

# A brief introduction to mixed models

Alencar Xavier

Quantitative Geneticist, DowDuPont

<http://alenxav.wix.com/home>

- **Part 1 – Concepts**

- Mixed models in plant breeding
- Fixed and Random terms
- History of mixed models
- Model notation
- Henderson's equation
- Variance decomposition

- **Part 2 – Applications**

- Phenotypic selection
- Practical example
  - Cunningham & Henderson (1986)
  - Robinson (1991)
  - Xavier et al. (2016)
- Models with genomic data
  - G-BLUP and RR-BLUP
  - GWAS

- **References**

# Outline

# PART I - Concepts

# Mixed models in plant breeding

- Variance components and heritability
- Genetic correlations
- Estimation of genetic values
- Estimation of breeding values
- Prediction of breeding values (unphenotyped material)
- Selection indexes
- Association analysis

# Fixed and Random terms

- **Fixed effect**

- Assumed to be invariable (you can't recollect the data)
- Inferences are made upon the parameters
- Results can not be extrapolated to other datasets
- Example: Overall mean and environmental effects

- **Random effects**

- You may not have all the levels available
- Inference are made on variance components
- Prior assumption – coefficients are normally distributed
- Results can not be extrapolated to other datasets
- Regularized (shrinkage)
- Example: Genetic effects

Robinson (1991) reciting Searle

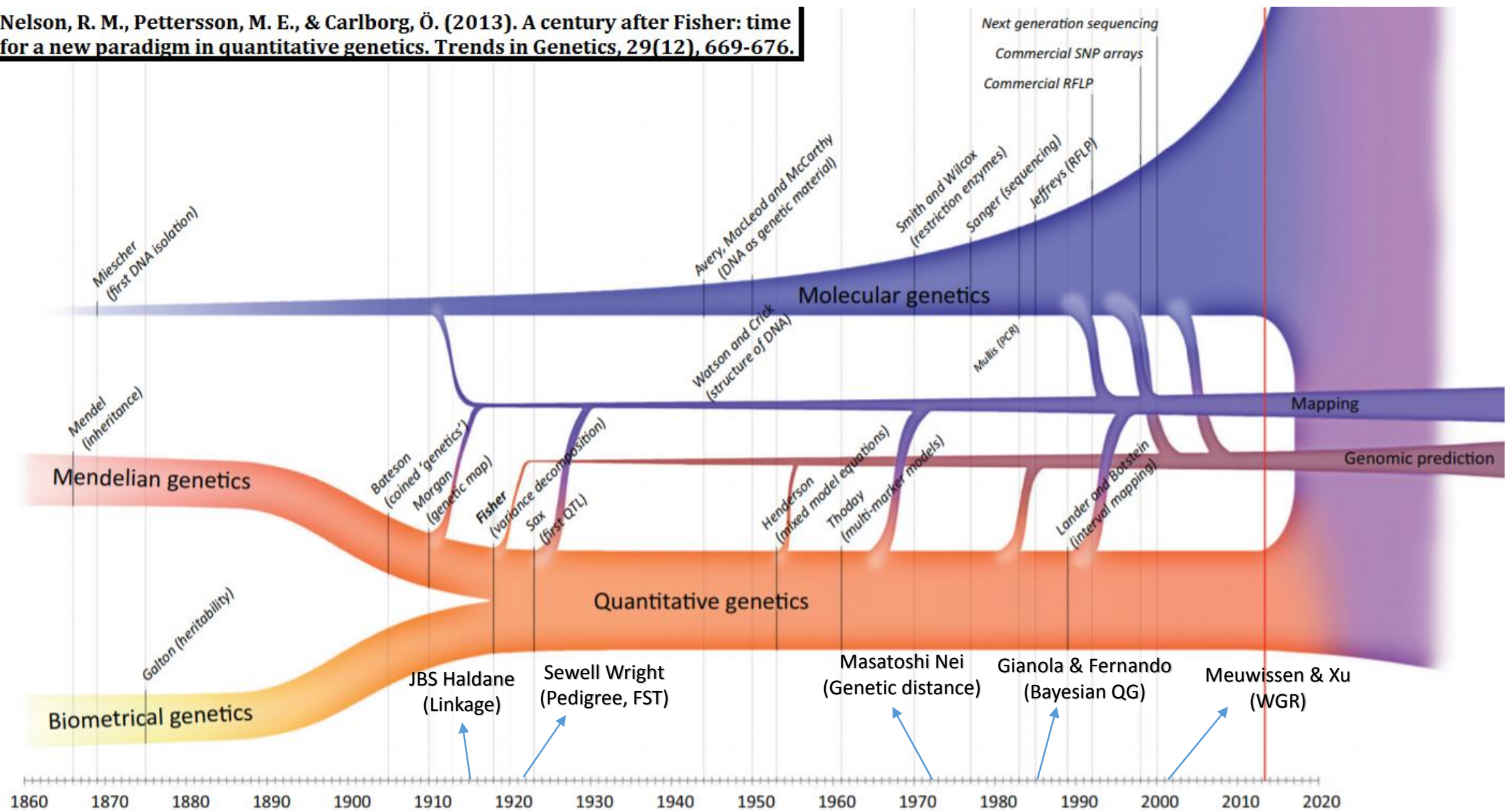
... when inferences are going to be confined to the effects in the model the effects are considered fixed; and when inferences will be made about a population of effects from which those in the data are considered to be a random sample then the effects are considered as random.

