

HOS 6236 Molecular Marker Assisted Plant Breeding Fall 2017

Last Class:

Methods in Genomic Selection

Today's Class

The G-Matrix, GBLUP, and beyond the additive model

Pedigree and Relationship Matrices

- Why worry about the pedigree in genetic analyses?

Statistically, random genetic effects (i.e. BLUPs) are not independent and their matrix of correlations or co-variances (**G** or **A**) needs to be specified.

Genetically, it is important to consider information about relatives as they will share some alleles, and therefore their response is correlated.

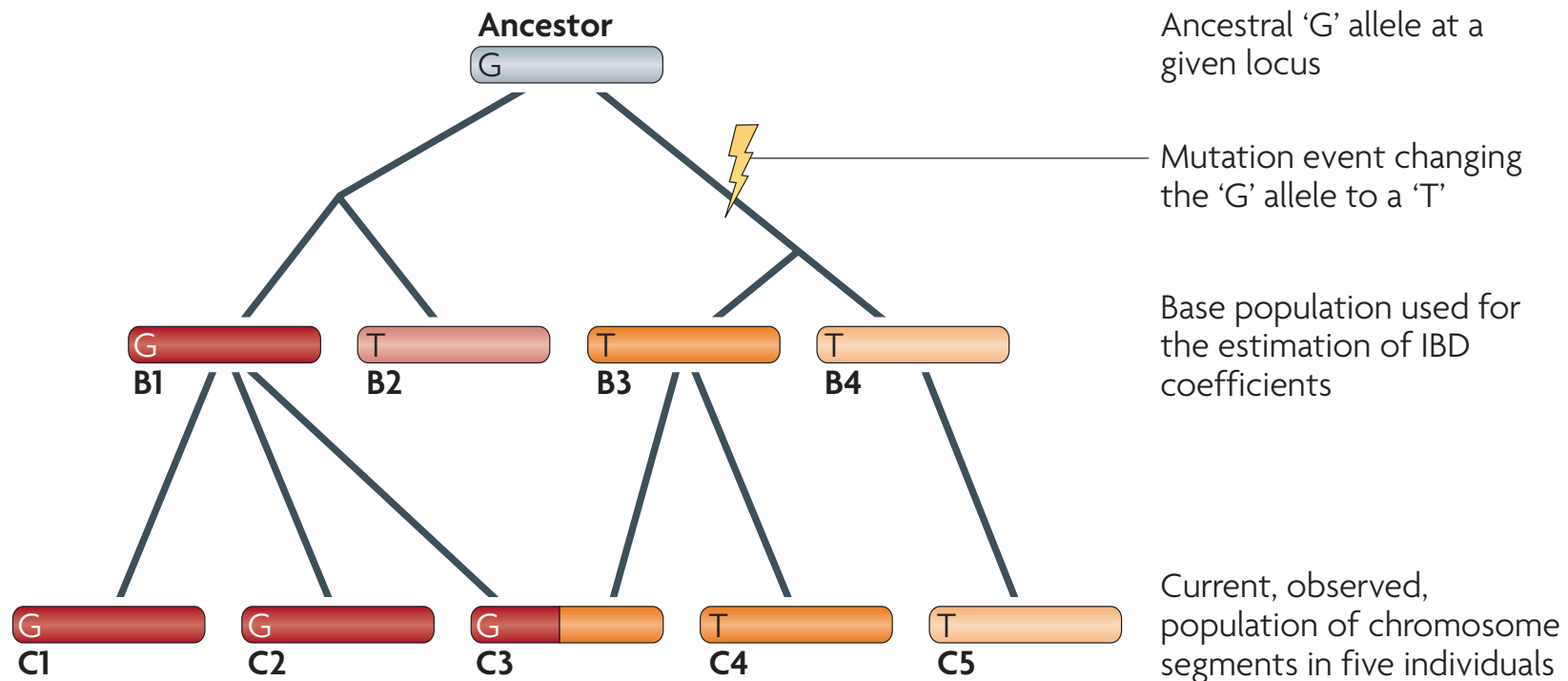
- How to incorporate this information?

Genetic relationships (pedigree) can be calculated using **genetic theory** (expected values) or **molecular information** (e.g. SNPs), and included into the linear mixed model by specifying a Relationship matrix

- Are there other benefits?

Many. It is a more **efficient** use of the information about individuals, but also genetic values of individual **not tested**, but with relatives tested, can be *predicted* and selected.

Pedigree and Relationship Matrices



Ancestor is the point of coalescent for current alleles C1-C5.

Identity by descent (IBD) of current alleles can be defined respect to B1-B4, thus G allele in C1-C3 are IBD

T allele in C4 and C5 are identity by state (IBS)

The chromosome segment C2 and C3 are IBS as well.

Powell et al 2010

Pedigree and Relationship Matrices

Construction / Check

- Pedigree information is associated with proper management and validation/check of data.
- Individuals need to be ordered by generation (e.g. parents need to be defined before progeny).
- All parents need to be defined in pedigree file (the inclusion of founder parents is optional).
- All individuals present in dataset (i.e. levels associated with pedigree file) need to be defined in pedigree file.
- Individuals can be defined as male or female parents (but this should be checked if is not biologically possible).

