

Creation and handling of genomic relationship matrices with preGSf90

Ignacio Aguilar

Instituto Nacional de Investigación Agropecuaria

INIA Las Brujas, Uruguay

iaguilar@inia.org.uy

Genomic Relationship Matrix - G

- $G = ZZ'/k$

- Z = matrix for SNP marker

- Dimension Z= n*p

- n animals,

- p markers

Genotype Codes

0 – Homozygous

1 – Heterozygous

2 – Homozygous

5 – No Call (Missing)

Data file with SNP marker

```

80  21101011002012011011010110111111211111210100
8014 21110101511101120221110111511112101112210100
516  21100101202252021120210121102111202212111101
181  21110111112201120550200020101022212211111100
  
```

HOWTO: Creation of Genomic Matrix

- Read SNP marker information => M
$$\begin{bmatrix} 2 & 1 & 2 & \dots \\ 0 & 1 & 0 & \dots \\ \dots & \dots & \dots & \dots \end{bmatrix}$$
- Get 'means' to center
 - Calculate allele frequency from observed genotypes (p_i)
 - $p_i = \text{sum}(\text{SNPcode}_i) / 2n$
- Matrix for center $W(3,p)$
$$\begin{matrix} 0 \\ 1 \\ 2 \end{matrix} \begin{bmatrix} 0-2p_1 & 0-2p_2 & \dots \\ 1-2p_1 & 1-2p_2 & \dots \\ 2-2p_1 & 2-2p_2 & \dots \end{bmatrix}$$
- Center matrix $Z = W(M)$

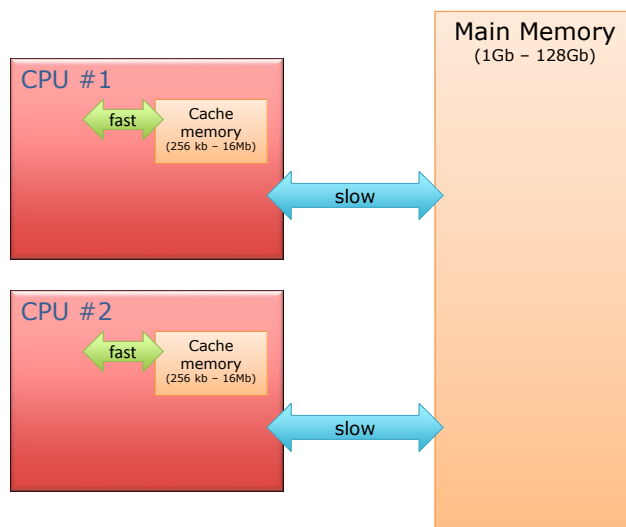
Creation of Genomic

- Issues
 - Large number of genotyped individuals
 - Large number of SNP markers
 - Matrix multiplication $\sim \text{cost } n^2 * p$
- Large amount of data put in (cache) memory for doing 'matmul' for each pair of animals and indirect memory access (center)
 - Memory hierarchy

Matrix multiplication

- Matrix multiplication
 - Several methods
 - Intrinsic matmul (good for small examples !!!)
 - “do-loops”
 - Packages (BLAS, LAPACK)
 - Non-optimized
 - Optimized (ATLAS, MKL, etc.)
 - Several Compilers
 - Perform automatic optimization
 - Vectorize loops
 - Detect permuted loops
 - Can use OpenMP directives for parallelization

Memory Hierarchy



Alternative codes to create G matrix

Original

```

Do i=1,n
  Do j=i,n
    S=0
    Do k=1,p
      S=S+Z(M(i,k),k)
      *Z(M(j,k),k)
    End do
    G(i,j)=S/sqrt(d(i)*d(j))
    G(j,i)=G(i,j)
  End do
End do

```

Optimize Indirect Memory Access -OPTM

```

Do k=1,p
  X(:,k)=Z(M(:,k),k)
End do
Do i=1,n
  Do j=i,n
    S=0
    Do k=1,p
      S=S+X(i,k)
      *X(j,k)
    End do
    G(i,j)=S/sqrt(d(i)*d(j))
    G(j,i)=G(i,j)
  End do
End do

```

Optimize Memory and Loops - OPTML

```

Do k=1,p
  X(:,k)=Z(M(:,k),k)
End do
Do i=1,n
  Do j=1,n
    Do k=1,p
      G(i,j)=G(i,j)
      +X(i,k)*X(j,k)
    End do
  End do
Do i=1,n
  Do j=1,n
    G(i,j)=G(i,j)/sqrt(d(i)*d(j))
  End do
End do

