

UNIVERSITY OF GEORGIA  
College of Agricultural & Environmental Sciences

# Introduction to genomics

Daniela Lourenco  
UGA USA

Andres Legarra  
INRA France

BLUPF90 TEAM, 02/2023

1

UNIVERSITY OF GEORGIA  
College of Agricultural & Environmental Sciences

## Genomic Information

Initial sequencing and analysis of the human genome

Mutation < 1% < SNP

2

UNIVERSITY OF GEORGIA  
College of Agricultural & Environmental Sciences

## What are SNP used for?

Theor Appl Genet (1993) 67:25-33

Genetic polymorphism in varietal identification and genetic improvement\*

M. Soller<sup>1</sup> and J. S. Beckmann<sup>2</sup>

<sup>1</sup> Department of Genetics, The Hebrew University of Jerusalem, 91904 Jerusalem, Israel  
<sup>2</sup> Institute of Field and Garden Crops, Agricultural Research Organization, The Volcani Center 50250 Bet Dagan, Israel

Received July 14, 1992; Accepted July 3, 1993  
Communicated by A. Robertson

**Summary.** New sources of genetic polymorphisms provide significant additions to the number of useful genetic markers in agricultural plants and animals, and prompt this review of potential applications of polymorphic genetic markers in plant and animal breeding. Two major areas of application can be distinguished. The first is based on the utilization of genetic markers to determine genetic relationships. These applications include varietal identification, protection of breeder's rights, and parentage determination. The second area of application is based on the use of genetic markers to identify and map loci affecting quantitative traits, and to monitor these loci during introgression or selection programs. A variety of breeding applications based on

Use of DNA polymorphisms as genetic markers

- Construct genetic relationships
- Parentage determination
- Identification of QTL

RFLP  
Expensive

3

UNIVERSITY OF GEORGIA  
College of Agricultural & Environmental Sciences

## Excitement about genomics

Copyright © 2001 by the Genetics Society of America

### Prediction of Total Genetic Value Using Genome-Wide Dense Marker Maps

T. H. E. Meuwissen,<sup>\*</sup> B. J. Hayes<sup>1</sup> and M. E. Goddard<sup>1,2</sup>

<sup>\*</sup> Research Institute of Animal Science and Health, 8200 AB Lelystad, The Netherlands, <sup>1</sup>Victorian Institute of Animal Science, Attwood 3049, Victoria, Australia and <sup>2</sup>Institute of Land and Food Resources, University of Melbourne, Parkville 3052, Victoria, Australia

Manuscript received August 17, 2000  
Accepted for publication January 17, 2001

- Genotyping will become cheap
  - Thousands of SNP
- Compute GEBV based on SNP
  - High accuracy
  - Animals with no phenotypes
  - Select the best animals earlier

4

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## Genotyping became cheaper in 2008

- First genomic evaluation for dairy and beef cattle in 2009
  - \$300 in 2009 vs. \$30 in 2022
  - 50,000 SNP

What about statistical methods able to fit genomic information?

5

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## Statistical methods before genomics

- BLUP (Henderson, 1949 - 1976)
  - Best: minimizes MSE
  - Linear: linear function of the data
  - Unbiased:  $E(u) = E(\hat{u})$
  - Prediction: for random effects

That BLUP is a Good Thing: The Estimation of Random Effects  
G. R. Wiggins

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W} \\ \mathbf{W}'\mathbf{X} & \mathbf{W}'\mathbf{W} + \mathbf{A}^{-1} \frac{\sigma_e^2}{\sigma_u^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{W}'\mathbf{y} \end{bmatrix}$$

6

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## Henderson's MME

- Model
 
$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{W}\mathbf{u} + \mathbf{e}$$
- Joint probability of phenotypes and EBV
 
$$p(\mathbf{y}, \mathbf{u}) = p(\mathbf{u}|\mathbf{y}) p(\mathbf{y}) = p(\mathbf{y}|\mathbf{u}) p(\mathbf{u})$$
- Joint probability density function of phenotypes and EBV
 
$$p(\mathbf{y}, \mathbf{u}) = p(\mathbf{y}|\mathbf{u}) p(\mathbf{u}) = \frac{1}{\sqrt{2\pi}|\mathbf{R}|} e^{-\frac{1}{2}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta}-\mathbf{W}\mathbf{u})'\mathbf{R}^{-1}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta}-\mathbf{W}\mathbf{u})} \frac{1}{\sqrt{2\pi}|\mathbf{G}|} e^{-\frac{1}{2}(\mathbf{u}-\mathbf{0})'\mathbf{G}^{-1}(\mathbf{u}-\mathbf{0})}$$

