# Main Scripts During Postgraduate

## 3.绘图

### 1、单倍型分析

### 2、单倍型聚类

#### ①NJ树

#### ②contree

#### ③豪哥的办法

#### ④VCF2Dis

### 3、单倍型网络图

## 曼哈顿图系列

### 1、TASSEL结果画曼哈顿图

### 2、注释显著位点

### 3、R语言CMplot包绘制曼哈顿图

### 5、Candihap画图

### 6、单倍型群体结构可视化-haplostrips

https://bitbucket.org/dmarnetto/haplostrips/src/master/

### 7、fastphase和beagle原理相同-提取候选区段

### 8、fst

manhadun:https://www.plob.org/article/21662.html

## 图片组合

## ！计算核苷酸多样性（π）

画图1：https://cloud.tencent.com/developer/article/1593265

## 折线图

https://www.jianshu.com/p/7b6701ae593b

# 1、GWAS前：文件准备工作（过滤按照第一次过滤，再进行填充）

## 1.加ID号

https://www.jianshu.com/p/2677dc3b1383 方法1

```
1
2  方法二
3  /data/apps/bcftools/1.9/bin/bcftools annotate --set-id +'%CHROM\_%POS' 938.raw.vcf.gz |/
```

```
1  gunzip
2  ## add variance id to vcf file（看vcf有无ID，需要加上）
3  perl lecture06_07_add_id.pl 516sample.vcf combine_data.vcf
4
```

## 最终过滤版本

```
1  gunzip
2  ## add variance id to vcf file（看vcf有无ID，需要加上）
3  perl lecture06_07_add_id.pl /w/00/g/g01/user305/Result/00.450_rawvcf/combine_450raw.vcf
4  ## filter by quality and depth（最小质量值、最小/最大深度、平均深度）
5  ## 可以用le log查看
6  ## plink对vcf过滤比较弱，如果要用plink先用vcftools进行过滤然后再用plink进行其他操作
7  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /w/00/g/g01/user305/Result/TResul
8
9  开始
10 ## keep only indels 分开indel和snp也可以使用gatk
11 ## 要输出vcf文件一定要加上--recode
12 ##vcftools --vcf filter1.vcf --keep-only-indels --recode --recode-INFO-all --out lecture0
13 ## keep only snps
14 /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf filter.vcf --remove-indels --reco
15
16
17 #1.SNP过滤
18 ## missing rate， maf，allele number filtering
19 /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf snp_only_filter.vcf --max-missing
```

```
20
```

## 查看错误位点

```
1  ## 用vcftools提取指定区间
2  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf new.filter.vcf  --chr At1 --from-
3
```

## 填充

/w/00/g/g01/user305/Result/TResult/01.filter_method2_mer5.vcf/01_afilter

```
1  cat snp_bi.filter.vcf | perl -pe "s/\s\.:/\t.\/.:/g" | /data/apps/htslib/1.9/bin/bgzip -
```

## 2、课题组过滤

```
1  ***重要***
2  每一步跑完，记得用bcftools输出一个stats文件，以便于掌握过滤之后位点的情况
3  /data/apps/bcftools/1.9/bin/bcftools stats new.filter.vcf.gz>new.filter.vcf.gz.stats
```

```
1  1.标记质量值小于30和indel5bp
2  /data/apps/bcftools/1.9/bin/bcftools filter -O z  -o genetalks-mihoutao.filter_variant.v
3  2.删除低质量和indel5bp
4  gzip -dc genetalks-mihoutao.filter_variant.vcf.gz |grep -v 'LOWQUAL'|grep -v 'SnpGap' |/
```

```
1
```

```
1  3.2比对个体
2  mkdir output.01.bam
3  $time -vo output.01.bam/Pd14008_04.01.bam.time bash -c "$bwa mem -t 16 -R '@RG\tID:Pd140(
4  $time -vo output.01.bam/Pd14008_04.02.bam.time $java -Xmx4G -jar  $picard SortSam COMPRES
5  $time -vo output.01.bam/Pd14008_04.03.bam.time $java -Xmx30G -jar $picard MarkDuplicates
```

```
6   $time -vo output.01.bam/Pd14008_04.04.bam.time $java -Xmx4G -jar  $picard SortSam COMPRE
7   $time -vo output.01.bam/Pd14008_04.05.bam.depth.time bash -c "$samtools depth output.01.
```

```
1   /data/apps/bcftools/1.9/bin/bcftools view -m2 -M2 -v snps genetalks-mihoutao.filter_vari
2
3   4.依据深度过滤（根据公司返回的报告）
4   过滤深度大于平均深度3倍的位点和小于平均深度1/3的位点
5   #####这里要注意，个体原始文件顺序和算深度的数据，个体顺序要一致，即脚本输出顺序与vcf文件里面的个
6   首先拿平均深度
7   过滤
8   注意输入输出都为gz文件
9   python filterstep.depth.py Indepth.txt Mp.cp.Low.Gap.2all.vcf.gz Mp.QUL30.indel5bp.alle2
10
```

```
1   5.个性化过滤
2   /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --gzvcf Mp.QUL30.indel5bp.alle2.depth.v
```

```
1   然后使用下面两个选项对vcf文件保留或者删除样品。
2   --keep <filename>    保留样品
3   --remove <filename>    删除样品
4   代码如下：
5   /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /w/00/g/g01/user305/Result/00.448_
```

```
1   6.画structure
2   ## missing rate and maf filtering
3   /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/01.gx_first
4   ## make bed and LD pruning
5   /w/00/g/g01/user305/soft/plink/plink --vcf snp.int0.8.maf0.05.vcf --indep-pairwise 100 50
6
7   ##extract prune in snps and recode to admixture format（转成ped格式）
```

```
8  /w/00/g/g01/user305/soft/plink/plink --bfile snp.int0.8.maf0.05 --extract snp.int0.8.maf0
9
10
11  #课题组方法无法完整的画图，先得到群体结构的结果
12  for i in {1..15}; do /data/apps/admixture/1.3.0/admixture --cv -j30 --seed=666 prunData.p
13  grep -h CV log*.out
```

```
At1:12861335 [*]
At1:12861804 [*]
At1:12848896 [*]
```

```
of A, C, T, G, or N characters at At1:10437381 [*]
```

## 3.过滤2/3再beagle再过滤1/3

```
1   source ~/.profile
2   conda
3
4   #0.下载下来的其实是.gz文件，后续命令可以直接调用gz文件，也可以解压之后再调用，但必须每个都先做in
5   gatk IndexFeatureFile -I /w/00/g/g01/user305/Result/TResult/mer5.vcf/combine_data.vcf
6
7   ## 过滤1: gatk VariantFiltration : --filter-expression 参考：https://www.omicsclass.com/ar
8
9   gatk --java-options "-Xmx4g" VariantFiltration -R /w/00/g/g01/user305/data/ref/Tibet.fast
10
11
12  gzip all.raw.gatk.vcf
13
14
15  zcat all.raw.gatk.vcf.gz |awk '$0~/#/ || ($7 =="PASS"){print $0}' |gzip - > all.raw.gatke
16
17  #过滤2: vcfutils.pl
18  # -w INT     SNP within INT bp around a gap to be filtered [3]
19  # -W INT     window size for filtering adjacent gaps [10]
20
21  /data/apps/bcftools/1.9/bin/vcfutils.pl varFilter -w 5 -W 10 "gzip -dc all.raw.gatked.vc
22
```

```
1
2  # 2. 补缺失impute，使用Beagle软件进行补缺失
3  gzip combine_data.vcf
4
5  #zcat input.rmUnusedAlt.vcf.gz | sed -i "/\tAt1__9166711\t/d" input.rmUnusedAlt.vcf.gz>
6
7
8  cat filt_done.vcf | perl -pe "s/\s\.:/\t.\/.:/g" >output.vcf
9  sed -e '/At1__10437381/d' filt_done.vcf > binput.vcf
10 sed -e '/At1__12861335/d' filt_done.vcf > binput.vcf
11 sed -e '/At1__12861804/d' filt_done.vcf > binput.vcf
12 sed -e '/At1__12848896/d' filt_done.vcf > binput.vcf
13 sed -e '/At1__10437381/d' filt_done.vcf > binput.vcf
14
15
16 https://zhuanlan.zhihu.com/p/109071456
17 #vcf转为 ped/map
18 /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/01.filter_
19
20 ## ped/map转换为bed/bim/fam
21 /w/00/g/g01/user305/soft/plink/plink --file tran --make-bed --out tran_test --allow-extra
22
23 /w/00/g/g01/user305/soft/plink/plink --bfile tran_test --list-duplicate-vars ids-only su
24 /w/00/g/g01/user305/soft/plink/plink --bfile tran_test --alleleACGT --snps-only just-acgt
25 ## bed/bim/fam 转为 vcf
26 /w/00/g/g01/user305/soft/plink/plink --bfile wbinput --export vcf --out bbinput --allow-
27
28
29
```

beagle填充REF field is not a sequence of A, C, T, G, or N characters at At1:10437381 - Google 搜索    v

```
#CHROM    POS      ID       REF     ALT     QUAL      FILTER  INFO
At1       10437381  .        .       A       T,*       6882.63 .
At1       10437384  .        .       A       AT,*      6890.79 .
```

