



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

Basics: Mean, Variances and Linkage Disequilibrium

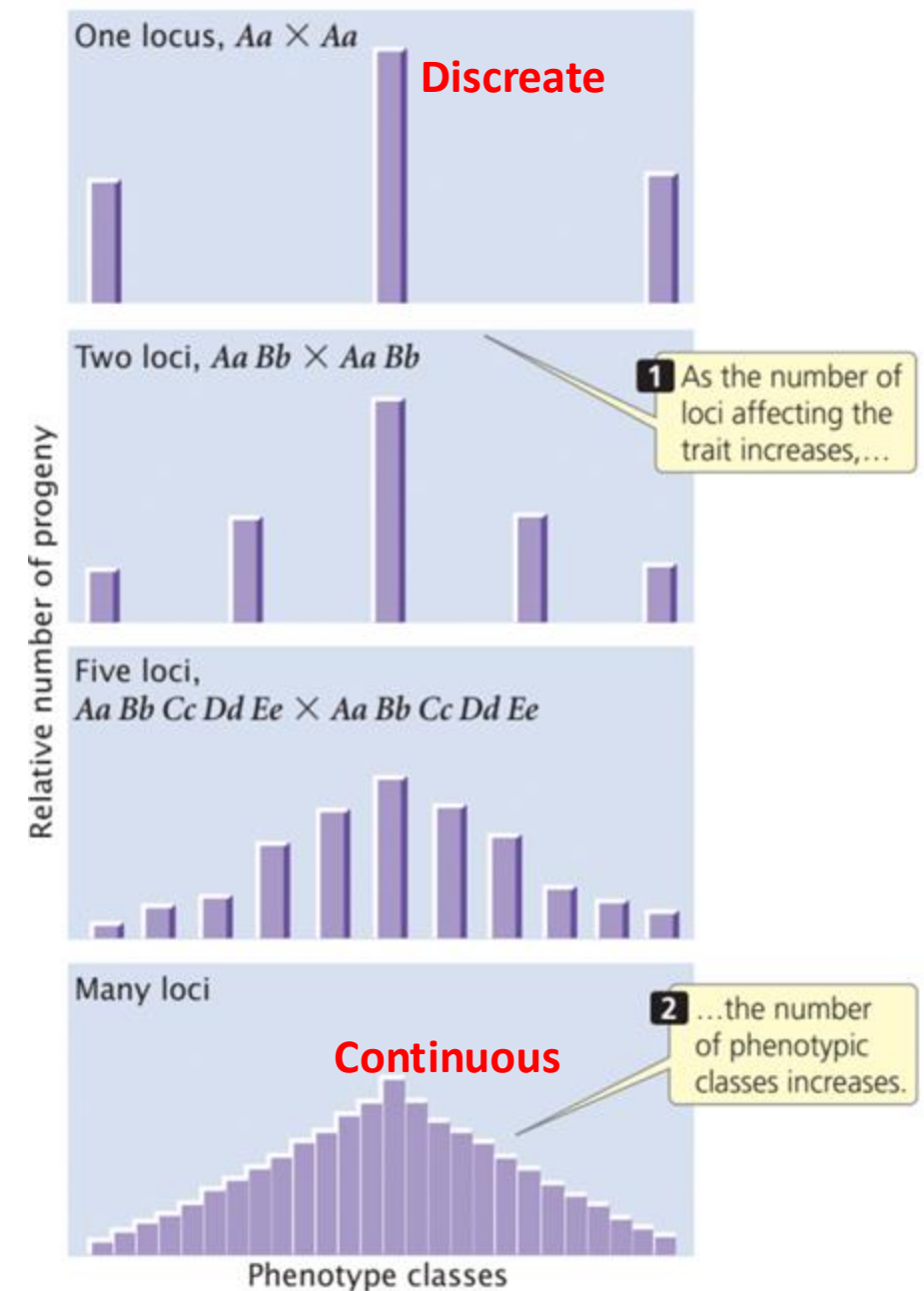
**Module 2
August 5, 2025**

Waseem Hussain and Mahender Anumalla
Rice Breeding Innovations Platform
IRRI

Quantitative Genetics

- Deals with complex traits determined by genes and environment
- Vary continuously and highly uninformative
- Effected by alleles at multiple loci
- Inherited as Mendelian Traits