$$\begin{cases} \mathbf{X}'\mathbf{R}^{-1}\mathbf{X}\boldsymbol{\beta} + \mathbf{X}'\mathbf{R}^{-1}\mathbf{W}\mathbf{u} = \mathbf{X}'\mathbf{R}^{-1}\mathbf{y} \\ \mathbf{W}'\mathbf{R}^{-1}\mathbf{X}\boldsymbol{\beta} + (\mathbf{W}'\mathbf{R}^{-1}\mathbf{W} + \mathbf{G}^{-1})\mathbf{u} = \mathbf{W}'\mathbf{R}^{-1}\mathbf{y} \end{cases} \quad \begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W} \\ \mathbf{W}'\mathbf{X} & \mathbf{W}'\mathbf{W} + \mathbf{A}^{-1} \frac{\sigma_e^2}{\sigma_u^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{W}'\mathbf{y} \end{bmatrix}$$

7

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## Henderson's MME for dairy in 1989

- BLUP (Henderson, 1949 - 1976)
- Implementation for dairy in 1989

Journal of Dairy Science  
Volume 72, Supplement 2, June 1989, Pages 54-69

Implementation of an Animal Model for Genetic Evaluation of Dairy Cattle in the United States  
G. R. Wiggins, I. Hovell, L.D. Van Vleet

National genetic improvement programs for dairy cattle in the United States  
G. R. Wiggins  
*J. Anim. Sci.* 1991, 69:3853-3860.

Challenges  
Genetic improvement programs are in a period of rapid change. Advances in computer capability enable adoption of sophisticated computational procedures. Advances in repro-

- 9.5 M animals
- 11 M lactations
- 23.5 M equations to solve
- 7.5 hours

ACKNOWLEDGMENTS  
This research was conducted using the Cornell National Supercomputer Facility, a resource of the

8

**From 1989 to 2009**

- How to add genomic information to the evaluation system in 2009?

**Multistep**

9

**Bayesian Alphabet**

- SNP effect models = outputs SNP effects
- BayesA (Meuwissen et al., 2001)
  - All SNPs have effect on the trait (few with large effect)  $a_i \sim N(\mu, \sigma_{a_i}^2)$
  - Different variances for each SNP
- BayesB (Meuwissen et al., 2001)
  - $p(a_i | \sigma_{a_i}^2, \pi) = \begin{cases} t(0, v, \sigma_{a_i}^2) & \text{with probability } (1 - \pi) \\ 0 & \text{with probability } \pi \end{cases}$
- When  $\pi = 0$ , BayesB becomes BayesA

10

**Bayesian Alphabet**

- BayesC (Habier et al., 2011)
  - $p(a_i | \sigma_a^2) = \begin{cases} N(0, \sigma_a^2) & \text{with probability } (1 - \pi) \\ 0 & \text{with probability } \pi \end{cases}$
- BayesR (Erbe et al., 2012)
  - $p(a_i | \pi, \sigma_a^2) = \pi_1 \times N(0, 0 \times \sigma_a^2) + \pi_2 \times N(0, 10^{-4} \times \sigma_a^2) + \pi_3 \times N(0, 10^{-3} \times \sigma_a^2) + \pi_4 \times N(0, 10^{-2} \times \sigma_a^2)$
- BayesRC (MacLeod et al., 2016)
  - BayesR using biological information to assign SNP to classes
- High computing cost and simple models
- After > 10 years, assumption of normality is good enough!

11

**SNP-BLUP (ridge regression)**

- SNP effect model = outputs SNP effects
- $a \sim N(0, \sigma_a^2)$

$$y = X\beta + Za + e$$

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + I \frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} \hat{\beta} \\ \hat{a} \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix}$$

**GEBV = Z $\hat{a}$**

- All SNP explain the same proportion of variance on the trait

12

**SNP-BLUP (ridge regression)**

- SNP effect model = outputs SNP effects
- All SNP explain the same proportion of variance on the trait

$$\mathbf{GEBV} = \mathbf{Z}\hat{\mathbf{a}}$$

$$\mathbf{u} = \mathbf{Z}\hat{\mathbf{a}}$$

$$\text{Var}(\mathbf{u}) = \mathbf{Z}\mathbf{Z}' \frac{\sigma_a^2}{2 \sum_{i=1}^{\text{SNP}} p_i(1-p_i)}$$

$$\text{Var}(\mathbf{u}) = \frac{\mathbf{Z}\mathbf{Z}'}{2 \sum_{i=1}^{\text{SNP}} p_i(1-p_i)} \sigma_a^2$$

Genomic relationship matrix  
VanRaden (2008)

$$\mathbf{G} = \frac{\mathbf{Z}\mathbf{Z}'}{2 \sum_{i=1}^{\text{SNP}} p_i(1-p_i)}$$

$$\text{Var}(\mathbf{u}) = \mathbf{G}\sigma_a^2 \quad \rightarrow \quad \text{GBLUP assumption!!!}$$

$$\sigma_a^2 = \frac{\sigma_u^2}{2 \sum_{i=1}^{\text{SNP}} p_i(1-p_i)}$$

13

**Understanding SNP variance**

$$\sigma_a^2 = \frac{\sigma_u^2}{2 \sum_{i=1}^{\text{SNP}} p_i(1-p_i)}$$

How do we get the variance of SNP effects,  $\sigma_a^2$  ?

