**Fundamentals of Genomic Prediction and Data-Driven Crop Breeding (August 4-8, 2025)**

**Understanding Prediction Models: Ridge Regression to Bayesian**

**Module 3**
**August 6, 2025**

**Waseem Hussain and Mahender Anumalla**
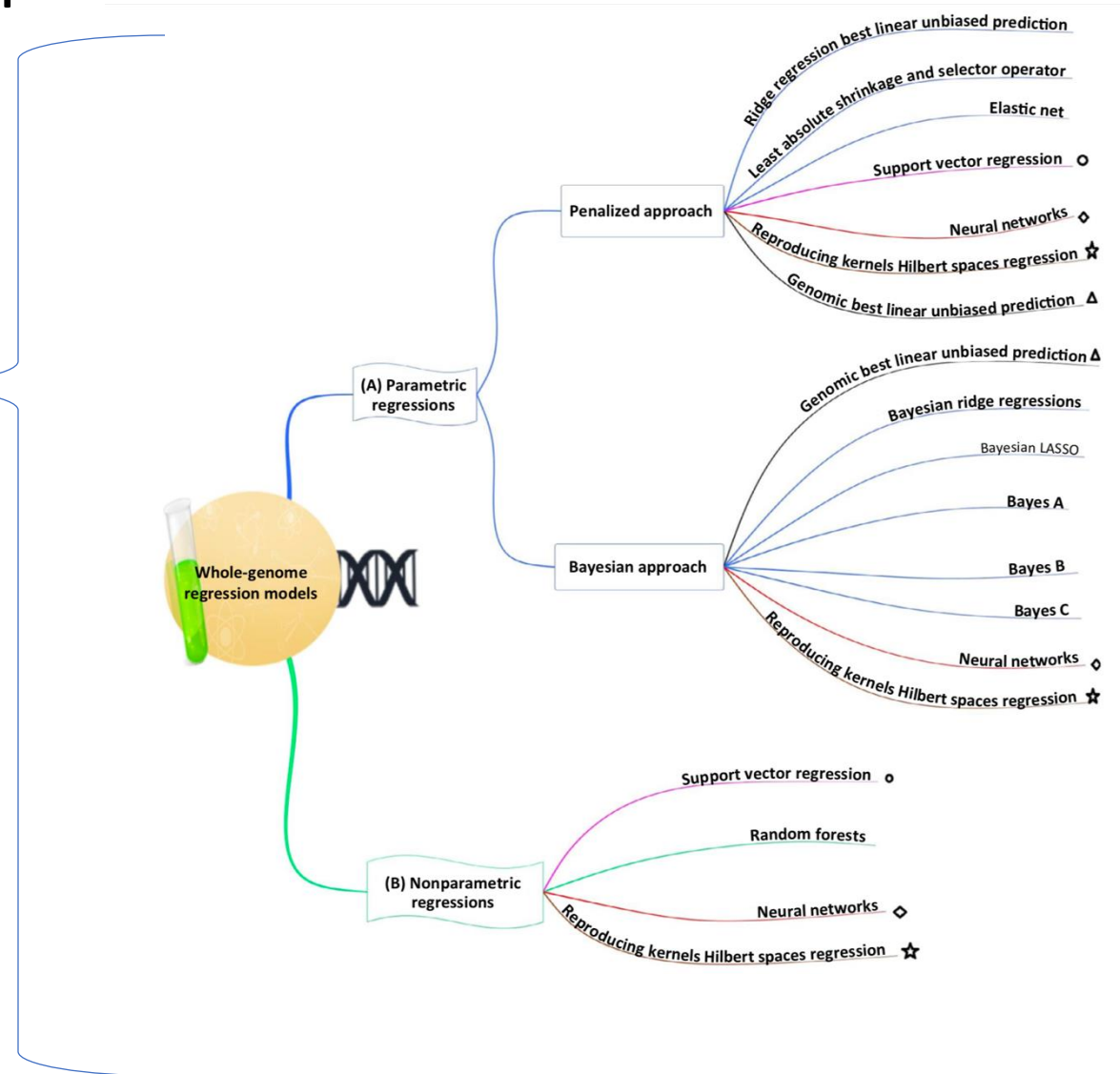**Rice Breeding Innovations Platform**
**IRRI**