- *Genes influencing Quantitative traits called as **QTLs***
- *Same principles of inheritance, however, more genes take part in the determination of quantitative characteristics.*



Measurement of Quantitative Traits

Quantitative Traits are Measured at the Population Level using Statistical Tools



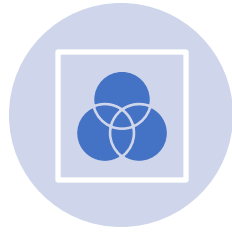
MEAN

Center of distribution
of data



VARIANCE

Spread of the data



CO-VARIANCES

Co-vary together



CORRELATIONS

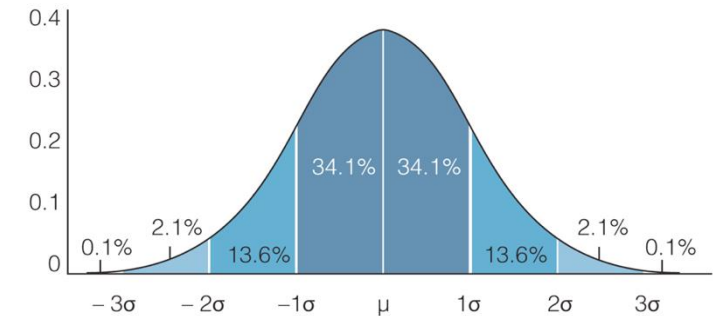
Relationship
between variables



REGRESSION

Relationship
between variables

Distribution of Variance



$$\text{Mean} = \frac{1}{n} \sum_i^n x_i$$

$$\text{Variance} = \frac{1}{n} \sum_i^n (x_i - \mu)^2$$

$$\text{COV}(X, Y) = \sum p_{xy} (x_i - \mu_x) (y_i - \mu_y)$$

Basic Quantitative Genetics Model

$$\text{Phenotype}(P) = \text{Genotype}(G) + \text{Environment}(E) + \text{Residual}(e)$$

or

$$y_{ijk} = \mu_i + G_i + E_j + e_{ijk}$$

Trait value observed in Individual

Average phenotypic value over a range of environments $G = E(P)$

Environments Individuals evaluated

Residual: Uncontrolled

The diagram illustrates the components of the Basic Quantitative Genetics Model. At the top, the general equation is given: $\text{Phenotype}(P) = \text{Genotype}(G) + \text{Environment}(E) + \text{Residual}(e)$. Below this, an alternative notation is shown: $y_{ijk} = \mu_i + G_i + E_j + e_{ijk}$. Arrows indicate the scope of each term: a dashed black arrow points from $\text{Phenotype}(P)$ to the text 'Trait value observed in Individual'; a dashed blue arrow points from G_i to 'Average phenotypic value over a range of environments $G = E(P)$ '; a dashed orange arrow points from E_j to 'Environments Individuals evaluated'; and a dashed blue arrow points from e_{ijk} to 'Residual: Uncontrolled'.

Dissecting Basic Quantitative Genetics Model

Genotype	A_2A_2	0	A_2A_1	A_1A_1
Genotype values	MP-a	MP	MP+d	MP+a
Code genotype	-a	0	d	a
Frequency	q^2	$2pq$		p^2

$$\text{Mid parent (MP)} = \frac{A_1A_1 + A_2A_2}{2} \text{ (Half difference between homozygotes)}$$

$$\text{Mean} = q^2(MP - a) + 2pq(MP) + p^2(MP + a)$$

a is additive effect

d is dominance

d=a, complete dominance

d=0, no dominance

0<d<a, partial dominance

d>a overdominance

$$\text{Mean} = MP + \text{Homozygote contribution} + \text{Heterozygote contribution}$$

