



## 简介

**blupADC** 是一个专注于分析动植物育种中的系谱数据、基因型数据及遗传评估的工具。在设计该工具时，我们对数据处理时可能遇到的各种问题均进行了详细的考量(**ps.如果您有好的建议，请积极联系作者!**)。此外，为了提高分析效率，**blupADC** 可支持并行计算(通过 `openMP`) 及大数据处理(通过 `bigmemory`)，并且 **blupADC** 中的核心函数均通过 C++ ( `Rcpp` and `RcppArmadillo` ) 进行编写。

**blupADC** 提供了许多有用的功能在整个动植物育种的流程中，包括 系谱分析(系谱追溯、重命名及纠错)，基因型数据格式转换(支持 **Hapmap, Plink, BLUPF90, Numeric, Haplotype** 及 **VCF** 格式)，基因型数据的质控填充，亲缘关系矩阵的构建(**系谱, 基因组及一步法亲缘关系矩阵**)以及遗传评估 (仅需几行代

码即可通过 DMU 和 BLUPF90 完成遗传评估)。

最后，为了进一步方便用户的使用(尤其是编程基础弱的用户)，我们创建了一个免费的在线网站 ([shinyapp](#))。相关的功能仍在开发中。但是，网站的一个缺点就是，不能处理大数据，请大家合理选择!

☺ 祝好运！如果你有建议或者问题，请联系: [hzau\\_qsmei@163.com](mailto:hzau_qsmei@163.com) !

## 新添加的功能

### 1.0.3

- 目前能够通过 DMU 自动分析 母性效应，永久环境效应，随机回归效应 及 社会遗传效应的模型 (2021.8.24)

### 1.0.4

- 支持单倍型格式转换，单倍型-数字矩阵及单倍型加性亲缘关系矩阵的构建(2021.10.8)
- 引入 `bigmemory` 对象支持大数据分析(2021.10.8)

## 开始

### 📦 安装

安装 **blupADC** 之前，用户首先需要安装如下3个包: `Rcpp`, `RcppArmadillo` and `data.table`.

```
install.packages(c("Rcpp", "RcppArmadillo", "data.table", "bigmemory"))
```

📖 **Note:** 在 DMU 和 BLUPF90 的分析中，我们通常需要提前下载好 DMU 软件 ([DMU 下载网站](#)) 和 BLUPF90 软件 ([BLUPF90 下载网站](#))。为了方便用户使用，我们已经将两款软件中基础模块封装进了 **blupADC**，请大家合理使用。

如果您想将DMU和BLUPF90用作商业用途，请务必联系 DMU 和 BLUPF90的作者！！！！

## 在 Linux 上安装 blupADC

```
packageurl <- "https://github.com/TXiang-  
lab/blupADC/releases/download/v1.0.4/blupADC_1.0.4_R_x86_64-pc-linux-gnu.tar.gz"  
install.packages(packageurl, repos=NULL, method="libcurl")
```

## 在 Windows 上安装 blupADC

```
packageurl <- "https://github.com/TXiang-  
lab/blupADC/releases/download/v1.0.4/blupADC_1.0.4.zip"  
install.packages(packageurl, repos=NULL)
```

📌 **Note:**针对github连接比较慢的地区，用户可以通过如下代码进行下载(国内用户推荐如下方式下载)：

## 在 Linux 上安装 blupADC

```
packageurl <-  
"https://gitee.com/qsmei/blupADC/attach_files/851170/download/blupADC_1.0.4_R_x8  
6_64-pc-linux-gnu.tar.gz"  
install.packages(packageurl, repos=NULL, method="libcurl")
```

## 在 Windows 上安装 blupADC

```
packageurl <-  
"https://gitee.com/qsmei/blupADC/attach_files/851169/download/blupADC_1.0.4.zip"  
install.packages(packageurl, repos=NULL)
```

安装成功后，我们输入如下代码即可加载R包：

```
library(blupADC)
```

## 🧑‍💻 功能

- 功能 1. 基因型数据间的格式转换
- 功能 2. 基因型数据的质控与填充
- 功能 3. 品种分析及基因型数据重复性检测
- 功能 4. 系谱追溯、重命名及纠错
- 功能 5. 系谱可视化
- 功能 6. 亲缘关系矩阵的构建(A, G, H)
- 功能 7. 利用DMU软件进行遗传评估
- 功能 8. 利用BLUPF90软件进行遗传评估

## 使用

为了方便用户使用，所有的文档均支持双语模式(中英文说明书)。

`blupADC` 内置了几个数据集对象，包括 `data_hmp` 及 `origin_pedigree`。

此外，`blupADC` 提供一些示例文件，这些文件存储在 `~/blupADC/extdata` 路径下。通过输出下面的代码，我们就能得到这些文件的绝对路径了：

```
system.file("extdata", package = "blupADC") # path of provided files
```

## 功能 1. 基因型数据间的格式转换 ([see more details](#))

```
library(blupADC)
format_result=geno_format(
  input_data_hmp=example_data_hmp, # provided data variable
  output_data_type=c("Plink","BLUPF90","Numeric"),# output data format
  output_data_path=getwd(), #output data path
  output_data_name="blupADC", #output data name
  return_result = TRUE, #save result in R environment
  cpu_cores=1 # number of cpu
)

#convert phased VCF data to haplotype format and haplotype-based numeric format
library(blupADC)
data_path=system.file("extdata", package = "blupADC") # path of example files
phased=geno_format(
  input_data_path=data_path, # input data path
  input_data_name="example.vcf", # input data name,for vcf data
  input_data_type="VCF", # input data type
  phased_genotype=TRUE, # whether the vcf data has been phased
  haplotype_window_nSNP=5, # according to nSNP define haplotype-
block,
  bigmemory_cal=TRUE, # format conversion via bigmemory
object
  bigmemory_data_path=getwd(), # path of bigmemory data
  bigmemory_data_name="test_blupADC", #name of bigmemory data
  output_data_type=c("Haplotype","Numeric"),# output data format
  return_result=TRUE, #save result in R environment
  cpu_cores=1 # number of cpu
)
```

## 功能 2. 基因型数据的质控与填充 ([see more details](#))

```
library(blupADC)
geno_qc_impute(
  input_data_hmp=example_data_hmp, #provided data variable
  data_analysis_method="QC_Imputation", #analysis method type,QC +
imputatoin
  output_data_path=getwd(), #output data path
  output_data_name="YY_data", #output data name
  output_data_type="VCF" #output data format
)
```

## 功能 3. 品种分析及基因型数据重复性检测 ([see more details](#))

```

library(blupADC)
check_result=geno_check(
    input_data_hmp=example_PCA_data_hmp,    #provided hapmap data
    object
    duplication_check=FALSE,                #whether check the duplication
    of genotype
    breed_check=TRUE,                       # whether check the record of
    breed
    breed_record=example_PCA_Breed, # provided breed record
    output_data_path=getwd(),              #output path
    return_result=TRUE                     #save result as a R
    environment variable
)

```

#### 功能 4. 系谱追溯、重命名及纠错 ([see more details](#))

```

library(blupADC)
pedigree_result=trace_pedigree(
    input_pedigree=example_ped1,    #provided pedigree data variable
    trace_generation=3,              # trace generation
    output_pedigree_tree=T           # output pedigree tree
)

```

#### 功能 5. 系谱可视化 ([see more details](#))

```

library(blupADC)
plot=ggped(
    input_pedigree=example_ped2,
    trace_id=c("121"),
    trace_sibs=TRUE    #whether plot the sibs of subset-id
)

```

#### 功能 6. 亲缘关系矩阵的构建(A,G, H) ([see more details](#))

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # path of example files
kinship_result=cal_kinship(
    input_data_path=data_path,          # input data path
    input_data_name="example.vcf",      # input data name,for vcf data
    input_data_type="VCF",              # input data type
    kinship_type=c("G_A","G_D"),        #type of kinship matrix
    dominance_type=c("genotypic"),      #type of dominance effect
    inbred_type=c("Homozygous"),        #type of inbreeding
    coefficients
    return_result=TRUE)                 #save result as a R
    environment variable

```

#### 功能 7. 利用DMU软件进行遗传评估 ([see more details](#))

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # path of example files

run_DMU(

```

```

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2"
, "Age"), # colnames of phenotype
      target_trait_name=list(c("Trait1")),          #trait name
      fixed_effect_name=list(c("Sex", "Herd_Year_Season")), #fixed effect
name
      random_effect_name=list(c("Id", "Litter")),    #random effect
name
      covariate_effect_name=NULL,                    #covariate
effect name
      phe_path=data_path,                             #path of phenotype file
      phe_name="phenotype.txt",                       #name of phenotype file
      integer_n=5,                                    #number of integer variable
      analysis_model="PBLUP_A",                      #model of genetic
evaluation
      dmua_module="dmuai",                            #module of estimating
variance components
      relationship_path=data_path,                    #path of relationship file
      relationship_name="pedigree.txt",               #name of relationship file
      output_result_path=getwd()                     # output path
)

```

## 功能 8. 利用BLUPF90软件进行遗传评估 ([see more details](#))

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # path of example files

run_BLUPF90(

  phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2"
, "Age"), # colnames of phenotype
      target_trait_name=list(c("Trait1")),          #trait name
      fixed_effect_name=list(c("Sex", "Herd_Year_Season")), #fixed effect
name
      random_effect_name=list(c("Id", "Litter")),    #random effect
name
      covariate_effect_name=NULL,                    #covariate
effect name
      phe_path=data_path,                             #path of phenotype file
      phe_name="phenotype.txt",                       #name of phenotype file
      analysis_model="PBLUP_A",                      #model of genetic
evaluation
      relationship_path=data_path,                    #path of relationship file
      relationship_name="pedigree.txt",               #name of relationship file
      output_result_path=getwd()                     # output path
)

```

## 功能1

### 简介

大家好,通过前一章节的学习,相信大家已经对 blupADC 有了一个初步的了解了。从本节开始,我们将对 blupADC 中的几个重要的函数——进行讲解。这一节主要给大家讲述的是如何使用 blupADC 中 geno\_format 函数来进行多种基因型格式数据的转换。

## 示例

### 格式转换-提供R中的变量名称