# Genomic Prediction Models

$$y = \mu + \sum_{k} x_k \beta_k + \varepsilon \longrightarrow \text{Error}$$

Trait     mean     Design matrix     Effect of $k$th marker

**Key: How we estimate and Assume Marker Variances.**

a. **Shrinkage Models**: RR BLUP and GBLUP

b. **Dimension Reduction Methods: Partial Least Squares (PLS) and Singular Value Decomposition (SDV)**

c. **Bayesian Approach: Variable Selection Models (priors)**

d. **Kernel and Machine Learning Methods**



Whole-genome regression models

(A) Parametric regressions

Penalized approach
- Ridge regression best linear unbiased prediction
- Least absolute shrinkage and selector operator
- Elastic net
- Support vector regression ○
- Neural networks ◇
- Reproducing kernels Hilbert spaces regression ☆
- Genomic best linear unbiased prediction Δ

Bayesian approach
- Genomic best linear unbiased prediction Δ
- Bayesian ridge regressions
- Bayesian LASSO
- Bayes A
- Bayes B
- Bayes C
- Neural networks ◇
- Reproducing kernels Hilbert spaces regression ☆

(B) Nonparametric regressions
- Support vector regression ○
- Random forests
- Neural networks ◇
- Reproducing kernels Hilbert spaces regression ☆

Adapted from manuscript Desta and Ortiz, 2014

# Ordinary Least squares (OLS)

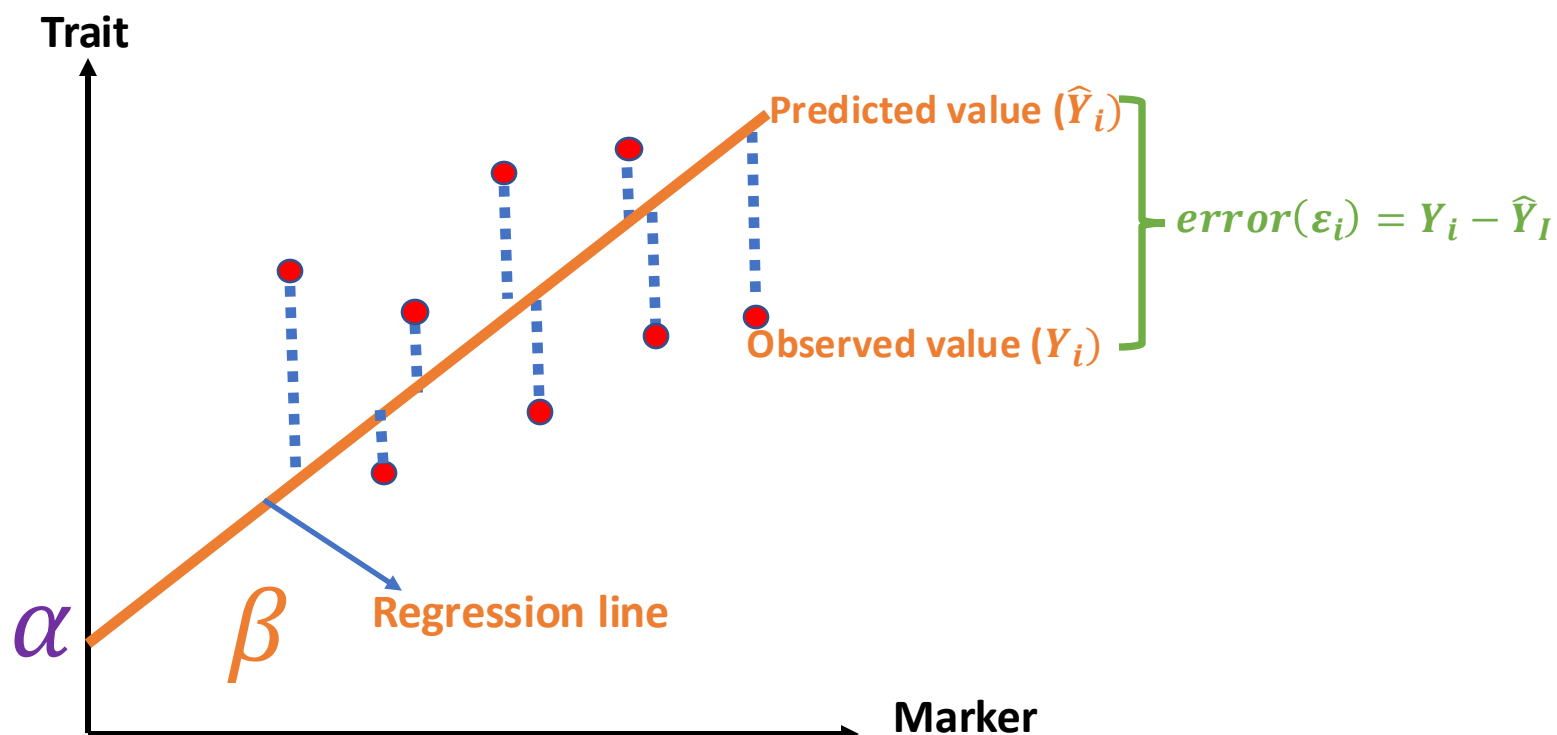Aim is to estimate $\alpha$ and $\beta$ parameters by minimize the squared errors