Genome Selection: Animal Model/ GBLUP

$$y_i = a_i + \varepsilon_i \quad (i=1,\dots,n)$$

$$a \sim MVN[0, G\sigma_a^2] \quad \varepsilon \sim MVN[0, I\sigma_\varepsilon^2]$$

Additive Relationship Matrix

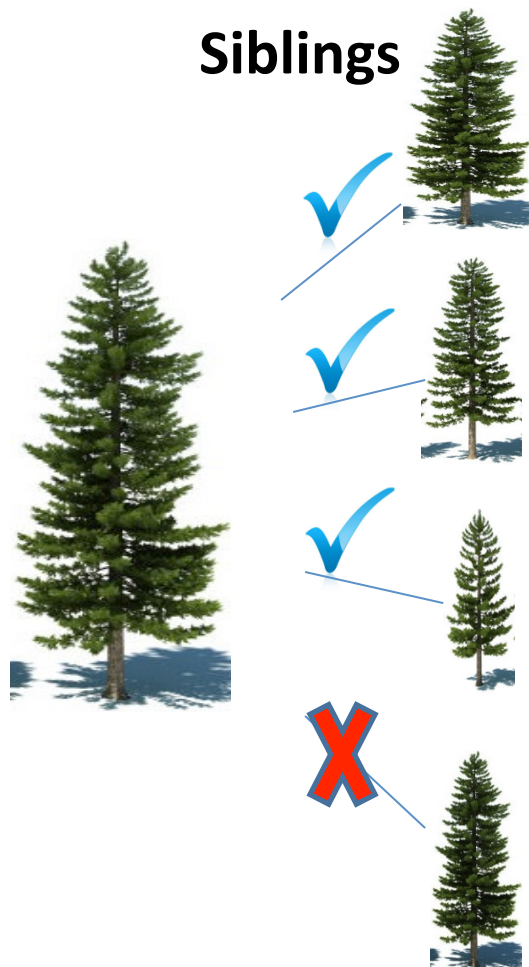
$$G_A \propto \textit{coancestry}$$

Genomic Relationship Matrix

$$G_M \propto XX'$$

$$\begin{array}{l} \text{TRN} \Rightarrow \\ \text{TST} \Rightarrow \end{array} \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \sim MVN \left[\mathbf{0}, \begin{pmatrix} \mathbf{G}_1 \mathbf{G}_{12} \sigma_a^2 & \mathbf{G}_{22} \sigma_a^2 & \mathbf{0} & \mathbf{0} \\ \mathbf{G}_{21} \sigma_a^2 & \mathbf{G}_{22} \sigma_a^2 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_1 \sigma_\varepsilon^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I}_2 \sigma_\varepsilon^2 \end{pmatrix} \right]$$

Use of relationships to estimate genetic effects



Pedigree Relationships	Marker Relationships	Benefit
0.5	0.38	Improve relationship estimation; By estimating Mendelian segregation
0.5	0.52	
0.5	0.63	
0.0	0.40	Estimate relationships inexistent by pedigree

From markers to G

Example:

$$\hat{a} = (X - P)\hat{g} \quad P = E(X) = 2p; \dots X - P = W$$

$$W = \begin{bmatrix} 1-p_i & 0-p_i & 1-p_i & 2-p_i \\ 2-p_i & 2-p_i & 0-p_i & 2-p_i \\ 2-p_i & 1-p_i & 1-p_i & 0-p_i \\ 0-p_i & 2-p_i & 2-p_i & 1-p_i \end{bmatrix} \quad \hat{g} = \begin{bmatrix} 0.24 \\ 0.02 \\ -0.08 \\ 0.14 \end{bmatrix} \quad \hat{a} = \begin{bmatrix} 0.44 \\ 0.80 \\ 0.42 \\ 0.02 \end{bmatrix}$$

- If the markers are capturing **all genetic variation**, then we can assume that (Van Raden, 2008):

$$a = Wg$$

- If we also assume:

$$V(g) = I\sigma_g^2$$

- Then we get:

$$V(a) = WW'\sigma_g^2$$

which is a covariance matrix for the individual breeding values “**a**”

From markers to G

- Ideally, we want to model this covariance using the same classical Linear Mixed Model framework, therefore, it would be desirable to have this matrix in terms of σ_a^2

$$\sigma_a^2 = \sum_{i=1}^{ALL_SNPs} 2p_i q_i \sigma_g^2 \longrightarrow \sigma_g^2 = \frac{\sigma_a^2}{\sum_{i=1}^{ALL_SNPs} 2p_i q_i}$$

- Gianola *et al.* (2009) showed that (under Hardy-Weinberg equilibrium):

$$V(a) = W W' \sigma_g^2$$

- If we recall then:

$$V(a) = \frac{W W' \sigma_a^2}{\sum_i 2p_i q_i} = G_A \sigma_a^2$$

by replacing σ_m^2 .

The Mixed Model Equation

Consider a model with block as fixed and variety as random effects.

$$\text{yield} = \mu + \text{block} + \text{variety} + \text{error}$$

$$y_{ij} = \mu + \alpha_i + g_j + e_{ij}$$

y_{ij} observation belonging to i^{th} treatment j^{th} block

α_i fixed effect of the i^{th} block

g_j random additive effect of the j^{th} variety, $E(g_j) = 0$, $V(g_j) = \mathbf{G}\sigma_g^2$

e_{ij} random error of the ij^{th} observation, $E(e_{ij}) = 0$, $V(e_{ij}) = \mathbf{I}\sigma^2 = \mathbf{R}$

$$\text{Cov}(g_j, e_{ij}) = 0$$

Where \mathbf{G} is the genomic relationship matrix calculated from the molecular markers (genomic data).

So, we keep using the same methods and models we already have!!!