## 提取指定区间

```
1  ## 用vcftools提取指定区间
2  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf combine_raw.vcf  --chr At1 --from
3
```

```
1
2  #下载件
3  wget https://faculty.washington.edu/browning/beagle/beagle.28Sep18.793.jar
4  # Pre-phasing
5  #该步骤为非必需步骤，对输入的VCF文件进行预phasing，理论上可以提高基因型推断质量
6  # 考虑换beagle最新款
7  -Xms2g -Xmx30g
8  /data/apps/jdk/jdk1.8.0_131/bin/java -Xms2g -Xmx30g -jar /w/00/g/g01/user305/soft/beagle
9
10
11 #注意：如果提示没有安装java，需要额外执行如下两句：
12 sudo apt update
13 sudo apt install default-jre     4373
```

```
1
2  #过滤3：vcftools
3  #--max-missing Exclude sites on the basis of the proportion of missing data
4  #(defined to be between 0 and 1, where 0 allows sites that are completely missing
5  #and 1 indicates no missing data allowed).
6
7  !!!!!#第3次过滤分离变异前（原始）
8  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --gzvcf all.varFilter.vcf.gz --recode --
```

```
9
```

```
10  #分离indel与SNP 到不同的文件
11  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --remove-indels --recode --recode-INFO-
12
13
14  zcat demo.filter.vcf.gz|wc -l
15  gunzip demo.filter.vcf
```

## 提取指定个体

```
1  然后使用下面两个选项对vcf文件保留或者删除样品。
2  --keep <filename>   保留样品
3  --remove <filename>    删除样品
4  代码如下：
5  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /data/apps/vcftools/vcftools-0.1.1
6
```

## 提取候选基因

```
1  awk '$3=="gene"' Tibet.gff3 |awk '$1=="At1"&& $4>20580546 && $5<20620546' > /w/00/g/g01/u
2
```

## 4.annovar注释

```
1  # 安装gffread和gtfToGenePred
2  #git clone https://github.com/gpertea/gffread
3  #下载annovar并解压缩: unzip zipped_file.zip -d unzipped_directory
4
5  #提取目标参考基因组
6  python /w/00/g/g01/user305/scripts/extract_seq.py
7
8
9  # build ANNOVAR index
10 /w/00/g/g01/user305/soft/gffread/gffread Tibet.gff3 -T -o Tibet.gtf
11
```

```
12  /w/00/g/g01/user305/soft/gtfToGenePred -genePredExt /w/00/g/g01/user305/soft/Annovar/tibe

13

14  perl retrieve_seq_from_fasta.pl --format refGene --seqfile tibetdb/Tibet.fasta  tibetdb/u

15

16

17  ##为了避免别的脚本权限问题，干脆全部777

18  #chmod 777 *  成功解决权限问题。

19  #一定要在当前目录下进行

20  perl table_annovar.pl /w/00/g/g01/user305/Result/TResult/1.gx_first_mer5.vcf/05.gx_indoo

21

22

23

24

25  #测试之前成功的:  perl table_annovar.pl /w/00/g/g01/user305/Result/TResult/first_filter_me
```

## 5.群体结构分析

https://zhuanlan.zhihu.com/p/109164405

## structure（改50/10/0.2

```
1   2.1 利用plink按照LD衰减过滤vcf文件（每100kb保留50个snp，筛选标准是R2为0.2）
2   /w/00/g/g01/user305/soft/plink/plink --vcf out_417individual.vcf --indep-pairwise 100 50
3   /w/00/g/g01/user305/soft/plink/plink --vcf out_417individual.vcf --recode vcf-iid --extr

4

5

6

7

8   2.2 利用plink转换为structure可以识别的格式
9   /w/00/g/g01/user305/soft/plink/plink --bfile LD_filter --extract LD_filter.prune.in --ou
10  2.3 structure运行，linux只需这个命令（参数都在文件里面）
11  structure -m test_mainpara1_1 -e extraparams
```

## 改

```
1   2.1 利用plink按照LD衰减过滤vcf文件（每100kb保留50个snp，筛选标准是R2为0.2）
2   /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/3.HJ/01.29
3   /w/00/g/g01/user305/soft/plink/plink --bfile snp --indep-pairwise 50 10 0.2 --out snp.ld
4   /w/00/g/g01/user305/soft/plink/plink --bfile snp --extract snp.ldprun.prune.in --out pru
```

```
5
6
7
8  sed -i 's/^At//g' prunData.bim
9  for i in {1..25}; do /data/apps/admixture/1.3.0/admixture --cv -j30 --seed=666 prunData.
10
11 grep -h CV log*.out
12
```

plink --bfile snp --maf 0.05 --geno 0.5 --indep-pairwise 100 10 0.05 --out snp.ldprun --allow-extra-chr

# extract prun in snps

# prune出来的以后转成bed文件; record有很多种格式

plink --bfile snp --extract snp.ldprun.prune.in --out prunData --make-bed --allow-extra-chr

# vcf 转成structur

--recode12

#admixture接受plink格式, bed格式

/w/00/g/g01/user305/soft/plink/plink --vcf admixture.in.vcf --make-bed --out demo.admixture.in

## 配置文件

```
Basic Program Parameters

#define MAXPOPS    13      // (int) number of populations assumed
#define BURNIN   50000   // (int) length of burnin period
#define NUMREPS  500000   // (int) number of MCMC reps after burnin
```

```
number of samples:     490
number of records:    720210
number of no-ALTs:     0
number of SNPs: 720210
```

```
1  # 查看位点信息
2  /data/apps/bcftools/1.9/bin/bcftools stats demo.impute.vcf>demo.impute.vcf.stats
3
4  #outfile
5  /w/00/g/g01/user305/Result/TResult/mer5.vcf/2.third_filter/03.structure/str_1_1
6  #infile
7  /w/00/g/g01/user305/Result/TResult/mer5.vcf/2.third_filter/03.structure/prunData.recode.s
```

```
 8  #样本量
 9  490
10  #标记数
11  720210
12
13
```