- You can estimate it (Bayes C, REML)
- You can « guess » from the genetic variance  $\sigma_g^2$

SNP 1 contributes  $2p_1q_1a_1^2$  to the genetic variance  
 SNP 2 contributes  $2p_2q_2a_2^2$  to the genetic variance  
 ...

Reversing the expression gives

$$\sigma_a^2 = 2 \sum p_i q_i a_i^2 \approx 2 \left( \sum p_i q_i \right) \times (\overline{a_i^2}) \approx 2 \left( \sum p_i q_i \right) \sigma_a^2$$

$$\sigma_a^2 \approx \frac{\sigma_g^2}{2 \left( \sum p_i q_i \right)}$$

14

**GBLUP: equivalent to SNP-BLUP**

- GBEV-based model = outputs genomic predictions
- $\mathbf{u} \sim N(0, \mathbf{G}\sigma_u^2)$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{W}\mathbf{u} + \mathbf{e}$$

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W} \\ \mathbf{W}'\mathbf{X} & \mathbf{W}'\mathbf{W} + \mathbf{G}^{-1} \frac{\sigma_e^2}{\sigma_u^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{a}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{W}'\mathbf{y} \end{bmatrix}$$

$$\mathbf{G} = \frac{\mathbf{Z}\mathbf{Z}'}{2 \sum p_i(1-p_i)}$$

Bernardo (1994)  
Nejati-Javaremi et al. (1997)

VanRaden (2008)

15

**Genomic relationship matrix**

$$\mathbf{G} = \frac{\mathbf{Z}\mathbf{Z}'}{2 \sum p_i(1-p_i)} = \frac{(\mathbf{M} - 2\mathbf{P})(\mathbf{M} - 2\mathbf{P})'}{2 \sum p_i(1-p_i)}$$

Genotypes {0,1,2}

Shifted to refer to the average of a population with allele frequencies  $p$

Scaled to refer to the genetic variance of a population with allele frequencies  $p$

16

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## What are genomic relationships?

- Relationships were conceived as standardized covariances (Fisher, Wright)
  - $Cov(u_i, u_j) = R_{ij} \sigma_a^2$
  - $R_{ij}$  "some" relationship
  - $\sigma_a^2$  genetic variance
- True relationships: two individuals are genetically identical (for a trait) if they carry the same genotype at the causal QTL or genes
- Genomic relationships: due to shared (Identical By State) alleles at *causal genes*
  - If I share the blood group A with someone, we are like twins!
  - Most of the genes are unknown
  - We use proxies (SNP markers)

17

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## Early use of markers to infer A

- A** = pedigree relationships: due to shared (Identical By Descent) alleles at *causal genes*
- In conservation genetics
- Gather markers, then reconstruct pedigrees, then construct **A**
  - Either estimates of  $A_{xy}$ , or estimates of « the most likely relation » (son-daughter, cousins, whatever)  
Li and Horvitz 1953, Cockerham 1969, Ritland 1996, Caballero & Toro 2002, and many others
- With abundant marker data we can do better than this

18

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## Pedigree vs. Genomic relationships

- Identical By Descent Relationships based on pedigree are average relationships which assume infinite loci
- « Real » IBD relationships are a bit different due to finite genome size (Hill and Weir, 2010)
- Therefore **A** is the expectation of realized or observed relationships
- SNPs more informative than **A**
  - Two full sibs might have a correlation of 0.4 or 0.6
- Many markers needed to better estimate relationships
  - Estimators of IBD

19

UNIVERSITY OF GEORGIA  
College of Agriculture & Environmental Sciences

## Pedigree vs. Genomic relationships

Adapted from Lourenco et al. (2015)

20

UNIVERSITY OF GEORGIA  
College of Agriculture and Environmental Sciences

## Genomic relationships

$$\mathbf{G} = \frac{\mathbf{ZZ}'}{2 \sum p_i(1-p_i)} = \frac{(\mathbf{M} - 2\mathbf{P})(\mathbf{M} - 2\mathbf{P})'}{2 \sum p_i(1-p_i)}$$

Genotypes {0,1,2}

Shifted to refer to the average of a population with allele frequencies  $p$

Scaled to refer to the genetic variance of a population with allele frequencies  $p$

If base allelic frequencies are used,  $\mathbf{G}$  is an unbiased and efficient estimator of IBD realized relationships