Homozygote contribution

For RILs, d=0, thus

$$\text{Mean} = MP + a(p - q)$$

Multiple loci

$$\text{Mean} = \sum_{i=1}^k a_i(p_i - q_i) + 2 \sum_{i=1}^k p_i q_i d_i$$

where mean is summation over all the loci and k is number of loci

Average Effect and Why

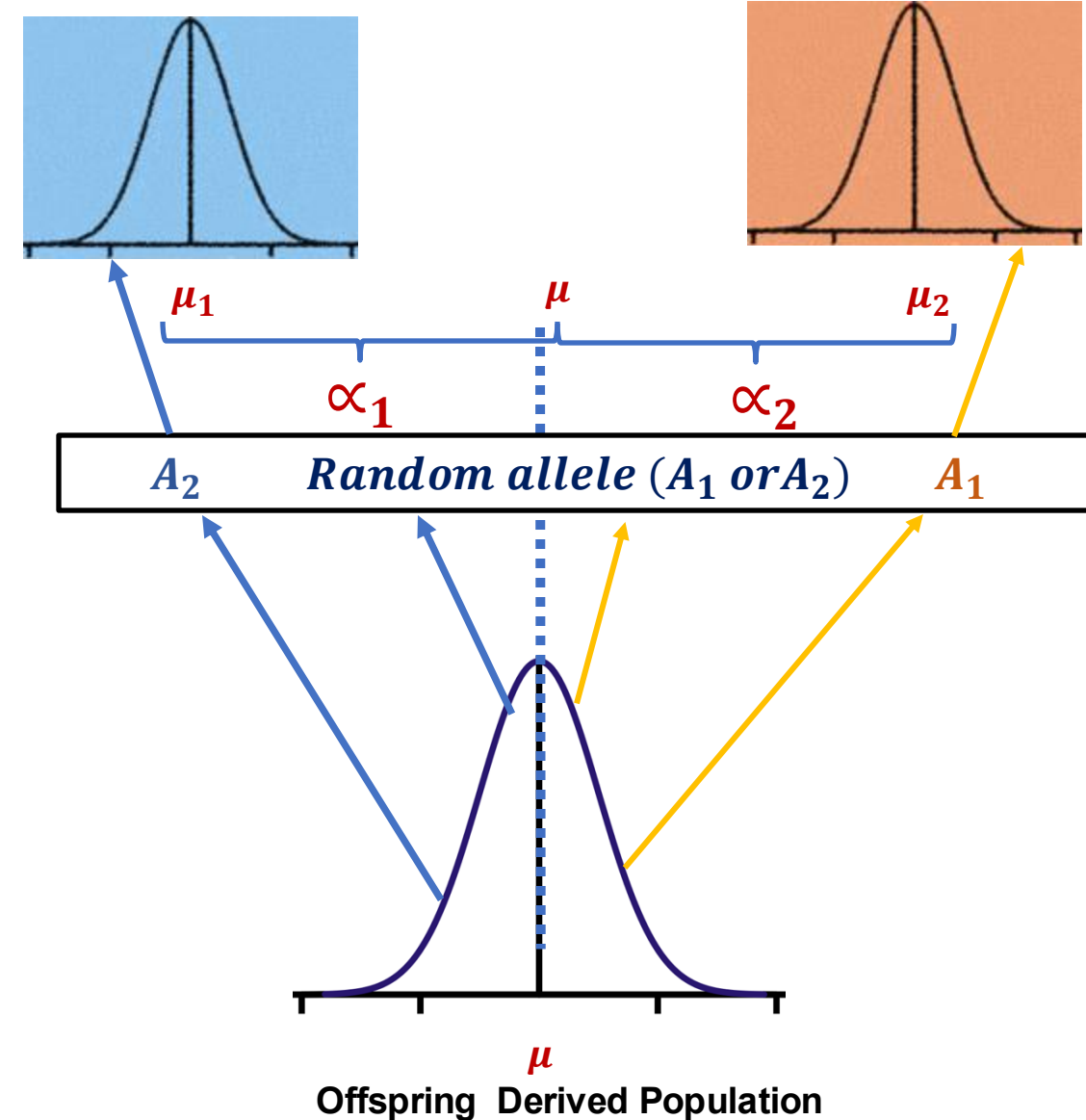
- Parents **pass a single allele** to offspring at a given locus, and hence **PART of genotypic value** is passed, which is determined as the **Average Effect**.
- It is a deviation of individuals from the population mean who received a particular allele from one parent, and the other allele comes as random