Simple linear regression:
$$\hat{Y}_i = \alpha + \beta x_i + \varepsilon_i$$

dependent variable | intercept | coefficient | Independent variable | error

Parameters (un-observed)

Trait

Predicted value $(\hat{Y}_i)$

$error(\varepsilon_i) = Y_i - \hat{Y}_I$

Observed value $(Y_i)$

$\alpha$   $\beta$   Regression line

Marker

# OLS in Relevance to Genomic Predictions

*Aim is to minimize residual squares*

$$argmin(\varepsilon`\varepsilon) = \text{argmin}(y - X\beta)`(y - x\beta)$$

$$\dots$$

$$\beta = (\acute{X}X)^{-1}\acute{X}Y$$

*Simple marker regression model*

$$Y_n = \mu_m + X_m\beta_n + \varepsilon$$

$\mu_m$ is mean of $n$th marker

|  | SNP$_1$ | SNP$_m$ |
|---|---|---|
| $y_{11}$ | 0 | 0 |
| $y_{21}$ | 2 | 2 |
| $y_{31}$ | 2 | 2 |
| $y_{41}$ | 0 | 0 |
| $y_{51}$ | 0 | 0 |
| . | . | . |
| . | . | . |
| . | . | . |
| $y_{m1}$ | 2 | 2 |

$y_1$
$y_2$
$y_3$
$y_4$
.
.
$y_n$

$\beta_n$ is the additive effect of $n$th marker

$Y_n$ are the phenotypic values of $n$th plants/individuals

$X_m$ is a vector or design matrix containing alleles (0,2 or 1) and connecting it to phenotypic values.

error for marker assumed $\varepsilon \sim N(0, \sigma_e^2)$, with mean 0 and marker variance $\sigma_e^2$

## Limitations

**1. m>>n, we cannot fit OLS (curse on dimensionality)**

$\acute{X}X$ is a singular, and determinate of $\acute{X}X$ is 0 so cannot take inverse, $\acute{X}X$ is not invertible

**2. m>>n, multi-collinearity produces singular matrix**

**3. m>>n, not perfect multi-collinearity, but predictors are highly correlated**

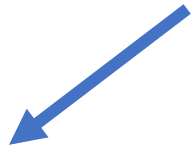*OLS estimates are not stable and have high variance*

## Need alternative strategy

# Ridge Regression

**Provides solution and overcomes problems of OLS**

*Aim is to minimize residual squares*

$$argmin(\varepsilon`\varepsilon) = \text{argmin}(y - X\beta)`(y - x\beta)$$
$$\ldots$$
$$\beta = (\acute{X}X + \lambda I)^{-1} \acute{X}Y$$