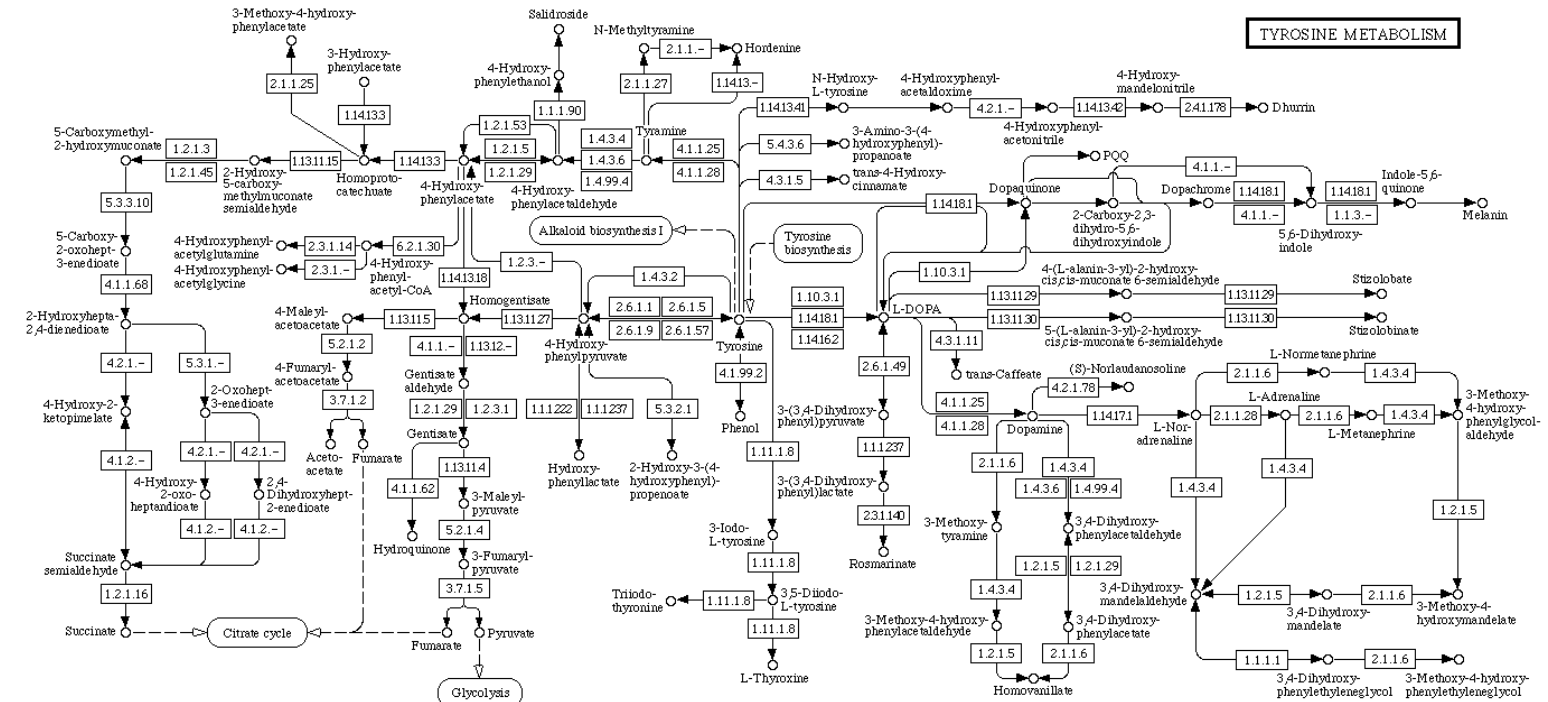
Non-Additive Effects Exist! - Evidence

Molecular evidence

Dominance at the gene level

Interaction of different genes

DO WE REALLY THINK EVERYTHING IS ADDITIVE!!