```

Gmatrix.f90 (VanRaden, 2009)

CPU time for alternative codes for G matrix and machines

- Testing
 - 6500 genotyped animals
 - 40k SNPs

		Algorithms		
Processor	Cache	Original	OPTM	OPTML
Xeon 3.5 GHz	6 MB	24 m	26 m	7 m
Opteron 3.0 GHz	1 MB	265 m	59 m	17 m

CPU time (m) with alternative codes and compilers

- Testing
 - 6500 genotyped animals
 - 40k SNPs
 - Opteron 3.02 GHz 1 MB Cache memory

Compiler	Original	OPTM	OPTML
Intel	265	59	17
Absoft	241	60	34
gfortran	213	63	>1day

PreGSf90 program

- From BLUPF90 package
- Uses a genomic module
- Creation and handling of genomic relationship matrices and relationship based on pedigree
- Different methods to optimize calculations using parallel processing

Input files

- Same parameter file as for all BLUPf90 programs
 - But with “OPTION SNP_file xxxx”
 - indicate to run genomic subroutines
- Pedigree file
- Marker information (SNP file)
- Cross Reference file for renumber ID
 - Links genotypes files with codes in pedigree, etc.

SNP map file (optional)

- For some genomic analyses or checks
- Format:
 - snp number
 - Index number of SNP in the sorted map
 - chromosome number
 - position
- First row corresponds to first column SNP in genotype file !!!

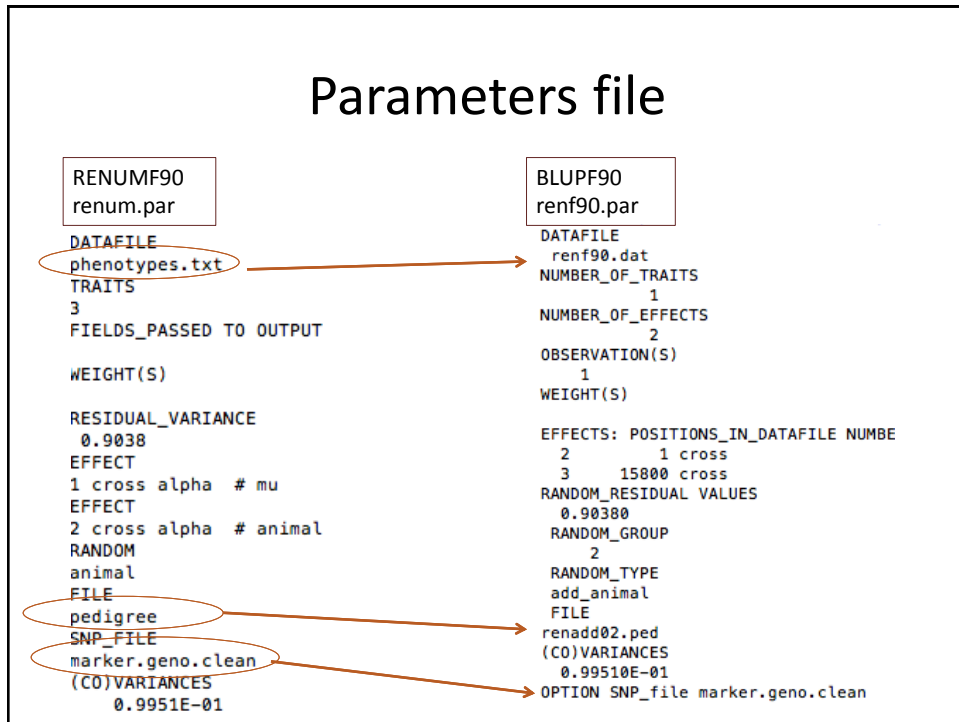
OPTIONS – BLUPF90 parameter file

- PreGSF90
 - controled by adding OPTIONS commands to the parameter file
 - `OPTION SNP_file marker.geno.clean`
 - Read 2 files:
 - `marker.geno.clean`
 - `marker.geno.clean.XrefID`

RENUMF90

- Add keyword to the “animal effect”
`SNP_FILE`
`marker_geno_clean`
- Renumber tool to prepares:
 - data
 - pedigree
 - genotypes
 - parameter files for BLUPF90 programs including PREGSF90
- Check wiki:
- <http://nce.ads.uga.edu/wiki/doku.php>