**Add constant, Penalty/Shrinkage Factor**

$$E(\hat{\beta}_{ridge}) = (\acute{X}X + \lambda I)^{-1} X\acute{E}(Y)$$

$$E(\hat{\beta}_{ridge}) = \beta - \lambda(\acute{X}X + \lambda I)^{-1}$$

**Biased, but with minimum variance**

*adding $\lambda$, we take inverse but we have biased estimates*

$$when, \lambda = 0, E(\widehat{\beta)} = \beta, i.e, OLS$$

# Ridge-Regression BLUP (RRBLUP)

$$y = \mu + \sum_k x_k \beta_k + \varepsilon$$

estimate $\beta$ by adding positive constraint

$$\beta_{ridge} = (X^T X)^{-1} + \lambda I)\, X^T y$$

$$\lambda = \frac{\sigma_e^2}{\sigma_\beta^2}$$

$$\beta_k \sim N(0, \sigma_e^2)$$

❖ **Ridge regression induces homogeneous shrinkage and it depends upon allele frequency**

❖ **Assumes all markers have same variance with small but non-zero effect.**

- We add constant to diagonal, thus it is invertible
- Degree of shrinkage depends upon $\lambda$, larger the $\lambda$ larger is the shrinkage

$$\frac{\beta_{OLS}}{1 + \lambda}$$

- We basically shrink OLS estimates towards 0
- $\lambda I$ term reduces collinearity and prevents the matrix $X^T X$ from becoming singular.

# Things to Know (RRBLUP)

$$y = \mu + \sum_k z_k \beta_k + \varepsilon$$

$$\beta_{ridge} = (Z^T Z)^{-1} + \lambda I)\, Z^T\, y$$

- We add constant to diagonal, thus it is invertible
- Degree of shrinkage depends upon $\lambda$, larger the $\lambda$ larger is the shrinkage

$$\frac{\beta_{OLS}}{1 + \lambda} \qquad \lambda = \frac{\sigma_\beta^2}{\sigma_\epsilon^2}$$

- We basically shrink $\beta$ (OLS) estimates towards 0
- $\lambda I$ term reduces collinearity and prevents the matrix $X^T X$ from becoming singular.

❖ Ridge regression induces homogeneous shrinkage, same for all markers

❖ Assumes all markers have the same variance with small but non-zero effect.

❖ Shrinkage depends upon allele frequency

$$\widehat{\beta}_{ridge} = \frac{2p_j(1-p_j)n}{2p_j(1-p_j)n + \lambda} * \widehat{\beta}_{OLS}$$

❖ Markers with extreme frequency are shrunk more.

# RR BLUP Model

$$\mathbf{y} = \mathbf{Z}g + \boldsymbol{\varepsilon}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ . \\ y_n \end{bmatrix} = \begin{bmatrix} \text{SNP1} & \text{SNP2} & \text{SNP3} \\ 1 & 2 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 2 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_p \end{bmatrix} + \boldsymbol{\varepsilon}$$

$n \times 1$  $n \times m$  $p \times 1$  $n \times 1$

**n is # of individuals**
**P is # of markers**

$$\begin{bmatrix} a \\ \varepsilon \end{bmatrix} \sim N\left(0, \begin{matrix} \sigma_a^2 & 0 \\ 0 & \sigma_\varepsilon^2 \end{matrix}\right)$$

**Solve mixed model equation (Henderson, 1989)**

$$\begin{bmatrix} \acute{X}X & \acute{X}Z \\ \acute{Z}X & \acute{Z}Z + \lambda I^{-1} \end{bmatrix} \begin{bmatrix} b \\ a \end{bmatrix} = \begin{vmatrix} \acute{X}y \\ \acute{Z}y \end{vmatrix}$$

$$\begin{bmatrix} \hat{a}_1 \\ \hat{a}_2 \\ \hat{a}_3 \\ \hat{a}_4 \\ \hat{a}_5 \\ \hat{a}_6 \end{bmatrix} = \begin{bmatrix} -0.35 \\ -0.25 \\ 1.35 \\ -1.15 \\ 1.45 \\ 0.45 \end{bmatrix}$$