# History of mixed models

- Francis Galton
  - **1888** – Regression and  $h^2$
- Ronald Fisher
  - **1918** - Infinitesimal model ( $P = G + E$ )
- Sewall Wright
  - **1922** - Pedigree matrix ( $A$ )
- Charles Henderson
  - **1950** - BLUP ( $u \sim A\sigma_a^2$ )



**Nelson, R. M., Pettersson, M. E., & Carlborg, Ö. (2013). A century after Fisher: time for a new paradigm in quantitative genetics. *Trends in Genetics*, 29(12), 669-676.**



# Model notation

$$y = 1\mu + g + e$$

$$\sigma_y^2 = \sigma_g^2 + \sigma_e^2$$

$$\mu = 50$$

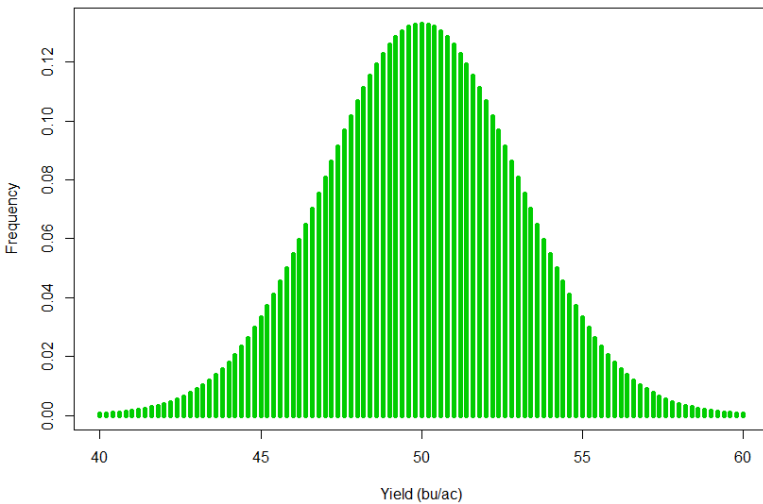
$$\sigma_y^2 = 9$$

$$\sigma_g^2 = 4$$

$$\sigma_e^2 = 5$$

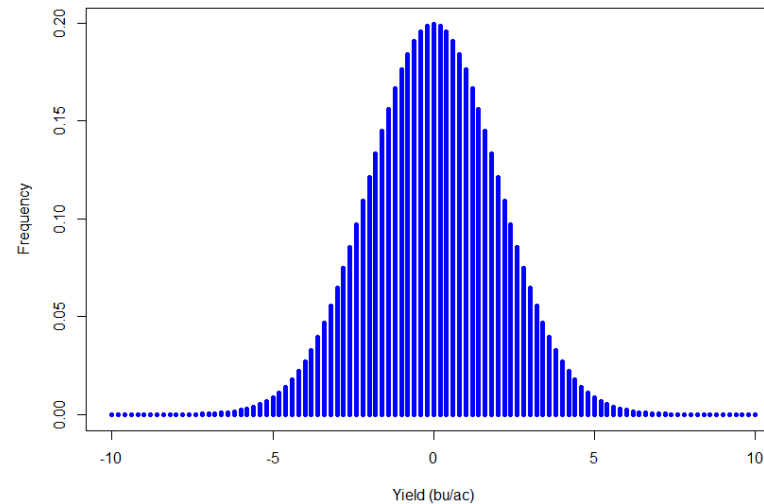
$$\text{cov}(g, e) = 0$$

Phenotype



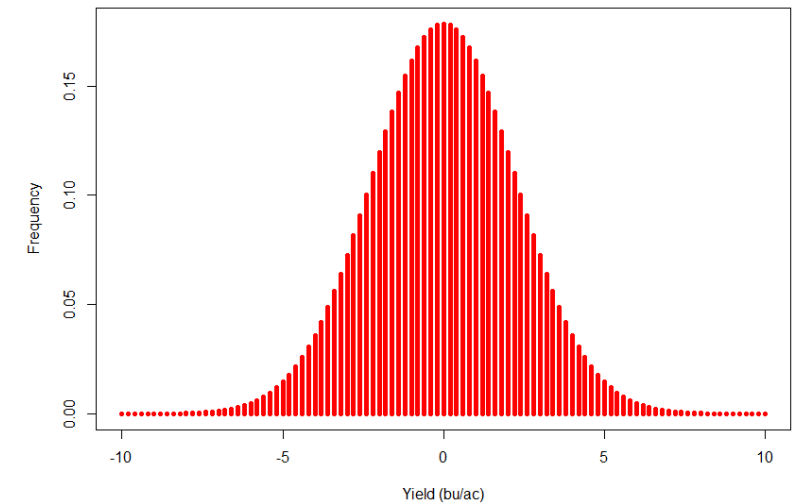
$$y \sim N(\mu, I\sigma_y^2)$$

Genotype



$$u \sim N(0, I\sigma_g^2)$$

Residual



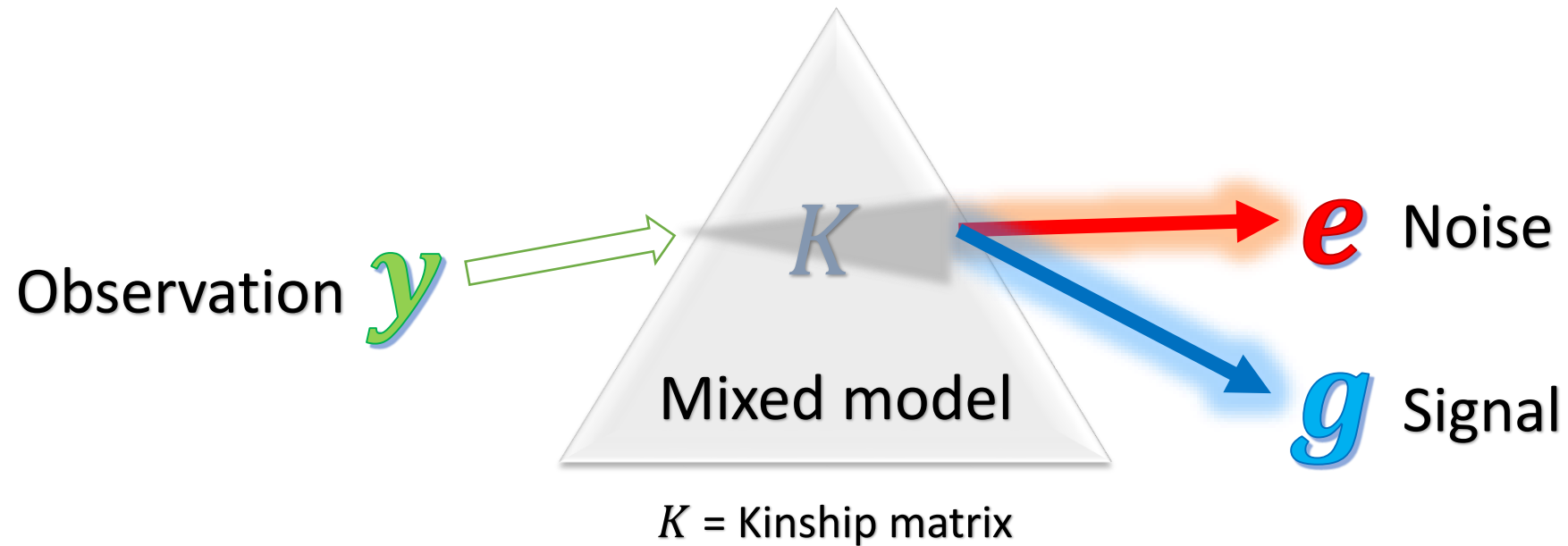