## ①自己研究的方法

Linkage disequilibrium based SNP pruning

```
--indep-pairwise 50 5 0.5
```

```
--indep-pairwise 50 5 0.2
```

a) consider a window of 50 SNPs

b) calculate LD between each pair of SNPs in the window

b) remove one of a pair of SNPs if the LD is greater than 0.5

c) shift the window 5 SNPs forward and repeat the procedure

```
 1  ## missing rate and maf filtering
 2  /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/mer5.vcf/3
 3  ## make bed and LD pruning
 4  /w/00/g/g01/user305/soft/plink/plink --vcf snp.int0.8.maf0.05.vcf --indep-pairwise 100 50
 5
 6  ##extract prune in snps and recode to admixture format（转成ped格式）
 7  /w/00/g/g01/user305/soft/plink/plink --bfile snp.int0.8.maf0.05 --extract snp.int0.8.maf0
 8
 9  #R需要先安装RSvgDeviece包
10  /data/apps/R/3.5.1/bin/R
11  >install.packages("RSvgDevice")
12  选的兰州
13
14  #for i in {1..30}; do /data/apps/admixture/1.3.0/admixture --cv -j30 --seed=666 /w/00/g/g
```

```
15  #perl plotAdmixtureCV.pl /data/apps/R/3.5.1/bin/Rscript output output.cv.pdf        #要把所有
16  #课题组方法无法完整的画图，先得到群体结构的结果
17  for i in {1..30}; do /data/apps/admixture/1.3.0/admixture --cv -j30 --seed=666 /w/00/g/g0
18  grep -h CV log*.out
19
20  30
```

```
1   ##修改bim的名称，否则会报错
2   #sed -i 's/^At//g' /w/00/g/g01/user305/Result/TResult/mer5.vcf/3.1_second_filter/03.stru
3   #for K in {1..50}; do /data/apps/admixture/1.3.0/admixture --cv  /w/00/g/g01/user305/Resu
4   #grep -h CV log*.out
5   选择CV error最低的对应K值，.Q文件为要使用到的文件
6
7   #群体结构绘图
8   #先生成一个csv文件，然后再第一列添加上分组信息
9   paste <(cut -f1 prunData.nosex) prunData.10.Q |sed 's/^/G\t/g' |sed 's/\s/,/g'>plotQ.csv
10  #plotQ.csv可以上传到以下网址进行可视化：
11  http://omicsspeaks.com/strplot2/
12
```

K in {1..36}

## ②课题组方法

```
1   #计算LD
2   /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --gzvcf demo.impute.vcf.gz --ld-window-l
3   删除LD位点
4   /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --gzvcf Mp.QUL30.indel5bp.alle2.depth.r
5
6
```

```
1   $bgzip = /data/apps/htslib/1.9/bin/bgzip
2   $vcftools = /data/apps/vcftools/vcftools-0.1.15/bin/vcftools
3   $plink = /data/apps/plink/20181202/plink
4   $admixture = /data/apps/admixture/1.3.0/admixture
```

```
5  $Rscript = /data/apps/R/3.5.1/bin/Rscript
6  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --gzvcf demo.impute.vcf.gz --stdout --th
```

## 1、过滤并生成bgzip文件

对于大型的VCF文件而言，如何快速访问其中的记录也是个难点。tabix可以对VCF文件构建索引，索引构建好之后，访问速度会快很多。tabix对VCF文件建立索引的用法如下

注意输入的VCF文件必须是使用bgzip压缩之后的VCF文件，生成的索引文件为view.vcf.gz.tbi，后缀为.tbi。

构建好索引之后，可以快速的获取指定区域的记录，示例如下

## https://www.jianshu.com/p/07640917764d Tabix— bgzip操作

```
1  #建立索引
2  tabix -p vcf view.vcf.gz
3  #1、获取位于11号染色体的SNP位点
4  tabix view.vcf.gz 11
5  #2、获取位于11号染色体上突变位置大于或者等于2343545的SNP位点
6  tabix view.vcf.gz 11:2343545
7  #3、获取位于11号染色体上突变位置介于2343540到2343596的SNP位点
8  tabix view.vcf.gz 11:2343540-2343596
9
```

```
1  1、重命名
2  perl vcfRenameChr.pl /data/apps/htslib/1.9/bin/bgzip variants.demo.impute.vcf.gz output/I
3  2、转换格式
4  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --gzvcf output/Lmar.rename.chr.vcf.gz --
5  /data/apps/plink/20181202/plink --noweb --file output/Lmar.rename.plink --recode12 --out
6  for i in 1 2 3 4 5 6 7 8 9; do /data/apps/admixture/1.3.0/admixture --cv -j30 --seed=666
7  perl plotAdmixture.pl /data/apps/R/3.5.1/bin/Rscript output/Lmar.rename.plink12.ped Lmar
8  perl plotAdmixtureCV.pl /data/apps/R/3.5.1/bin/Rscript output output.cv.pdf     #要把所有的
9
10
11
```

## 6、亲缘关系分析-GCTA

```
1  ## convert to plink bed（GCTA接受bed格式，此处使用过滤后的vcf文件）
2  /w/00/g/g01/user305/soft/plink/plink --vcf demo.impute.vcf --make-bed --out snp --allow-
3
4
5  ## gcta grm matrix
6  ①输出亲缘关系矩阵②和③为文件前缀④常染色体的数目此处应该为5
7  /w/00/g/g01/user305/soft/gcta_1.91.6beta/gcta_1.91.6beta/gcta64 --make-grm-gz --out snp.g
8
9
```

```
1  ## convert to normal format
2  setwd("D:\\work_dir\\GWAS\\00.project\\000000.hifi\\01.gx\\02.leaves\\11.kinship")
3  library(reshape2)
4  tmp <- read.table(gzfile("snp.gcta.grm.gz"), header = F, stringsAsFactors = F)
5  ids <- read.table("snp.gcta.grm.id", header = F, stringsAsFactors = F)
6  tmp <- tmp[,c(1,2,4)]
7  result_matrix <- acast(tmp, V1~V2, value.var="V4", drop = F)
8  makeSymm <- function(m) {
9    m[upper.tri(m)] <- t(m)[upper.tri(m)]
10   return(m)
11 }
12
```

```
1  #转换为全矩阵
2  result_full <- makeSymm(result_matrix)
3  diag(result_full) <- 2    #这一步选做
4  result_df <- as.data.frame(result_full)
5  row.names(result_df) <- ids$V2
6  colnames(result_df) <- ids$V2
7  #生成的gcta.kinship.txt可以直接用作GWAS分析
8  write.table(result_df, file = "gcta.kinship.txt", row.names = T, col.names = NA, sep = "
9
```

```
 3
 4  ##修改bim的名称，否则会报错
 5  sed -i 's/^At//g' snp.bim
 6  ## gcta grm matrix（PCA分析用的--make-grm结果，而非gz后的，所以重新跑）
 7  /w/00/g/g01/user305/soft/gcta_1.91.6beta/gcta_1.91.6beta/gcta64 --make-grm --out /w/00/g,
 8  ## gcta pca（保留前20个PCA）
 9  /w/00/g/g01/user305/soft/gcta_1.91.6beta/gcta_1.91.6beta/gcta64 --grm snp.gcta --pca 20 -
10  ## snp.gcta.eigenval is eigenvalue file
11  ## snp.gcta.eigenvec is eigenvector file
12
13
```

## 7、PCA绘制

```
1  ## convert results to normal data frame
2  setwd("F:/works/developing/course/gwas/data/lecture08")
3  eigvec <- read.table("snp.gcta.eigenvec", header = F, stringsAsFactors = F)
4  colnames(eigvec) <- c("FID", "Sample", paste0("PC", 1:20))
5  write.table(eigvec[2:ncol(eigvec)], file = "gcta.eigenvector.xls", sep = "\t", row.names
6
```

```
1  eigval <- read.table("snp.gcta.eigenval", header = F)
2  pcs <- paste0("PC", 1:nrow(eigval))
3  eigval[nrow(eigval),1] <- 0
4  percentage <- eigval$V1/sum(eigval$V1)*100
5  eigval_df <- as.data.frame(cbind(pcs, eigval[,1], percentage), stringsAsFactors = F)
6
```

```
1  names(eigval_df) <- c("PCs", "variance", "proportion")
2  #str(eigval_df)为字符型，所以下面要换成数字型
3
4  eigval_df$variance <- as.numeric(eigval_df$variance)
5  eigval_df$proportion <- as.numeric(eigval_df$proportion)
6  write.table(eigval_df, file = "gcta.eigenvalue.xls", sep = "\t", quote = F, row.names =
7
```

