

## GDA II, §1

- This book is about analyzing DNA data

→ A T G C (A) T T T  
→ A T G C (C) T T T

↳ Humans have 2 copies

↳ Some plants have more than 2

↳ Most loci are the same in the same species

↳ This individual has one difference

- "SNP" = "Single Nucleotide Polymorphism"  
= Locus on genome where there is a single difference that varies in the population

- "Allele" = a difference  
↳ A allele and C allele above

- There are many ways DNA can differ at a locus

- INDEL (insertion/deletion)

A T T G C T ← insert here?

A T T - C T ← delete here?

- Large structural changes

A T T T C Diff 1 G G C C A

A T T T C Diff 2 G G C C A

- This first chapter describes assays to see these DNA differences.

- Phenotype: observed consequence of an allele

↑  
Round peas versus wrinkled peas

- Homozygote: 2 of the same allele

- Heterozygote: Contains different alleles.

- "Dominant Phenotypes" - AA observe round

Aa observe round

aa observe wrinkled

↪ "recessive"



expressed if homozygote or heterozygote

↑  
only expressed if homozygote

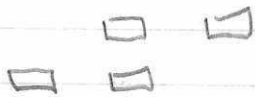
- "Codominant" - can see presence/absence of allele no matter what

- "Allozyme" - old fashioned

DNA → RNA → protein

↳ how far a protein migrates

↳ changes based on DNA sequence



aa    Aa    AA

↑  
heterozygote

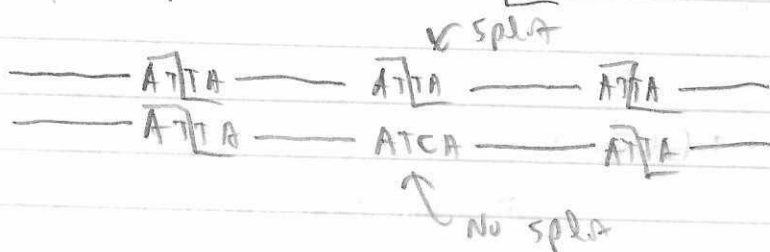
- Proteins: Sequences of 20 Amino Acids  
↳ can sequence these

Structures  
still used

## • RFLP (Restriction Fragment Length Polymorphisms)

- Restriction enzymes: Proteins that split DNA wherever they find a certain sequence

Eg. Split at ATTA



- So set of fragment sizes indicates presence/absence of allele

- Steps:
  - 1.) Apply enzyme to DNA
  - 2.) Let migrate in gel
  - 3.) smaller molecules migrate (faster?) differently
  - 4.) Use a probe to determine region of interest  
↳ SR gives you presence/absence

Structures  
still used

• Microsatellites = short tandem repeats (STR)  
= simple sequence repeats (SSR)  
= region of DNA that has lots of repeats

ATC ATC ATC ATC ...

↳ Length of this repeat is an allele

- All of these methods try to find properties of DNA,
- If we can see DNA, don't need them
- DNA sequencing is possible now and is gold standard

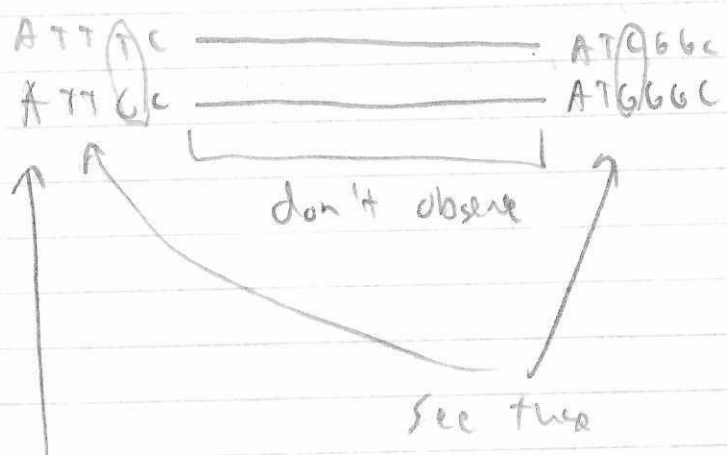
↳ but is sometimes really noisy, so sometimes older technologies are used

- Sequencing data:
 

ATTGC	A	TTT
ATTGC	C	TTT

↑ SNP, we see it!

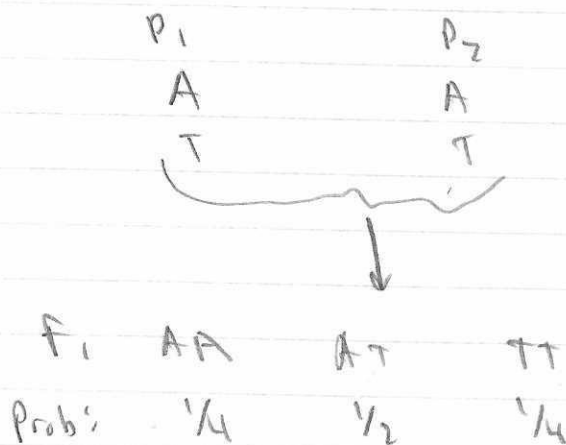
- Sometimes, have sequences at different regions of genome



↳ Is T on same molecule as C or as G?  
 ↳ can't tell sometimes

## Genetic and Statistical Sampling:

- Statistical Sampling: There is a population of size  $N$ , but we only observe  $n$  of them
- Genetic Sampling: From generation to generation, random alleles will segregate to offspring



↑ This shows up in later chapters (not 2 and 3)  
↓ we are mostly interested in present population

## Mutation:

$N$ : Population size

$n$ : Sample size

$A$ : reference allele

$a$ : alternative allele

If more than 2 alleles,  $A_1, A_2, A_3, \dots, A_n$

$p_A$ : "Allele frequency"

↑ Proportion of  $A$ 's in the population

• If 2 loci, alleles are A/a and B/b

$P_{AB}$  = proportion of genomes w/ allele A at first locus and allele B at second locus.