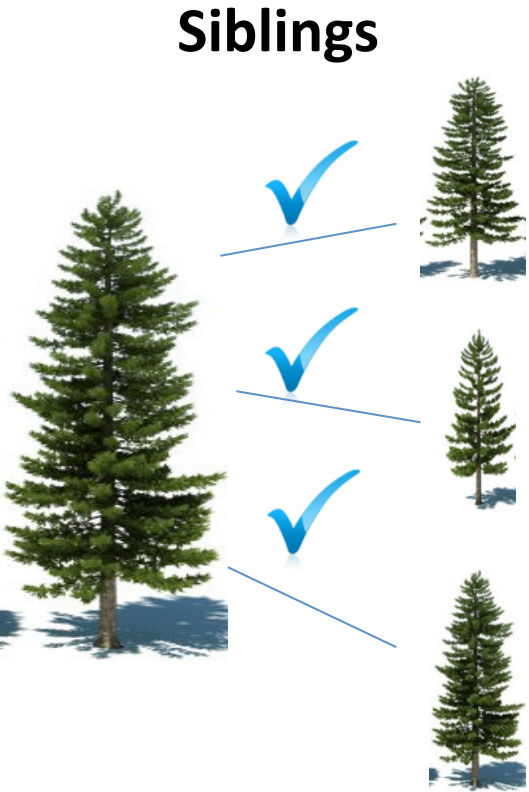
00350 10/9/98

Hybrid Vigor

<http://www.genome.jp/kegg/pathway.html>

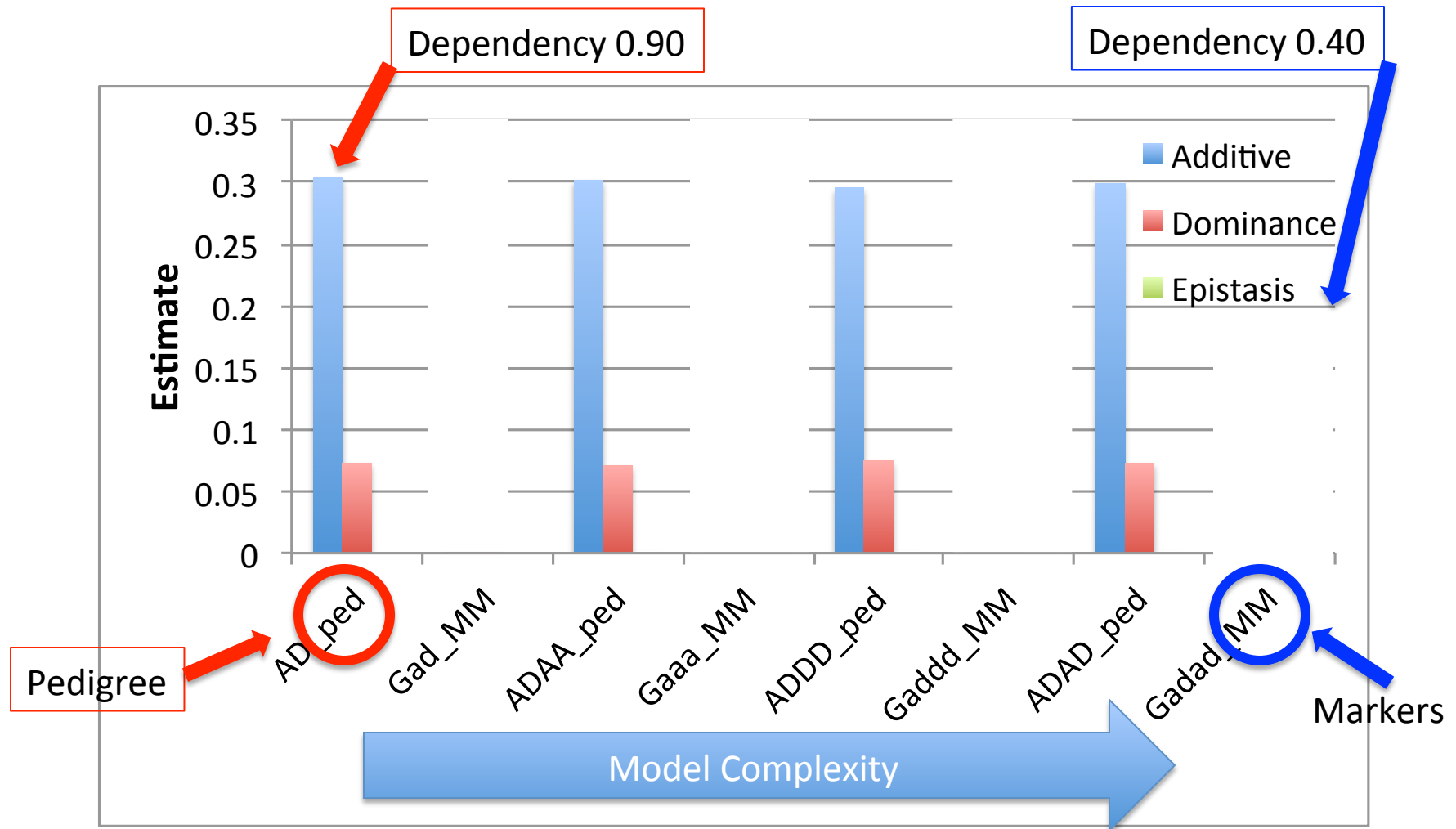
McKeand et al 2007

Markers-derived Relationships



Additive Pedigree Relationships	Dominance Pedigree Relationships
0.5	0.25
0.5	0.25
0.5	0.25

Estimates of Variance Components



Even when complexity of the model increases partition is the same with pedigree.

A different architecture is observed with marker-derived matrices

Prediction in 10-fold Cross Validation

Model	Cor(RBV,PBV)	MSE(RBV,PBV)	Top10%RankCor
A_ped	0.640	1335.800	0.17
Ga_MM	0.670	1291.800	0.34
ADAD_ped	0.727	657.258	0.16
Gadad_MM	0.872	108.240	0.37
ADDD_ped	0.732	638.464	0.18
Gaddd_MM	0.873	151.199	0.32

RBV → Breeding Value using all data within model

PBV → Breeding Value in cross validation

Top10%RankCor → correlation between R-ranking and P-ranking
for top 10% performance

Estimates of Variance Components

GENOMIC SELECTION

Unraveling Additive from Nonadditive Effects Using Genomic Relationship Matrices

Patricio R. Muñoz,^{*,†,1,2} Marcio F. R. Resende, Jr.,^{†,‡,1} Salvador A. Gezan,[†] Marcos Deon Vilela Resende,^{§,}
Gustavo de los Campos,^{††} Matias Kirst,^{†,‡†} Dudley Huber,[†] and Gary F. Peter^{†,‡†,2}**

^{*}Agronomy Department, [†]School of Forest Resources and Conservation, [‡]Genetics and Genomics Graduate Program, and
^{‡†}University of Florida Genetics Institute, University of Florida, Gainesville, Florida 32611, [§]EMBRAPA Forestry, Estrada da Ribeira,
Colombo, PR 83411-000 Brazil, ^{**}Department of Forest Engineering, Universidade Federal de Viçosa, Viçosa, MG 36571-000
Brazil, and ^{††}Department of Biostatistics, University of Alabama at Birmingham, Birmingham, Alabama 35294

The Mixed Model Equation

Consider a model with block as fixed, additive effect and dominance effect as random effects.

$$y_{ij} = \mu + \alpha_i + a_j + d_j + e_{ij}$$

y_{ij} observation belonging to i^{th} treatment j^{th} block

α_i fixed effect of the i^{th} block

a_j random additive effect of the j^{th} variety, $E(a_j) = 0$, $V(a_j) = \mathbf{G}\sigma_a^2$

d_j random dominance effect of the j^{th} variety, $E(d_j) = 0$, $V(d_j) = \mathbf{D}\sigma_d^2$

e_{ij} random error of the ij^{th} observation, $E(e_{ij}) = 0$, $V(e_{ij}) = \mathbf{I}\sigma^2 = \mathbf{R}$

$$\text{Cov}(a_j, e_{ij}) = 0 ; \text{Cov}(d_j, e_{ij}) = 0 ; \text{Cov}(a_j, d_j) = 0$$

Where \mathbf{G} is the additive genomic relationship matrix calculated from the molecular markers (genomic data).