**Marker effects (BLUP)**

$$\text{GEBVs}(\widehat{u}) = \mathbf{Z}\widehat{a}$$

|       | $a_1$ | $a_2$ | $a_3$ | $a_4$ |
|-------|-------|-------|-------|-------|
| $y_1$ | 0     | 0     | 0     | 2     |
| $y_2$ | 2     | 1     | 0     | 2     |
| $y_3$ | 2     | 2     | 2     | 2     |
| $y_4$ | 2     | 2     | 0     | 0     |

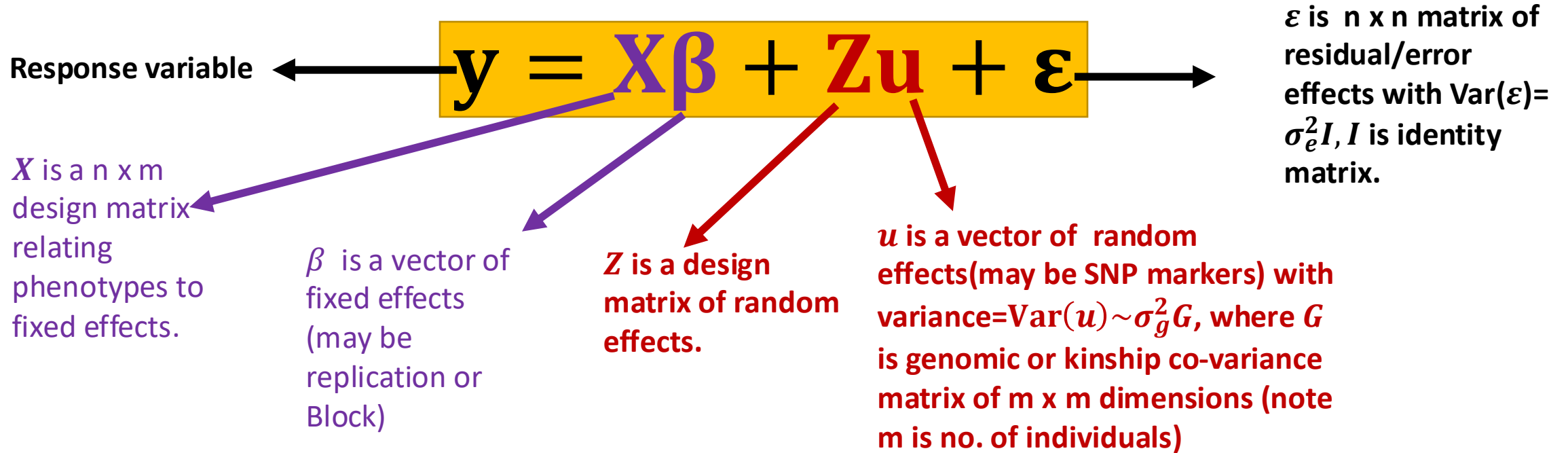$$GEBV(y_1) = [0*(-0.35) + [0*(-0.25)] + \cdots. + [2*(-1.15)]$$

What if LD exists between markers?

*RRBLUP or Ordinary least squares (OLS) do not consider:*

**LD Between Markers and account for the relationship between the genotypes**

*Better to compute and Integrate the Relationship matrix in the model*

# Mixed Effect Model

$$y = X\beta + Zu + \varepsilon$$

**Response variable**

$\varepsilon$ is n x n matrix of residual/error effects with Var($\varepsilon$)= $\sigma_e^2 I$, $I$ is identity matrix.

$X$ is a n x m design matrix relating phenotypes to fixed effects.

$\beta$ is a vector of fixed effects (may be replication or Block)

$Z$ is a design matrix of random effects.