•  $P_{AA}$  = genotypic frequencies  
↑  
capital P

$P_{AA}$  = proportion 2 A's @ a locus

$P_{Aa}$  = proportion 1 A and 1 a @ a locus

$P_{aa}$  = proportion 2 a's @ a locus

•  $P_{ab}^{AB}$  = proportion AB on one chromosome and ab on other

$P_{aB}^{Ab}$  = proportion Ab on one chromosome and aB on other

$P_{AaBb}$  = proportion w/ 1 A and 1 a at locus 1 and 1 B and 1 b at locus 2

I but don't know whether A pairs w/ B or b

$n_{AA}$  = # AA genotypes

$n_{Aa}$  = # Aa genotypes

$n_{aa}$  = # aa genotype

$n = n_{AA} + n_{Aa} + n_{aa} = \text{total \# people}$

$n_A = 2n_{AA} + n_{Aa} = \text{\# A genomes}$

- Sample frequencies w/ tildes

$$\tilde{p}_A = \frac{1}{n} n_A$$

$$\tilde{p}_{AA} = \frac{1}{n} n_{AA}$$

- Estimates of parameters are hats.

$$\hat{p}_A = \tilde{p}_A$$

↑ But for other parameters, will not equal sample frequencies

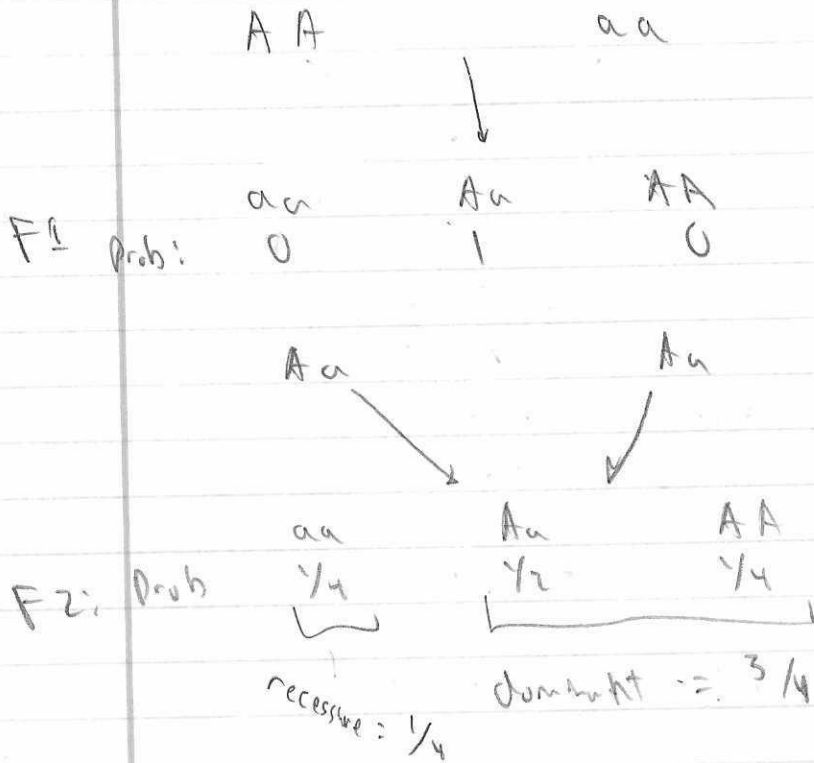
- Parameters: numbers that summarize population

- Statistics: numbers that summarize sample

- $E[X]$  = average of  $X$  over all possible samples  
= "expected value"

## Mendel and Fiske

- Mendel looked at lots of phenotypes that were controlled by single loci and were dominant/recessive



- Table 1.8: Rows 1 and 2: all F2 individuals  
 Rows 3 and 22: specific F2 individuals from the same cross of F1 individuals

Seed Characteristics

- Look at Row 3:

$$n_{AA} + n_{Aa} = 45, \quad n_{aa} = 12, \quad n = 57$$

- Does 45 and 12 closely follow the  $\frac{3}{4}$ ,  $\frac{1}{4}$  proportions we expect?



chi-square  
 $\chi^2$  goodness-of-fit test

If a category has observed count  $o$  and expected count  $e$  under some hypothesis, we can test this hypothesis w/ statistic

$$\chi^2 = \sum_{\text{categories}} \frac{(o-e)^2}{e}$$

If hypothesis is true,  $o$  and  $e$  should be close

$\chi^2 \sim \chi^2_{df}$  where  $df = \text{degrees of freedom}$   
(depends on hypothesis being tested)  
 $(\# \text{ cells}) - (\# \text{ parameters under null}) - 1$

Expected  $Aa$  or  $AA = n \cdot \frac{3}{4} = 57 \cdot \frac{3}{4} = 42.75$

Expected  $aa = n \cdot \frac{1}{4} = 57 \cdot \frac{1}{4} = 14.25$

$$\chi^2 = \frac{(45 - 42.75)^2}{42.75} + \frac{(12 - 14.25)^2}{14.25} = 0.47$$

Compare to  $\chi^2_1$  distribution