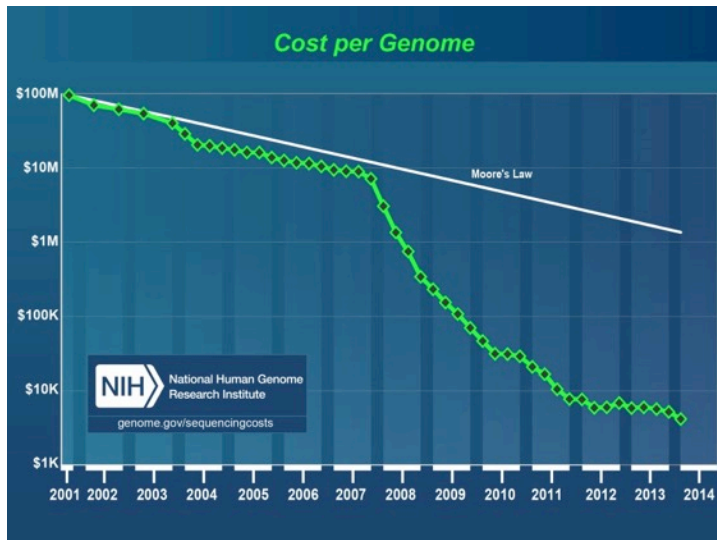
Where \mathbf{D} is the dominance genomic relationship matrix calculated from the molecular markers (genomic data).

AND we keep using the same methods and models we already have!!!

Family-Bulk Genome-wide Prediction

Goal: perform predictions using bulk genotyping and phenotyping

In many species “cultivars” corresponded to a population of genotypes (e.g. forages).



Even when genotyping cost has decreased significantly using NGS, it is still too high for most breeding programs.

Model Species and Method

Pine breeding (CCLONES) population with 71 families and genotyped with Chip for ~4,700 SNPs markers