$$e \sim N(0, I\sigma_e^2)$$

DOW RESTRICTED



# Model notation

$$y = \mu + g + e$$



# Model notation

$n$  = number of observations

$p$  = number of parameters

$q$  = number of individuals

$$y = Xb + Zu + e$$

$$\left\{ \begin{array}{l} y \sim N(Xb, V) \\ y \sim N(Xb, ZAZ'\sigma_a^2 + R\sigma_e^2) \\ y \sim N(Xb + Zu, R\sigma_e^2) \end{array} \right.$$

$$u \sim N(0, A\sigma_a^2)$$

$$e \sim N(0, R\sigma_e^2)$$

$$\text{cov}(u, e) = 0$$

$y$  = vector of observations ( $n$ )

$X$  = design matrix of fixed effects ( $n \times p$ )

$b$  = vector of fixed effect coefficients ( $p$ )

$Z$  = incidence matrix of random effects ( $n \times q$ )

$u$  = vec. of random effects – genetics values ( $q$ )

$e$  = vector of residuals ( $n$ )

$\sigma_a^2$  = random effect variance (1)

$\sigma_e^2$  = residual variance (1)

$A$  = random effect correlation matrix ( $q \times q$ )

$R$  = residual correlation matrix ( $n \times n$ )

$\lambda = \sigma_e^2 : \sigma_a^2$  = regularization parameter (1)

# Henderson's equation

Model statement

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \quad \mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}) \quad \mathbf{V} = \mathbf{Z}\mathbf{D}\mathbf{Z}' + \mathbf{R} \quad \begin{matrix} D = A\sigma_a^2 \\ R = I\sigma_e^2 \end{matrix}$$

Generalized Mixed  
Linear Model Equation

$$\begin{bmatrix} \mathbf{X}'\mathbf{R}^{-1}\mathbf{X} & \mathbf{X}'\mathbf{R}^{-1}\mathbf{Z} \\ \mathbf{Z}'\mathbf{R}^{-1}\mathbf{X} & \mathbf{D}^{-1} + \mathbf{Z}'\mathbf{R}^{-1}\mathbf{Z} \end{bmatrix} \begin{bmatrix} \tilde{\boldsymbol{\beta}} \\ \tilde{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{R}^{-1}\mathbf{y} \\ \mathbf{Z}'\mathbf{R}^{-1}\mathbf{y} \end{bmatrix}$$