21

UNIVERSITY OF GEORGIA  
College of Agriculture and Environmental Sciences

## Some “interesting” properties of $\mathbf{G}$

- If  $p$  are computed from the data  
This implies that  $E(\text{Breeding Values})=0$
- Positive and negative inbreeding  
Some individuals are more heterozygous than the average of the population (OK, no biological problem)
- Positive and negative genomic relationships  
Individuals  $i$  and  $j$  are more distinct than an average pair of individuals in the data  
Fixing negative estimates of relationships to 0 is a wrong praxis

22

UNIVERSITY OF GEORGIA  
College of Agriculture and Environmental Sciences

## Some “interesting” properties of $\mathbf{G}$

- VanRaden (2008)
  - $\mathbf{G}$  can be singular if few SNP or identical genotypes (twins)
  - $\mathbf{G}$  must be singular if number of individuals > number of SNP
- Strandén and Christensen (2011)
  - $\mathbf{G}$  is singular if  $p$ 's are averages across the sample

$$\mathbf{G} = 0.95 \frac{\mathbf{ZZ}'}{2 \sum p_i(1-p_i)} + 0.05\mathbf{I} \quad \text{OR} \quad \mathbf{G} = 0.95 \frac{\mathbf{ZZ}'}{2 \sum p_i(1-p_i)} + 0.05\mathbf{A} \quad \rightarrow \quad \mathbf{G} = \alpha \mathbf{G}_0 + \beta \mathbf{A}$$

- Blending  $\approx$  Adding a residual polygenic effect

23

UNIVERSITY OF GEORGIA  
College of Agriculture and Environmental Sciences

## Some “interesting” properties of $\mathbf{G}$

- For all matrices of the kind  $\mathbf{G} = \frac{\mathbf{ZZ}'}{2 \sum p_i(1-p_i)} = \frac{(\mathbf{M} - 2\mathbf{P})(\mathbf{M} - 2\mathbf{P})'}{2 \sum p_i(1-p_i)}$ 
  - We don't need to put the same  $p$ 's in the upper and in the lower part
- Changing allele frequencies in  $\mathbf{P}$  shifts EBV's by a constant
  - Irrelevant if there is an overall mean or fixed effect in the model (Strandén and Christensen, 2011)
- Changing allele frequencies in  $\frac{1}{2 \sum p_i q_i}$  “scales”

24

**Not all individuals are genotyped**

- Genomic evaluation would be simpler if all individuals were genotyped
- What to do when there are genotyped and non-genotyped individuals?
  - SNPs are capturing relationships
  - Pedigrees give information about relationships
  - Genomic and pedigree relationships can be combined in a single matrix!

Non-genotyped

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$$

Genotyped

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{G} \end{bmatrix}$$

~~$$\mathbf{H} = \mathbf{A} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G} - \mathbf{A}_{22} \end{bmatrix}$$~~

Misztal et al., 2009 25

25

**Not all animals are genotyped**

- Genomic info can be extended to non-genotyped animals
  - joint distribution of EBV for non-genotyped ( $u_1$ ) and genotyped ( $u_2$ )

$$p(u_1, u_2) = p(u_2)p(u_1|u_2)$$

Legarra et al., 2009

$$\mathbf{H} = \begin{pmatrix} \text{var}(u_1) & \text{cov}(u_1, u_2) \\ \text{cov}(u_2, u_1) & \text{var}(u_2) \end{pmatrix} = \begin{pmatrix} \mathbf{A}_{11} + \mathbf{A}_{12}\mathbf{A}_{22}^{-1}(\mathbf{G} - \mathbf{A}_{22})\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G} \\ \mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{G} \end{pmatrix}$$

Error in the prediction

Variance of prediction of genotypes for non-genotyped animals

Prediction generates a covariance

$$\mathbf{H} = \begin{bmatrix} \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} + \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G} \\ \mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{G} \end{bmatrix}$$

Relationships from genotypes

26

26

**Understanding H**

- It is a projection of **G** matrix on the rest of individuals "so that" **G** matrix makes sense
  - e.g. parents of two animals related in **G** should be related in **A**
- It is a Bayesian update of the pedigree matrix based on new information from genotypes
- Typically
  - **A** in the millions
  - **G** and **A<sub>22</sub>** in the thousands
  - Leads to a very efficient method of genomic evaluation:
    - **Single Step GBLUP**

27

27

**Some properties of H**

- Always semi-positive definite
  - eigenvalues are always positive or zero
- Positive definite & invertible if **G** is invertible
- In practice, if **G** is too different from **A<sub>22</sub>** (wrong pedigree or genotyping), this gives lots of numerical problems
- If no one is genotyped, Single-step is BLUP
- If everyone is genotyped, Single-step is GBLUP

28

28

UNIVERSITY OF GEORGIA  
College of Agriculture and Forestry  
Department of Statistics