$u$ is a vector of random effects(may be SNP markers) with variance=Var($u$)$\sim\sigma_g^2 G$, where $G$ is genomic or kinship co-variance matrix of m x m dimensions (note m is no. of individuals)

$$\begin{bmatrix} u \\ \varepsilon \end{bmatrix} \sim \left( N\left( 0, \begin{matrix} \sigma_u^2 K & 0 \\ 0 & \sigma_\varepsilon^2 R \end{matrix} \right) \right)$$

Note: cov($u$, $\varepsilon$)=0

# Solving Mixed model Equation

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

$$\begin{bmatrix} \acute{X}X & \acute{X}Z \\ \acute{Z}X & \acute{Z}Z + \lambda G^{-1} \end{bmatrix} \begin{bmatrix} \beta \\ u \end{bmatrix} = \begin{bmatrix} \acute{X}y \\ \acute{Z}y \end{bmatrix}$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\textbf{LHS}} \qquad\qquad \underbrace{\qquad}_{\textbf{RHS}}$$

$$\tilde{u} = DZ'V^{-1}(y - X\beta) = BLUP(u)$$

$$\begin{bmatrix} \widehat{\beta} \\ \widehat{u} \end{bmatrix} = \begin{bmatrix} \acute{X}X & \acute{X}Z \\ \acute{Z}X & \acute{Z}Z + \lambda G^{-1} \end{bmatrix}^{-1} \cdot \begin{bmatrix} \acute{X}y \\ \acute{Z}y \end{bmatrix}$$

*We decompose the matrices and get solutions through iterative methods*

# Genomic BLUP (gBLUP)

$$\mathbf{y} = \mathbf{X\beta} + \mathbf{Zu} + \mathbf{\varepsilon}$$

$$\hat{u} = \left[\mathbf{I} + \mathbf{G}^{-1}\frac{\sigma_e^2}{\sigma_u^2}\right]\mathbf{y}$$

GRM to account for mendelian sampling

**Equivalence between rrBLUP and gBLUP**

For gBLUP the $Var(y) = \mathbf{ZGZ'}\sigma_u^2 + \mathbf{I}\sigma_e^2$

For rrBLUP the $Var(y) = \mathbf{XX'}\sigma_\beta^2 + \mathbf{I}\sigma_e^2$

# Genomic BLUP (gBLUP Model)

$$\mathbf{y} = \mathbf{Z}a + \boldsymbol{\varepsilon}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ . \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & 0.5 & 0.3 & 0.9 \\ 0.5 & 1 & -0.2 & 0.8 \\ 0.3 & -0.2 & 1 & 0.4 \\ 0.9 & 0.8 & .4 & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_p \end{bmatrix} + \boldsymbol{\varepsilon}$$

**n x 1**   **n x g**   **g x 1**   **n x 1**

**n is # of individuals**
**g is # of breeding values**

$$\begin{bmatrix} g \\ \varepsilon \end{bmatrix} \sim \left( MNV \left( 0, \begin{matrix} G\sigma_a^2 & 0 \\ 0 & \sigma_\varepsilon^2 \end{matrix} \right) \right)$$

GRM accounts for Mendelian sampling

**Solve mixed model equation (Henderson, 1989)**

$$\begin{bmatrix} \acute{X}X & \acute{X}Z \\ \acute{Z}X & \acute{Z}Z + \lambda G^{-1} \end{bmatrix} \begin{bmatrix} b \\ g \end{bmatrix} = \begin{vmatrix} \acute{X}y \\ \acute{Z}y \end{vmatrix}$$

$$\begin{bmatrix} \hat{g}_1 \\ \hat{g}_2 \\ \hat{g}_3 \\ \hat{g}_4 \\ g \\ \hat{g}_6 \end{bmatrix} = \begin{bmatrix} -0.35 \\ -0.25 \\ 1.35 \\ -1.15 \\ 1.45 \\ 0.45 \end{bmatrix}$$

**(BLUP of Breeding values)**

**GEBVs($\hat{g}$)**

## Construction of G matrix

3 steps
1. Create a centered Z matrix
2. Create the cross product
3. Divide it by $\sum_{i=1}^{m} 2p_j (1-p_j)$

$$\mathbf{G} = \frac{\mathbf{Z}\mathbf{Z}^{\mathbf{T}}}{\sum_{i=1}^{m} 2p_j (1-p_j)}$$