BLUE & BLUP  
solutions

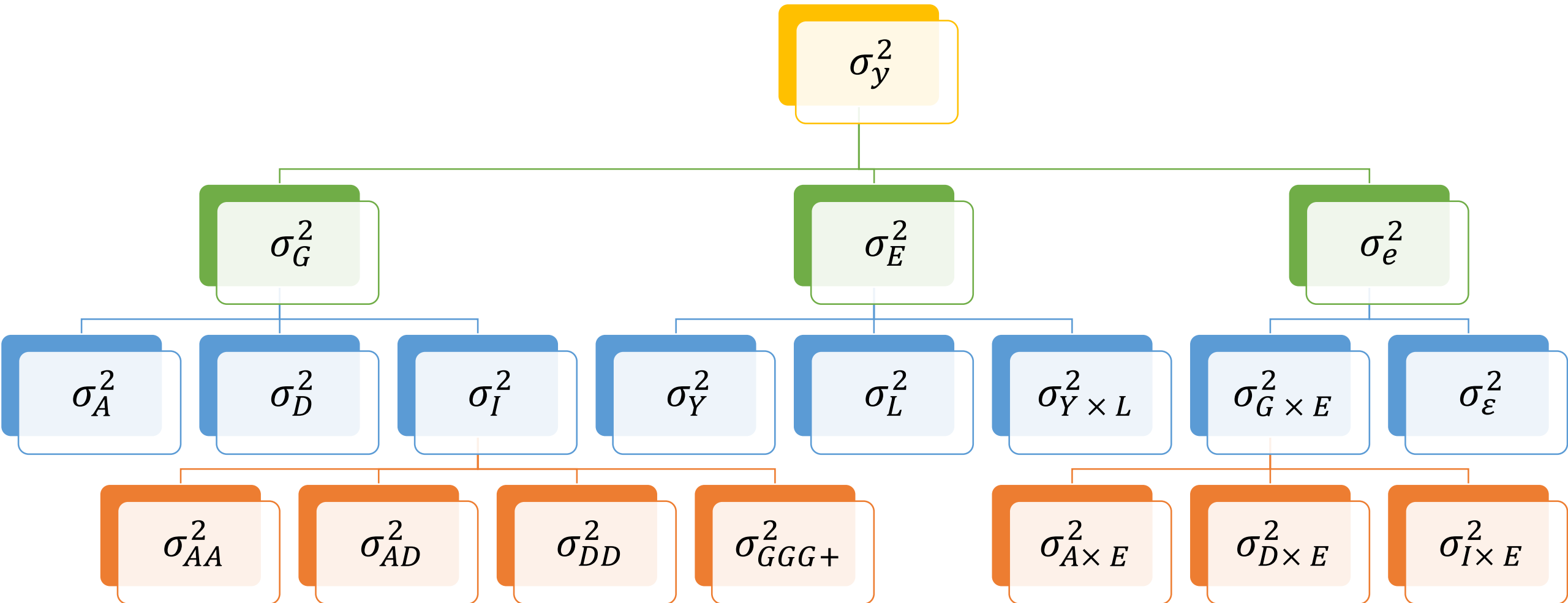
$$\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \mathbf{X}'\mathbf{V}^{-1}\mathbf{y}, \quad \text{with } \text{BLUE}(\mathbf{X}\boldsymbol{\beta}) = \mathbf{X}\tilde{\boldsymbol{\beta}}$$

$$\tilde{\mathbf{u}} = \mathbf{D}\mathbf{Z}'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) = \text{BLUP}(\mathbf{u})$$

Reduce Model  
( $R = I$ )

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \left\{ \lambda_i \mathbf{I}_{q_i} \right\}_{i=1}^r \end{bmatrix} \begin{bmatrix} \tilde{\boldsymbol{\beta}} \\ \tilde{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{Z}'\mathbf{y} \end{bmatrix}$$