$$\alpha_x = \mu - \mu_x$$

$$\alpha_1 = q(a + d(q - p))$$

$$\alpha_2 = -p(a + d(q - p))$$

$$\alpha = -p(a + d(q - p))$$

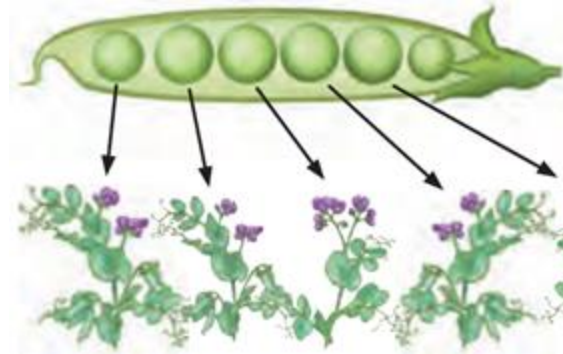


Average Effect and Why

- Alleles, not genotypes, are inherited
- Link the genotype values from one generation to another
- Average effect is population, depends upon allele frequency
- Average effect of random allele is 0?

Mendelian Trait

*an allele causes
certain phenotype*



Quantitative Trait

*an allele has a large
average effect*



Note

We don't measure individual loci average effects; we use **Regression** to get the effects.

Breeding value (Additive Genetic Value)

Deviations expected from the offspring of a particular genotype when it is mated with another individual of the same genetic worth

Genotype = *mean* + *additive* + *dominance* + *epistasis* + *error*

$$G_{ij} = \alpha_i + \alpha_j \quad \text{A part of G (additive) is only transmitted}$$

$$A = \sum_{k=1}^n \left(\alpha_i^{(k)} \alpha_j^{(k)} \right) = \sum_{k=1}^n a_{ik}$$

$\alpha_i^{(k)}$ and $\alpha_j^{(k)}$ = is the effect of allele *i* and *j* and locus *k*

A is called Breeding Value or Additive Genetic Value and

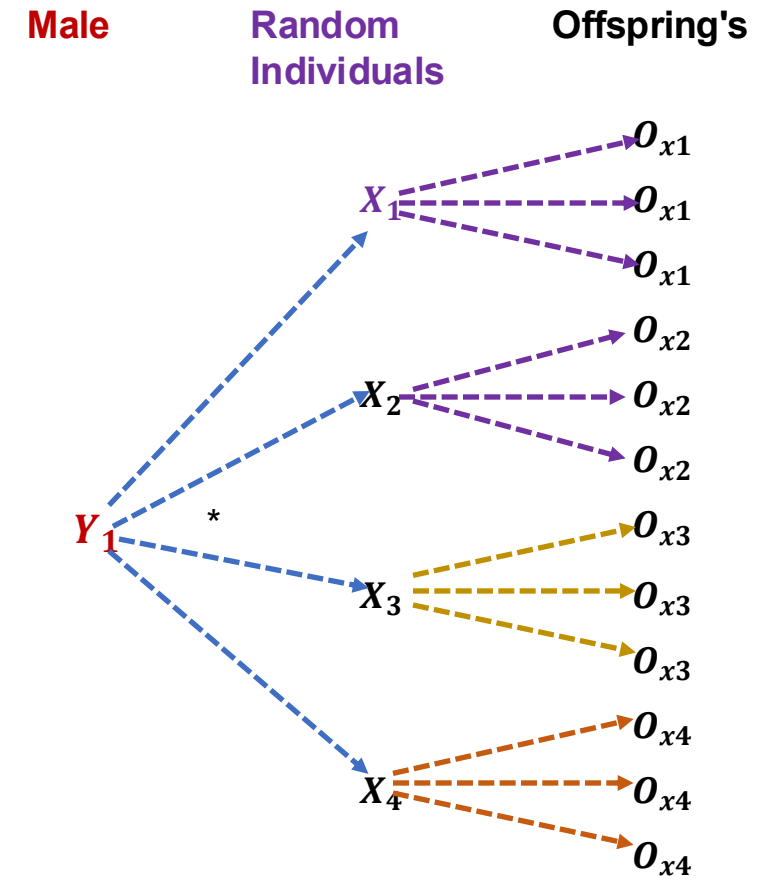
a_{ik} is additive genetic effect for *k* loci

Genetic value (Additive) = Sum of average effect of alleles

$$\begin{aligned} \text{Mean BV of all genotypes} \\ = 2pq(p\alpha + q\alpha - p\alpha - q\alpha) = 0 \end{aligned}$$