## Realized relationship matrix (H)

Animal	Sire	Dam
1	0	0
2	0	0
3	1	2
4	1	2

Pedigree Relationship Matrix (A)

$$\begin{bmatrix} 1.0 & 0.0 & 0.5 & 0.5 \\ . & 1.0 & 0.5 & 0.5 \\ . & . & 1.0 & 0.5 \\ . & . & . & 1.0 \end{bmatrix}$$

Genomic Relationship Matrix (G) for animals 3 and 4

$$\begin{bmatrix} 1.0 & 0.52 \\ . & 1.0 \end{bmatrix}$$

Realized Relationship Matrix (H)

$$\begin{bmatrix} 1.004 & 0.0 & 0.507 & 0.507 \\ . & 1.004 & 0.507 & 0.507 \\ . & . & 1.0 & 0.52 \\ . & . & . & 1.0 \end{bmatrix}$$

29

UNIVERSITY OF GEORGIA  
College of Agriculture and Forestry  
Department of Statistics

## Single-step Genomic BLUP (ssGBLUP)

- Because not all animals are genotyped
  - 5% to 10% in large populations

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \mathbf{H}^{-1} \frac{\sigma_e^2}{\sigma_u^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{Z}'\mathbf{y} \end{bmatrix}$$

$$\mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix}$$

Aguilar et al., 2010  
Christensen and Lund, 2010

30

29

30

UNIVERSITY OF GEORGIA  
College of Agriculture and Forestry  
Department of Statistics

## Combining two sources of relationships

$$\mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} \mathbf{A}_{12}\mathbf{A}_{22}^{-1}(\mathbf{G} - \mathbf{A}_{22})\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{A}_{12}\mathbf{A}_{22}^{-1}(\mathbf{G} - \mathbf{A}_{22}) \\ (\mathbf{G} - \mathbf{A}_{22})\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & (\mathbf{G} - \mathbf{A}_{22}) \end{bmatrix}$$

- A**
  - Contains expected relationships
  - Is limited by the pedigree depth and completeness
  - Depends on accuracy of recording pedigrees
- G**
  - Contains number of alleles shared between animals weighted by heterozygosity
  - No limitations regarding to the number of past generations
  - Depends on allele frequency and quality of genomic data

31

31

UNIVERSITY OF GEORGIA  
College of Agriculture and Forestry  
Department of Statistics

## Combining two sources of relationships

$$\mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix}$$

Computed using Henderson-Quaas' algorithm with inbreeding

Computed using VanRaden's formula, which considers inbreeding

Computed using Colleau's algorithm, which considers inbreeding

- Tuning**
  - Base of **G** is *genotyped* animals
  - Base of **A** is *founders of the pedigree*
  - For SSGBLUP, Vitezica et al. 2011 modeled a mean in genotyped animals:

$$p(\mathbf{u}_2) = N(\mathbf{1}\boldsymbol{\mu}, \mathbf{G})$$

Integrate  $\boldsymbol{\mu} : \mathbf{G}^* = a + b\mathbf{G}$

$\boldsymbol{\mu} = (\text{Pedigree base}) - (\text{Genomic base})$

Tries to put G and A on the same scale

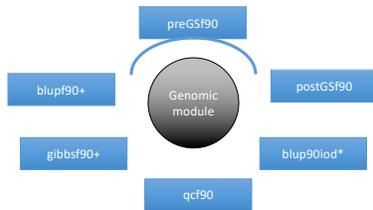
32

32



## preGSf90

- Interface program to the genomic module to process the genomic information in the BLUPF90 family of programs



37

## preGSf90

- Performs Quality Control of SNP information
- Creates the genomic relationship matrix (**G**)
  - and relationships based on pedigree (**A<sub>22</sub>**)
  - Inverse of relationship matrices



38

## preGSf90

- Same parameter file as for all BLUPF90 programs
- Needs an extra OPTION in renf90.par
  - OPTION SNP\_file *marker.geno*
- Reads 2 extra files (besides data and pedigree):
  - *marker.geno*
  - *marker.geno\_XrefID* (created by renumf90)

*\_XrefID* has 2 columns: Renumbered ID Original ID

39

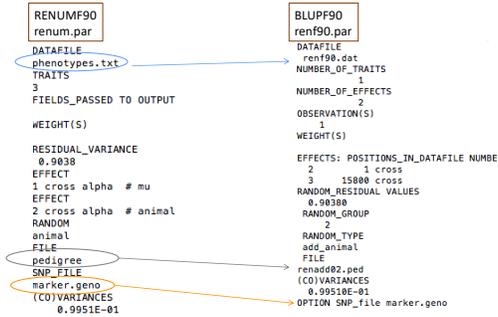
## Run renumf90 before preGSf90

- Use renumf90 for renumbering data and creating XrefID and files

```
EFFECT
1 cross alpha
RANDOM
animal
FILE
ped3.txt
FILE_POS
1 2 3 0 0
SNP_FILE
marker.geno
PED_DEPTH
0
(CO) VARIANCES
0.30
```