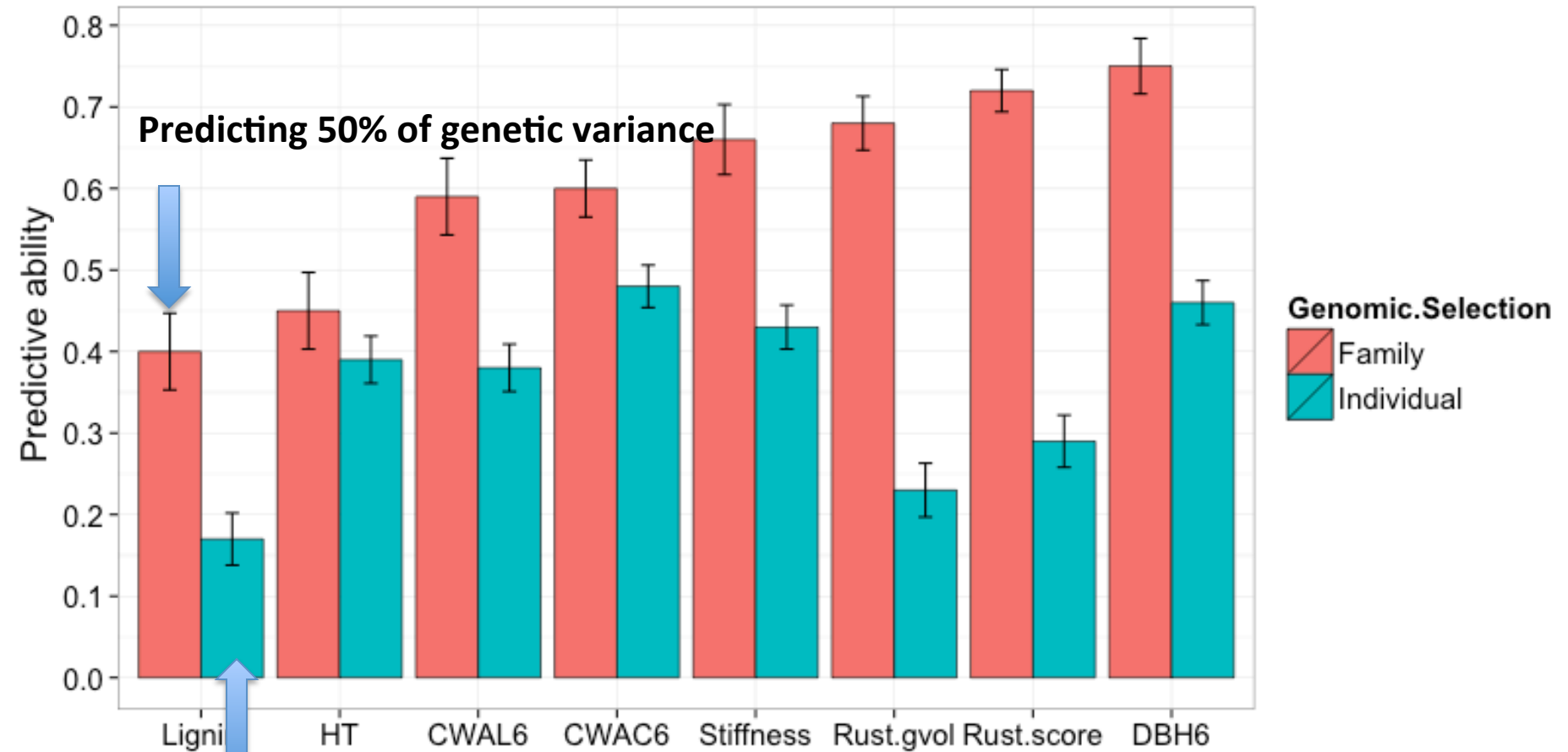


Phenotype: family phenotype mean

Genotype: family allele frequency

Family-Bulk Genotyping vs Traditional GWS

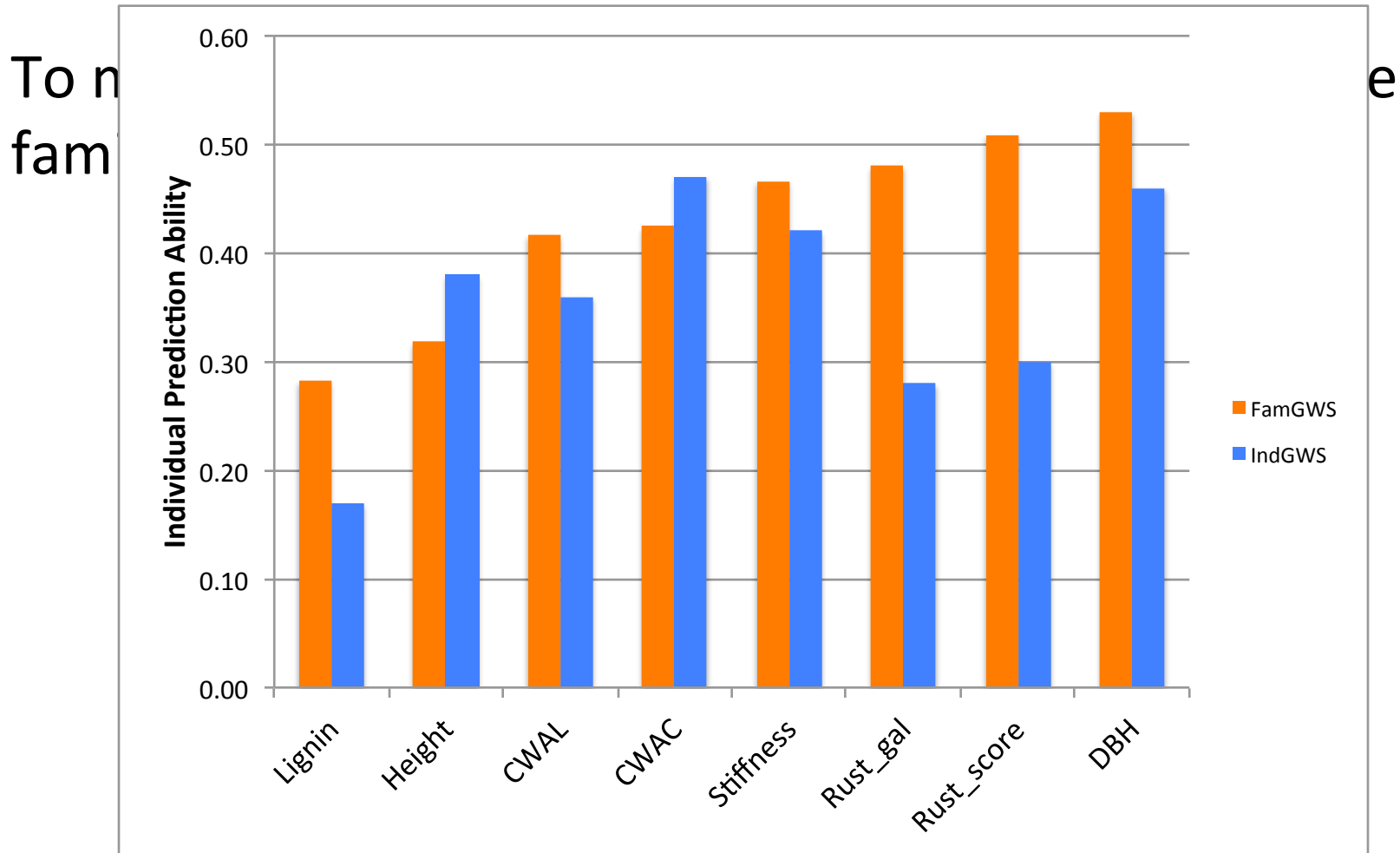
Pine



Capitalize in information at the mean level, which is more precise than at the individual level. But...