Parameters file



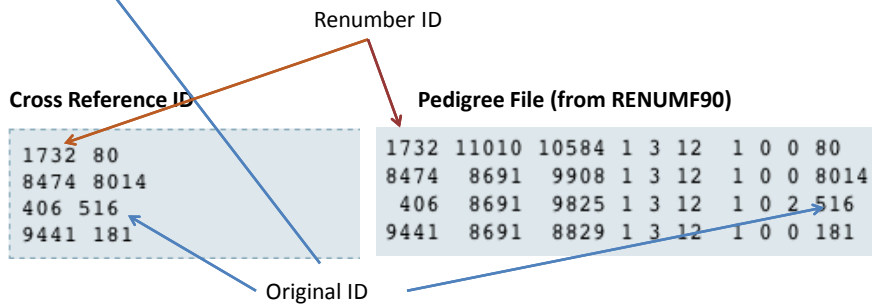
Pedigree file from RENUMF90

- 1 - **animal number**
- 2 - parent 1 number or UPG
- 3 - parent 2 number or UPG
- 4 - 3 minus number of known parents
- 5 - known or estimated year of birth
- 6 - number of known parents;
 if animal is genotyped 10 + number of known parents
- 7 - number of records
- 8 - number of progenies as parent 1
- 9 - number of progenies as parent 2
- **10 - original animal ID**

SNP file & Cross Reference Id

SNP File First col: Identification, could be alphanumeric
 Second col: SNP markers {codes: 0,1,2 and 5 for missing}

```
80  211010110020120110110101101111
8014 211101015111011202211101115111
516  211001012022520211202101211021
181  211101111122011205502000201010
```



Genomic Matrix default options

- $G^* = ZZ'/k$ as in VanRaden, 2008
- With:
 - Z center using allele frequencies estimated from the genotyped individuals
 - $k = 2 \sum (p * (1-p))$
- $G = G*0.95 + A*0.05$ (to invert)
- Tuning of G (see Z. Vitezica work)
 - Adjust G to have mean of diagonals and off-diagonals equal to A

Genomic Matrix Options

- OPTION whichfreq x
 - 0: read from file *freqdata* or other specified
 - 1: 0.5
 - 2: current calculated from genotypes (default)
- OPTION FreqFile *file*
 - Reads allele frequencies from a file
- OPTION maxsnps x
 - Set the maximum length of string for reading marker data from file => BovineHD chip

Options for Blending G and A

- OPTION AlphaBeta alpha beta
 - $G = \alpha * G^r + \beta * A$
- OPTION tunedG
 - 0: no adjustment
 - 1: $\text{mean}(\text{diag}(G))=1, \text{mean}(\text{offdiag}(G))=0$
 - 2: $\text{mean}(\text{diag}(G))=\text{mean}(\text{diag}(A)),$
 $\text{mean}(\text{offdiag}(G))=\text{mean}(\text{offdiag}(A))$ (default)
 - 3: $\text{mean}(G)=\text{mean}(A)$
 - 4: Use Fst adjustment. Powell et al. (2010) & Vitezica et al. (2011)

Creation of 'raw' genomic matrix

- Tricks:
 - Use dummy pedigree
 - 1 0 0
 - 2 0 0
 - ...
 - Change blending parameters
 - OPTION AlphaBeta 0.99 0.01
 - No adjustment for compatibility with A
 - OPTION tunedG 0

$$G = 0.99 * G + 0.01 * I$$

Storing and Reading Matrices

- PreGSF90:
 - Facilitate the implementation of single-step
 - Matrix A is replaced by H with:

$$\mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix}$$
 - Default output is the matrix GimA22i, to be included in application programs (BLUPF90, REMLF90..)
- BUT: intermediate matrices could be stored for examination, use in application programs, etc.

Storing and Reading Matrices

- Matrices that can be stored:
 - A22, inv(A22), G, inv(G), GmA22, inv(GmA22), inv(H)
- All matrices are stored in same format:
 - upper triangle
 - By default in binary format
 - But to store in text (Ascii) format:
 - Use: `OPTION saveAscii`
- Values
 - `i j val`
 - `i & j` refers to the row number in the genotype file !!!!!
 - Renumber ID could be obtained from the XrefID file

Storing and Reading Matrices

To save our 'raw' genomic matrix:

- `OPTION saveG [all]`
 - If the optional `all` is present all intermediate G matrices will be saved!!!
- or it inverse
- `OPTION saveGInverse`
 - Only the final matrix G, after blending, scaling, etc. is inverted !!!
- Look in wiki for keywords for other matrices

