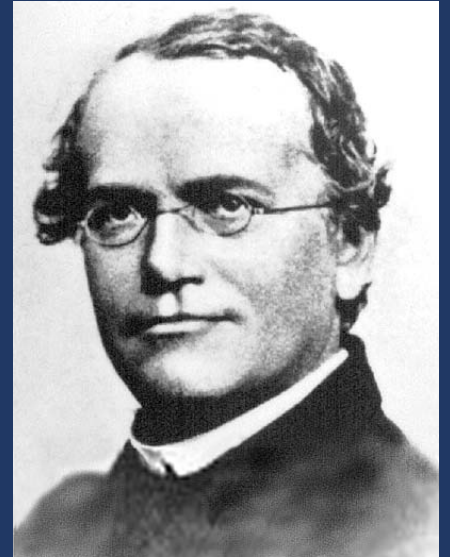# SnpArrays.jl, VCFTools.jl, and ADMIXTURE.jl

OpenMendel Workshop

ASHG Annual Meeting 2020

Hua Zhou (huazhou@ucla.edu)

# Motivation

SnpArrays.jl and VCFTools.jl are two gateway packages that

- Translate (potentially big) genotype data files to vectors/matrices for statistical analysis

- Perform common manipulations such as subsetting, summary statistics, filtering, quality control, GRM calculation, PCA calculation, and so on in a **high-level language** environment

I'll demonstrate functionalities of SnpArrays.jl here. Those for VCFTools.jl are similar.

# Read in Plink binary format

Given a set of Plink files

```
1  cd("/Users/huazhou/.julia/dev/SnpArrays/data/")
2  readdir(glob"mouse.*")
```

```
3-element Array{String,1}:
 "mouse.bed"
 "mouse.bim"
 "mouse.fam"
```

we initialize a SnpArray

```
1  mouse = SnpArray("mouse.bed")
```

```
1940×10150 SnpArray:
 0x02  0x02  0x02  0x02  0x03  0x02  …  0x03  0x03  0x03  0x03  0x03  0x03
 0x02  0x02  0x03  0x02  0x02  0x02     0x03  0x03  0x03  0x03  0x03  0x03
 0x03  0x03  0x03  0x03  0x03  0x03     0x03  0x03  0x03  0x03  0x03  0x03
 0x02  0x02  0x02  0x02  0x02  0x02     0x03  0x03  0x03  0x03  0x03  0x03
 0x03  0x03  0x03  0x03  0x03  0x03     0x02  0x02  0x02  0x02  0x02  0x02
```

# Gzipped Plink files?

No problem. Use the same command

```
1  # requires corresponding `.fam.gz` file
2  SnpArray("mouse.bed.gz")
```

```
1940×10150 SnpArray:
 0x02  0x02  0x02  0x02  0x03  0x02  …   0x03  0x03  0x03  0x03  0x03  0x03
 0x02  0x02  0x03  0x02  0x02  0x02      0x03  0x03  0x03  0x03  0x03  0x03
 0x03  0x03  0x03  0x03  0x03  0x03      0x03  0x03  0x03  0x03  0x03  0x03
```

Other supported compressed formats are

```
1  SnpArrays.ALLOWED_FORMAT
```

```
6-element Array{String,1}:
 "gz"
 "zlib"
 "zz"
 "xz"
 "zst"
 "bz2"
```

# Translate genotypes to numbers

Codebook:

| Genotype | SnpArray | model=ADDITIVE_MODEL | model=DOMINANT_MODEL | model=RECESSIVE_MODEL |
|---|---|---|---|---|
| A1,A1 | 0x00 | 0 | 0 | 0 |
| missing | 0x01 | NaN | NaN | NaN |
| A1,A2 | 0x02 | 1 | 1 | 0 |
| A2,A2 | 0x03 | 2 | 1 | 1 |