The breeding value is population specific
(depends upon allele frequency)

$$\begin{aligned} y &\sim N(\mu, \sigma_a^2 + \sigma_e^2) \\ a &\sim N(0, \sigma_a^2) \\ e &\sim N(0, \sigma_e^2) \end{aligned}$$



$$BV = 2 (\text{mean progeny } A_1A_1 - \text{Mean (Population)})$$

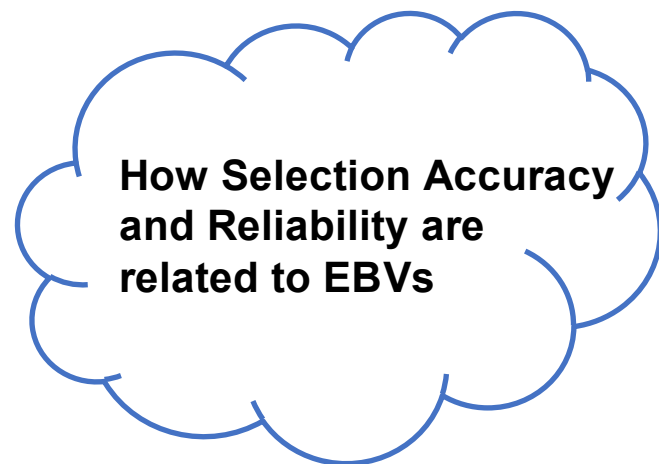
$$E(BV) = \frac{p_1 + p_2}{2} + \delta$$

δ is Mendelian Sampling (variation between offspring between same parent)

Breeding value (Additive Genetic Value)

Expresses the value transmitted from parent to offspring

- We don't observe breeding values, we estimate it
- Breeding values (EBVs) are estimated from phenotypic (True values) and relative information using a mixed model (Regression); BLUP
- Conditional on Phenotypic value ($E(a|y)$)



$$E(a|y) = \underbrace{\text{cov}(a, y)[\text{var}(y)]^{-1}}_{\text{Regression coefficient (b)}} \underbrace{(y - \mu)}_{\text{Adjusted mean}}$$

Breeding value is adjusted mean weighted by the coefficient that is heritability

Under a simple genetic model

$$y_i = \mu + a_1 + \varepsilon_i$$

$$b_{(a,y)} = \frac{\text{Cov}[y_i, a_i]}{\text{Var}[y_i]} = \frac{\sigma_a^2}{\sigma_y^2} = h^2$$

$$\hat{a}_i = h^2(y_i - \mu)$$

EBV is adjusted mean multiplied by the heritability

Variances and Co-variances

$$\text{Phenotype}(P) = \text{Genotype}(G) + \text{Environment}(E) + \text{Residual}(e)$$

$$V_P = V_g + V_e + V_{ge}$$

$$V_a + V_d + V_i$$

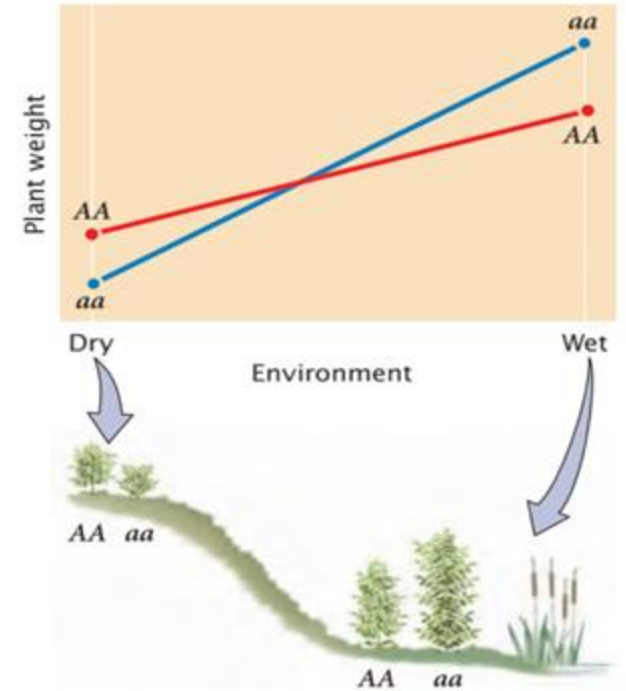
Heritability is used to measure the proportion of variance
The proportion of variation that is due to genetics is
heritability

$$\text{Broad Sense heritability}(H^2) = \frac{V_g}{V_p}$$

$$\text{Broad Sense heritability}(h^2) = \frac{V_a}{V_g}$$

$$\text{Cov}(a_x, a_y) = A_{xy}\sigma_a^2$$

where, a_x and a_y are the breeding BV's and σ_a^2 is the additive variance
And A is the additive genetic relationship related to IBD (*identity-by-descent*)