*First G matrix, VanRanden (2008)*

$\sigma^2 a$ is additive variance = $\sum_{i=1}^{m} 2p_j (1-p_j$

# Difference Between RRBLUP and gBLUP

| RR BLUP | gBLUP |
|---|---|
| Marker as Predictors (BLUPs) | Predict Additive breeding Values directly |
| Markers are directly used | Markers are used to compute covariance matrix (Relationship Matrix) |
| Dimension of genetic effects is p x p | Dimension of effects is n x n, computationally easy |
| Does not account for the relationships between genotypes. Assuming genotypes are independent | Accounts relationship between genotypes, thus more appropriate for predictions or dissecting G x E (WHY?) |
| $var(y) = Z\acute{Z}\,\sigma_u^2 + I\sigma_e^2$ | $var(y) = ZG\acute{Z}\,\sigma_u^2 + I\sigma_e^2$ |
| $blup(\hat{u}) = (Z\acute{Z} + I\lambda)^{-1}Z(Y\acute{} - \hat{\mu})$ | $blup(\hat{u}) = (Z\acute{Z} + G^{-1}\lambda)^{-1}Z(Y\acute{} - \hat{\mu})$ |

# Pedigree BLUP (pBLUP Model)

$$\mathbf{y} = \mathbf{Z}a + \mathbf{\epsilon}$$
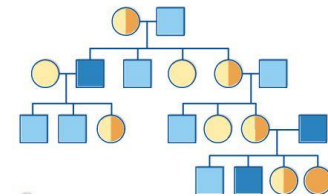
Solve mixed model equation (Henderson, 1989)

$$\begin{bmatrix} \acute{X}X & \acute{X}Z \\ \acute{Z}X & \acute{Z}Z + \lambda A^{-1} \end{bmatrix} \begin{bmatrix} b \\ g \end{bmatrix} = \begin{vmatrix} \acute{X}\,y \\ \acute{Z}y \end{vmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ . \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & 0.5 & 0.3 & 0.9 \\ 0.5 & 1 & -0.2 & 0.8 \\ 0.3 & -0.2 & 1 & 0.4 \\ 0.9 & 0.8 & .4 & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_p \end{bmatrix} + \mathbf{\epsilon}$$

n x 1          n x g          g x 1          n x 1

**n is # of individuals**
**g is # of breeding values**

$$\begin{bmatrix} g \\ \epsilon \end{bmatrix} \sim \left( MNV \left( 0, \begin{pmatrix} A\sigma_a^2 & 0 \\ 0 & \sigma_\epsilon^2 \end{pmatrix} \right) \right)$$

$$\begin{bmatrix} \hat{g}_1 \\ \hat{g}_2 \\ \hat{g}_3 \\ \hat{g}_4 \\ g \\ \hat{g}_6 \end{bmatrix} = \begin{bmatrix} -0.35 \\ -0.25 \\ 1.35 \\ -1.15 \\ 1.45 \\ 0.45 \end{bmatrix}$$

**(BLUP of Breeding values)**

$$\text{GEBVs}(\widehat{g})$$

## Construction of A matrix



Accounts for relationships (Coefficient of co-ancestry , $r = 2\theta_{xy}$)

$$y = \mathbf{Z}a + \epsilon$$

**A matrix**

# Dimension Reduction methods

$$y = \mu + \sum_k x_k \beta_k + \varepsilon$$

- ❖ Reason is to avoid inverse
- ❖ Suitable when predictors are correlated
- ❖ And calculations are tedious

$$x_{k = UDV^T}$$

*U = n x m orthogonal matrix*
*D = n x m diagonal matrix with singular values*
*V = n x n orthogonal matrix*

$$\beta_{OLS} = V\,D^{-1}U^T\,y$$

Regress phenotypes directly on eigenvectors or principle components

# Bayesian Approaches
## (Relax the Distribution Assumptions)

*posterior*

*Sampling distr.*

*prior*

$$p(\theta|y) = \frac{p(\theta, y)}{p(y)} = \frac{p(y|\theta)p(\theta)}{p(y)}$$

$y$ is observed data $y \sim p(\theta|y)$
$\theta$ is unknown parameter (random)

Bayesian Models differ with respect to the *Prior*

➢ The even distribution of genetic causation is not satisfactory

➢ The assumption of common variance does not imply that the effects of all markers