40

## Parameter files



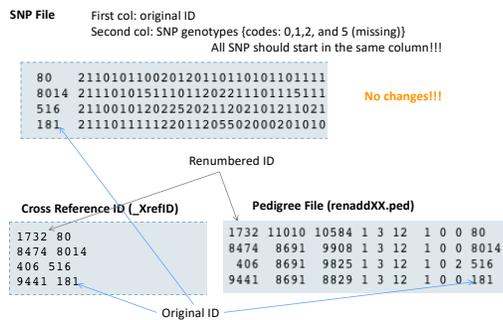
41

## New pedigree file from RENUMF90

- 1 - renumbered animal ID
- 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

42

## SNP file, XrefID, and ped from renumf90



43

## preGSf90

- Same parameter file as for all BLUPF90 programs
- Needs an extra OPTION in renf90.par
  - OPTION SNP\_file marker.geno
- Reads 2 extra files (besides data and pedigree):
  - marker.geno
  - marker.geno\_XrefID (created by renumf90)

\_XrefID has 2 columns: Renumbered ID Original ID

44

## SNP map file – new default

- `OPTION chrinfo <file>`
- `OPTION map_file <file>`
  - For GWAS and QC
- Format:
  - A header must be provided
    - Names for SNP, chromosome, and physical position are mandatory
  - SNPID for SNP
  - CHR for chromosome
  - POS for position

```
SNPID CHR POS REFID MARK
21420 14 7921199 AAG-BPGL-BAC-1003 2
22050 14 32183746 AAG-BPGL-BAC-10045 3
21371 14 4131029 AAG-BPGL-BAC-10046 4
21075 14 34424444 AAG-BPGL-BAC-10067 8
21950 14 31287768 AAG-BPGL-BAC-10068 9
21096 10 18881218 AAG-BPGL-BAC-10082 10
21550 10 20680250 AAG-BPGL-BAC-10086 11
21646 10 21228390 AAG-BPGL-BAC-10075 12
21612 10 24527357 AAG-BPGL-BAC-10086 13
24705 10 79223200 AAG-BPGL-BAC-10098 14
24712 10 79220023 AAG-BPGL-BAC-11000 15
24724 10 804810073 AAG-BPGL-BAC-11009 16
24741 10 80783719 AAG-BPGL-BAC-11007 17
24827 10 86518019 AAG-BPGL-BAC-11020 18
25065 11 21274134 AAG-BPGL-BAC-11039 21
```

45

## Saving 'clean' files

- SNP excluded from QC are set to missing (i.e., Code=5)
  - 5 is replaced by 0 in calculations
- `OPTION saveCleanSNPs`
- Save clean genotype data without excluded SNP and individuals
  - For example, for a SNP\_file named *marker.geno*
  - Clean files will be:
    - *marker.geno\_clean*
    - *marker.geno\_clean\_XrefID*
  - Removed SNP/animals will be output in files:
    - *marker.geno\_SNPs\_removed*
    - *marker.geno\_Animals\_removed*

46

## Only QC in preGSf90

- Quality control
- Genomic relationship matrices and inverses
  - Inverse is costly
- How to do only QC avoiding the inverses:
  - `OPTION SNP_file marker.geno`
  - `OPTION saveCleanSNPs`
  - `OPTION createGInverse 0`
  - `OPTION createA22Inverse 0`
  - `OPTION createGimA22i 0`

47

## No QC in the application programs

- **ONLY use:**
  - If QC was performed in a previous run
  - and "clean" genotype file is used
- `OPTION SNP_file marker.geno_clean`
- `OPTION no_quality_control`

48

## Use in application programs

- Use `renumf90` for renumbering and creation of XrefID and files  
`SNP_FILE`  
`marker.geno`

```
EFFECT
1 cross alpha
RANDOM
animal
FILE
ped1.txt
FILE_POS
1 2 3 0 0
SNP_FILE
marker.geno
PED_DEPTH
0
(CO)VARIANCES
0.30
```

- Run `preGSf90` with quality control, saving clean files
- Run further programs with clean files as needed
  - `blupf90+`, `gibbs2f90+`, ...