Storing with Original IDs

- Some matrices could be stored in text files with the original IDs extracted from *renaddxx.ped* created by the RENUMF90 program (col #10)
- For example:
 - OPTION saveGOrig
 - OPTION saveDiagGOrig
 - OPTION saveHinvOrig
- Values
 - origID_i, origID_j, val

Printout: Same heading as other programs

Options read from parameter file:

```
* SNP file: marker.geno.clean
* SNP Xref file: marker.geno.clean_XrefID
* Matrix in Ascii format(default=binary)
```

```
*-----*
*              Genomic Library: Version 1.110              *
*-----*
* Modified relationship matrix (H) created for effect:  2  *
*-----*
```

```
Read 18600 animals from pedigree file: "renadd02.ped"
Number of Genotyped Animals: 1500
```

All options that were enter in the parameter file should be here !!.
IF not check that keywords are correct (upper and lower case)

Check number of animals and individuals with genotypes

Printout

```

Creating A22
  Extracting subset of: 4634 pedigrees from: 18600 elapsed time:    0.0019
  Calculating A22 Matrix by Colleau ...elapsed time  1.250464

Reading SNP file
  Column position in file for the first marker: 7
  Format to read SNP file: (6x,400000i1)
  Number of SNPs: 3000
  Number of Genotyped animals: 1500
  Reading SNP file elapsed time: .41

Statistics of alleles frequencies in the
  N:          3000
  Mean:       0.500
  Min:        0.101
  Max:        0.898
  Var:        0.016

```

Information from genotype file.
The format is detected from
the first line !!!

So all genotypes should start in
the same column !!!

Number of SNP is also
determined by the first line!!

Looking stored matrices

- Avoid open with text editors, huge files !!!
- For example:
- 1500 genotyped individuals => 1,125,750 rows
- Inspection could be done by Unix commands:
 - head G => first 10 lines
 - tail G => last 10 lines
 - less G => scroll document by line/page
 - wc -l G => count number of lines

good for checks with the number of
genotypes $(n) = (n*(n+1)/2)$

head G

```
1 1 .999382118619
1 2 .355052761478
2 2 1.014521277458
1 3 -.048184197960
2 3 -.057513012886
3 3 .976558921904
1 4 -.101734083083
2 4 -.007644724611
3 4 .196757165096
4 4 1.018165021903
```

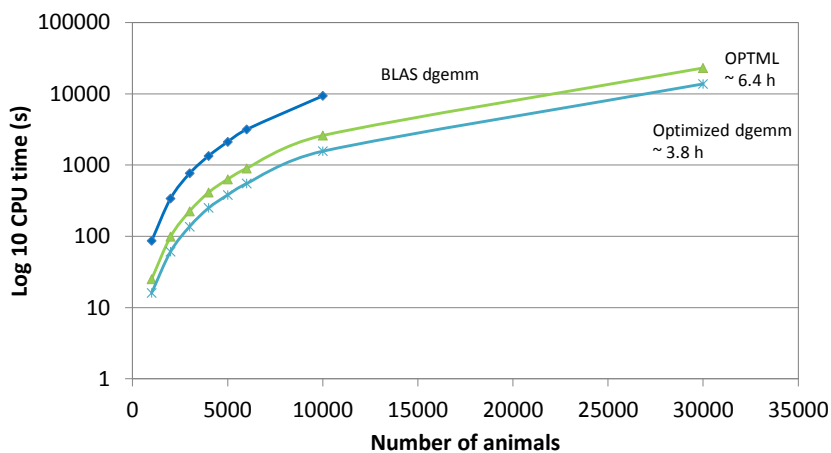
PreGSf90 inside BLUPF90 ??

- Almost all programs from package support creation of genomic relationship matrices, Hinv, etc.
- OPTION SNP_file xxxx
- Why preGSF90 ?
 - Same genomic relationship matrix for several models, traits, etc. Just do it once and store.
 - Uses of optimized subroutines for efficient matrix multiplications, inversion and with support for parallel processing