```
library(blupADC)
format_result=geno_format(
  input_data_hmp=example_data_hmp, #provided hapmap data object
  output_data_type=c("Plink","BLUPF90","Numeric","VCF"),# output data
  format
  return_result = TRUE,           # return result
  cpu_cores=1                     # number of cpu
)
```

### 格式转换-提供本地文件的路径和名称

```
#convert phased VCF data to haplotype format and haplotype-based numeric format
library(blupADC)
data_path=system.file("extdata", package = "blupADC") # path of example files
phased_result=geno_format(
  input_data_path=data_path,      # input data path
  input_data_name="example.vcf", # input data name,for vcf data
  input_data_type="VCF",         # input data type
  phased_genotype=TRUE,          # whether the vcf data has been phased
  haplotype_window_nSNP=5,       # according to nSNP define block,
  output_data_type=c("Haplotype","Numeric"),# output data format
  return_result=TRUE,            #save result as a R environment
  variable
  cpu_cores=1                    # number of cpu
)
```

### 格式转换-通过bigmemory方法

```
library(blupADC)
data_path=system.file("extdata", package = "blupADC") # path of example files
phased_result=geno_format(
  input_data_path=data_path,      # input data path
  input_data_name="example.vcf", # input data name,for vcf data
  input_data_type="VCF",         # input data type
  phased_genotype=TRUE,          # whether the vcf data has been phased
  haplotype_window_nSNP=5,       # according to nSNP define block,
  bigmemory_cal=TRUE,            # format conversion via bigmemory
  object
  bigmemory_data_path=getwd(),    # path of bigmemory data
  bigmemory_data_name="test_blupADC", #name of bigmemory data
  output_data_type=c("Haplotype","Numeric"),# output data format
  return_result=TRUE,            #save result in R environment
  cpu_cores=1                    # number of cpu
)
```

## 输出

输出结果主要分为6个部分(长度为6的列表)，分别为：

- **hmp** : Hapmap 格式的基因型数据

第1列为SNP，第3列为染色体，第4列为物理位置，第12列开始为基因型数据

rs#	alleles	chrom	pos	strand	assembly	center	protLSID	assayLSID	panelLSID	QCcode	3098	3498	3297	2452
SNP1	NA	1	224488	NA	NA	NA	NA	NA	NA	NA	CC	AC	AC	CC
SNP2	NA	1	293696	NA	NA	NA	NA	NA	NA	NA	GG	TG	TG	GG
SNP3	NA	1	333333	NA	NA	NA	NA	NA	NA	NA	GG	TT	TT	GG
SNP4	NA	1	464830	NA	NA	NA	NA	NA	NA	NA	CC	CC	CC	CC
SNP5	NA	1	722623	NA	NA	NA	NA	NA	NA	NA	AA	GG	GG	AA
SNP6	NA	1	838596	NA	NA	NA	NA	NA	NA	NA	CC	TC	TT	CC

- **ped** : Plink 格式的基因型数据

第1列为家系，第2列为个体号，第7列开始为基因型数据。

3098	3098	0	0	0	0	C	C	G	G
3498	3498	0	0	0	0	A	C	T	G
3297	3297	0	0	0	0	A	C	T	G
2452	2452	0	0	0	0	C	C	G	G
4255	4255	0	0	0	0	A	C	G	G
2946	2946	0	0	0	0	C	C	T	G

- **map**: Plink 格式的基因型数据

第1列为染色体，第2列为SNP，第3列为遗传距离(cM)，第4列为物理位置。

1	SNP1	0.224488	224488
1	SNP2	0.293696	293696
1	SNP3	0.333333	333333
1	SNP4	0.464830	464830
1	SNP5	0.722623	722623
1	SNP6	0.838596	838596

- **blupf90**: BLUPF90 格式的基因型数据

第1列为个体号，第2列为基因型。

3098	200000
3498	112021
3297	112022
2452	200000
4255	102011
2946	212000

- **numeric:** Numeric 格式的基因型数据

行名为个体，列名为SNP，0,1,2表示的是个体在某个SNP位点的基因型数据

2	0	0	0	0	0
1	1	2	0	2	1
1	1	2	0	2	2
2	0	0	0	0	0
1	0	2	0	1	1
2	1	2	0	0	0

- **haplotype\_hap:** Haplotype 格式的基因型数据

行表示的为SNP,列表式的是个体，每个个体占两列。

0	0	0	1	1	0	0	0
0	0	1	0	0	1	0	0
1	1	0	0	0	0	1	1
0	0	1	1	1	1	0	0
0	0	0	1	1	1	0	0

- **haplotype\_sample:** Haplotype 格式的基因型数据

基因型数据的个体名称。

3098
3498
3297
2452
4255
2946

- **haplotype\_map:** Haplotype 格式的基因型数据

--	--	--	--	--



1	SNP1	224488	C	A
1	SNP2	293696	G	T
1	SNP3	333333	T	G
1	SNP4	464830	A	G
1	SNP5	722623	C	T
1	SNP6	838596	C	A

- **vcf**: VCF 格式的基因型数据

##fileformat=VCFv4.2										
##source="beagle.29May21.d6d.jar"										
##INFO<ID=AF,Number=A,Type=Float>										
##INFO<ID=IMP,Number=0,Type=Flag>										
##FORMAT<ID=GT,Number=1,Type=String>										
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	3498	3297
1	6260	M2	T	A	.	PASS	.	GT	1 0	0 1
1	15289	M17	A	T	.	PASS	.	GT	0 0	0 0

## 参数

### Basic

- **参数1:input\_data\_plink\_ped**

用户提供的 Plink-ped格式的数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数2:input\_data\_plink\_map**

用户提供的 Plink-map格式的数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数3:input\_data\_hmp**

用户提供的 Hapmap格式的数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数4:input\_data\_blupf90**

用户提供的 BLUPF90格式数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数5:input\_data\_numeric**

用户提供的 Numeric格式的数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数6:input\_data\_haplotype\_hap**

用户提供的 Haplotype 格式的数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数7:input\_data\_haplotype\_sample**

用户提供的 Haplotype 格式的数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数8:input\_data\_haplotype\_map**

用户提供的 Haplotype 格式的数据， `data.frame` or `matrix` 类型。具体格式见结果部分

- **参数9:input\_data\_vcf**

用户提供的 VCF 格式的数据, `data.frame` or `matrix` 类型。具体格式见结果部分

**Note:** `input_data_numeric` 数据应该包含行名和列名。

此外, 为了方便用户使用, 用户还可以直接通过提供本地数据的路径、名称与数据类型即可完成数据提供这一步骤, 而无需将数据读入到R里面。

- **参数10:input\_data\_type**

用户提供的本地数据的格式, `character` 类型。数据格式包括:

- Hapmap
- Plink
- BLUPF90
- Numeric
- Haplotype
- VCF

- **参数11:input\_data\_path**

用户提供的本地数据的文件路径, `character` 类型。

- **参数12:input\_data\_name**

用户提供的本地数据的文件名称, `character` 类型。

**Note:** 如果提供的数据类型为Plink, 那么本地文件名称不需要带后缀, eg. 本地文件名为test1.map test1.ped, 我们提供文件名称为: `input_data_name="test1"`。除了Plink格式的数据外, 其他数据格式必须提供完整的名称(带后缀)。

- **参数13:output\_data\_path**

输出的基因型数据保存到本地的路径, `character` 类型。

- **参数14:output\_data\_name**

输出的基因型数据保存到本地的文件名称, `character` 类型。

- **参数15:output\_data\_type**

用户提供的本地数据的格式, `character` 类型。数据格式包括:

- Hapmap
- Plink
- BLUPF90
- Numeric
- Haplotype
- VCF

- **参数16:return\_result**

是否在R中返回输出的结果, `logical` 类型。默认为FALSE。

- **参数17:bigmemory\_cal**

是否使用bigmemory方式进行计算. `logical` 类型. 默认为 FALSE.

- **参数18:bigmemory\_data\_path**

bigmemory数据保存的路径. `character` 类型.

- **参数19:bigmemory\_data\_name**

bigmemory数据保存的文件名称. `character` 类型.

- **参数20:phased\_genotype**

是否基因型数据已经经过定向. `logical` 类型.默认为 `FALSE`.

- **参数21:haplotype\_window\_nSNP**

根据连续的SNP数目来定义单倍型block. `numeric` 类型.

- **参数22:haplotype\_window\_kb**

根据物理位置信息来定义单倍型block. `numeric` 类型.

- **参数23:haplotype\_window\_block**

根据用户自定义的信息来定义单倍型block. `data.frame` or `matrix` 类型.

第一列是window起始位置, 第二列是window结束位置

1	5
6	10
11	15
16	20
21	25
26	30

## 🔑 Advanced

- **参数24:cpu\_cores**

函数调用的cpu个数, `numeric` 类型。默认调用1个

- **参数25:miss\_base**

缺失值在原基因型数据中所表示的的字符, `character` 类型。默认为"NN".

- **参数26:miss\_base\_num**

数字化格式转换中缺失值转换成的数字, `numeric` 类型。默认为 5。

## 功能2

### 简介

🔍 通常来讲, 我们公司拿到的原始芯片数据大都是包含缺失值且未经过质控的。而在实际的数据分析中, 我们一般都要求数据是经过质控和填充。为此, `blupADC` 中提供了 `geno_qc_impute` 函数, 可以方便我们在R中调用**Plink**(用于质控)和**Beagle**(用于填充)软件进行基因型数据的质控和填充。

📌 **Note:** 为了方便用户使用, `blupADC` 已经事先封装好了 **Plink**(用于质控) `version-1.9` 和 **Beagle**(用于填充) `version-5.2` 软件, 用户无需再重新下载. 如果用户想自行指定软件的版本, 可以通过更改相关的参数(在进阶参数部分里).

老规矩, 我们还是用一个小例子来看下函数的用法:

## 示例