# Variance decomposition: Multiple random effects

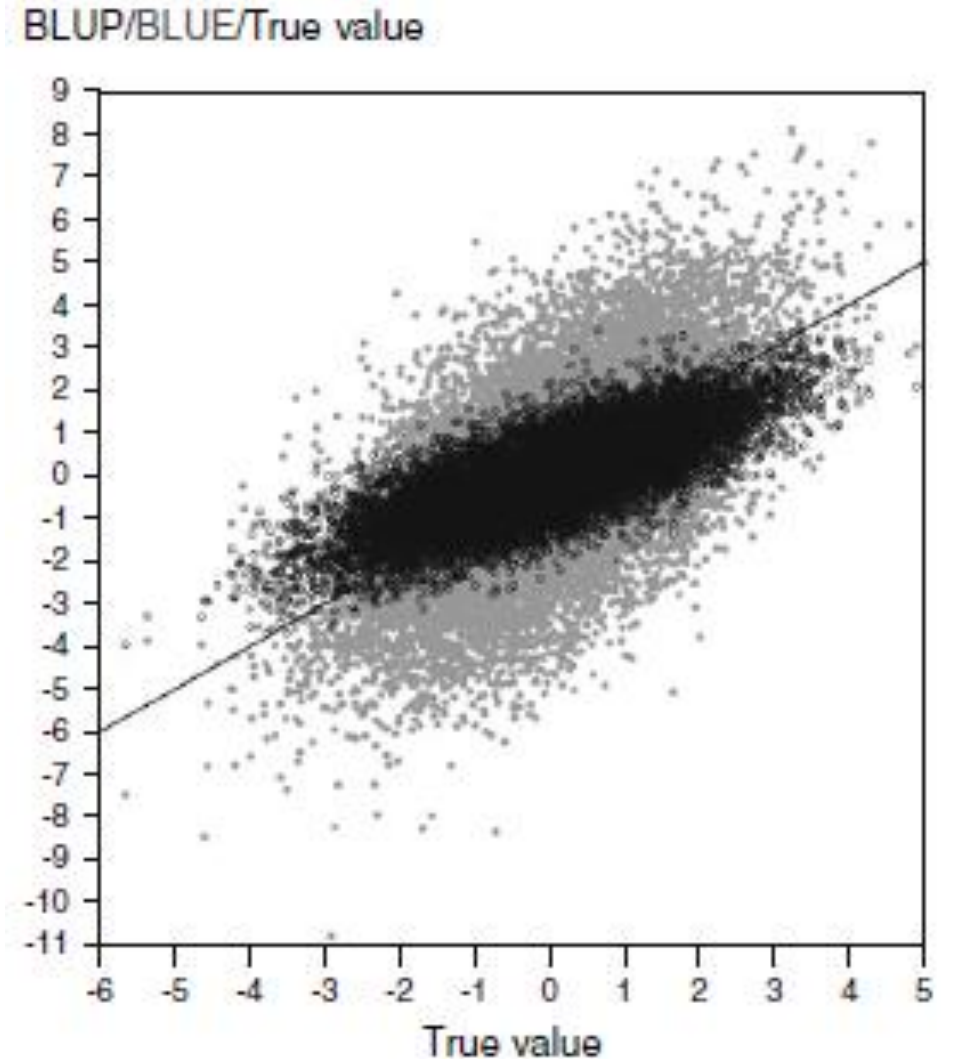


DOW RESTRICTED

# PART II - Applications

# Phenotypic selection

1. Genetic values (image on right)
  - From BLUP or BLUE
  - Non-additive
  - Require replicated trials
2. Breeding values
  - Use pedigree information
  - Additive genetics
  - Not necessarily replicated
3. Genomic Breeding values
  - Genotypes replace pedigree
  - Not necessarily additive