## 2.GWAS关联分析

### fast-lmm(51.19)

```
1  ## convert vcf to tped format【先用plink将vcf格式转换为tped格式，不能接受vcf格式但是ped，bed
2  plink --vcf lecture09_sort.vcf --recode 12 transpose --output-missing-genotype 0 --out fa
3  ## create fastlmm input trait file in R【创造特有的表型文件】
4  setwd("F:/works/developing/course/gwas/data/lecture09")
5  tfam <- read.table("fastlmm_in.tfam", header = F, stringsAsFactors = F)
6  tr <- read.table("lecture09.trait.txt", header = T, check.names = F, stringsAsFactors = I
7  tr <- tr[match(tfam$V1, tr$`<Trait>`),] ## sort by tfam
8  tr[tr == -999] <- NA
9  tre <- cbind(tr[,1], tr)
10 for(i in 3:ncol(tre)){
11   file_name <- paste0("fastlmm.trait.", names(tre)[i], ".txt")
12   write.table(tre[, c(1, 2, i)], file = file_name, col.names = F, row.names = F, sep = "`
13 }
14 ## run fast-lmm gwas with kinship only
15 fastlmmc -tfile fastlmm_in  -pheno fastlmm.trait.Trait04.txt -tfileSim fastlmm_in -simOu
16 ## prepare covariant in R
17 setwd("F:/works/developing/course/gwas/data/lecture09")
18 tfam <- read.table("fastlmm_in.tfam", header = F, stringsAsFactors = F)
19 covs <- read.table("lecture09.Q.txt", header = T, check.names = F, stringsAsFactors = F,
20 covs <- covs[match(tfam$V1, covs$`<Trait>`), ]
21 covs <- cbind(covs[,1], covs[,1], covs[,2:ncol(covs)]) ## for fast-lmm
22 write.table(covs[, -ncol(covs)], file = "fastlmm.cov.txt", col.names = F, row.names = F,
23 ## run fast-lmm gwas with kinship and Q-matrix
24 fastlmmc -tfile fastlmm_in -pheno fastlmm.trait.Trait04.txt -tfileSim fastlmm_in -simOut
25
```

# EMMAX

```r
## convert vcf to tped format
/w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/01.gx_first
## kinship estimate using emmax-kin-intel64【软件自己算的kinship】
emmax-kin-intel64 emmax_in -v -d 10 -o emmax_in.BN.kinf
## create emmax input trait file in R【脚本同上】
setwd("F:/works/developing/course/gwas/data/lecture09")
tfam <- read.table("emmax_in.tfam", header = F, stringsAsFactors = F)
tr <- read.table("lecture09.trait.txt", header = T, check.names = F, stringsAsFactors = 
tr <- tr[match(tfam$V1, tr$`<Trait>`),] ## sort by tfam
tr[tr == -999] <- NA
tre <- cbind(tr[,1], tr)
for(i in 3:ncol(tre)){
  file_name <- paste0("emmax.trait.", names(tre)[i], ".txt")
  write.table(tre[, c(1, 2, i)], file = file_name, col.names = F, row.names = F, sep = "
}
## run emmax gwas with kingship only
emmax-intel64 -t emmax_in -o Trait04_emmax.k -p emmax.trait.Trait04.txt -k emmax_in.BN.k
## prepare covariant in R
setwd("F:/works/developing/course/gwas/data/lecture09")
tfam <- read.table("emmax_in.tfam", header = F, stringsAsFactors = F)
covs <- read.table("lecture09.Q.txt", header = T, check.names = F, stringsAsFactors = F,
covs <- covs[match(tfam$V1, covs$`<Trait>`), ]
covs <- cbind(covs[,1], covs[,1], rep(1, nrow(covs)), covs[,2:ncol(covs)]) ## for emmax
write.table(covs[, -ncol(covs)], file = "emmax.cov.txt", col.names = F, row.names = F, s
## run emmax gwas with kinship and Q-matrix
emmax-intel64 -t emmax_in -o Trait04_emmax.qk -p emmax.trait.Trait04.txt -k emmax_in.BN.
```

# GEMMA

https://www.jianshu.com/p/372ec8585b49

https://cloud.tencent.com/developer/article/1677009

https://cloud.tencent.com/developer/article/1638988 !!!!!!

https://www.codenong.com/cs105808876/

https://cloud.tencent.com/developer/article/1593313

对此, Zhou等[29]提出了GEMMA方法, 是EMMA的精确算法, 计算速度比EMMA大大提高。近年来, 快速运算得到长足发展, 不断涌现新方法, 如FaST-LMM[30]、FaST-LMM-Select[31]和BOLT-LMM[32]。最近, Wang等[33]将CMLM和FaST- LMM两种

算法整合，提出了运算速度更快的SU

```
1  ## 下载GEMMA
2  wget -c https://github.com/genetics-statistics/GEMMA/releases/download/0.98.1/gemma-0.98
3  ## 解压
4  gzip -d gemma-0.98.1-linux-static.gz
5  ## 设置权限
6  chmod a+x gemma-0.98.1-linux-static
```

```
1  ## 转换格式
2  /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/01.gx_first
3  /w/00/g/g01/user305/soft/plink/plink --file snp --make-bed --out gemma_input --allow-ext
```

```
1  接下来计算kinship矩阵。
2
3  /w/00/g/g01/user305/soft/gemma-0.98.1-linux-static -bfile gemma_input -gk 2 -o gemma
```

**-bfile：输入Plink二进制格式文件的前缀。**
<mark>**-gk：**指定生成的kinship矩阵类型。-gk 1 为centered matrix，-gk 2 为standardized matrix。</mark>
**-o：**输出文件前缀。

## 3.绘图

## 1、单倍型分析

```
1  1.先生成det文件找到强连锁区域
2  ## whole genome block
3  ## 需要加上no-pheno-req
4
5  /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/mer5.vcf/de
6  #该参数来源于--allow-extra-chr --blocks no-pheno-req --blocks-max-kb 1000 --geno 0.2 --bl
7  #基于选择信号和关联分析鉴定温带和热带玉米中优势单倍型和关联基因
```

```
 8
 9   ## subset vcf and convert to Haploview format
10   ## 用vcftools提取指定区间
11   /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf demo.impute.vcf  --chr At1 --from
12
13   2.beagle进行phase
14   /data/apps/jdk/jdk1.8.0_131/bin/java -Xms2g -Xmx8g -jar /w/00/g/g01/user305/soft/beagle.
15
16   20362400-20602400= 211.605kb
17
```

## 2、单倍型聚类

### ①NJ树

```
1   /data/apps/plink/20181202/plink --vcf /w/00/g/g01/user305/Result/TResult/mer5.vcf/3.1_se
2
3   /data/apps/plink/20181202/plink --bfile plink --allow-extra-chr --distance-matrix
4
5   perl ~/script/NJTree/04.mdist2phylip.pl
6
7   fastme建树
8   /data/apps/fastme/2.1.5/bin/fastme -i plink.mdist.rename.phylip -m NJ -T 10 -o fastme.tr
```

### ②contree

```
 1   #在正式运行phylip子程序之前，需要将VCF格式转换成phylip所需的格式。
 2   # github中已经有开源的转换代码：edgardomortiz/vcf2phylip
 3   # 查看文件权限是否是可执行的
 4   ## 用vcftools提取指定区间
 5   /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf demo.impute.vcf  --chr At1 --from
 6
 7   python vcf2phy.py -i chr1_20304000-20604000.region.vcf
 8
 9
10   wget http://evolution.gs.washington.edu/phylip/download/phylip-3.69.tar.gz ./ #下载软件
11   tar zxvf phylip-3.697.tar.gz #解压
12   cd phylip-3.697/src
```

```
13  make -f Makefile.unx install
14  #以上几步即安装完软件，文件夹中的exe目录里为可执行程序
```

```
1   1、构建四个子程序的.par文件# seqboot.par
2   chr1.region.min4.phy
3   R #选择bootstrap
4   100 #设置bootstrap的值，即重复的replicate的数目，通常使用1000或者100，注意此处设定好后，后续两
5   Y #yes确认以上设定的参数
6   9 #设定随机参数，输入奇数值。
7
8
9   #2、 dnadist.par
10  seqboot.out
11  T  #选择设定Transition/transversion的比值
12  2.3628   #比值大小
13  M  #修改M值
14  D  #修改M值
15  100   #设定M值大小
16  2   #将软件运行情况显示出来
17  Y   #确认以上设定的参数
18
```

```
1   # 3、neighbor.par（注意不要有空格）
2   dnadist.out
3   M
4   100    #设定M值大小
5   9   #设定随机数，输入奇数值
6   Y   #确认以上设定的参数
7
8   #4、consense.par
9   nei.tree
10  Y #确认以上设定的参数
```

```
1   #/w/00/g/g01/user305/soft/phylip-3.697/exe
2   /w/00/g/g01/user305/soft/phylip-3.697/exe/seqboot < ./seqboot.par && mv ./outfile ./seqb
```

```
3   /w/00/g/g01/user305/soft/phylip-3.697/exe/dnadist < ./dnadist.par &&  mv ./outfile ./dna
4   /w/00/g/g01/user305/soft/phylip-3.697/exe/neighbor < ./neighbor.par && mv  ./outfile ./ne
5   /w/00/g/g01/user305/soft/phylip-3.697/exe/consense < ./consense.par && mv ./outfile ./co
6
```