```
library(blupADC)
geno_qc_impute(
  input_data_hmp=example_data_hmp,      #provided hapmap data object
  data_analysis_method="QC_Imputation",  #analysis method type,QC +
  imputatoin
  output_data_path=getwd(),              #output data path
  output_data_name="YY_data",            #output data name
  output_data_type="VCF"                 #output data format
)
```

通过上述代码，我们即可对本地的**Hapmap**格式的基因型数据进行质控和填充，并将其以**VCF**格式并保存到本地。

进行质控和填充时，我们必须事先提供基因型数据，这部分的参数与 `geno_format` 函数中的参数用法一致。具体大家可参阅之前的介绍: [基因型数据间的格式转换](#)。

完成了基因型数据的提供后，我们可以通过以下参数进行质控填充的相关分析。

## 参数详解

### ♥ 基础参数

- **参数1: data\_analysis\_method**

指定基因型数据的处理方法，`character` 类型。可选方法包括：

- "QC" :进行质控
- "Imputation" :进行填充
- "QC\_Imputation" :先质控，后填充

- **参数2: qc\_snp\_rate**

使用Plink进行质控时，SNP call rate 的阈值，`numeric` 类型，默认为0.1。

- **参数3: qc\_ind\_rate**

使用Plink进行质控时，IND call rate 的阈值，`numeric` 类型，默认为0.1。

- **参数4: qc\_maf**

使用Plink进行质控时，MAF 的阈值，`numeric` 类型，默认为0.05。

- **参数5: qc\_hwe**

使用Plink进行质控时，哈代温伯格平衡的阈值，`numeric` 类型，默认为1e-6。

### ♥ 进阶参数

- **参数6: plink\_software\_path**

Plink软件的路径，`character` 类型。

- **参数7: plink\_software\_name**

Plink软件的名称，`character` 类型。

- **参数8: beagle\_software\_path**

Beagle软件的路径，`character` 类型。

- **参数9: beagle\_software\_name**

Beagle软件的名称, `character` 类型。

- **参数10: beagle\_ref\_data\_path**

使用beagle进行填充时, 提供的reference data的路径, `character` 类型。

- **参数11: beagle\_ref\_data\_name**

使用beagle进行填充时, 提供的reference data的名称, `character` 类型。

- **参数12: beagle\_ped\_path**

使用beagle进行填充时, 提供的系谱数据的路径, `character` 类型。

- **参数13: beagle\_ped\_name**

使用beagle进行填充时, 提供的系谱数据的名称, `character` 类型。

## 功能3

### Overview

🧐 品种成分分析一直以来都是数据分析中的一个难题。blupADC 为用户提供了 `geno_check` 函数, 使得用户能够方便的解决这个问题。此外, 用户还可以用这个函数进行基因型数据的重复性检测。

### 示例

#### 品种成分分析

```
library(blupADC)
check_result=geno_check(
  input_data_hmp=example_PCA_data_hmp,    #provided hapmap data
  object
  duplication_check=FALSE,                 #whether check the duplication
  of genotype
  breed_check=TRUE,                       # whether check the record of
  breed
  breed_record=example_PCA_Breed,         # provided breed record
  return_result=TRUE                      #return result
)
```

#### 重复性检测

```
library(blupADC)
check_result=geno_check(
  input_data_hmp=example_data_hmp,    #provided hapmap data
  object
  duplication_threshold=0.95, #threshold of duplication
  duplication_check=TRUE,     #whether check the duplication of
  genotype
  breed_check=FALSE,          # whether check the record of breed
  return_result=TRUE          #return result
)
```

输出

输出的结果主要包括以下两个部分，如下：

- duplicated\_genotype

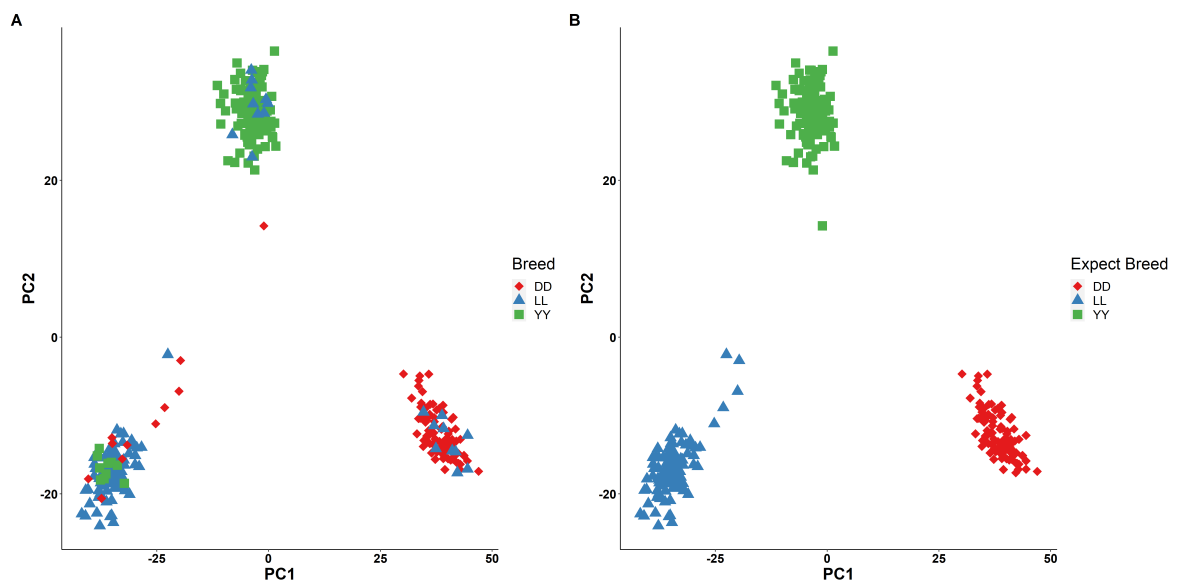
IND1	IND1	1
IND2	IND2	1
IND3	IND3	1
IND4	IND4	1

第一列和第二列为个体名称，第三列为重复的比例

- pca\_outlier

Id	Breed	Expeced_Breed
IND100	LL	YY
IND233	DD	YY
IND91	LL	YY
IND92	LL	YY
IND93	LL	YY
IND94	LL	YY

图A是进行品种分析前的品种记录，图B是进行品种分析后(可以理解为对错误的品种记录数据进行纠正)的品种记录



参数

geno\_check 函数中的许多参数均与 geno\_format 函数中一致。故此，接下来我们将主要介绍 geno\_check 函数中独有的参数[see more details](#)).

- 1: selected\_snps

进行基因型数据重复性检测时,所选用的SNP数目, `numeric` 类型. 默认为 1000.

- **2: overlap\_threshold**

判定两个个体为重复的阈值, `numeric` 类型. 默认为 0.95.

- **3: duplication\_check**

是否进行基因型数据重复性检测, `logical` 类型. 默认为 TRUE.

- **4: breed\_check**

是否进行品种分析, `logical` 类型. 默认为 FALSE.

- **5: ind\_breed**

个体的品种记录数据, `data.frame` 类型.

`ind_breed` 数据格式如下所示:

Id	Breed
IND1	YY
IND2	YY
IND3	YY
IND4	YY
IND5	YY
IND6	YY

## 功能4

### 简介

大家好,这一节主要给大家讲述的是如何使用 `b1upADC` 中的函数来进行系谱数据处理。`b1upADC` 提供的 `trace_pedigree` 函数,可以帮助我们非常方便的的对系谱数据进行多种处理:包括系谱重命名、系谱追溯及系谱纠错等。

### 示例

同样的,我们还是用一个小例子来简单的看下该函数的用法

```
library(b1upADC)
pedigree_result=trace_pedigree(
  input_pedigree=origin_pedigree,    #provided pedigree data object
  output_pedigree_tree=TRUE          # output pedigree tree
)
```

我们可以通过 `str` 查看函数的输出结果, 如下所示:

```
str(pedigree_result)
## List of 5
```

```
## $ ped      : chr [1:15945, 1:3] "DD19348310" "DD19386807" "DD19119705"
"DD16007415" ...
## ..- attr(*, "dimnames")=List of 2
## .. ..$ : NULL
## .. ..$ : chr [1:3] "Offspring" "Sire" "Dam"
## $ rename_ped : 'data.frame':  15945 obs. of  6 variables:
## ..$ offspring : chr [1:15945] "DD19348310" "DD19386807" "DD19119705"
"DD16007415" ...
## ..$ Generation : num [1:15945] 0 0 0 0 0 0 0 0 0 0 ...
## ..$ offspring_Id: int [1:15945] 1 2 3 4 5 6 7 8 9 10 ...
## ..$ Sire_Id      : num [1:15945] 0 0 0 0 0 0 0 0 0 0 ...
## ..$ Dam_Id       : num [1:15945] 0 0 0 0 0 0 0 0 0 0 ...
## ..$ Order        : int [1:15945] 1 2 3 4 5 6 7 8 9 10 ...
## $ pedigree_tree: chr [1:15945, 1:15] "DD19348310" "DD19386807" "DD19119705"
"DD16007415" ...
## ..- attr(*, "dimnames")=List of 2
## .. ..$ : NULL
## .. ..$ : chr [1:15] "Offspring" "Sire" "Dam" "SireSire" ...
## $ error_id_set :List of 4
## ..$ error_duplicated_id: chr [1:24] "DD19119705" "DD20488904" "DD20153801"
"DD20376912" ...
## ..$ error_sex_id: chr "DD13006182"
## ..$ error_breed_id: NULL
## ..$ error_birth_date_id: NULL
```

可以很明显的看到,输出结果包括以下几个部分:

- **ped:** 经过处理后(纠错、追溯等)的原始系谱数据且未进行重命名
- **rename\_ped:** 经过处理(纠错、追溯等)且重命名的系谱数据。第1列为原始系谱ID,第2列为个体在系谱中的代数,第3-5列为重命名后的系谱数据
- **pedigree\_tree:** 个体的系谱树矩阵。可以通过设置 `pedigree_tree_depth` 指定系谱树包含的代数,默认不输出系谱树(节省时间)
- **error\_id\_set:** 系谱记录错误个体数据集。根据错误的种类可以分为以下4个子集
  - `error_duplicated_id`: 相同的个体,父母却不相同
  - `error_sex_id`: 个体同时出现在父亲列和母亲列
  - `error_breed_id`: 个体和父母的品种不相同(仅针对特殊格式)
  - `error_birth_date_id`: 个体的出生日期要早于父母的出生日期(需在系谱的第四列提供个体的出生日期)

下面,我们将具体讲解 `trace_pedigree` 函数中各种参数的含义:

## 参数详解

### ✧ 基础参数

- **参数1: input\_pedigree**

用户提供的系谱数据, `data.frame` 或 `matrix` 类型。

📁 用户提供系谱数据需为以下几种格式中的一种,包括:

- 3列系谱格式:



Offspring	Sire	Dam
DD19575312	DD18768902	DD16376015
DD19513112	DD18768902	DD17111017
DD20348012	DD19132207	DD19234510
DD20361110	DD19331001	DD19293112
DD20471212	DD19331001	DD19320808
DD20564818	DD19331001	DD19311009

- 4列系谱格式:

Offspring	Sire	Dam	Birth_Date
DD19575312	DD18768902	DD16376015	20200101
DD19513112	DD18768902	DD17111017	20200102
DD20348012	DD19132207	DD19234510	20200103
DD20361110	DD19331001	DD19293112	20200104
DD20471212	DD19331001	DD19320808	20200105
DD20564818	DD19331001	DD19311009	20200106

- 多列系谱格式:

Offspring	Sire	Dam	SireSire	DamSire	SireSireSire
DD20231905	DD19581002	DD18750810	DD16785512	DD15507717	DD14008512
DD20458701	DD19564302	DD18925809	DD15349017	DD15245411	DD16771212
DD20324707	DD19232903	DD18571012	DD16794714	DD16744412	DD16714516
DD19288609	DD18713408	DD18552609	DD15180015	DD15479214	DD15243711
DD16222012	DD15145005	DD15378812	DD14110014	DD15501518	DD15206217
DD17684713	DD16672107	DD15122311	DD15505715	DD15347415	DD16383111

**Note:**需要注意的是，当系谱为多列时，用户必须要设置 `multi_col=TRUE`，并且系谱的列名需要标注为特殊形式,e.g. SireSire:父亲的父亲, SirSireSire:父亲的父亲的父亲.

系谱数据中缺失值可以设置为: **NA或0**。

同样的，为了便于用户操作，用户还可以直接提供本地系谱数据的路径和名称

- **参数2: input\_pedigree\_path**

本地系谱数据的路径, `character` 类型。

- **参数3: input\_pedigree\_name**

本地系谱数据的名称, `character` 类型。

- **参数4: multi\_col**

是否将提供的多列系谱转换到3列, `logical`类型。如果用户提供的系谱数据包含多列, 那么用户必须设置 `multi_col` 。

- **参数5: trace\_id**

在追溯系谱时, 所追溯的个体集, `character` 类型. 默认为 `NULL` , 即追溯系谱中的所有个体

- **参数6: trace\_generation**

在追溯系谱时, 所追溯的最大代数, `numeric` 类型。默认为 `NULL` , 即追溯系谱中的全部代数。

- **参数7: trace\_birth\_date**

追溯出生日期不早于指定日期的个体, `numeric` 类型。个体出生日期早于用户提供的出生日期将会被排除在系谱追溯过程中。

- **参数8: output\_pedigree\_path**

系谱输出到本地的路径, `character` 类型。

- **参数9: output\_pedigree\_name**

系谱输出到本地的名称, `character` 类型。

## 🔗 进阶参数

- **参数10: dup\_error\_check**

检查相同个体的父母却不相同的错误, `logical` 类型, 默认为TRUE。

- **参数11: sex\_error\_check**

检查个体同时出现在父亲列和母亲列的错误, `logical` 类型, 默认为TRUE。

- **参数12: birth\_date\_error\_check**

检查个体出生日期早于父母的错误, `logical` 类型, 默认为FALSE。

- **参数13: output\_pedigree\_tree**

是否输出系谱树, `logical` 类型, 默认为FALSE。

- **参数14: pedigree\_tree\_depth**

系谱树的深度(系谱代数), `numeric` 类型, 默认为3。

## 功能5

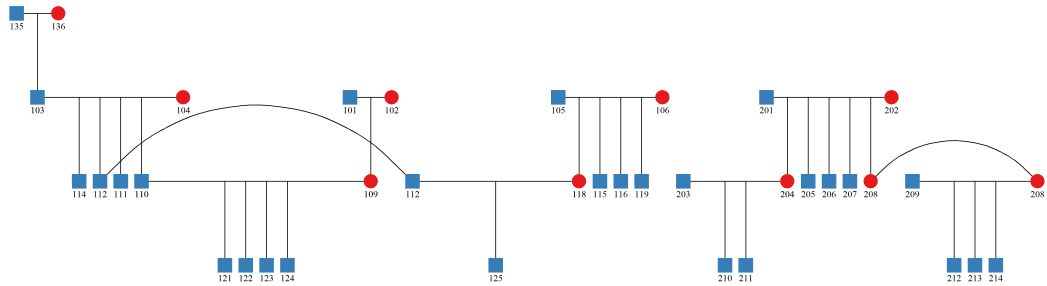
### 简介

🖼️ 一个直观和清楚的系谱结构图能够帮助育种家和研究者做出更好的育种规划。通过使用 `ggped` 函数, 用户即能非常简单的绘制出所需的系谱图。

### 示例

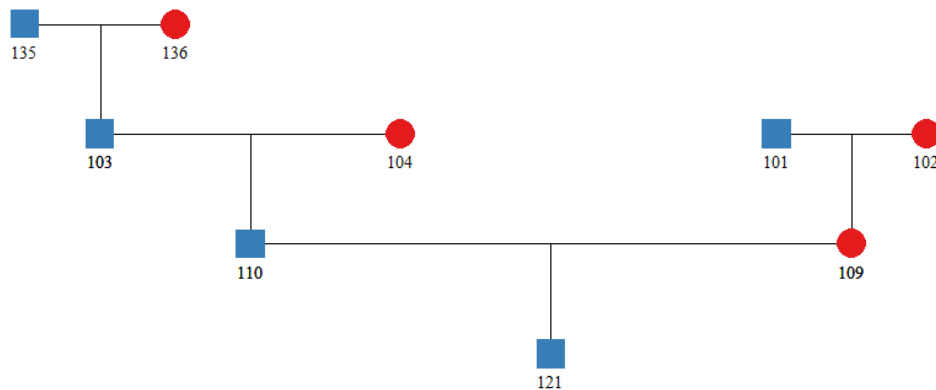
#### 绘制系谱中所有个体

```
library(blupADC)
pedigree_result=ggped(
  input_pedigree=example_ped2
)
```



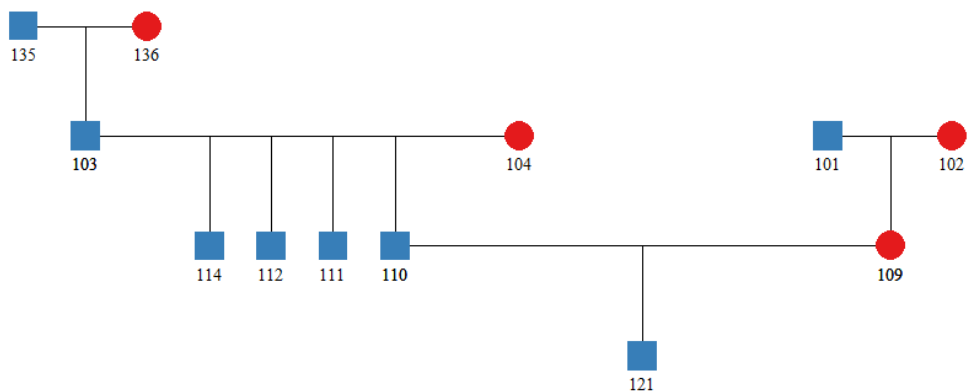
## 绘制系谱中的子集