# Convert SnpArray to a numeric matrix

```
1  convert(Matrix{Float64}, mouse;
2      model=ADDITIVE_MODEL, center=false, scale=false, impute=false)
```

```
1940×10150 Array{Float64,2}:
1.0  1.0  1.0  1.0  2.0  1.0  2.0  1.0  …   2.0   2.0   2.0   2.0   2.0
1.0  1.0  2.0  1.0  1.0  1.0  1.0  2.0      2.0   2.0   2.0   2.0   2.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  1.0  1.0  1.0  1.0  1.0  2.0      2.0   2.0   2.0   2.0   2.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      1.0   1.0   1.0   1.0   1.0
1.0  1.0  1.0  1.0  2.0  1.0  2.0  1.0  …   2.0   2.0   2.0   2.0   2.0
1.0  1.0  1.0  1.0  2.0  1.0  2.0  1.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  2.0  1.0  1.0  1.0  1.0  2.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  2.0  1.0  1.0  1.0  1.0  2.0      2.0   2.0   2.0   2.0   2.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      1.0   1.0   1.0   1.0   1.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0  …   0.0   0.0   0.0   0.0   0.0
1.0  1.0  1.0  1.0  2.0  1.0  2.0  1.0      2.0   2.0   2.0   2.0   2.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      0.0   0.0   0.0   0.0   0.0
⋮                              ⋮          ⋱   ⋮
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      2.0   2.0   2.0   2.0   2.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  1.0  1.0  1.0  1.0  1.0  2.0  …   2.0   2.0   2.0   2.0   2.0
1.0  1.0  1.0  1.0  1.0  1.0  2.0  1.0      2.0   2.0   2.0   2.0   2.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  1.0  1.0  2.0  1.0  2.0  1.0      2.0   2.0   2.0   2.0   2.0
2.0  2.0  2.0  2.0  2.0  2.0  2.0  2.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  1.0  1.0  2.0  1.0  2.0  1.0  …   2.0   2.0   2.0   2.0   2.0
1.0  1.0  2.0  1.0  1.0  1.0  1.0  2.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  2.0  1.0  1.0  1.0  1.0  2.0      2.0   2.0   2.0   2.0   2.0
1.0  1.0  1.0  1.0  1.0  1.0  1.0  2.0      NaN   NaN   NaN   NaN   NaN
0.0  0.0  0.0  0.0  2.0  0.0  2.0  0.0      2.0   2.0   2.0   2.0   2.0
```

# Convert a column of genotypes

```julia
1  # convert(Vector{Float64}, view(mouse, :, 1)) # alternative syntax
2  # @views convert(Vector{Float64}, mouse[:, 1]) # alternative syntax
3  convert(Vector{Float64}, @view(mouse[:, 1]))
```

```
1940-element Array{Float64,1}:
 1.0
 1.0
 2.0
 1.0
 2.0
 1.0
 1.0
 1.0
 1.0
 2.0
```

# Summary statistics

- Minor allele frequencies (MAF)

```
1  maf(mouse)
```

```
10150-element Array{Float64,1}:
 0.4434984520123839
 0.4438144329896907
 0.359504132231405
 0.4439855446566856
 0.13119834710743805
 0.44404332129963897
 0.1412706611570248
```

# Summary statistics

Missing rate by SNP

```
1  missingrate(mouse, 1)
```

```
10150-element Array{Float64,1}:
 0.0010309278350515464
 0.0
 0.0020618556701030930
```

Missing rate by individual

```
1  missingrate(mouse, 2)
```

```
1940-element Array{Float64,1}:
 0.00019704433497536947
 0.0
 0.018423645320197045
```

# Empirical kinship matrix

SnpArrays.jl implements commonly used methods for calculating empirical kinship matrix.

```
1  # grm(mouse, method=:MoM)
2  # grm(mouse, method=:Robust)
3  grm(mouse, method=:GRM, minmaf=0.05)
```

```
1940×1940 Array{Float64,2}:
  0.478556    -0.0331783    0.013541    …   -0.0348225   -0.0129761
 -0.0331783    0.422993    -0.0389741        0.0457975    0.00554753
  0.013541    -0.0389741    0.50952         -0.0357183   -0.0609305
  0.0203209    0.00777944  -0.00887047      -0.0297696   -0.00972836
  0.0567523   -0.0163798   -0.00498406      -0.0413874   -0.0416146
 -0.0166009   -0.0191523   -0.0112531    …   0.0176939   -0.0193442
```

