

Genomic models and BLUP

Ole F. Christensen

Center for Quantitative Genetics and Genomics, Aarhus University

BovReg course - September 2023

Genetic markers

- ▶ Fragments of DNA associated with certain locations in genome.
- ▶ Must be polymorphic
- ▶ Many types: Microsatellites, insertions/deletions, **Single Nucleotide Polymorphisms (SNP)**, etc

SNP data

- ▶ SNP's have two alleles (**diallelic**)

	SNP_1	SNP_2	SNP_2	
ID_1	AB	BB	BB	...
ID_2	AA	AB	BB	...
ID_3	AA	AA	AB	...
ID_4	BB	BB	AA	...
...	

- ▶ Chose a reference allele
- ▶ **Coding**: 0,1,2 copies of the reference allele

SNP data

100011122002001211101111211101111001121100020122002220111
1202101200211122110021112001111001011011010220011002201101
1200201101020222121122102010011100011220221222112021120120
2010020220200002110001120201122111211102201111000021220200
0221012020002211220111012100111211102112110020102100022000
2201000201100002202211022112101121110122220012112122200200
020020202012221100222222002212111121002111120011011101120
0202220001112011010211121211102022100211201211001111102111
2110211122000101101110202200221110102011121111011202102102
1211011022122001211011211012022011002220021002110001110021
1021101110002220020221212110002220102002222121221121112002
0110202001222222112212021211210110012110110200220002001002
0001111011001211021212111201010121202210101011111021102112
2111112121121011012001111021111011111220121012121101022
202021211222120222002121210121210201100111222121101

SNP-BLUP

- ▶ **Model:** $\mathbf{y} = \mathbf{Xb} + \sum_i \mathbf{Z}_i g_i + \mathbf{e}$
- ▶ g_j 's are SNP effects (genetic)
- ▶ Centering of marker codes: $\mathbf{Z}_j = \mathbf{M}_j - 2p_j\mathbf{1}$; relate SNP effects to individuals
- ▶ Compact notation: $\mathbf{y} = \mathbf{Xb} + \mathbf{Zg} + \mathbf{e}$
- ▶ **SNP-BLUP:** assumes SNP effects are normally distributed and independent

Mixed model equations for SNP-BLUP

$$\begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{Z} \\ \mathbf{Z}^T \mathbf{X} & \mathbf{Z}^T \mathbf{Z} + \mathbf{I}\alpha \end{bmatrix} \begin{bmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{g}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{Z}^T \mathbf{y} \end{bmatrix},$$

where $\alpha = \sigma_e^2 / \sigma_g^2$.

- ▶ MME provide BLUP SNP effects $\hat{\mathbf{g}}$
- ▶ GEBV on individual i is $\hat{a} = \sum_j \mathbf{Z}_{\cdot,j} \hat{g}_j$

Equivalent model - GBLUP

- ▶ Genetic effects $\mathbf{a} = \mathbf{Zg}$
- ▶ Genetic variance-covariance:

$$\text{Var}(\mathbf{a}) = \mathbf{Z}\text{Var}(\mathbf{g})\mathbf{Z}^T = \sigma_g^2 \mathbf{Z}\mathbf{Z}^T$$

- ▶ Model

$$\mathbf{y} = \mathbf{Xb} + \mathbf{Wa} + \mathbf{e}$$

where

$$\text{Var}(\mathbf{a}) = \sigma_a^2 \mathbf{G},$$

$$\mathbf{G} = \mathbf{Z}\mathbf{Z}^T / \left(\sum_j 2p_j(1 - p_j) \right)$$

- ▶ Easy to compute
- ▶ Genetic variance $\sigma_a^2 = \sigma_g^2 \sum_j 2p_j(1 - p_j)$

GBLUP versus Animal model (pedigree **A**)

- ▶ GBLUP is an animal model
- ▶ Genomic versus Pedigree relationships
- ▶ Realised (marker genotypes) versus expected (pedigree) relationships
- ▶ More accurate estimate of proportion of chromosome segments shared between individuals
- ▶ GBLUP is an improved version of pedigree BLUP

Mixed model equations for GBLUP

$$\begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{W} \\ \mathbf{W}^T \mathbf{X} & \mathbf{W}^T \mathbf{W} + \mathbf{G}^{-1} \alpha \end{bmatrix} \begin{bmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{a}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{W}^T \mathbf{y} \end{bmatrix},$$

where $\alpha = \sigma_e^2 / \sigma_a^2$.