```
library(blupADC)
pedigree_result=ggped(
  input_pedigree=example_ped2,
  trace_id=c("121") #provided subset-id
)
```



## 绘制系谱中的子集 (考虑子集的同胞)

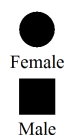
```
library(blupADC)
pedigree_result=ggped(
  input_pedigree=example_ped2,
  trace_id=c("121"),
  trace_sibs=TRUE #whether plot the sibs of subset-id
)
```



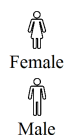
## 系谱图的样式

通过修改 `shape_type` 参数， 用户即可改变系谱图的样式。默认的 `shape_type` 为 4。

Shape type=1



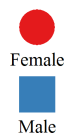
Shape type=2



Shape type=3



Shape type=4



Shape type=5



Shape type=6



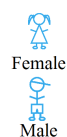
Shape type=7



Shape type=8



Shape type=9



## 参数详解

### ✳️ 基础参数

- 参数1: `input_pedigree`

用户提供的系谱数据, `data.frame` 或者 `matrix` 类型.

📁 提供的系谱数据类型应为3列如下的3列结构:

Offspring	Sire	Dam
DD19575312	DD18768902	DD16376015
DD19513112	DD18768902	DD17111017
DD20348012	DD19132207	DD19234510
DD20361110	DD19331001	DD19293112
DD20471212	DD19331001	DD19320808
DD20564818	DD19331001	DD19311009

- **参数2: trace\_id**

追溯子集的系谱, `character` 类型. 默认为 `NULL` (绘制系谱中所有个体)

- **参数3: trace\_sibs**

追溯子集的系谱过程中, 是否追溯子集的同胞, `logical` 类型. 默认为 `FALSE` .

- **参数4: ind\_sex**

个体的性别记录, `data.frame` 类型.

数据格式如下:

Individual	Sex
101	Male
102	Female
103	Male
104	Female

- **参数5: plot\_family**

绘制系谱过程中, 是否划分家系结构, `logical` 类型. 默认为 `FALSE` .

- **参数6: shape\_type**

系谱图的样式, `numeric` 类型. 默认为 4.

## 🔗 进阶参数

- **参数7: shape\_size**

系谱图的大小, `numeric` 类型. 默认为 8 .

- **参数8: ind\_text\_size**


个体名称的文本大小, `numeric` 类型. 默认为 4` .

- **参数9: ind\_text\_vjust**

个体名称的文本垂直距离, `numeric` 类型. 默认为 3 .

## 功能6

### 简述

 在动植物育种中，亲缘关系矩阵的构建是其中的关键步骤。在本章，我们将主要介绍如何利用 blupADC 中的 cal\_kinship 函数完成各种亲缘关系矩阵的构建，包括：**加性亲缘关系矩阵**(系谱，基因组，一步法)及**显性亲缘关系矩阵**(系谱，基因组，一步法)等。此外，cal\_kinship函数还能方便的计算各种类型的**近交系数**，包括：系谱近交系数，基因组近交系数(Homozygous,Digonal)及一步法近交系数(Digonal)。

在构建基因组亲缘关系矩阵及一步法亲缘关系矩阵的时候，我们必须事先提供基因型数据，这部分的参数与 geno\_format 函数中的参数用法一致。具体大家可参阅之前的介绍:[基因型数据间的格式转换](#)。

### 示例

#### 构建系谱亲缘关系矩阵

```
library(blupADC)
kinship_result=cal_kinship(
  object              input_pedigree=example_ped3,          #provided hapmap data
                      kinship_type=c("P_A"),              #type of kinship matrix
                      inbred_type=c("Pedigree"),            #type of inbreeding coefficients
                      return_result=TRUE)                   #return result
```

#### 构建基因组亲缘关系矩阵

```
library(blupADC)
kinship_result=cal_kinship(
  object              input_data_hmp=data_hmp,              #provided hapmap data object
                      kinship_type=c("G_A","G_D"),          #type of kinship matrix
                      dominance_type=c("genotypic"),         #type of dominance effect
                      inbred_type=c("Homozygous"),           #type of inbreeding
                      coefficients
                      return_result=TRUE)                   #return result
```

#### 构建一步法亲缘关系矩阵

```
library(blupADC)
kinship_result=cal_kinship(
  object              input_data_hmp=example_data_hmp,      #provided hapmap data
                      input_pedigree=example_ped3,
                      kinship_type=c("H_A"),                #type of kinship matrix
                      inbred_type=c("H-diag"),              #type of inbreeding coefficients
                      return_result=TRUE)                   #return result
```

#### 构建一步法亲缘关系矩阵(via bigmemory method)

```
library(blupADC)
phased_kinship_result=cal_kinship(
  object          input_data_hmp=example_data_hmp,          #provided hapmap data
  input_pedigree=example_ped3,
  bigmemory_cal=TRUE,
  bigmemory_data_path=getwd(),
  bigmemory_data_name="blupADC",
  kinship_type=c("H_A"),          #type of kinship matrix
  inbred_type=c("H_diag"),        #type of inbreeding coefficients
  return_result=TRUE)            #return result
```

## 参数详解

### 基础参数

- **参数1: kinship\_type**

指定构建亲缘关系矩阵的类型，`character` 类型。可选关系矩阵类型：

- "G\_A" :基因组加性亲缘关系矩阵"
- "G\_Ainv" :基因组加性亲缘关系逆矩阵"
- "G\_D" :基因组显性亲缘关系矩阵"
- "G\_Dinv" :基因组显性亲缘关系逆矩阵"
- "P\_A" :系谱加性亲缘关系矩阵"
- "P\_Ainv" :系谱加性亲缘关系逆矩阵"
- "P\_D" :系谱显性亲缘关系矩阵"
- "P\_Dinv" :系谱显性亲缘关系逆矩阵"
- "H\_A" :一步法加性亲缘关系矩阵"
- "H\_Ainv" :一步法加性亲缘关系逆矩阵"
- "H\_D" :一步法显性亲缘关系矩阵"
- "H\_Dinv" :一步法显性亲缘关系逆矩阵"

**Note:**如果计算系谱及一步法亲缘关系矩阵，必须要提供系谱数据。关于如何提供系谱数据，我们会在后面进行讲解。

- **参数2: dominance\_type**

指定计算的显性亲缘关系矩阵的类型，`character` 类型，可选类型包括：

- "genotypic" : 按照 $(0 - 2pq, 1 - 2pq, 0 - 2pq)$ 方式编码方式构建显性亲缘关系矩阵
- "classical" : 按照 $(-2q^2, 2pq, -2p^2)$ 方式编码方式构建显性亲缘关系矩阵

关于二者的区别，具体可阅读该文献：[On the Additive and Dominant Variance and Covariance of Individuals Within the Genomic Selection Scope](#)

- **参数3: inbred\_type**

指定计算的近交系数的类型，`character` 类型。可选近交系数类型包括：

- "Homozygous" :根据纯合子位点计算
- "G\_Diag" :G矩阵对角线-1
- "H\_diag" :H矩阵对角线-1
- "Pedigree" :A矩阵对角线-1

- **参数4: input\_pedigree**

用户提供的系谱数据，`data.frame` 或 `matrix` 类型。具体的系谱数据格式可以参阅之前的介绍：[系谱追溯、重命名及纠错](#)。

- **参数5: IND\_rename**

是否根据系谱的重命名结果对基因型数据中的个体进行重命名, `logical` 类型, 默认为FALSE(不进行重命名)。

- **参数6:bigmemory\_cal**

是否使用bigmemory方式进行计算. `logical` 类型. 默认为 FALSE.

- **参数7:bigmemory\_data\_path**

bigmemory数据保存的路径. `character` 类型.

- **参数8:bigmemory\_data\_name**

bigmemory数据保存的文件名称. `character` 类型.

- **参数9: output\_matrix\_type**

输出亲缘关系矩阵的格式, `character` 类型。可选参数包括:

1. "col\_all" :按照n\*n的格式输出亲缘关系矩阵
2. "col\_three" : 按照3列矩阵的格式输出亲缘关系矩阵, 第1列和第2列为个体号, 第3列为亲缘系数。DMU和BLUPf90软件均需提供这种格式的亲缘关系矩阵。形如:

1001	1001	0.989
1001	1002	0.421
1001	1003	0.567

默认参数为 "col\_all"

- **参数10: output\_matrix\_path**

亲缘关系矩阵保存到本地的路径, `character` 类型。

- **参数11: output\_matrix\_name**

亲缘关系矩阵保存到本地的名称, `character` 类型。

## 进阶参数

- **参数12: cpu\_cores**

调用的cpu个数, `numeric` 类型, 默认为1。

- **参数13: kinship\_base**

是否按照基础群的方式构建基因组亲缘关系矩阵( $p=q=0.5$ ), `logical` 类型, 默认为FALSE。

- **参数14: kinship\_trace**

是否按照矩阵迹和的方式对基因组亲缘关系矩阵进行标准化, `logical` 类型, 默认为FALSE。

- **参数15: Metafounder\_algorithm**

是否按照metafounder的方法计算一步法亲缘关系矩阵, `logical` 类型, 默认为FALSE。

- **参数16: APY\_algorithm**

是否按照APY的方法计算亲缘关系矩阵的逆矩阵, `logical` 类型, 默认为FALSE。



- **参数17: APY\_eigen\_threshold**

特征值所能解释的遗传变异的比例的阈值, `numeric` 类型. 默认为 0.95.

- **参数18: APY\_n\_core**

核心个体数, `numeric` 类型. 默认为 NULL.

- **参数19: SSBLUP\_omega**

构建一步法亲缘关系矩阵时G矩阵和A矩阵的比例, `numeric` 类型, 默认为0.05。

- **参数20: gene\_dropping**

是否使用 gene dropping 的方法构建系谱显性亲缘关系矩阵, `logical` 类型, 默认为FALSE。

- **参数21: gene\_dropping\_iteration**

gene dropping方法的迭代次数, `numeric` 类型, 默认为1000。

## 功能7

### 简述

☺在讲述完各种各样的数据预处理方法后, 我们正式进入到育种数据的分析层面。在目前的动植物育种领域, 主要的育种软件包括但不限于以下两种: **DMU**和**BLUPf90**。这两款软件均于是于上世纪80-90年代开发的, 并且一直处于维护中。截至目前, 两款软件的引用次数均已接近千次(实际可能更多), 这也足见这两款软件的认可度。

但是, 这两款软件均存在一个共同的缺点, 就是使用起来较为麻烦(需要提供准备好的参数文件)。笔者当时学习参数文件的配置时, 前前后后花费了近一个月的时间, 足以见这两款软件的不友好性☹。

为此, 我们在R中编写了相应的函数, 使得用户可以轻松的完成两款软件参数文件的配置及对应的数据分析。本章我们主要讲述如何通过 `BLUP_ADC` 中的 `run_DMU` 函数, 在R中调用**DMU**软件并完成数据分析。在下一章, 我们将会讲述如何在R中调用**BLUPf90**软件, 并完成数相应的据分析。

📖 **Note:** 为了方便用户使用, `blupADC` 已经封装了DMU中的几个基本模块(`dmu1`, `dmuai`, and `dmu5`), 更多的模块可以从DMU官网进行下载([DMU下载网站](#))。

如果您想将DMU用作商业用途, 请务必联系 DMU的作者!!!

### 示例

#### 单性状模型-系谱