➢ RR BLUP over shrinkage of large marker effects

➢ Large effect and small effect QTLs (natural in breeding populations)

# Bayesian Ridge-Regression (BRR)

$$\widehat{\boldsymbol{\beta}} = (\boldsymbol{XX^T} + \boldsymbol{K^{-1}}\frac{\sigma_e^2}{\sigma_\beta^2})^{-1}(\boldsymbol{yX^T} + \boldsymbol{K^{-1}}\frac{\sigma_e^2}{\sigma_\beta^2})\beta_o$$

$$\lambda = \frac{\sigma_e^2}{\sigma_\beta^2}$$

$$\beta \sim N(0, \sigma_\beta^2)$$
$$\sigma_\beta^2 \sim \mathcal{X}^{-2}(df, s)$$

If we assume prior mean=0, $\beta$=0, then $\boldsymbol{K}$=I

$$\boldsymbol{\beta_{ridge}} = (\boldsymbol{X^TX})^{-1} + \lambda\boldsymbol{I})\, \boldsymbol{yX^T}$$

Prior for markers is given by Gaussian distribution, which differs from LASSO
Scaled inverted chi-square distribution
The only difference in prior information

# LASSO (Regularization)

$$\beta_{lasso} = argmin\left[\sum(y_i - \sum_{j=1}^{m} x_{ij}\beta_j)^2 + \lambda\sum_{i=1}^{j}\beta_j^2\right]$$

$$in\ LASSO\ |\beta_j|\ is\ absolute\ value\ (L_1\ Penality)$$
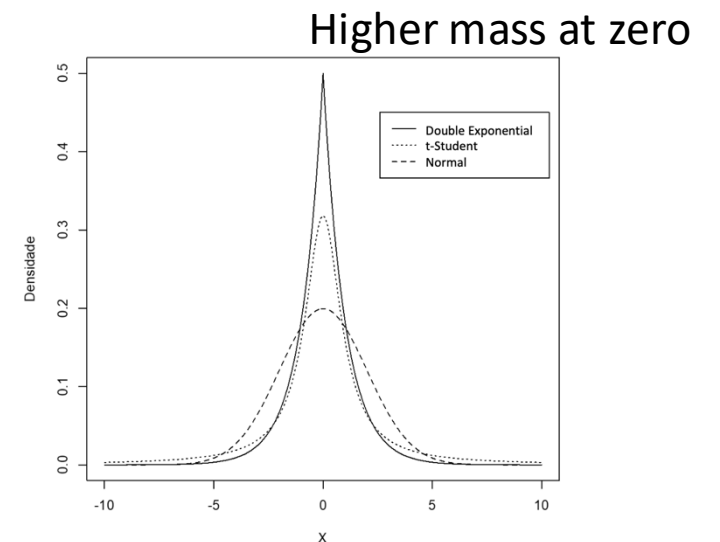$$however, in\ Rideg\ regression\ value\ is\ \beta_j^2\ (L_2\ Penality)$$

Type equation here.
In LASSO, shrinkage is stronger-sets some marker effects equal to 0

## Bayesian LASSO

$$\beta_{lasso} = argmin[(y - x\beta)^T(y - x\beta) + \lambda|\beta_j|]$$
$$= exp\left(-\frac{1}{2\sigma_e^2}(y - x\beta)^T(y - x\beta)\right)exp(\lambda|\beta_j I|)$$

Gaussian Sampling        Double exponential

In Bayesian Lasso, assumption of marker effects is Double exponential
Double exponential shrink more marker effects close to zero

Higher mass at zero

# Bayesian Approaches
## (Relax the Distribution Assumptions)

**Bayes A**

- Assumes marker-specific variance

- Utilizes an inverse chi-square on marker variances

- Shrinks tiny marker effects (variances) towards zero, and larger values survive.

**Bayes B**

- Fraction of markers with zero effect

- More realistic prior because we expect that some regions of the genome will carry no QTL

**Bayes C**

- Assumes t-distribution one with large variance for SNP fraction and other with small variance

Thank You

Questions

**Genomic Selection is Simply a tool to Supplement the Breeding Pipeline**