

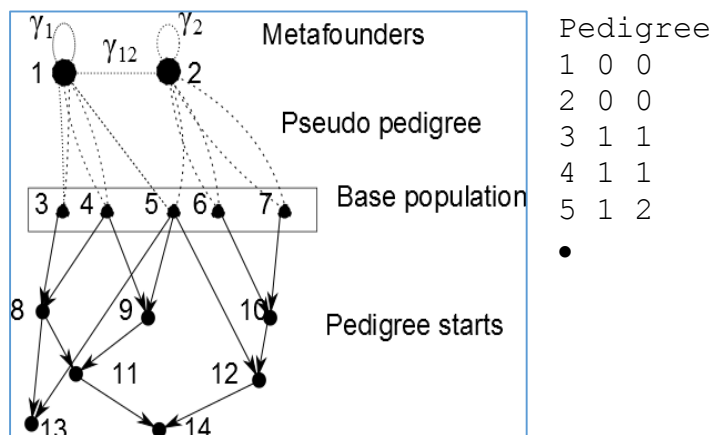
# Use of createHmf to compute H matrices with metafounders

Software /save/alegarra/createHmetafounders/createHmf makes two different tasks.

On input:

a recoded pedigree and genotypes are needed.

- Genotypes must be coded as {0,1,2} and the reference allele must be taken at random so that, in the observed data, the average  $p$  across loci is 0.5.
- Pedigree must be renumbered so that parents precede progeny, and at the beginning there are the metafounders. Example from the Genetics paper:



## A file with gamma coefficients

This file is required unless  $\Gamma$  is estimated from the data. This file contains the gamma coefficients, then the s parameter. For instance:

```
0.70 0.52
0.52 0.61
23876.4
```

## Estimate the parameters $\Gamma$

Three methods are used. In all cases  $\Gamma = 8\Omega = 8Cov(\mathbf{P})$  where  $\Omega = Cov(\mathbf{P})$  is the covariance of allelic frequencies across loci in the base populations.  $\mathbf{P}$  is a matrix with as many rows as markers and as many columns as populations,  $p_{ij}$  is the base allelic frequency at locus  $i$  and population  $j$ .  $\Omega_{ii}$  is the variance of allelic frequencies in base  $i$  and,  $\Omega_{ij}$  is the covariance of allelic frequencies in population  $i$  with allele frequencies in population  $j$ . The virtue of using  $\Gamma$  instead of using  $\mathbf{P}$  explicitly is that whereas  $\mathbf{P}$  is poorly estimated  $\Gamma$  is much better estimated. Methods are:

- Naif (crude estimate of marker allelic frequencies ignoring pedigree structure)
- Generalized least squares (GLS), although this is actually Gengler et al. (2007, 2008) BLUP method for prediction of gene content. This method, extended to several populations and their crosses, is, for one locus:

$$\begin{pmatrix} Q'Q & Q'W \\ W'Q & W'W + A^{-1}\lambda \end{pmatrix} \begin{pmatrix} \mu \\ u \end{pmatrix} = \begin{pmatrix} Q'm \\ W'm \end{pmatrix}$$

- This can be rewritten as:

$$\hat{\mu}_i = (Q'A_{22}^{-1}Q)^{-1}Q'A_{22}^{-1}m_i$$

Where  $\mu_i$  is the mean of the gene content,  $p_i = \mu_i/2$  is a vector containing allelic frequencies across the different populations,  $m_i$  is a vector containing genotypes coded as  $\{-1,0,1\}$  and  $Q$  is a matrix of population fractions for each individual. Use of  $A_{22}^{-1}$  ignores the fact that variances of gene content are heterogeneous across populations, but this is likely to be negligible

(NOTE: it is very easy to extend this to quality control of h2 of gene content!!)

- Maximum likelihood (ML), where  $P$  is supposed to come from a normal distribution,
  - in the case of a single population this is  $p \sim N(0, I\sigma_p^2)$ . ML is done by an EM algorithm.
  - For several populations, this is  $P \sim N(0, I \otimes P_0)$ ,  $P_0 = \begin{pmatrix} \sigma_{\mu^1\mu^1}^2 & \sigma_{\mu^1\mu^2}^2 & \dots \\ \dots & \sigma_{\mu^2\mu^2}^2 & \dots \\ \dots & \dots & \dots \end{pmatrix}$  and  $\hat{\Gamma} = 2\hat{P}_0$ . In this case, the covariance of  $m$  is more complex and involves partial relationship matrices (as in Garcia-Cortes and Toro 2006). Here we take a simpler approach in which we assume that for locus  $i$   $var(m_i) = A_{22}2\bar{p}_i\bar{q}_i$  where  $\bar{p}_i$  is the average allelic frequency across populations.

○

In Garcia-Baccino simulations with a single population, GLS and ML are almost identical and very accurate.

Files with  $\hat{\Gamma}$  are `estimatedGamma_EM` and `estimatedGamma_GLS`.

## Estimate the parameter $s$

This parameter is such that  $G_{0.5} = ZZ'/s$ .

In practice its computation is trivial: is just half the number of markers including in the analysis (regardless of whether they are polymorphic or not; this is unimportant because fixed markers add constants to  $G$ ).

## Construct the matrix $A$ inverse

The program constructs  $A^{(r)-1}$  into file `AiWithMetafounders`.

## Construct the matrix H inverse

The program constructs

$$\mathbf{H}^{(F)-1} = \mathbf{A}^{(F)} + \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}_{0.5}^{-1} - \mathbf{A}_{22}^{(F)-1} \end{pmatrix}$$

in file **HiWithMetafounders** . This matrix is in the format (row, col, val), lower stored, so that it can be read by the blupf90 software. At the beginning the program asks some questions.

If  $\mathbf{F}$  is not estimated, **a file with gamma coefficients** needs to be provided. This file contains the gamma coefficients, then the s parameter. For instance:

```
0.70 0.52
0.52 0.61
23876
```

## Computation of inbreeding

This is achieved by sparse inversion (using FSPAK) of  $\mathbf{H}^{(F)-1}$ , followed by extraction of the diagonal elements. Thus, it may become prohibitive for large matrices. No doubt smarter strategies exist. In file `diagHWithMetafounders` there are 5 columns:

Id, kind (metafounder, genotyped or ungenotyped),  $(H_{ii}^F \text{ rescaled})-1$  , 4th column  $F^F$  rescaled, 5th column  $F$  (normal, gamma=0)'.