```
library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2",
"Age"), # colnames of phenotype
  target_trait_name=list(c("Trait1")), #性状名称
  fixed_effect_name=list(c("Sex", "Herd_Year_Season")), #固定效应名称
  random_effect_name=list(c("Id", "Litter")), #随机效应名称
  covariate_effect_name=NULL, #协变量效应名称
  genetic_effect_name="Id", #遗传效应名称
```

```

phe_path=data_path,
phe_name="phenotype.txt",
integer_n=5,
analysis_model="PBLUP_A",
dmu_module="dmuai",
relationship_path=data_path,
relationship_name="pedigree.txt",
output_result_path=getwd()
)
#表型文件路径
#表型文件名
#整型变量数
#遗传评估模型
#方差组分估计使用的DMU模块
#亲缘关系文件路径
#亲缘关系文件名
#结果输出路径

```

## 单性状模型-GBLUP

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2",
"Age"), # colnames of phenotype
target_trait_name=list(c("Trait1")),
fixed_effect_name=list(c("Sex", "Herd_Year_Season")),
random_effect_name=list(c("Id", "Litter")),
covariate_effect_name=NULL,
genetic_effect_name="Id",
phe_path=data_path,
phe_name="phenotype.txt",
integer_n=5,
analysis_model="GBLUP_A",
dmu_module="dmuai",
relationship_path=data_path,
relationship_name="G_Ainv_col_three.txt",
output_result_path=getwd()
)
#性状名称
#固定效应名称
#随机效应名称
#协变量效应名称
#遗传效应名称
#表型文件路径
#表型文件名
#整型变量数
#遗传评估模型
#方差组分估计使用的DMU模块
#亲缘关系文件路径
#亲缘关系文件名
#结果输出路径

```

## 单性状模型-Single-step(一步法)

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2",
"Age"), # colnames of phenotype
target_trait_name=list(c("Trait1")),
fixed_effect_name=list(c("Sex", "Herd_Year_Season")),
random_effect_name=list(c("Id", "Litter")),
covariate_effect_name=NULL,
genetic_effect_name="Id",
phe_path=data_path,
phe_name="phenotype.txt",
integer_n=5,
analysis_model="SSBLUP_A",
dmu_module="dmuai",
relationship_path=data_path,
relationship_name=c("pedigree.txt", "G_A_col_three.txt"),
)
#性状名称
#固定效应名称
#随机效应名称
#协变量效应名称
#遗传效应名称
#表型文件路径
#表型文件名
#整型变量数
#遗传评估模型
#方差组分估计使用的DMU模块
#亲缘关系文件路径
#亲缘关系文件

```

```
output_result_path=getwd()          #结果输出路径
)
```

细心的同学应该能发现，我们仅需改变 `analysis_model` 及 `relationship_name` 这两个参数即可完成系谱、GBLUP及SSBLUP的分析，极大的简化了我们的分析步骤(PS: `G_Ainv_col_three.txt` 和 `G_A_col_three.txt` 文件 均可通过 `cal_kinship` 函数得到 [了解更多](#))。

上面我们介绍的都是单性状模型(只包括了一个目标性状)，而在实际育种分析中，多性状模型也是非常常见。在上面代码的基础上稍作修改，我们就能够非常方便的完成多性状模型的运算，如下：

## 多性状模型-系谱

```
library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2",
"Age"), # colnames of phenotype
  target_trait_name=list(c("Trait1"), c("Trait2")),
  #性状名称

fixed_effect_name=list(c("Sex", "Herd_Year_Season"), c("Herd_Year_Season")), #
固定效应名称
  random_effect_name=list(c("Id", "Litter"), c("Id")), #随机效应名称
  covariate_effect_name=NULL, #协变量效应名称
  genetic_effect_name="Id", #遗传效应名称
  phe_path=data_path, #表型文件路径
  phe_name="phenotype.txt", #表型文件名
  integer_n=5, #整型变量数
  analysis_model="PBLUP_A", #遗传评估模型
  dmua_module="dmuai", #方差组分估计使用的DMU模块
  relationship_path=data_path, #亲缘关系文件路径
  relationship_name="pedigree.txt", #亲缘关系文件名
  output_result_path=getwd() #结果输出路径
)
```

## 单性状模型-系谱 (用户提供方差组分先验文件)

```
library(blupADC)
data_path=system.file("extdata", package = "blupADC") # path of provided files

run_DMU(phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter",
"Trait1", "Trait2", "Age"), # colnames of
phenotype
  target_trait_name=list(c("Trait1")), #trait name
  fixed_effect_name=list(c("Sex", "Herd_Year_Season")), #fixed effect
name
  random_effect_name=list(c("Id", "Litter")), #random effect
name
  covariate_effect_name=NULL, #covariate
effect name
  genetic_effect_name="Id", #遗传效应名称
  phe_path=data_path, #path of phenotype file
  phe_name="phenotype.txt", #name of phenotype file
```

```

        provided_prior_file_path=data_path,          #path of user-provided
prior file
        provided_prior_file_name="PAROUT",          #name of user-provided
prior file
        integer_n=5,                                #number of integer variable
        analysis_model="PBLUP_A",                  #model of genetic
evaluation
        dmu_module="dmuai",                          #module of estimating
variance components
        relationship_path=data_path,                #path of relationship file
        relationship_name="pedigree.txt",           #name of relationship file
        output_result_path=getwd()                 # output path
    )

```

## 单性状模型-系谱 (包含母性效应)

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(

    phe_col_names=c("Herd", "B_month", "D_age", "Litter", "Sex", "HY", "ID", "DAM", "L_Dam"
,
                    "W_birth", "W_2mth", "W_4mth", "G_0_2", "G_0_4", "G_2_4"), #
colnames of phenotype
    target_trait_name=list(c("W_birth")),          #trait
name
    fixed_effect_name=list(c("B_month", "D_age", "Litter", "Sex", "HY")),
#fixed effect name
    random_effect_name=list(c("ID", "L_Dam")),      #random effect name
    maternal_effect_name=list(c("DAM")),
    genetic_effect_name="ID",                       #遗传效应名称
    covariate_effect_name=NULL,                     #covariate effect name
    phe_path=data_path,                             #path of phenotype file
    phe_name="maternal_data",                       #name of phenotype file
    integer_n=9,                                     #number of integer variable
    analysis_model="PBLUP_A",                       #model of genetic
evaluation
    dmu_module="dmuai",                             #module of estimating
variance components
    relationship_path=data_path,                     #path of relationship file
    relationship_name="maternal_pedigree",           #name of relationship file
    output_result_path=getwd()                      # output path
)

```

## 单性状模型-系谱 (包含永久环境效应)

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(
    phe_col_names=c("id", "year_grp", "breed", "time", "t_dato",
                    "age", "L1", "L2", "L3", "gh"), # colnames of
phenotype
    target_trait_name=list(c("gh")),                #trait name
    fixed_effect_name=list(c("year_grp", "breed", "time")), #fixed effect name

```

```

name      random_effect_name=list(c("id","t_dato")),          #random effect
name      covariate_effect_name=list(c("age")),              #covariate effect
name      genetic_effect_name="id",                          #遗传效应名称
effect    included_permanent_effect=list(c(TRUE)),           #whether include permant
effect    phe_path=data_path,                                #path of phenotype file
           phe_name="rr_data",                                #name of phenotype file
           integer_n=5,                                        #number of integer variable
           analysis_model="PBLUP_A",                          #model of genetic
evaluation dmu_module="dmuai",                                #modeule of estimating
variance components relationship_path=data_path,             #path of relationship file
           relationship_name="rr_pedigree",                   #name of relationship file
           output_result_path=getwd()                         # output path
           )

```

## 单性状模型-系谱 ( 包含随机回归效应)

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(
  phe_col_names=c("id","year_grp","breed","time","t_dato",
                  "age","L1","L2","L3","gh"),           # colnames of
phenotype target_trait_name=list(c("gh")),              #trait name
           fixed_effect_name=list(c("year_grp","breed","time")), #fixed effect name
           random_effect_name=list(c("id","t_dato")),        #random effect
name      covariate_effect_name=list(c("age")),          #covariate effect
name      genetic_effect_name="id",                      #遗传效应名称
effect    included_permanent_effect=list(c(TRUE)),        #whether include permant

  random_regression_effect_name=list(c("L1&id","L1&pe_effect","L2&id","L2&pe_effec
ct")), #random regression effect name
           phe_path=data_path,                             #path of phenotype file
           phe_name="rr_data",                             #name of phenotype file
           integer_n=5,                                     #number of integer variable
           analysis_model="PBLUP_A",                        #model of genetic
evaluation dmu_module="dmuai",                             #modeule of estimating
variance components relationship_path=data_path,           #path of relationship file
           relationship_name="rr_pedigree",                 #name of relationship file
           output_result_path=getwd()                       # output path
           )