49

## PreGSf90 wiki

50

## preGSf90

- Performs Quality Control of SNP information
- Creates the genomic relationship matrix (**G**)
  - and relationships based on pedigree (**A<sub>22</sub>**)
  - Inverse of relationship matrices



51

## BLUP-based models

$$\begin{bmatrix} X'X & X'W \\ W'X & W'W + A^{-1} \frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} \beta \\ \bar{a} \end{bmatrix} = \begin{bmatrix} X'y \\ W'y \end{bmatrix} \quad \text{BLUP} \quad \text{Henderson, 1963}$$

$$\begin{bmatrix} X'X & X'W \\ W'X & W'W + G^{-1} \frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} \beta \\ \bar{a} \end{bmatrix} = \begin{bmatrix} X'y \\ W'y \end{bmatrix} \quad \text{GBLUP} \quad \text{Nejati-Javaremi et al., 1997; Fernando, 1998; VanRaden, 2008}$$

$$\begin{bmatrix} X'X & X'W \\ W'X & W'W + H^{-1} \frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} \beta \\ \bar{a} \end{bmatrix} = \begin{bmatrix} X'y \\ W'y \end{bmatrix} \quad \text{ssGBLUP} \quad \text{Miztal et al. (2009); Legarra et al. (2009); Aguilar et al. (2010); Christensen & Lund (2010)}$$

$$H^{-1} = \begin{bmatrix} A_{11}^{-1} & A_{12}^{-1} \\ A_{21}^{-1} & A_{22}^{-1} \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & G^{-1} - A_{22}^{-1} \end{bmatrix} \quad H^{-1} = A^{-1} + \begin{bmatrix} 0 & 0 \\ 0 & G^{-1} - A_{22}^{-1} \end{bmatrix}$$

52

## PreGSf90

- Created to construct the matrices using in ssGBLUP

$$\mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{bmatrix}$$

$$\mathbf{G} \quad \mathbf{G}^{-1}$$

$$\mathbf{A}_{22} \quad \mathbf{A}_{22}^{-1}$$

$$\mathbf{G}^{-1} - \mathbf{A}_{22}^{-1}$$

53

## Genomic Relationship Matrix - G

$$\mathbf{G} = \frac{\mathbf{Z}\mathbf{Z}'}{2 \sum p_i(1-p_i)} \quad (\text{VanRaden, 2008})$$

- $\mathbf{Z}$  = matrix for SNP marker
- Dimension of  $\mathbf{Z} = n \times i$
- $n$  animals
- $i$  markers

**Genotype Codes**  
 0 – Homozygous  
 1 – Heterozygous  
 2 – Homozygous  
 5 – No Call (Missing)

SNP file

```

80  21101011002012011011010110111111211111210100
8014 2111010151110112022110111511112101112210100
516  2110010120225202112021012110211120221211101
181  2111011111220112055020002010102221221111100
    
```

54

## PreGSf90

- Efficient methods
  - create the genomic relationship matrix and the relationship matrix based on pedigree
  - Invert the relationship matrices
- Computes statistics for the matrices
  - Means, Var, Min, Max
  - Correlations between diagonals
  - Correlations for off-diagonals
  - Correlations for the full matrices
  - Regression coefficients

55

## Genomic Matrix default options

$$\mathbf{G}_0 = \frac{\mathbf{Z}\mathbf{Z}'}{2 \sum p_i(1-p_i)} \quad (\text{VanRaden, 2008})$$

- With:
  - $\mathbf{Z}$  centered using current allele frequencies
  - Current genotyped animals

56

## Genomic Matrix Options

- **OPTION whichfreq  $x$** 
  - 0: read from file *freqdata* or other specified name (needs OPTION FreqFile)
  - 1: 0.5
  - 2: current calculated from genotypes (default)
- **OPTION FreqFile *file***
  - Reads allele frequencies from a file

57

## Genomic Matrix default options

- **Blending** - to avoid singularity problems  

$$\mathbf{G} = 0.95 * \mathbf{G}_0 + 0.05 * \mathbf{A}_{22}$$
  - OPTION AlphaBeta 0.95 0.05 #(default)
  - Beta may vary from 0.2 to 0.01

58

## Genomic Matrix default options

- **Tuning**
  - Adjust  $\mathbf{G}$  to have mean of diagonals and off-diagonals equal to  $\mathbf{A}_{22}$
  - OPTION tunedG 2 #(default) Chen et al. (2011)
    - Base of GBLUP is *genotyped* animals
    - Base of pedigree is *founders of the pedigree*
    - For SSGBLUP modelled as a mean for genotyped animals
      - $p(\mathbf{u}_2) = N(\mathbf{1}\mu, \mathbf{G})$
      - Integrate  $\mu$ :  $\mathbf{G}^* = 11'\lambda + (1 - \lambda/2)\mathbf{G}$
      - $\mu = (\text{Genomic base}) - (\text{Pedigree base})$
      - Vitezica et al. 2011