Predicting 100% of genetic variance

Family-Bulk Genotyping vs Traditional GWS Pine

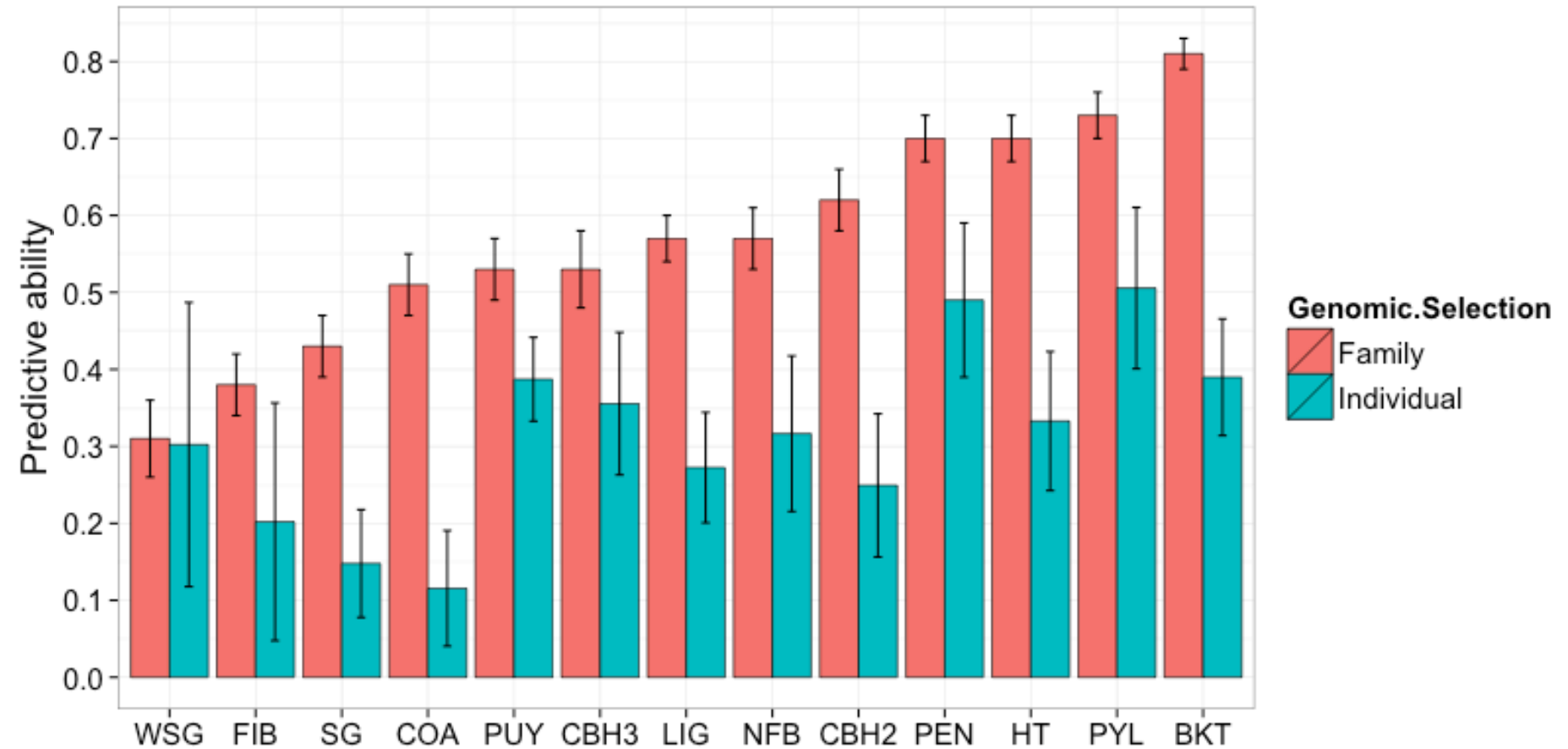


Still better in 6 out of 8 traits

Family-Bulk Genotyping vs Traditional GWS

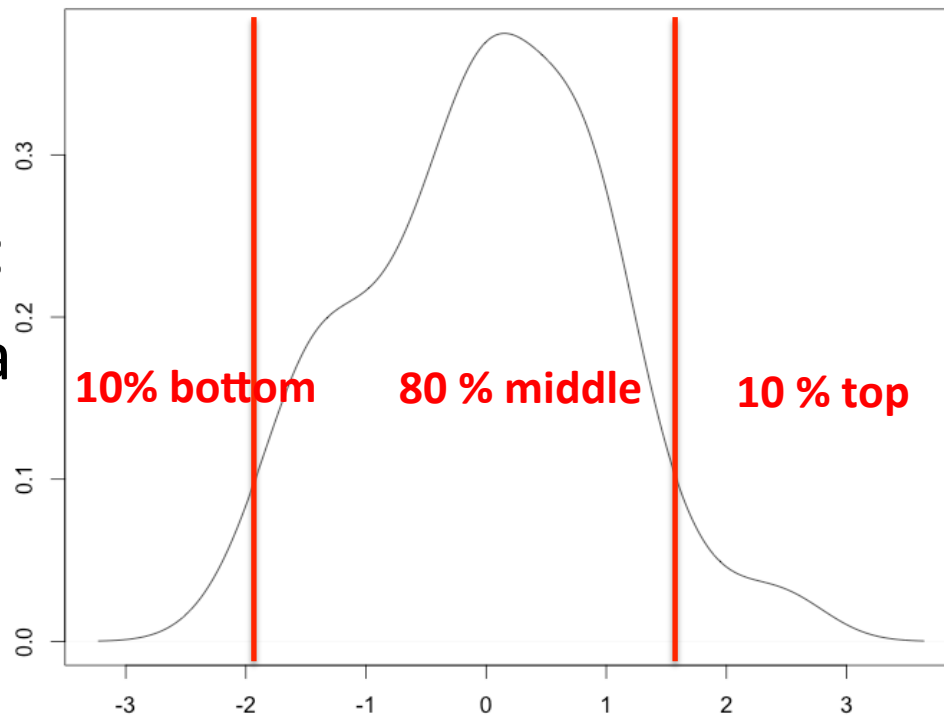
Eucalyptus

Eucalyptus breeding population with 68 families and genotyped with sequence capture for 500,000 SNPs markers.