```

## 单性状模型-系谱( 包含 社会遗传效应)

用户提供的表型文件不需要包含 最大群体大小相关的列

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(
  phe_col_names=c("Id", "Group", "Sex", "Phe"), # colnames of phenotype
  target_trait_name=list(c("Phe")),           #trait name
  fixed_effect_name=list(c("Sex")),           #fixed effect name
  random_effect_name=list(c("Id", "Group")),  #random effect name
  covariate_effect_name=NULL,                 #covariate effect name
  genetic_effect_name="Id",                   #遗传效应名称
  include_social_effect=list(c(TRUE)),
  group_effect_name="Group",
  phe_path=data_path,                         #path of phenotype file
  phe_name="raw_social_data",                 #name of phenotype file
  integer_n=3,                                #number of integer variable
  analysis_model="PBLUP_A",                  #model of genetic
evaluation
  dmu_module="dmuai",                         #module of estimating
variance components
  relationship_path=data_path,                #path of relationship file
  relationship_name="socail_pedigree",        #name of relationship file
  output_result_path=getwd() # output path
)

```

## 单性状模型-系谱( 包含 社会遗传效应)

用户提供的表型文件需要包含 最大群体大小相关的列

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_DMU(phe_col_names=c("Id", "Group", "Sex", "Gr_id1", "Gr_id2", "Gr_id3", "Gr_id4", "
Gr_id5",

  "Phe", "Status_Gr_id1", "Status_Gr_id2", "Status_Gr_id3", "Status_Gr_id4", "Status_G
r_id5"), # colnames of phenotype
  target_trait_name=list(c("Phe")),           #trait name
  fixed_effect_name=list(c("Sex")),           #fixed effect name
  random_effect_name=list(c("Id", "Group")),  #random effect name
  covariate_effect_name=NULL,                 #covariate effect name
  genetic_effect_name="Id",                   #遗传效应名称
  include_social_effect=list(c(TRUE)),        #whether include social genetic
effect
  integer_group_names=c("Gr_id1", "Gr_id2", "Gr_id3", "Gr_id4", "Gr_id5"),
  #integer variable name of max group size
  real_group_names=
c("Status_Gr_id1", "Status_Gr_id2", "Status_Gr_id3", "Status_Gr_id4", "Status_Gr_id5
"), #real variable name of max group size
  phe_path=data_path,                         #path of phenotype file
  phe_name="social_data",                     #name of phenotype file
  integer_n=8,                                #number of integer variable
  analysis_model="PBLUP_A",                  #model of genetic
evaluation
  dmu_module="dmuai",                         #module of estimating
variance components
  relationship_path=data_path,                #path of relationship file
  relationship_name="socail_pedigree",        #name of relationship file
)

```

```
output_result_path=getwd() # output path
)
```

我们将对 `run_DMU` 函数中的参数——进行讲解。

## 参数详解

### 基础参数

- **参数1: phe\_path**

本地表型数据文件的路径, `character` 类型。

- **参数2: phe\_name**

本地表型数据文件的名称, `character` 类型。

- **参数3: phe\_col\_names**

表型数据的列名, `character` 类型。

- **参数4: integer\_n**

整型变量的数目, `numeric` 类型。

- **参数5: genetic\_effect\_name**

遗传效应的名称(一般为个体号), `character` 类型。

- **参数6: target\_trait\_name**

目标性状的名称, `list` 类型。每个列表均代表一个性状。

通过添加多个性状的名称, 我们即可完成多性状模型的分析, e.g.

`target_trait_name=list(c("Trait1"),c("Trait2"))` 即可完成 `Trait1` 和 `Trait2` 的双性状模型

- **参数7: fixed\_effect\_name**

目标性状的固定效应名称, `list` 类型。在多性状模型中, `fixed_effect_name` 为每个性状的固定效应名称向量组成的列表结构, 性状的顺序需与 `target_trait_name` ——对应。

e.g. 第一个性状的固定效应名称为: `c("Sex","Herd_Year_Season")`

第二个性状的固定效应名称为: `c("Sex")`

那么 `fixed_effect_name=list(c("Sex","Herd_Year_Season"),c("Sex"))`

- **参数8: random\_effect\_name**

目标性状的随机效应名称, `list` 类型。在多性状模型中, `random_effect_name` 为每个性状的随机效应名称向量组成的列表结构, 性状的顺序需与 `target_trait_name` ——对应。

e.g. 第一个性状的随机效应名称为: `c("Id","Litter")`

第二个性状的随机效应名称为: `c("Id")`

那么 `random_effect_name=list(c("Id","Litter"),c("Id"))`

- **参数9: covariate\_effect\_name**

目标性状的协变量效应名称, `list` 类型。在多性状模型中, `random_effect_name` 为每个性状的协变量效应名称向量组成的列表结构, 性状的顺序需与 `target_trait_name` ——对应。



e.g. 第一个性状的协变量效应名称为: `c("Age")`

第二个性状的协变量效应名称为: `NULL` (意味着无协变量)

那么 `covariate_effect_name=list(c("Age"),NULL)`

- **参数10: maternal\_effect\_name**

母性效应名称(一般为母亲名称), `list` 类型.

在多性状模型中, `maternal_effect_name` 为每个性状的母性效应名称向量组成的列表结构, 性状的顺序需与 `target_trait_name` 一一对应.

eg. `target_trait_name=list(c("Trait1"),c("Trait2"))`

`maternal_effect_name=list(c(NULL),c("Dam"))`

- **参数11: random\_regression\_effect\_name**

随机回归效应名称, `list` 类型.

在多性状模型中, `random_regression_effect_name` 为每个性状的随机回归效应名称向量组成的列表结构, 性状的顺序需与 `target_trait_name` 一一对应.

eg. `target_trait_name=list(c("Trait1"),c("Trait2"))`

`random_regression_effect_name=list(c("L1&id", "L1&pe_effect", "L2&id", "L2&pe_effect"),  
c("L1&id", "L1&pe_effect", "L2&id", "L2&pe_effect"))`

在每个列表中, `&` 左边 代表的是多项式系数名称, `&` 右边 代表的是嵌套在多项式里的相应的随机效应名称. 如果用户想将 永久环境效应也嵌套在多项式里, `&` 右边 代表的随机效应名称应设置为 "pe\_effect", 并且需要设置 `included_permanent_effect` 参数为 `TRUE`.

- **参数12: included\_permanent\_effect**

是否包括永久环境效应在分析中, `list` 类型.

在多性状模型中, `included_permanent_effect` 为每个逻辑向量组成的列表结构, 性状的顺序需与 `target_trait_name` 一一对应.

eg. `target_trait_name=list(c("Trait1"),c("Trait2"))`

`included_permanent_effect=list(c(TRUE),c(TRUE))`

- **参数13: include\_social\_effect**

是否包括社会遗传效应在分析中, `list` 类型.

在多性状模型中, `include_social_effect` 为每个逻辑向量组成的列表结构, 性状的顺序需与 `target_trait_name` 一一对应.

eg. `target_trait_name=list(c("Trait1"),c("Trait2"))`

`include_social_effect=list(c(TRUE),c(TRUE))`

- **参数14: group\_effect\_name**

Group效应的名称在社会遗传效应分析中, `character` 类型.

当用户提供的表型数据中不包含最大群体大小相关的列时, 用户需要提供 `group_effect_name` 参数. 当用户提供了 Group效应的名称后, 软件将会自动生成包含 最大群体大小相关的列的表型并进行后续的社会遗传分析.



- **参数15: integer\_group\_names**

最大群体大小相关的整型列的变量名称, `character` 类型。

当用户提供的表型数据中包含最大群体大小相关的列时, 用户需要指定 最大群体大小相关的整型列的变量名称。

- **参数16: real\_group\_names**

最大群体大小相关的实型列的变量名称, `character` 类型。

当用户提供的表型数据中包含最大群体大小相关的列时, 用户需要指定 最大群体大小相关的实型列的变量名称。

- **参数17: analysis\_model**

遗传评估的分析模型, `character` 类型。可选模型包括:

- `"PBLUP_A"`: 系谱-加性效应模型
- `"GBLUP_A"`: 基因组-加性效应模型
- `"GBLUP_AD"`: 基因组-加性-显性效应模型
- `"SSBLUP_A"`: 一步法加性效应模型
- `"User_define"`: 用户自定义模型

- **参数18: dmu\_module**

DMU分析时使用的模块, `character` 类型。可选模块包括:

- `"dmuai"`
- `"dmu4"`
- `"dmu5"`

- **参数19: DMU\_software\_path**

DMU软件在本地的路径, `character` 类型。

- **参数20: relationship\_path**

提供的亲缘关系文件的路径, `character` 类型。

- **参数21: relationship\_name**

提供的亲缘关系文件的名称, `character` 类型。

针对不同的遗传评估分析模型, 我们需要提供不同类型的亲缘关系文件。

针对 `"PBLUP_A"` 模型, 我们需要提供系谱文件, e.g. `relationship_name="pedigree.txt"`;

针对 `"GBLUP_A"` 或 `"GBLUP_AD"` 模型, 我们需要提供3列格式的基因组亲缘关系矩阵的逆矩阵文件, e.g. `relationship_name=c("G_A_inv_matrix.txt", "G_D_inv_matrix.txt")`;

