test (`--ppc 1e-3`)
  - Clusters contain 1 case and 1 control (`--cc`)



## Clustering

```
plink --bfile wgas3
      --read-genome ibs1.genome
      --cluster
      --ppc 1e-3
      --cc
      --mds-plot 2
      --out strat1
```

### Output: strat1.genome

```
CH18526 NA18526 0
CH18524 NA18524 0
CH18529 NA18529 0
CH18558 NA18558 1
CH18532 NA18532 1
CH18561 NA18561 0
....
JA18987 NA18987 3
JA18990 NA18990 2
JA18991 NA18991 2
JA18994 NA18994 2
JA18992 NA18992 4
JA18997 NA18997 2
...
```

## Clustering

```
plink --bfile wgas3
      --read-genome ibs1.genome
      --cluster
      --ppc 1e-3
      --cc
      --mds-plot 2
      --out strat1
```

### Output: strat1.mds

FID	IID	SOL	C1	C2
CH18526	NA18526	0	-0.0245884	0.00917367
CH18524	NA18524	0	-0.0242271	-0.0278564
CH18529	NA18529	0	-0.0182157	0.00810387
CH18558	NA18558	1	-0.0233174	-0.0318428
CH18532	NA18532	1	-0.023291	-0.000359224
CH18561	NA18561	0	-0.0219817	-0.00855256
CH18562	NA18562	1	-0.0203434	-0.0379835
CH18537	NA18537	0	-0.0202644	-0.0320186
CH18603	NA18603	1	-0.0247671	0.0113186
CH18540	NA18540	1	-0.0236313	0.00147393
CH18605	NA18605	0	-0.0223527	-0.00800508

...

# Make MDS Plot

- Open R
- Navigate to your working directory

```
> d <- read.table("strat1.mds", header=TRUE)
> plot( d$C1, d$C2, pch=20, cex=2, col = d$SOL+1)

> pop = numeric(nrow(d))
> pop[which(substr(d[,1],1,1)=="C")] = 1
> pop[which(substr(d[,1],1,1)=="J")] = 2
> plot( d$C1, d$C2, pch=20, cex=2, col = pop)
```



## ***Association using MDS***

```
plink --bfile wgas3  
      --mh  
      --within strat1.cluster2  
      --out mdsaccoc1
```

```
plink --bfile wgas3  
      --logistic  
      --covar strat1.mds  
      --covar-number 1,2  
      --out mdsaccoc2
```