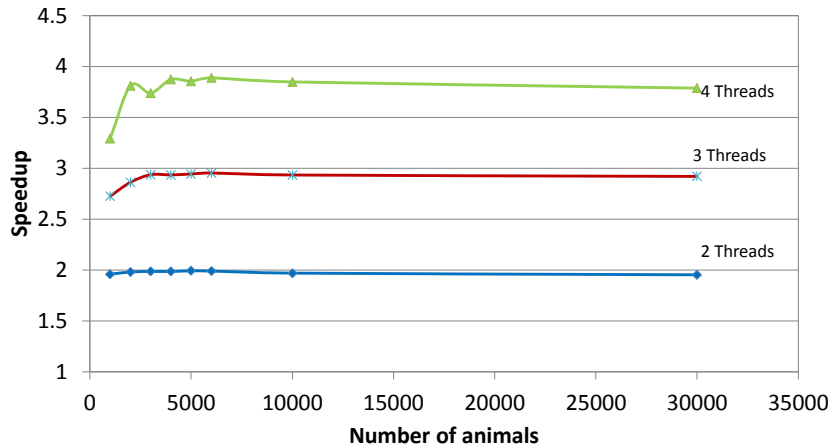
Matrix multiplication subroutines

- Optimized memory and loops (compiler optimization)
 - *dgemm* subroutine from BLAS
 - Optimized *dgemm* (ATLAS or MKL libraries*)
 - Serial
 - Parallel (Automatic use of OpenMP)
- * Intel Fortran Compiler

Matrix multiplication using 40k SNPs



Speedup for matrix multiplications



Speedup = time using one thread/time using n threads

Efficient methods to construct genomic relationship matrices

Elapsed time for different number of individuals
BLADE INIALB 24 cpu

Number of genotypes	Genomic Relationship Matrix	
	Creation	Inversion
10k	0.6 m	0.1 m
30k	5.4 m	3 m
50k	15 m	14 m
70k	30 m	36.4 m
100k	60 m	106 m

Creation a subset of relationship matrix (A₂₂)

- Create a relationship matrix for only genotyped animals (~ thousands)
- Full pedigree (~millions)
- Trace only ancestors of genotyped (reduce but still large number for A matrix)

Relationship Matrix of Genotyped Animals

- Colleau's algorithm to creates A₂₂
- No need to have explicit A matrix
- Method uses "matrix-vector" multiplication with a decomposition of A matrix

$$\mathbf{v} = \mathbf{A}\mathbf{r} = (\mathbf{I} - \mathbf{P})^{-1}\mathbf{D}(\mathbf{I} - \mathbf{P})^{-1}\mathbf{r}$$

Example A times a vector

Pedigree

```
[,1] [,2] [,3]
[1,] 1 0 0
[2,] 2 0 0
[3,] 3 1 2
```

Matrix P

```
[,1] [,2] [,3]
[1,] 0.0 0.0 0.0
[2,] 0.0 0.0 0.0
[3,] 0.5 0.5 0.0
```

Matrix (I-P)⁻¹

```
[,1] [,2] [,3]
[1,] 1.0
[2,] 0.0 1.0
[3,] 0.5 0.5 1.0
```

$$\mathbf{v} = \mathbf{A}\mathbf{r} = (\mathbf{I} - \mathbf{P})^{-1}\mathbf{D}(\mathbf{I} - \mathbf{P})^{-1}\mathbf{r}$$

Matrix (I-P)⁻¹

```
[,1] [,2] [,3]
[1,] 1.0
[2,] 0.0 1.0
[3,] 0.5 0.5 1
```

Matrix D

```
[,1] [,2] [,3]
[1,] 1
[2,] 1
[3,] 0.5
```

Vector q

```
[,1]
[1,] 25
[2,] 35 = [3,]
[3,] 30
```

Matrix (I-P)⁻¹

```
[,1] [,2] [,3]
[1,] 1 0 0.5
[2,] 1 0.5
[3,] 1.0
```

Vector r₂

```
[,1]
[1,] 10
[2,] 20
[3,] 30
```

```
Do i=1,n
  vi = qi*di+(qsi+qdi)/2
End do
```

```
Do i=n,1
  qi = qi+r2i
  qsi = qi+q/2
  qdi = qi+q/2
End do
```

Relationship Matrix of Genotyped Animals

- For each genotyped animal in A₂₂

$$\mathbf{v} = \mathbf{A}\mathbf{r}_2 = (\mathbf{I} - \mathbf{P})^{-1}\mathbf{D}(\mathbf{I} - \mathbf{P})^{-1}\mathbf{r}_2$$



Tabular method vs. Colleau algorithm

- Testing
 - 6,500 genotyped Holsteins
 - 57,000 pedigrees

	Tabular*	Colleau method
CPU Time	311 s	45 s
Memory	12.1GB	322MB

* Gmatrix.f90 (VanRaden, 2009)