Dominance Deviation

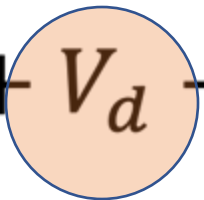
$$V_P = V_g + V_e + V_{gxe}$$

V_P : total variance

V_g : Genotype variance

V_e : environmental

V_{gxe} : variance due to interactions

$$V_g = V_a + V_d + V_i$$


Dominance deviations are obtained by taking the difference between the genotypic value and the breeding value for each genotype.

Points to Note:

- When $d=0$, a genotype has a breeding value that is identical to its genotypic value
- Without dominance, genotypic values in progeny as measured by the breeding value can be predicted perfectly from the combination of average effects.

Covariance Between Individuals

$$Cov(X, Y) = 2f_{xy}\sigma_a^2 (\text{assuming Dominance} = 0)$$

Where f_{xy} is the coefficient of co-ancestry between two individuals and σ_a^2 is the additive genetic variance

Example of Variance and Co-variance

Genotypes	1	2	3	4
1	Var (1)	cov(1,2)	cov(1,3)	cov(1,4)
2		Var (2)	cov(2,3)	cov(2,4)
3			Var (1)	cov(3,4)
4				Var (4)

Note:

- The diagonal represents variance and off-diagonal covariance
- If genotypes are independent off-diagonal elements=0
- Non-zero elements of off-diagonal reflect use of **relative information for BLUP estimation**

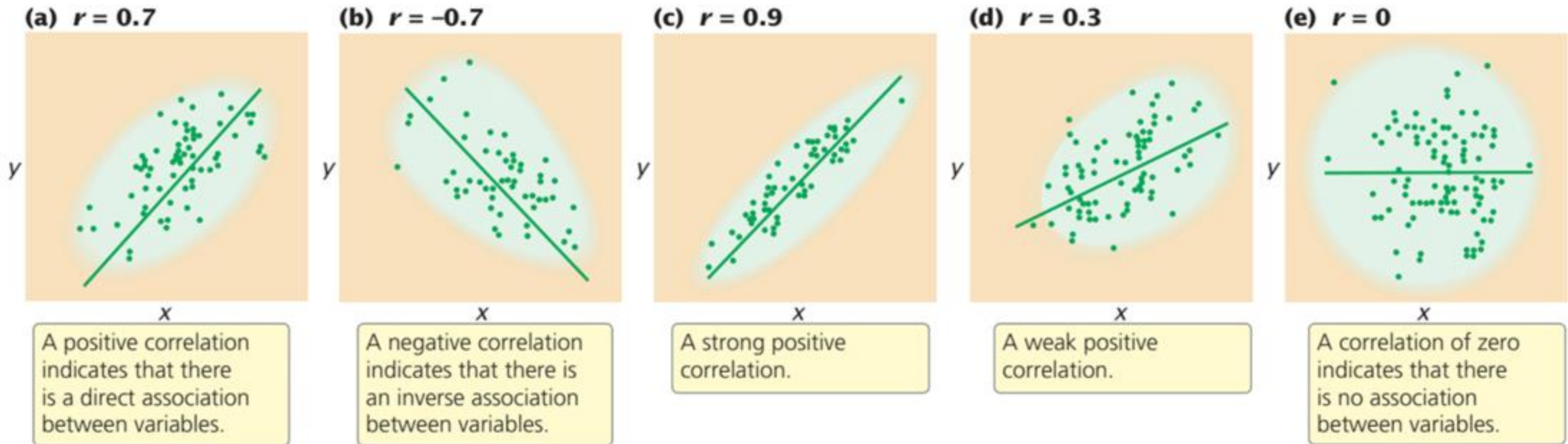
$$\text{var}(u) = A * \sigma_a^2$$

Elements of **the additive relationship matrix (A)** has the elements given by the $u(x, y) = 2f_{xy}$

Correlations

(with more than one character)

Relationship between two variables



24.11 The correlation coefficient describes the relation between two or more variables.

Linkage disequilibrium is simply correlation between the SNPs

Figure adopted from Book *Conceptual approaches to Genetics* by Benjamin A. Pierce

Regression

- Correlation provides information only about the strength and direction of association between variables.
- Whether two variables are associated.
- Basically, to predict the value of one variable, given a value of the other.