针对 `"SSBLUP_A"` 模型, 我们需要同时提供系谱文件及3列格式的基因组亲缘关系矩阵的文件, e.g. `relationship_name=c("pedigree.txt", "G_A_matrix.txt")`。

- **参数22: output\_result\_path**

DMU运行结果的保存路径, `character` 类型。

- **参数23: output\_ebv\_path**

输出的育种值、残差及校正表型文件的保存路径, `character` 类型。

- **参数24: output\_ebv\_name**

输出的育种值、残差及校正表型文件的名称, `character` 类型。

## 进阶参数

- 参数25: `provided_effect_file_path`

性状效应记录文件的路径，`character` 类型。为了方便用户输入固定效应、随机效应及协变量效应，用户可以在本地直接提供相应的文件，格式如下所示：

V1	V2	V3	V4	V5	V6	V7	V8	V9
Trait1	*	Sex	Herd_Year_Season	*	Id	Litter	*	*
Trait2	*	Sex	*	Id	*	Age	*	

每类效应都用\* 隔开，每一列的间隔均为制表符间隔。每个性状所在的行均有4个，第1-2个\*之间的效应代表的是固定效应，第2-3个\*之间的效应代表的是随机效应，第3-4个\*之间的效应代表的是协变量效应。

- 参数26: `provided_effect_file_name`

性状效应记录文件的名称，`character` 类型。

- 参数27: `provided_DIR_file_path`

用户自己提供的DIR文件的路径，`character` 类型。

- 参数28: `provided_DIR_file_name`

用户自己提供的DIR文件的名称，`character` 类型。

- 参数29: `included_permanent_effect`

是否进行永久环境效应分析，`logical` 类型，默认为FALSE。

- 参数30: `dmu_algorithm_code`

DMU模块内的算法代码，`numeric` 类型。

- 参数31: `provided_prior_file_path`

用户提供的方差组分-PRIOR文件的路径，`character` 类型。

- 参数32: `provided_prior_file_name`

用户提供的方差组分-PRIOR文件的名称，`character` 类型。

- 参数33: `missing_value`

表型数据的缺失值，`numeric` 类型，默认为 -9999。

- 参数34: `iteration_criteria`

方差组分迭代收敛的标准，`numeric` 类型，默认为 1.0e-7。

- 参数35: `genetic_effect_number`

SOL文件中，遗传效应所代表的数字，`numeric` 类型，默认为4。

- 参数36: `residual_cov_trait`

残差协方差约束为0的性状，`list` 类型。e.g. 将Trait1和Trait2的残差协方差约束为0，  
`residual_cov_trait=list(c("Trait1","Trait2"))`

- 参数37: `selected_id`

只输出这部分个体的育种值、残差及校正表型，`character` 类型。

- **参数38: cal\_debv**

是否计算DEBV，`logical` 类型，默认为FALSE。

- **参数39: debv\_pedigree\_path**

计算DEBV时，提供的系谱文件的路径，`character` 类型。

- **参数40: debv\_pedigree\_name**

计算DEBV时，提供的系谱文件的名称，`character` 类型。

## 功能8

### 简述

在前面的章节，我们已经讲述了如何在R中调用DMU软件并完成相应的分析。本章，我们将讲述如何通过BLUP\_ADC中的run\_BLUPF90函数，在R中调用BLUPF90软件并完成数据分析。为了方便用户的使用，我们已经尽可能地将run\_BLUPF90函数中的参数和run\_DMU函数中的参数进行了统一。

**Note:** 为了方便用户使用，blupADC已经封装了BLUPF90中的几个基本模块(renumf90, remlf90, airemlf90, 和 blupf90), 更多的模块可以从BLUPF90官网进行下载([BLUPF90下载网站](#))。

**如果您想将BLUPF90用作商业用途，请务必联系 BLUPF90的作者!!!**

接下来，我们还是用几个简单的例子看看该函数的用法:

### 示例

#### 单性状模型-系谱

```
library(blupADC)
data_path=system.file("extdata", package = "blupADC") # 示例文件的路径

run_BLUPF90(

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2",
"Age"), # 表型数据的列名(ps.表型文件无列名)
  target_trait_name=list(c("Trait1")), #性状名称
  fixed_effect_name=list(c("Sex", "Herd_Year_Season")), #固定效应名称
  random_effect_name=list(c("Id", "Litter")), #随机效应名称
  covariate_effect_name=NULL, #协变量效应名称
  genetic_effect_name="Id", #遗传效应名称
  phe_path=data_path, #表型文件路径
  phe_name="phenotype.txt", #表型文件名
  analysis_model="PBLUP_A", #遗传评估模型
  relationship_path=data_path, #亲缘关系文件路径
  relationship_name="pedigree.txt", #亲缘关系文件名
  output_result_path=getwd() #结果输出路径
)
```

#### 单性状模型-GBLUP

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC")    # 示例文件的路径

run_BLUPF90(

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2",
"Age"), # 表型数据的列名(ps.表型文件无列名)
    target_trait_name=list(c("Trait1")),                #性状名称
    fixed_effect_name=list(c("Sex", "Herd_Year_Season")), #固定效应名称
    random_effect_name=list(c("Id", "Litter")),          #随机效应名称
    covariate_effect_name=NULL,                          #协变量效应名称
    genetic_effect_name="Id",                            #遗传效应名称
    phe_path=data_path,                                  #表型文件路径
    phe_name="phenotype.txt",                            #表型文件名
    analysis_model="GBLUP_A",                            #遗传评估模型
    relationship_path=data_path,                          #亲缘关系文件路径
    relationship_name="blupf90_gnumeric",                 #亲缘关系文件名
    output_result_path=getwd()                          #结果输出路径
)

```

## 单性状模型-Single-step(一步法)

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC")    # 示例文件的路径

run_BLUPF90(

phe_col_names=c("Id", "Mean", "Sex", "Herd_Year_Season", "Litter", "Trait1", "Trait2",
"Age"), # 表型数据的列名(ps.表型文件无列名)
    target_trait_name=list(c("Trait1")),                #性状名称
    fixed_effect_name=list(c("Sex", "Herd_Year_Season")), #固定效应名称
    random_effect_name=list(c("Id", "Litter")),          #随机效应名称
    covariate_effect_name=NULL,                          #协变量效应名称
    genetic_effect_name="Id",                            #遗传效应名称
    phe_path=data_path,                                  #表型文件路径
    phe_name="phenotype.txt",                            #表型文件名
    analysis_model="SSBLUP_A",                          #遗传评估模型
    relationship_path=data_path,                          #亲缘关系文件路径
    relationship_name=c("pedigree.txt", "blupf90_gnumeric"), #亲缘关系文
件名
    output_result_path=getwd()                          #结果输出路径
)

```

同样的，与DMU使用类似，我们仅需改变 `analysis_model` 及 `relationship_name` 这两个参数即可完成系谱、GBLUP及SSBLUP的分析(PS: blupf90\_gnumeric 文件 均可通过 `genotype_data_format_conversion` 函数得到 [了解更多](#))。

## 多性状模型-系谱

上面我们介绍的都是单性状模型(只包括了一个目标性状)。与上一节的介绍的 `run_DMU` 函数类似，如果我们想完成双性状模型的计算，只需要在上面的脚本的基础上稍作修改就能实现目的，具体代码如下：

```

library(blupADC)
data_path=system.file("extdata", package = "blupADC")    # 示例文件的路径

```

```
run_BLUPF90(

phe_col_names=c("Id","Mean","Sex","Herd_Year_Season","Litter","Trait1","Trait2",
"Age"), # 表型数据的列名(ps.表型文件无列名)
      target_trait_name=list(c("Trait1"),c("Trait2")), #性状名称

fixed_effect_name=list(c("Sex","Herd_Year_Season"),c("Herd_Year_Season")), #固定效应名称

random_effect_name=list(c("Id","Litter"),c("Id")), #随机效应名称

covariate_effect_name=list(NULL,"Age"), #协变量效应名称

      genetic_effect_name="Id", #遗传效应名称
      phe_path=data_path, #表型文件路径
      phe_name="phenotype.txt", #表型文件名
      analysis_model="PBLUP_A", #遗传评估模型
      relationship_path=data_path, #亲缘关系文件路径
      relationship_name="pedigree.txt", #亲缘关系文件名
      output_result_path=getwd() #结果输出路径
)
```

## 参数详解

接下来，我们将对 `run_BLUPF90` 中特有的参数进行讲解，剩余的参数大家可移步[DMU软件的交互使用](#)进行查看，相同参数的用法和含义均是一模一样的。

### 🔗 特有参数

- **参数1: blupf90\_algorithm**

BLUPF90进行方差组分估计时选用的算法， `character` 类型。可选算法包括：

- `"AI_REML"`
- `"EM_REML"`
- `"BLUP"`：无需估计方差组分，根据提供的先验直接求解混合线性模型方程组。

默认算法为：`"AI_REML"`

- **参数2: provided\_blupf90\_prior\_file\_path**

用户提供的BLUPF90格式的方差组分-PRIOR文件的路径， `character` 类型。

**Note:**需要注意的一点是，**BLUPF90**格式的PRIOR文件和**DMU**格式的PRIOR文件是不相同的。关于二者的差异，以后有时间会再出一篇文章进行讲解，这里就不再赘述了。

- **参数3: provided\_blupf90\_prior\_file\_name**

用户提供的BLUPF90格式的方差组分-PRIOR文件的名称， `character` 类型。

- **参数4: provided\_blupf90\_prior\_effect\_name**

用户提供的PRIOR文件中，与方差组分对应的各个随机效应名称， `character` 类型。

- **参数5: BLUPf90\_software\_path**

BLUPF90软件在本地的路径， `character` 类型。