**Fig. 4** Plot of simulated BLUP (black circles), BLUE (grey dots), and true genetic values (solid line) versus true values

Piepho et al. 2008

# Practical examples

# Example from Cunningham & Henderson 1968

$$y = \mu + Xa + Zb + e.$$

DATA AND INCIDENCE MATRICES						
$y$	$\mu$	$a_1$	$a_2$	$b_1$	$b_2$	$b_3$
3	1	1	0	1	0	0
2	1	1	0	0	1	0
3	1	1	0	0	0	1
2	1	1	0	1	0	0
3	1	1	0	0	1	0
5	1	1	0	0	1	0
6	1	1	0	0	1	0
7	1	1	0	0	1	0
2	1	0	1	1	0	0
8	1	0	1	0	1	0
4	1	0	1	0	0	1
3	1	0	1	1	0	0
8	1	0	1	0	1	0
4	1	0	1	0	0	1
9	1	0	1	0	1	0
3	1	0	1	0	0	1
2	1	0	1	0	0	1
5	1	0	1	0	0	1

$$y = Xa + Zb + e$$

The least squares equations (ignoring  $\mu$ ) are

$$\begin{bmatrix} 8 & 0 & 2 & 5 & 1 \\ 0 & 10 & 2 & 3 & 5 \\ \hline 2 & 2 & 4 & 0 & 0 \\ 5 & 3 & 0 & 8 & 0 \\ 1 & 5 & 0 & 0 & 6 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \hline b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} 31 \\ 48 \\ \hline 10 \\ 48 \\ 21 \end{bmatrix}$$

In algebraic terms, these equations are

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix}$$

$$\lambda = 0.5721$$

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + I\lambda \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix}$$

$$\begin{bmatrix} 8 & 0 & 2 & 5 & 1 \\ 0 & 10 & 2 & 3 & 5 \\ \hline 2 & 2 & 4.5721 & 0 & 0 \\ 5 & 3 & 0 & 8.5721 & 0 \\ 1 & 5 & 0 & 0 & 6.5721 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \hline b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} 31 \\ 48 \\ \hline 10 \\ 48 \\ 21 \end{bmatrix}$$

```
> solve(C, g)
[1] 2.9371 4.8684 -1.2272 2.1826 -0.9554
```

a1            a2            b1            b2            b3



# Example from Robinson 1991

## Data

Herd	Sire	Yield
1	A	110
1	D	100
2	B	110
2	D	100
2	D	100
3	C	110
3	C	110
3	D	100
3	D	100

## Design matrices

$$X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix} \quad Z = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

## Solving

$$\left( \begin{array}{ccc|cccc} 2 & 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 3 & 0 & 0 & 1 & 0 & 2 \\ 0 & 0 & 4 & 0 & 0 & 2 & 2 \\ \hline 1 & 0 & 0 & 11 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 11 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 12 & 0 \\ 1 & 2 & 2 & 0 & 0 & 0 & 15 \end{array} \right) \begin{pmatrix} \hat{h}_1 \\ \hat{h}_2 \\ \hat{h}_3 \\ \hat{s}_A \\ \hat{s}_B \\ \hat{s}_C \\ \hat{s}_D \end{pmatrix} = \begin{pmatrix} 210 \\ 310 \\ 420 \\ 110 \\ 110 \\ 220 \\ 500 \end{pmatrix}$$

which has solution

$$(1.4) \quad \begin{aligned} \hat{\beta} &= (105.64, 104.28, 105.46)^T, \\ \hat{u} &= (0.40, 0.52, 0.76, -1.67)^T. \end{aligned}$$

## MME

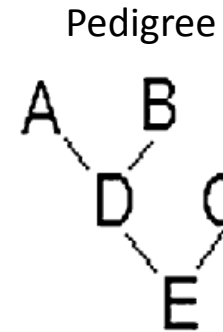
$$\begin{bmatrix} X'X & Z'X \\ X'Z & Z'Z + \lambda K^{-1} \end{bmatrix} \begin{bmatrix} b \\ u \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix}$$