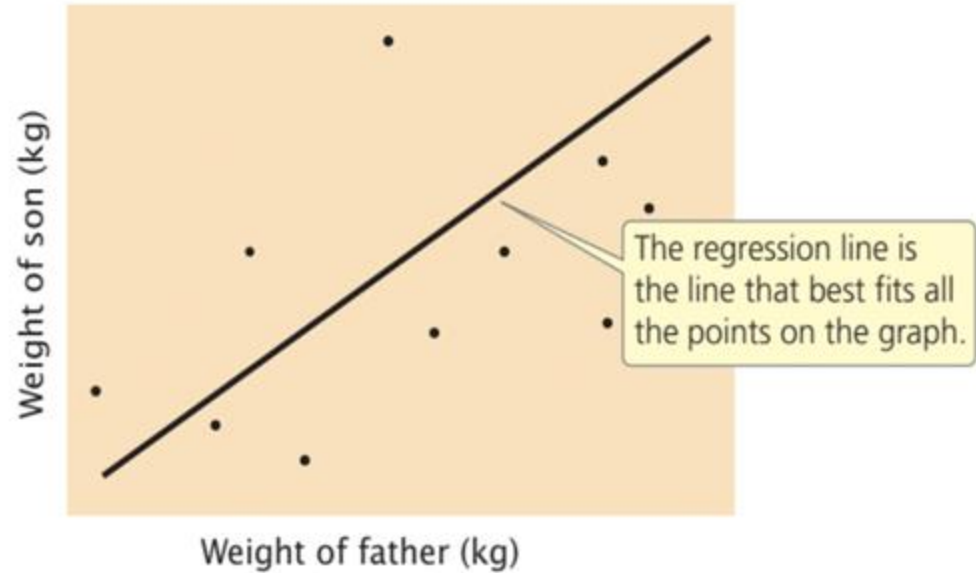
Dissect the Regression line

$$y = a + bx$$

- x and y represent the x and y variables
- a is intercept
- B is slope or regression coefficient.

$$b = \frac{\text{cov}(x, y)}{\text{variance}(x)}$$

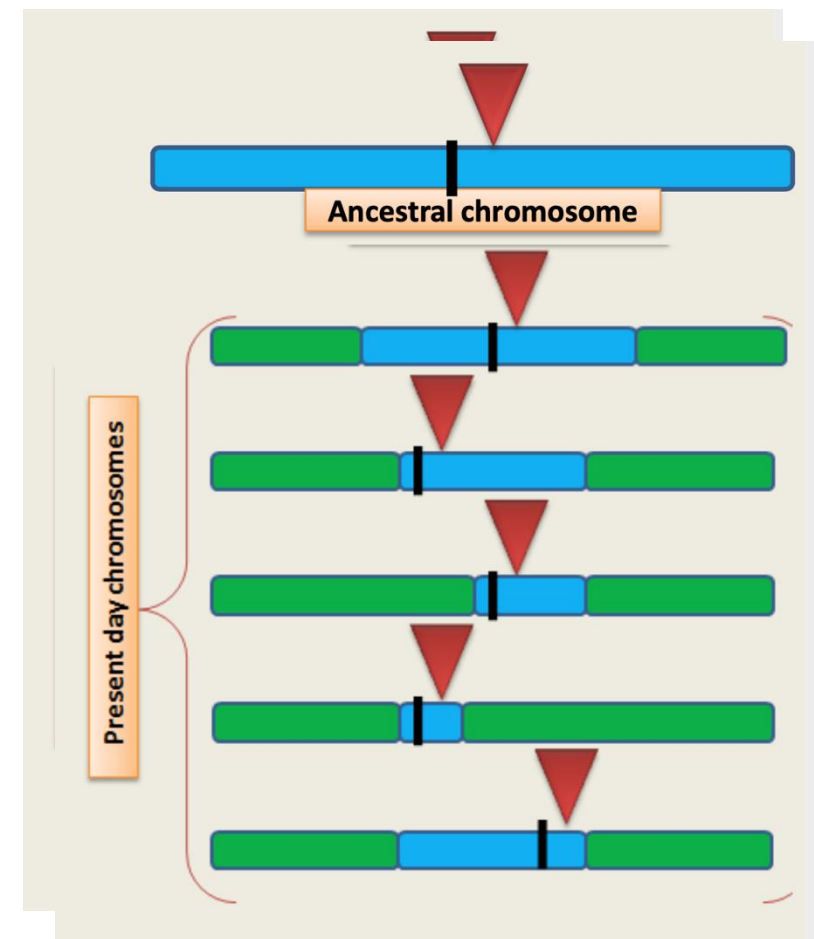
Regression coefficient indicates how much increases, on average, per increase in x