## ③豪哥的办法

```
1   /data/apps/iqtree/2.0.5/iqtree2 -s chr1_20304000-20604000.region.min4.phy -nt AUTO -m GTI
```

# ④VCF2Dis

## 1、an Example of nj-tree with no boostrap

```
1   /w/00/g/g01/user305/soft/VCF2Dis-1.44/bin/VCF2Dis -InPut  /w/00/g/g01/user305/Result/TRe
2
```

## 2、Run  PHYLIP

```
1   #   After p_distance done , software PHYLIPNEW 3.69 (http://evolution.genetics.washingto
2   #  Install PHYLIPNEW:https://github.com/BGI-shenzhen/VCF2Dis/blob/master/Install.NJ.en.m
3          PHYLIPNEW-3.69.650/bin/fneighbor  -datafile p_dis.matrix  -outfile tree.out1.
4
5      #    3.2 Run  MEGA
6      #    The MEGA6 (http://www.megasoftware.net/) was used to present the phylogenetic
```

## 3) an Example of nj-tree with boostrap

```
1   /w/00/g/g01/user305/soft/VCF2Dis-1.44/bin/VCF2Dis -InPut /w/00/g/g01/user305/Result/TResu
```

## 3、单倍型网络图

https://www.jianshu.com/p/f9371b782f05

https://www.cnblogs.com/chenwenyan/p/9877377.html

https://weibo.com/p/2304188741e6c10102x7aw?

pids=PI_Official_CardMixFeedv6__4&feed_filter=1 Popart使用方法

phylo_tools/phy2fasta.py at master · harmsm/phylo_tools · GitHub phy文件to fasta

method0:

python phy2vcf.py phased.min4.phy > input.fas

```python
1  #!/usr/bin/env python
2  __description__ = \
3  """
4  Converts the .phy file used as input for a phyml-ss calculation (with its
5  stockholm-style structural annotations) and converts it into a simple,
6  paml-friendly fasta file with no structural information.
7  """
8  __author__ = "Michael J. Harms"
9  __date__ = "130426"
10 __usage__ = "
11 phy2fasta.py phy_file > fasta_output_file"
12
13 import sys
14
15 f = open(sys.argv[1],'r')
16 lines = f.readlines()
17 f.close()
18
19 lines = [l for l in lines[2:] if l.strip() != "" and not l.startswith("#")]
20
21 for l in lines:
22     c = l.split()
23     print ">%s\n%s" % tuple(c)
```

```
1  #用plink过滤低频位点
2  plink --vcf xxxx.vcf --maf 0.05 --geno 0.2 --recode vcf-iid -out xxxx-maf0.05 --allow-ext
3  #用plink筛选连锁区
4  plink --vcf xxxx-maf0.05.vcf --indep-pairwise 100 50 0.2 -out xxxx-maf0.05-LD --allow-ext
5  #用plink找到强连锁区
6  plink --noweb --bfile xxxx-maf0.05-LD --blocks no-pheno-req --allow-extra-chr
```

Method1:

```
1  ## subset vcf and convert to Haploview format
2  ## 用vcftools提取指定区间
3  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /w/00/g/g01/user305/Result/TResul
4
5  /data/apps/jdk/jdk1.8.0_131/bin/java -Xms2g -Xmx10g -jar /w/00/g/g01/user305/soft/beagle
6  #VCF转换phy
7  vcf2phy.py -i phased.vcf.gz
8  #有了phy文件拿到Win操作系统下，用Bioedit就可以变成fasta格式
```

## Method2:

1）将vcf转化为plink格式，假定输入的vcf文件名为：17893893-17898893.vcf

```
1  /vcftools --vcf 17893893-17898893.vcf --plink-tped --out 17893893-17898893
2  /plink --tfile 17893893-17898893 --recode --out 17893893-17898893
```

2) 用PLINK确定要研究的位点是否处于连锁的状态；生成blocks和blocks.det两种后缀格式文件；

```
1  /plink --file 17893893-17898893 --blocks no-pheno-req --out 17893893-17898893
```


17893893-17898893.blocks
17893893-17898893.blocks.det

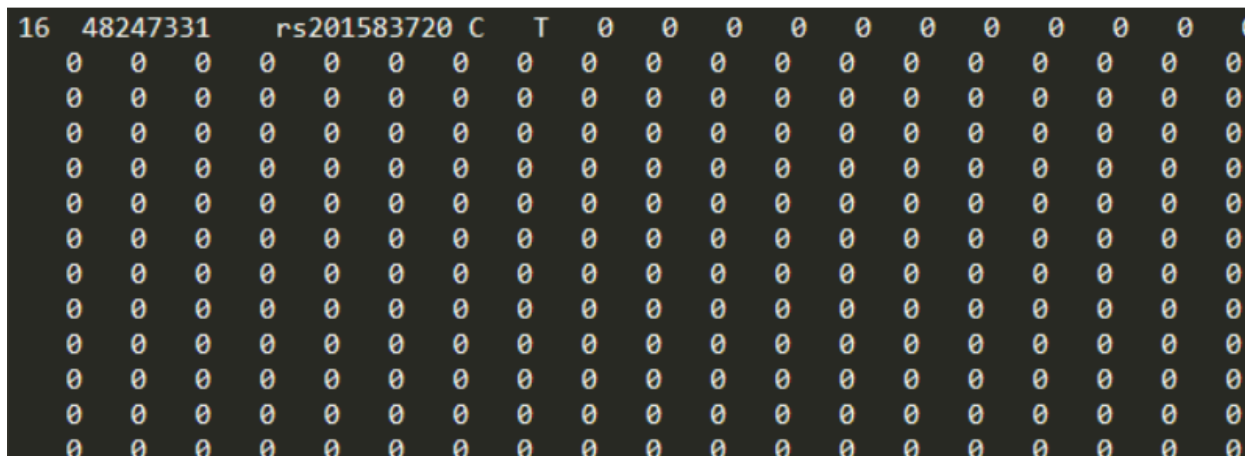3) 提取这35个位点作为接下来的单倍型网络构建；去掉vcf的头文件，保存为txt格式，见如下图，17893893-17898893.txt:



4）准备17893893-17898893_singstring.txt文件，该文件其实就是去掉 0|0，0|1，1|0，1|1的"|"。

```
16    48247331      rs201583720 C   T   0   0   0   0   0   0   0   0   0   0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
     0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
```

## 曼哈顿图系列

## 1、TASSEL结果画曼哈顿图

```
1  awk '$1=="Relativedryweight" {print $3, $4, $7}' MLM_result.txt > input_file.txt
2  sed 's/^AT//g' input_file.txt > new_file.txt
```

## 2、注释显著位点

https://zhuanlan.zhihu.com/p/96653959

```
1  library(ggrepel)
2  #准备数据
3  data <- Snp_pos %>%
4  # 添加高亮和注释信息：snpsOfInterest中的rs编号和P值大于6的点
5   mutate( is_highlight=ifelse(SNP %in% snpsOfInterest, "yes", "no")) %>%
6   mutate( is_annotate=ifelse(-log10(P)>6, "yes", "no"))
7  # 绘图
8  p1 <- ggplot(data, aes(x=BPcum, y=-log10(P))) +
9    geom_point( aes(color=as.factor(CHR)), alpha=0.8, size=1.3) +
10    scale_color_manual(values = rep(c("grey", "skyblue"), 22 )) +
11    scale_x_continuous( label = X_axis$CHR, breaks= X_axis$center ) +
12    scale_y_continuous(expand = c(0, 0) ) +
13    # 添加高亮点
14    geom_point(data=subset(data, is_highlight=="yes"), color="orange", size=2) +
15    # 添加高亮label，且防止重叠
16    geom_label_repel( data=subset(data, is_annotate=="yes"), aes(label=SNP), size=2) +
17    theme_bw() +
18  theme(
19      legend.position="none",
```

```
20        panel.border = element_blank(),
21        panel.grid.major.x = element_blank(),
22        panel.grid.minor.x = element_blank()
23     )
```