- ▶ MME provide GEBVs on individuals $\hat{\mathbf{a}}$

Polygenic effects

- ▶ If Markers are not capturing all the genetic effects, then add a residual polygenic effect

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{W}\mathbf{u} + \mathbf{W}\mathbf{a} + \mathbf{e}$$

- ▶ \mathbf{u} : residual polygenic effect, $\text{Var}(\mathbf{u}) = \sigma_u^2 \mathbf{A}$
- ▶ Mixed model equations

$$\begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{W} & \mathbf{X}^T \mathbf{W} \\ \mathbf{W}^T \mathbf{X} & \mathbf{W}^T \mathbf{W} + \mathbf{A}^{-1} \alpha_1 & \mathbf{W}^T \mathbf{W} \\ \mathbf{W}^T \mathbf{X} & \mathbf{W}^T \mathbf{W} & \mathbf{W}^T \mathbf{W} + \mathbf{I} \alpha_2 \end{bmatrix} \begin{bmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{u}} \\ \hat{\mathbf{a}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{W}^T \mathbf{y} \\ \mathbf{W}^T \mathbf{y} \end{bmatrix},$$

where $\alpha_1 = \sigma_e^2 / \sigma_u^2$ and $\alpha_2 = \sigma_e^2 / (\sigma_a^2 - \sigma_u^2)$.

- ▶ Genetic variance is $\sigma_a^2 = \sigma_u^2 + (\sigma_a^2 - \sigma_u^2)$.

ssGBLUP

- ▶ A general approach to handle the situation that some animals are genotyped whereas others are not.
- ▶ **Essential idea:** Construct a combined relationship matrix across all animals, using both marker genotypes and pedigree, and use a model for all animals.
- ▶ Special case where no animals are genotyped: Animal model based on pedigree relationships.
- ▶ Special case where all animals are genotyped: GBLUP

Combined relationship matrix

- ▶ Pedigree-based relationship matrix

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

where split is into non-genotyped and genotyped animals.

- ▶ Example

ID	FATHER	MOTHER	
1	0	0	
2	0	0	genotyped
3	1	2	
4	1	2	genotyped

Combined relationship matrix

Example:

ID	FATHER	MOTHER	
1	0	0	
2	0	0	genotyped
3	1	2	
4	1	2	genotyped

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0.5 & 0.5 \\ 0 & 1 & 0.5 & 0.5 \\ 0.5 & 0.5 & 1 & 0.5 \\ 0.5 & 0.5 & 0.5 & 1 \end{bmatrix}$$

Combined relationship matrix

- ▶ Pedigree-based relationship matrix

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

where split is into non-genotyped and genotyped animals
(note: animals reordered).

- ▶ A first naive attempt:

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

- ▶ Doesn't adjust relationships within non-genotyped and between genotyped and non-genotyped animals depending on \mathbf{G}
- ▶ Invalid (not a positive definite matrix).

Combined relationship matrix

- ▶ Pedigree-based relationship matrix

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

where split is into non-genotyped and genotyped animals

- ▶ Second attempt

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

- ▶ Modify relationships depending on \mathbf{G}

Combined relationship matrix

▶ Model: $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{W}\mathbf{a} + e$, where $\mathbf{a} \sim N(\mathbf{0}, \sigma_a^2 \mathbf{H})$

▶ Matrix \mathbf{H} has sparse inverse

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

▶ Matrix \mathbf{H} is valid (positive definite)

▶ Easy to compute, but not based on tracing of inheritance (sub-optimal)

▶ \mathbf{H}^{-1} is provided as input to genetic evaluation software

Combined relationship matrix

▶ Model: $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{W}\mathbf{a} + \mathbf{e}$ where $\mathbf{a} \sim N(\mathbf{0}, \sigma_a^2 \mathbf{H})$

▶ Matrix \mathbf{H} has sparse inverse

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

▶ Special case: all animals genotyped, then $\mathbf{H} = \mathbf{G}$

▶ Special case: no animals genotyped, then $\mathbf{H} = \mathbf{A}$

Compatibility of **G** and **A**

- ▶ Relationships and inbreeding in **A** are relative to the base population
- ▶ Relationships and inbreeding in **G** are relative to allele frequencies
- ▶ Adjust $\mathbf{G}_a = b\mathbf{G} + a\mathbf{1}\mathbf{1}^T$ such that \mathbf{G}_a and \mathbf{A}_{22} are compatible (ad hoc).
- ▶ Metafounders in pedigree and $\mathbf{G}_{0.5}$ with allele frequencies 0.5 (more elegant)

Genomic models and BLUP

- ▶ SNPBLUP, GBLUP, ssGBLUP presented here
- ▶ Other genomic models:
 - ▶ Non-additive genomic relationship matrices (Dominance, G^*G , G^*E)
 - ▶ Breed-specific (partial) genomic relationship matrices.
- ▶ Non-genetic similarity matrices (e.g. metabolomics).
- ▶ Standard genetic evaluation software with BLUP is powerful

BLUP vs Bayes+McMC - pros/cons

- ▶ Pros (compared to Bayes + McMC)
 - ▶ Usually computationally fast
 - ▶ No Monte Carlo noise in (important for routine genetic evaluation)
 - ▶ No need to assess convergence and mixing of McMC

- ▶ Cons (compared to Bayes + McMC)
 - ▶ Parameter uncertainty not incorporated in prediction (important?)
 - ▶ Not always sufficiently flexible (important for research)