59

## Options for matching $\mathbf{G}$ to $\mathbf{A}_{22}$

- **OPTION tunedG  $x$** 
  - 0: no adjustment
  - 1: mean(diag(G))=1, mean(offdiag(G))=0
  - 2: mean(diag(G))=mean(diag(A<sub>22</sub>)), mean(offdiag(G))=mean(offdiag(A<sub>22</sub>)) (default)
  - 3: mean(G)=mean(A<sub>22</sub>)
  - 4: Use Fst adjustment. Powell et al. (2010) & Vitezica et al. (2011)

$$\lambda = \frac{1}{n^2} \left( \sum_i \sum_j A_{22ij} - \sum_i \sum_j G_{ij} \right) \quad \mathbf{G}^* = 11'\lambda + (1 - \lambda/2)\mathbf{G}$$

60

## Storing and Reading Matrices

- preGSf90 saves  $\mathbf{G}^{-1} - \mathbf{A}_{22}^{-1}$  by default (file: GimA22i)

To save 'raw' genomic matrix:

- OPTION saveG [all]
  - If the optional *all* is present all intermediate  $\mathbf{G}$  matrices will be saved!!!

To save  $\mathbf{G}^{-1}$

- OPTION saveGInverse
  - Only the final  $\mathbf{G}$ , after blending, scaling, etc. is inverted !!!

To save  $\mathbf{A}_{22}$  and inverse

- OPTION saveA22 and OPTION saveA22Inverse

61

## Storing and Reading Matrices

- OPTION saveG [all] , OPTION saveGInverse, ...
  - Saves in binary format
  - "Dumped" format to save space and time
  - To save as row, column, value:
    - OPTION no\_full\_binary
    - Still binary, but can be easily read and converted to text

62

## 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

63

## Genomic Matrix - Population structure

```
OPTION plotpca
```

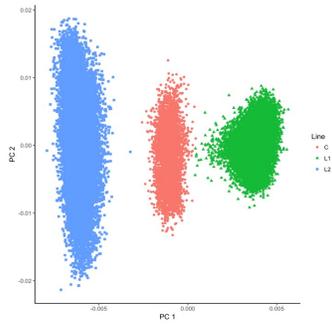
Plot first two principal components to look for stratification in the population.

```
OPTION extra_info_pca file col
```

Reads from *file* the column *col* to plot with different colors for different classes.

64

## Genomic Matrix - Population structure



65

## Tricks to setup **G** for GBLUP #1

- Tricks are needed because preGSf90 is set up for ssGBLUP

1) Use a dummy pedigree

```
1 0 0
2 0 0
...
```

2) Use PED\_DEPTH 1 in renumf90

3) Change blending parameters

- OPTION AlphaBeta 1.00 0.00 →  $G = 1.00 * G + 0.00 * I$

- OPTION AlphaBeta 0.95 0.05 →  $G = 0.95 * G + 0.05 * I$

4) No adjustment for compatibility with  $A_{22}$

- OPTION tunedG 0

66

## Tricks to setup **G** for GBLUP #2

1) In renum.par, remove any information about the pedigree. Example:

```
FILE
pedigree.txt
FILE_POS
1 2 3 0 0
PED_DEPTH
3
```

3) Change blending parameters

- OPTION AlphaBeta 1.00 0.00 →  $G = 1.00 * G + 0.00 * I$

- OPTION AlphaBeta 0.95 0.05 →  $G = 0.95 * G + 0.05 * I$

4) No adjustment for compatibility with  $A_{22}$

- OPTION tunedG 0

67

## PreGSf90 inside BLUPF90 ??

- Almost all programs from BLUPF90 support creating genomic relationship matrices

- OPTION SNP\_file xxxx

- Why preGSf90 ?

- Same genomic relationship matrix for several models, traits, etc.

- Just do it once and store GimA22i or Gi and A22i separate

68

## Use in application programs

- Use renumf90 for renumbering and creation of XrefID and files  
SNP\_FILE  
marker.geno
- Run preGSf90 with quality control, saving clean files
- Option 1:  
run blupf90 with clean files
- Option 2:  
run preGSf90 with clean files (program saves **GimA22i**)  
run blupf90 with option to read **GimA22i** from the file

69

## Reading external matrices

- BLUPF90 programs accept external matrices created outside
- [http://nce.ads.uga.edu/wiki/doku.php?id=user\\_defined\\_files\\_for\\_covariances\\_of\\_random\\_effects](http://nce.ads.uga.edu/wiki/doku.php?id=user_defined_files_for_covariances_of_random_effects)
- File should be row, column, value in plain text format (lower OR upper triangular)

```
renf90.par
RANDOM_GROUP
# genomic
2
RANDOM_TYPE
user_file
FILE
# matrix file
GI
```

Valid format	Non-valid format
<pre>1 1 1 1 2 0.5 2 2 1</pre>	<pre>1 1 1 1 2 0.5 2 1 0.5 2 2 1</pre>

- user\_file: if providing the inverse of the covariance structure
- user\_file\_inv: if the program has to invert the covariance structure

70