## 3、R语言CMplot包绘制曼哈顿图

https://cloud.tencent.com/developer/article/1481862

```
1
```

## 5、Candihap画图

https://mp.weixin.qq.com/s/1leAY8a-wYAFISFHgmu93Q

https://github.com/xukaili/CandiHap/tree/master/Sanger_ab1_Linux

```
1
2  perl  /w/00/g/g01/user305/soft/CandiHap/vcf2hmp.pl  /w/00/g/g01/user305/Result/TResult/me
```

## 6、单倍型群体结构可视化-haplostrips

## https://bitbucket.org/dmarnetto/haplostrips/src/master/

```
1  #命令行下使用：win+R / CMD
2  haplostrips -v chr22.1000genomes.example.vcf.gz -R -V chr22.archaic_hominids.example.vcf
```

## 7、fastphase和beagle原理相同-提取候选区段

https://mp.weixin.qq.com/s/YwXq_MA4SEXJhNzm-lF5-g

https://mp.weixin.qq.com/s/BT6wHdIA70ltuBusxzD-wA

```
1  #1.解压缩
2  tar -xzvf  phase.2.1.1.linux.tar.gz
3  ##直接使用beagle的结果
4  6.画structure
5  ## missing rate and maf filtering
6  /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/mer5.vcf/3
7  ## make bed and LD pruning
8  /w/00/g/g01/user305/soft/plink/plink --vcf /w/00/g/g01/user305/Result/TResult/mer5.vcf/2
9
```

```
10   ##extract prune in snps and recode to admixture format（转成ped格式）
11   /w/00/g/g01/user305/soft/plink/plink --bfile snp.int0.8.maf0.05 --extract snp.int0.8.maf(
12
13
14   #课题组方法无法完整的画图，先得到群体结构的结果
15   for i in {1..10}; do /data/apps/admixture/1.3.0/admixture --cv -j30 --seed=666 /w/00/g/g(
16   grep -h CV log*.out
```

## 8、fst(10k windows,2k 步长)

https://www.jianshu.com/p/7b6701ae593b

https://www.jianshu.com/p/b73a8d6233be

https://www.jianshu.com/p/98e56862347f

```
1   ##对每一个SNP变异位点进行计算
2   vcftools --vcf output.10.06filterSNPrmOutGroup.hwe0.0001.vcf --weir-fst-pop Cultivar01.t
3
4   ##按照区域来计算
5   vcftools --vcf output.10.06filterSNPrmOutGroup.hwe0.0001.vcf --weir-fst-pop Cultivar01.t
```

给的极端值txt文件只包含ID

```
文件(F)  编辑(E)  格式(
0-14A
0-15A
0-17A
0-18A
0-20A
0-23A
0-27A
```

```
1   /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /w/00/g/g01/user305/Result/TResul
2
3   ##图1
4   library(ggplot2)
5   data<-read.table("p_1_2_bin.windowed.weir.fst",header=T)
6   sc3 = subset(data,CHROM=="At1")
7   p <- ggplot(sc3,aes(x=BIN_END/1000000,y=WEIGHTED_FST)) + geom_point(size=0.5, colour="blu
8   p + theme_bw()

9
```

```
10  ##图2
11  library(ggplot2)
12  data<-read.table("test.out.weir.fst",header=T)
13  sc3 = subset(data,CHROM=="Gm01")
14  p <- ggplot(sc3,aes(x=POS,y=WEIR_AND_COCKERHAM_FST)) + geom_point(size=0.5, colour="blue"
15  p + theme_bwt1
```

```
1   library(ggplot2)
2   data<-read.table("p_1_2_bin.windowed.weir.fst",header=T)
3   sc3 = subset(data,CHROM=="At1")
4   p1 = ggplot(sc3,aes(x=BIN_END/1000000,y=WEIGHTED_FST)) + geom_point(size=0.5, colour="blu
5   p1 + theme_bw()
6
7   sc3 = subset(data,CHROM=="At2")
8   p2 = ggplot(sc3,aes(x=BIN_END/1000000,y=WEIGHTED_FST)) + geom_point(size=0.5, colour="blu
9   p2 + theme_bw()
10
11
12  sc3 = subset(data,CHROM=="At3")
13  p3 = ggplot(sc3,aes(x=BIN_END/1000000,y=WEIGHTED_FST)) + geom_point(size=0.5, colour="blu
14  p3 + theme_bw()
15
16
17  sc3 = subset(data,CHROM=="At4")
18  p4 = ggplot(sc3,aes(x=BIN_END/1000000,y=WEIGHTED_FST)) + geom_point(size=0.5, colour="blu
19  p4 + theme_bw()
20
21  sc3 = subset(data,CHROM=="At5")
22  p5 = ggplot(sc3,aes(x=BIN_END/1000000,y=WEIGHTED_FST)) + geom_point(size=0.5, colour="blu
23  p5 + theme_bw()
24  install.packages("ggpubr")
25  library(ggpubr)
26  ggarrange(p1, p2, p3, p4, p5, ncol = 2, nrow = 3)
```

manhadun:**https://www.plob.org/article/21662.html**

1、下面脚本需要改变染色体的名称

```
1  sed "1d" p_1_2_bin.windowed.weir.fst|awk '{if($5<0)print $1"\t"$2"\t0";else print $1"\t"$
2
3  ### 加载要使用的R包
4  library(qqman)
5  library(Cairo)
6  #z注意修改染色体编号
7  Fstfile<-read.table("TM_4500vsYJ_800_20k_5k_plot.txt", header=F, stringsAsFactors=F)
8  SNP<-paste(Fstfile[,1],Fstfile[,2],sep = ":")
9  Fstfile=cbind(SNP,Fstfile)
10 colnames(Fstfile)<-c("SNP","CHR","POS","Fst")
11 outfile<-"TM_4500vsYJ_800_20k_5k"
12 filePNG<-paste(outfile,"manhattan.png",sep=".")
13 CairoPNG(file=filePNG,width=1500,height=500)
14 colorset<-c("#FF0000","#FFD700","#2E8B57","#7FFFAA","#6495ED","#0000FF","#FF00FF")
15
16 ### 调用qqman,画出全部染色体
17 manhattan(Fstfile,chr="CHR",bp="POS",p="Fst",snp="SNP", col=colorset,logp=F,suggestivelin
18   dev.off()
19 ### 只绘制单条染色体
20 if (FALSE){
21     manhattan(subset(Fstfile,CHR=="10"),chr="CHR",bp="POS",p="Fst",snp="SNP", col=colors
22 }
23
24 dev.off()
25
```

```
1  library("qqman")
2  mydata=read.csv("p_c_t_bin.windowed.weir.fst",header = T,sep=",")
3  manhattan(mydata)
4  manhattan(mydata,main="p_c_t",ylim=c(0,1),cex=0.8,cex.axis=0.8,col = c("gray75","red","b
5  ###  main="GWAS-Manhattan"这个命令是给图片取名字；highlight = GL这个命令是标记出高亮的snp的点
6  suggestiveline = F，genomewideline = 4.596这两个命令是设定阈值用的可以按照自己的需要设定，如是
7  col = c("gray75","red","black")这个命令就是规定曼哈顿图的颜色的，可以根据自己的需要而定。
8  qq(mydata$P)
```

## 图片组合

```
1  # 绘制第一个图像
2  > p1 = ggplot(cars, aes(x = speed, y = dist)) +
3      geom_line(size = 1) +
4      labs(x = "speed", y = 'dist')
5
6  # 绘制第二个图像
7  > p2 = ggplot(cars, mapping = aes(x = speed, y = dist)) +
8      geom_point(size = 1, col = 'darker') +
9      labs(x = "speed", y = 'dist')
10
11 # 整合两张图
12 > ggarrange(p1, p2, ncol = 2, nrow = 1)
```

## ！计算核苷酸多样性（π）

非常简单，把某群体的全部样本ID放在一个TXT文件里就行，注意要每一行一个ID。

计算核苷酸多样性（π）

```
1  #局部和整体
2  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /w/00/g/g01/user305/Result/flower
3
4  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /w/00/g/g01/user305/Result/TResul
```