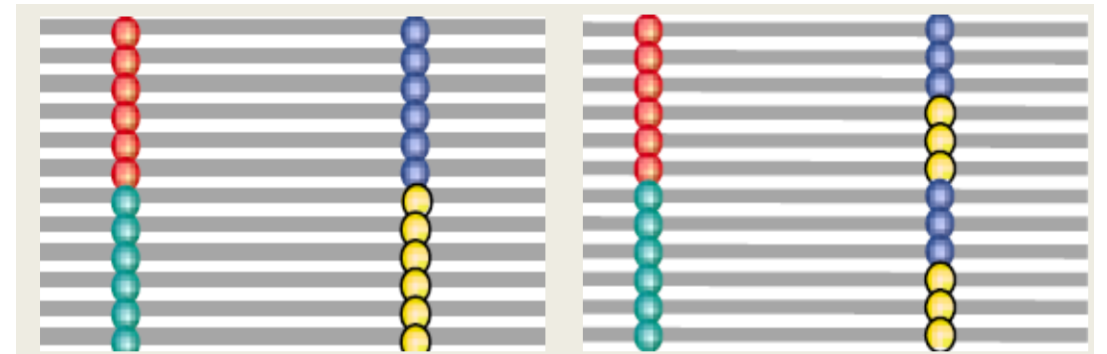
24.13 A regression line defines the relation between two variables. Illustrated here is a regression of the weights of fathers against the weights of sons. Each father–son pair is represented by a point on the graph: the x value of a point is the father's weight and the y value of the point is the son's weight.

Linkage Disequilibrium

- Non-random association of alleles at adjacent locus
- Closer the markers higher is the LD
- The resolution with which the QTL can be mapped is function of LD.
- Very Important in Genomic Selection to determine the prediction accuracy and number of markers.



LD across the Historical Combinations



How to Measure LD

Mathematically,

Linkage equilibrium **$PAB = P_A \times P_B$**

Linkage disequilibrium **$PAB \neq P_A \times P_B$**

where A and B are alleles at two different loci,

PAB is the frequency of haplotypes having both alleles at the two loci,

P_A and P_B are the frequency of haplotypes having only A allele and B allele, respectively.

D ranges from 0 -1 (At equilibrium, D= 0)

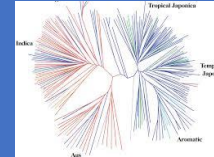
$$D' = \frac{D}{|D|}$$

$$r^2 = \frac{D^2}{p_A p_B (1-p_A)(1-p_B)} \quad (\text{Hill and Robertson, 1968})$$

D is determined by the range of allele frequency

r² is Pearson's correlation coefficient, and is the most relevant LD measurement (0-1)

Factors Affecting the LD



01

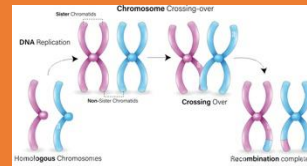
Germplasm

Longer in narrow base germplasm. 1 kb in cultivated and 2kb in landraces.

Recombination

Recombination decreases LD.
LD decays high with each generation or mating

02

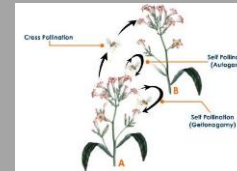
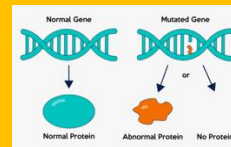


$$D_{t+1} = (1 - c)D_t$$

Mutations

Creation of LD and is eroded by Recombination.

04



03

Mating System

Longer in self-pollinated crops.
Less recombination in Self pollinated
LD decays faster in Cross Pollinated crops