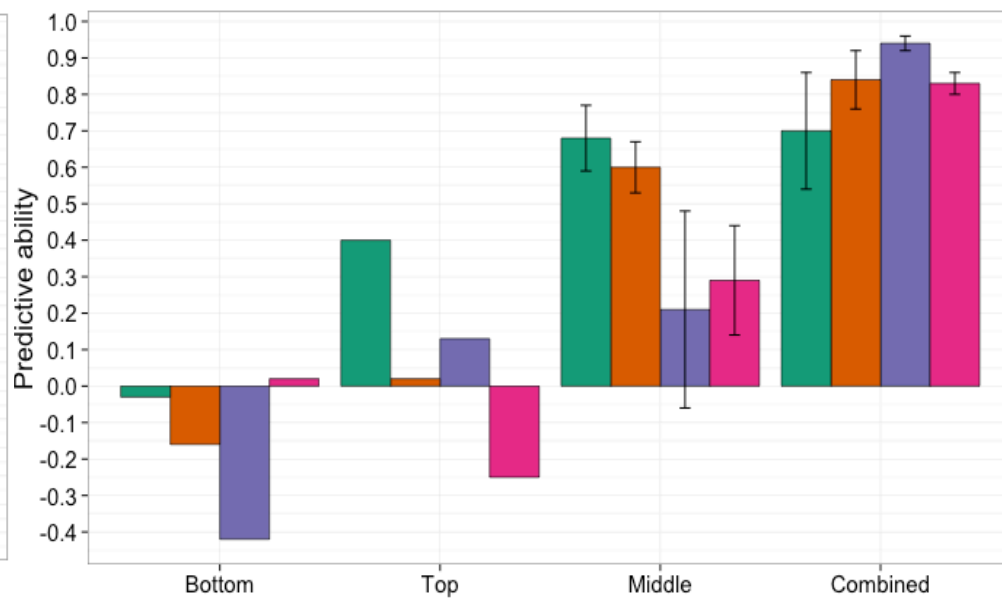
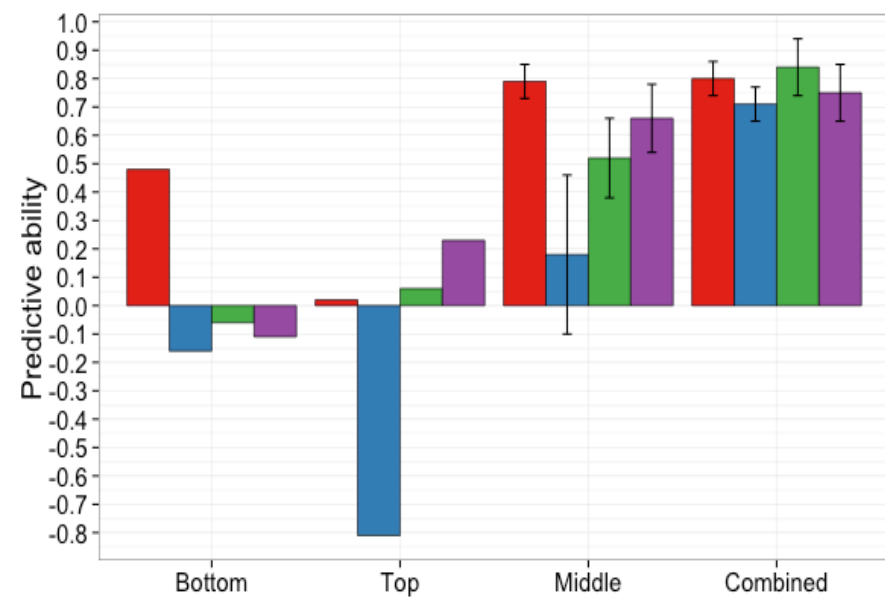


Still better in 8 out of 13 traits at the individual level

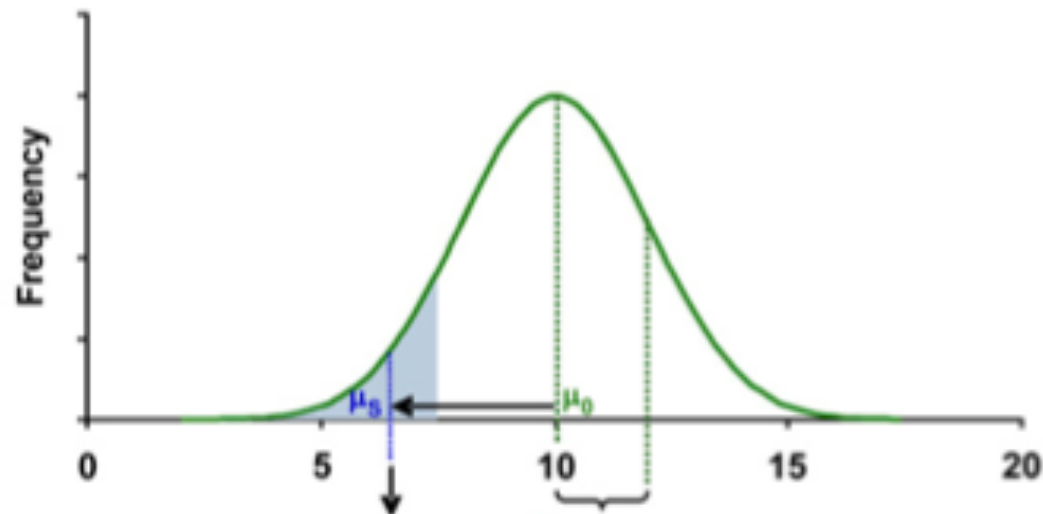
Challenges:
Very contra



ation folds



The Goal!!



Genetic Gain = $\Delta G = h^2 \sigma_P i / L$ → length of cycle interval (usually 1 generation)

heritability

Selection Intensity

proportion of population selected to produce the next generation

phenotypic variability in population

$$\sigma^2_{\text{Genotype}} - \sigma^2_{\text{Dominance}} - \sigma^2_{\text{Epistasis}} = \sigma^2_{\text{Additive}}$$

$$\sigma^2_{\text{Genotype}} + \sigma^2_{\text{Environment}} + \sigma^2_{\text{GxE}} + \sigma^2_{\text{error}} = \sigma^2_P$$

Breeder's Equation

$$\text{Genetic_Gain} = \Delta G = \frac{h^2 * \sigma_p^2 * i}{L}$$

h^2 = **Narrow-sense heritability** = is the portion of phenotypic variance due to additive genetic variation.

σ_p = **phenotypic standard deviation** = phenotypic variability in the population

i = **intensity of selection**

L = length of breeding cycle interval

Now, think what of the above parameters is impacted by GS!!