## 画图1：https://cloud.tencent.com/developer/article/1593265

箱线图

```
1  bb<-read.table("../../vcf_handling/Fish_Populations/benthic_pi.windowed.pi",header=T)
2  ll<-read.table("../../vcf_handling/Fish_Populations/littoral_pi.windowed.pi",header=T)
3  dim(bb)
4  head(bb)
5  bb$indiv<-"benthic"
6  dim(ll)
7  head(ll)
8  ll$indiv<-"littoral"
9  df<-rbind(bb,ll)
10 colnames(df)
11 dim(df)
12 library(ggplot2)
```

```
13   options(scipen=200)
14   ggplot(df,aes(x=indiv,y=PI,fill=indiv))+
15     geom_boxplot()+theme_bw()
```

## 折线图

```
1   ggplot(df,aes(x=BIN_START,y=PI,group=indiv,color=indiv))+
2     geom_line()+theme_bw()+
3     theme(legend.position = "top",
4           legend.title = element_blank(),
5           axis.text.y = element_blank(),
6           axis.ticks.y = element_blank(),
7           axis.title = element_blank())
```

## https://www.jianshu.com/p/7b6701ae593b

```
1   #读入数据；
2   dt1<- read.delim("pi.txt",sep="\t", header = T, check.names = F)
3   # 加载ggplot2包；
4   library(ggplot2)
5   #定义染色体位置。
6   #我分析的物种有8条染色体，我将所有的染色体都串起来作图，因此需要标出每条染色体的中间位置在哪。
7   br = c(440000, 1370000, 2410000, 3515000, 4610000,5800000,7095000,8399791)
8   la = c("1","2","3","4","5","6","7","8")
9   #自定义图表主题，对图表做精细调整；
10  mytheme<-theme(panel.grid.major =element_blank(),
11                 panel.grid.minor = element_blank(),
12                 panel.background = element_blank(),
13                 panel.border = element_blank(),
14                 axis.line.y = element_line(color = "black"),
15                 axis.line.x = element_line(color = "black"),
16                 #axis.title.x = element_text(size = rel(1.2)),
17                 axis.title.y = element_text(size = rel(1.2)),
18                 axis.text.y = element_text(size=rel(1.2),color="black"),
19                 #axis.text.x = element_text(size=rel(1.2),color="black"),
20                 plot.margin=unit(x=c(top.mar,right.mar,bottom.mar,left.mar),units="inches
21  #绘制
22  pa<-ggplot(data=dt1, mapping = aes(x=No,y=pop1))+geom_line(color=dt1$Color1,size=1)
23  #设置x轴范围，避免点的溢出绘图区；
```

```
24  pa<-pa+scale_x_continuous(limits = c(-1000, 9230000),breaks = br,labels = la)
25  #设置y轴范围
26  pa<-pa+scale_y_continuous(limits = c(-0.0005,0.01),breaks = c(0,0.002,0.004,0.006,0.008,
27  #设置图例、坐标轴、图表的标题;
28  pa<-pa+labs(y="Pi (10e-3)",x=NULL)
29  #自定义图表主题，对图表做精细调整;
30  pa<-pa+mytheme
31  #出图
32  pa
33
```

## 9、0|0转换-聚类热图

## https://www.cnblogs.com/Mao1518202/p/11269595.html 最重要的热图指南

RectChr 之多样品某一区域 genotype/gene/单倍型 热度图 - 知乎 (zhihu.com)

https://github.com/BGI-shenzhen/RectChr/tree/main/Example

**方法一：**

**https://mp.weixin.qq.com/s/8tchKCRVsBs5f5fQtR-d1A** ！！！！

**方法二：plink**

https://mp.weixin.qq.com/s/hX0XjOpTugTEIMEI3rz1_g

使用–bfile 、 --file 和 --tfile读取文件类型不一样：

–bfile 读取二进制文件，bed、bim和fam格式

–file 读取文本文件，ped和map格式

```
1
2  ## 用vcftools提取指定区间
3  /data/apps/vcftools/vcftools-0.1.15/bin/vcftools --vcf /w/00/g/g01/user305/Result/TResul
4
5  方法一:
6  #对数据进行清洗，将其转化为0,1,2的形式
7  #转ped、map格式
8
9  /w/00/g/g01/user305/soft/plink/plink --vcf chr4_10kb.vcf --recode --out region_10k --all
```

```
10   #首先，使用plink命令，将基因型数据转化为012的raw格式：
11   /w/00/g/g01/user305/soft/plink/plink --file region_10k --recodeA --allow-extra-chr
12
13
14   方法二：
15   /w/00/g/g01/user305/soft/plink/plink --vcf chr1_20304000-20604000.region.vcf --recode --
16   ## ped/map转换为bed/bim/fam
17   /w/00/g/g01/user305/soft/plink/plink --file region_300k --make-bed --out snp_region_300k
18
19
20
```

## extract_sample

```bash
#!/bin/bash
for i in `cat hap1.txt`
do
grep $i trait.txt >> result1.txt

done
```

```
file1（id）file2（id+数据）
grep -wFf file1 file2
```

[https://www.jianshu.com/p/5b6bdbebcd61](https://www.jianshu.com/p/5b6bdbebcd61) **重要热图**

[https://ww2.mathworks.cn/help/matlab/ref/matlab.graphics.chart.heatmapchart-properties.html](https://ww2.mathworks.cn/help/matlab/ref/matlab.graphics.chart.heatmapchart-properties.html) **帮助文档**

```r
setwd("D:\\work_dir\\10.40KB_picture")

install.packages("gplots")
```

```
4
5   test<- read.table("test1.xls", head=T)
6
7   library(pheatmap)
8   library(grid)
9   library(ComplexHeatmap)
10
11  data = as.matrix(test)
12  is.matrix(data)
13
14  Heatmap(data, cluster_rows = T, cluster_columns = F,show_row_names = F, show_column_names
15
16  pheatmap(data, cluster_cols = F, cluster_rows = T, fontsize = 1)
```

## 先将dataframe换为矩阵

https://blog.csdn.net/weixin_41938903/article/details/84983804

```
1   # R读入csv文件
2   data1<-read.csv('d:/ko_wt.csv')   #参数默认为TRUE
3   # 将dataframe转化为matrix，目前找到了2种方法
4   ①
5   y = data.metrix(data1)
6   ②
7   y <- apply(data1, as.numeric)
```

## 简单自己的总结

```
1   BiocManager::install("ComplexHeatmap")
2   setwd("D:\\pheatmap")
3   test<- read.csv("test1.csv", head=T)
4
5   library(pheatmap)
6   library(grid)
7   library(ComplexHeatmap)
8
9   data1 = as.matrix(data)
10  is.matrix(data1)
11
12  Heatmap(data1,  cluster_rows = F, cluster_columns = F,show_row_names = F, show_column_nam
```

## complexheatmap-https://www.jianshu.com/p/06705a7812ee

```
1   #设置工作路径
2   getwd()
3   setwd("D:\\pheatmap")
4   test<- read.table("genotype_transpose.csv", head=T)
5   rownames(test)<-test[,1]
6   test1<-test[,c(1:490)]
7
8   library(pheatmap)
9   pheatmap(test1)
10  ?pheatmap
11  vignette()
12  devtools::install_github('satijalab/seurat-data')
13  library(SeuratData)
14  InstallData("pbmc3k")
15  BiocManager::install("ComplexHeatmap")
16
```

```
1   最后用 awk 进行矩阵转置：
2   awk '{i=1;while(i <= NF){col[i]=col[i] $i " ";i=i+1}} END {i=1;while(i<=NF){print col[i]
3   awk -F, '{for(i=1;i<=NF;i=i+1){a[NR,i]=$i}}END{for(j=1;j<=NF;j++){str=a[1,j];for(i=2;i<=N
4
5
```

## 热图2https://mp.weixin.qq.com/s/INswPaJP1P4caEnAE_X--A

### 1.准备

```
1   rm(list=ls()) # 清除环境变量
2   setwd("C:/Rdata") # 设置工作目录
3   getwd() # 查看当前的工作目录
```

### 2.安装包

```
1  1.2 安装和加载R包
2  安装我们绘图需要用的包，如果已安装，直接加载就好了。
3  install.packages("pheatmap") # 安装包
4  install.packages("ggplot2") # 安装包
5  library(pheatmap) # 加载包
6  library(ggplot2) # 加载包
```