# Not covered here

- Filtering
- Linear algebra on SnpArray directly without converting to numerical matrix
- Linear algebra on SnpArray directly on GPU
- Split, concatenate, merge, and reordering of Plink files
- VCF to Plink

Learn from

- Package documentation: https://openmendel.github.io/SnpArrays.jl/latest/
- Workshop tutorial (Binder link): https://mybinder.org/v2/gh/OpenMendel/ASHG-OpenMendelWorkshop-2020-Oct/master?filepath=04-SnpArrays-HZhou%2FSnpArraysTutorial.ipynb

# Demo: empirical kinship from admixed samples

Step 1: Run Admixture program (available on Mac and Linux) to calculate ancestry fractions (per sample) and population specific allele frequencies (per SNP)

```
1  using ADMIXTURE, SnpArrays
2
3  # Step 1: Run Admixture with K=3 populations and j=4 threads
4  P, Q = admixture("mouse.bed", 3, j=4)
```

# Demo: empirical kinship from admixed samples

Step 2: Use the P (allele frequencies) and Q (ancestry fractions) matrix from Admixture to compute the kinship coefficients

```
1  grm_admixture(mouse, P', Q')
```

```
convert genotype: 0.38 seconds
Φ = GG': 0.63 seconds
convert G to {0,1} matrix: 0.03 seconds
S = GG': 0.63 seconds

1940×1940 Array{Float64,2}:
   0.459157     -0.0156932    -0.00323859  …  -0.0241032     0.0122942
  -0.0156932     0.382539     -0.00891408      0.00213638    0.0256826
  -0.00323859   -0.00891408    0.488814       -0.0171976    -0.046833
   0.00309657    0.0202297    -0.0243852      -0.0176487     0.00056089
   0.0332112     0.00312463   -0.0272922      -0.0327539    -0.0122725
  -0.0358108    -0.00853747   -0.0271514   …   0.0243008    -0.00138405
   0.121113     -0.0392917     0.00165833      0.0108531    -0.0377001
```

# Demo: calculate empirical kinship from admixed samples

Timing results on the ACCORD data with ~6,000 individuals from 4 populations, each with ~800,000 SNPs

- SnpArrays.jl takes about **6 minutes**
- REAP (C Program) takes about **10 hours**

# Is that easy to write Julia code?

```julia
73  function grm_admixture(
74      s::AbstractSnpArray,
75      P::AbstractMatrix,
76      Q::AbstractMatrix,
77      ::Type{T} = Float64
78      ) where T <: AbstractFloat
79      m, n = size(s)
80      # convert genotype
81      tic = time()
82      G = Mmap.mmap(Matrix{T}, m, n) # slightly faster than G = Matrix{T}(undef, m, n)
83      Base.copyto!(G, s, P, Q, model=ADDITIVE_MODEL, center=true, scale=true)
84      @printf("convert genotype: %.2f seconds\n", time() - tic)
85      # GG'
86      tic = time()
87      Φ = G * transpose(G) # m x m
88      @printf("Φ = GG': %.2f seconds\n", time() - tic)
89      # convert G to {0,1} matrix
90      tic = time()
91      @inbounds for (idx, g) in enumerate(G)
92          G[idx] = ifelse(iszero(g), T(0), T(1))
93          # iszero(g) || (G[idx] = T(1))
94      end
95      @printf("convert G to {0,1} matrix: %.2f seconds\n", time() - tic)
96      # Sij = # SNPs observed in both individuals i and j
97      tic = time()
98      S = G * transpose(G)
99      @printf("S = GG': %.2f seconds\n", time() - tic)
100     # Φ /= 2S
101     @inbounds for j in 1:m, i in 1:j
102         Φ[i, j] /= 2S[i, j]
103     end
104     copytri!(Φ, 'U')
105 end # function grm_admixture
```