( $\lambda = \sigma_e^2 / \sigma_a^2$ )

DOW RESTRICTED

# Example from Xavier et al. 2016

Field map (line and its yield)		
A = 27	Missing	E = 21
E = 27	C = 20	B = 27
B = 21	A = 25	Missing



INPUT

$$\mathbf{y} = \begin{bmatrix} \text{Yield} \\ 25 \\ 27 \\ 27 \\ 21 \\ 20 \\ 21 \\ 27 \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} \text{Intercept} \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\mathbf{Z} = \begin{bmatrix} & \text{A} & \text{B} & \text{C} & \text{D} & \text{E} \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} & \text{A} & \text{B} & \text{C} & \text{D} & \text{E} \\ \text{A} & 1 & 0 & 0 & 0.5 & 0.25 \\ \text{B} & 0 & 1 & 0 & 0.5 & 0.25 \\ \text{C} & 0 & 0 & 1 & 0 & 0.5 \\ \text{D} & 0.5 & 0.5 & 0 & 1 & 0.5 \\ \text{E} & 0.25 & 0.25 & 0.5 & 0.5 & 1 \end{bmatrix}$$

OUTPUT

$$\mathbf{b} = [23.812] \quad \mathbf{u} = \begin{bmatrix} 1.191 \\ 0.172 \\ -1.291 \\ 0.799 \\ -0.060 \end{bmatrix} \quad \sigma_a^2 = 4.004 \quad \sigma_e^2 = 6.987$$

# Models with genomic data

# G-BLUP and RR-BLUP

- **Purpose:** Estimate or predict breeding values using genomic data
- **Method 1 (G-BLUP):** MME using genomic relationship matrix ( $\mathbf{K}=\mathbf{G}$ )
- **Method 2 (RR-BLUP):** MME using centralized markers as the random design matrix ( $\mathbf{Z}=\mathbf{M}$ )

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

- Environment
- Block effect
- Covariates

- Breeding values (1):  $\mathbf{u} \sim N(0, \mathbf{G}\sigma_a^2)$
- Breeding values (2):  $\mathbf{u} = \mathbf{M}\boldsymbol{\alpha}, \boldsymbol{\alpha} \sim N(0, \mathbf{I}\sigma_u^2)$


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

$$\sigma_a^2 = \sigma_u^2 2 \sum_{j=1}^P p_j(1-p_j)$$

# Genome-Wide Association Studies (GWAS)

- **Purpose:** Identify markers associated to the trait of interest
- **Method:** Likelihood Ratio Test (LRT) between of model with and without markers

$$y = Xb + Zu + e$$

- 
- Intercept
  - Sub-population (**Q**)
  - Covariates
  - Marker (**with** vs. **without**)
- Breeding values:  $\mathbf{u} \sim N(0, \mathbf{K}\sigma_a^2)$

# REFERENCES

1. Cunningham, E. P., & Henderson, C. R. (1968). An iterative procedure for estimating fixed effects and variance components in mixed model situations. *Biometrics*, 13-25.
2. Robinson, G. K. (1991). That BLUP is a good thing: the estimation of random effects. *Statistical science*, 15-32.
3. Searle, S. R. (1997). The matrix handling of BLUE and BLUP in the mixed linear model. *Linear algebra and its applications*, 264, 291-311.
4. Piepho, H. P., Möhring, J., Melchinger, A. E., & Büchse, A. (2008). BLUP for phenotypic selection in plant breeding and variety testing. *Euphytica*, 161(1-2), 209-228.
5. Xavier, A., Muir, W. M., Craig, B., & Rainey, K. M. (2016). Walking through the statistical black boxes of plant breeding. *Theoretical and Applied Genetics*, 129(10), 1933-1949.