## 3.查看数据

```
1  data <- read.table("test.FPKM.txt",  # 读取的数据文件名称，这里文件是放在工作目录下
2                       header=T, # 数据集第一行为变量名
3                       row.names=1, # 第一列为行名
4                       sep="\t") # 指定分隔符号
5  dim(data) # 查看变量有多少行多少列
```

## 4.绘制热图

```
1  先绘制个基本款的热图。
2  p <- pheatmap(data)
3
4  3.1 热图横向标准化
5  设置标准化方向scale，对其横向标准化。
6  p <- pheatmap(data, scale="row")
```

```
1  3.2 图形外观调整
2  设置边框为白色，去掉横向、纵向聚类。
3  p <- pheatmap(data, scale="row",
4            border="white", # 设置边框为白色
5            cluster_cols = F, # 去掉横向、纵向聚类
6            cluster_rows = F)
```

```
1  3.3 去除横纵坐标标签
2  去掉横纵坐标中的id。
3  p <- pheatmap(data,scale="row",
4               border="white",  # 设置边框为白色
```

```
5              cluster_cols = F, # 去掉横向、纵向聚类
6              cluster_rows = F,
7              show_rownames = F, #去掉横、纵坐标id
8              show_colnames = F)
```

```
1  3.4 修改图例
2  3.4.1 不显示右上角图例
3  p <- pheatmap(data,scale="row",
4              border="white",  # 设置边框为白色
5              cluster_cols = F, # 去掉横向、纵向聚类
6              cluster_rows = F,
7              show_rownames = F, #去掉横、纵坐标id
8              show_colnames = F,     #此处参数可以选择对横纵坐标去聚类
9              legend = F) # 去掉图例
```

```
1  3.4.2 设置图例范围
2  p <- pheatmap(data,scale="row",
3              border="white",  # 设置边框为白色
4              cluster_cols = F, # 去掉横向、纵向聚类
5              cluster_rows = F,
6              show_rownames = F, #去掉横、纵坐标id
7              show_colnames = F,
8              legend = T, # 添加图例
9              legend_breaks=c(-1,0,1)) # 设置图例范围
10
11 # 也可以设置legend_breaks=c(-2,0,2)试试
```

## 5、颜色改变——使用heatmap

https://www.jianshu.com/p/06705a7812ee

https://www.jianshu.com/p/325107a11492

https://mp.weixin.qq.com/s/mESHydcbzyMc2p4rpBE-MA 格式转换，在微信公众号可以找到

```
1  #格式判断和格式转换
2  is.matrix(data)
3  data1 = as.matrix(data)
4  is.matrix(data1)
```

## 基本显示

```r
1  ht <- Heatmap(mat)
2  draw(ht)
```

## 颜色的改变

- 在热图当中如果需要改变颜色的话，**一定要使用**`circlize::colorRamp2()`来对颜色进行赋值。这个函数接受两个参数:**颜色分割的位置**以及**相应分割点上的颜色**。最后我们的`Heatmap`当中通过`col`来制定颜色

**PS**：`colorRamp2`默认的使用的配色方案是`LAB`的。如果想使用`RGB`的则可以通过其中的`space`参数进行修改。

```r
1  library(circlize)
2  col_fun <- colorRamp2(c(-2, 0, 2), c("navy", "white", "firebrick3"))
3  Heatmap(mat, col = col_fun)
```

# 聚类

## http://showteeth.tech/posts/31202.html

## https://www.yuque.com/bioitee/complexheatmap-reference/dimension-names

## 基本设置

无监督的聚类属于热图的可视化的一个重要组成部分。

1. `cluster_rows/columns`来设置**是否进行聚类**
2. `show_column/row_dend`设置**是否显示聚类树**(会发生聚类)
3. `column/row_dend_side`设置**聚类图绘制的位置**
4. `column/row_dend_height`设置**聚类树的高度**

```r
1  Heatmap(mat, cluster_rows = F, show_column_dend = T,
2          column_dend_side = "bottom",  column_dend_height = unit(2, "cm"))
```

## 去掉横纵坐标

```r
1  Heatmap(data1, cluster_rows = F, cluster_columns = F,show_row_names = F, show_column_name
2
```

```
3
```

## 注释

```
1  ha = HeatmapAnnotation(foo = anno_points(matrix(runif(20), nc = 2),
2      pch = 1:2, gp = gpar(col = 2:3)))
```

## 条形块

```
1  # color mapping for -log10(pvalue)
2  pvalue_col_fun = colorRamp2(c(0, 2, 3), c("green", "white", "red"))
3  ha = HeatmapAnnotation(
4      pvalue = anno_simple(-log10(pvalue), col = pvalue_col_fun, pch = pch),
5      annotation_name_side = "left")
6  ht = Heatmap(matrix(rnorm(100), 10), name = "mat", top_annotation = ha)
```

**如果想颠倒横纵坐标画图，可以直接使用t(matrix)进行转置。**

http://showteeth.tech/posts/31202.html

**R语言画热图时图例（legend）过多超出画图边界**

https://cloud.tencent.com/developer/article/1628471

## 三、解决方案——将整个矩阵转化为数值型变量

```
1  test=apply(test,2,as.numeric)
2  Warning message:
3  In apply(test, 2, as.numeric) : NAs introduced by coercion
```

```
1  是rt1的expression列不是数值型向量
2  > rt1 <- as.data.frame(rt1)
3  > rt1[,1] <- as.numeric(rt1[,1])
4  添加这两句进去就行
```

## 热图注释

https://www.jianshu.com/p/8ed637722b0b

```
1
```

## 曼哈顿图优化

数据准备

```
1  awk '{print $2, $3, $4, $7}' input1.txt > input_qqman.txt
2  sed -i "s/^At//g" input_qqman.txt
```

## ①qqman

```
1   #安装方法：
2   #====设置安装源为清华大学TUNA镜像站点====================
3   options("repos" = c(CRAN="https://mirrors.tuna.tsinghua.edu.cn/CRAN/"))
4   options(BioC_mirror="https://mirrors.tuna.tsinghua.edu.cn/bioconductor")
5   install.packages('qqman')
6   #=======================================================
7
8   library(qqman)
9   setwd('F:/PlantTech/training/GWAS/qqman')
10  #输入文件格式如下：
11  #   SNP    CHR    BP    P
12  # 其中SNP列为SNP名称，这一列可以没有
13  # CHR，BP分别为SNP所在的染色体和位置
14  # P表示该SNP的关联显著性p值
15  data<-read.table('input.txt',header=T)
16
17  # 取bonferroni矫正阈值
18  fdr=0.01/nrow(data)
19
20  manhattan(data, main = "Manhattan Plot",
21            ylim = c(0, 10), cex = 0.6,cex.axis = 0.9,
22            col = c("blue4", "orange3"),
23            suggestiveline = -log10(1e-05), genomewideline = -log10(fdr))
24
25  qq(data$P, main = "Q-Q plot of GWAS p-values",
26     xlim = c(0, 7), ylim = c(0,12),
27     pch = 18, col = "blue4", cex = 1.5)
```

## LD衰减-可以根据LD衰减距离来决定候选基因上下游的范围

https://cloud.tencent.com/developer/article/1677015

https://zhuanlan.zhihu.com/p/109164405

```
1  ## calculate ld decay using poplddecay全部
2  /w/00/g/g01/user305/soft/PopLDdecay/bin/PopLDdecay -InVCF /w/00/g/g01/user305/Result/TRe
3
4
5  #画图
6  perl /w/00/g/g01/user305/soft/PopLDdecay/bin/Plot_OnePop.pl -inFile demo.impute.pld.stat
7
```

**POPhelper**

**BLUP分析**

https://www.jianshu.com/p/c72b3835be27

# QTL定位
## ①文件格式

popt = RI7

```
name · = · RiceDemo
popt · = · F2
nind · = · 197
nloc · = · 439
```

## ②软件
https://github.com/bmansfeld/QTLseqr