

Advanced Algorithms and Models for Computational Biology -- a machine learning approach

Population Genetics: Pedigree and linkage analysis

Eric Xing

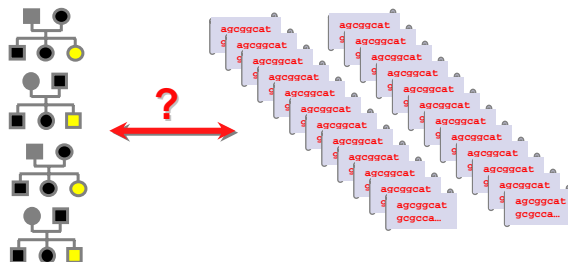
Lecture 16, March 20, 2006

Reading: DTW book, Chap 13



A crime or mass-disaster scene

- Given genetic fingerprints of F family pedigrees for alleged victims and genetic fingerprints of S samples found at a disaster site:
 - Who can you confirm died at the site? (legal)
 - Who died at the site that is outside the alleged set? (law enforcement)
 - Cluster the remains for burial. (closure)



Royal pedigree example

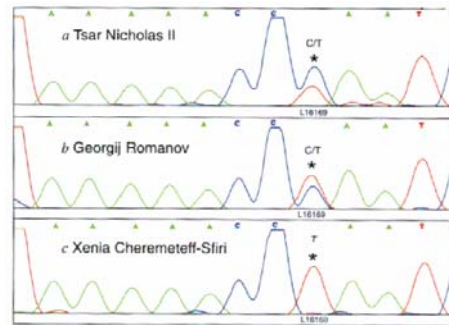
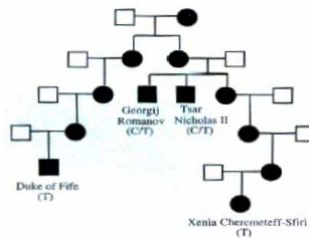
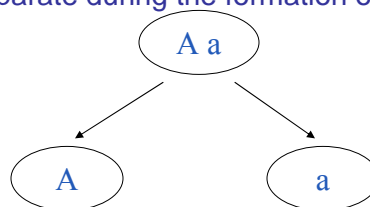


Fig. 2 Automated sequence chromatograms comparing mtDNA sequences of position 16189. a, Sequence from bones of putative Tsar Nicholas II, showing heteroplasmy with cytosine predominating thymine; b, sequence from bones of Grand Duke Georgij Romanov, showing heteroplasmy with thymine predominating cytosine; c, sequence from Countess Xenia Cheremeteff-Sfiri, homozygous for thymine.

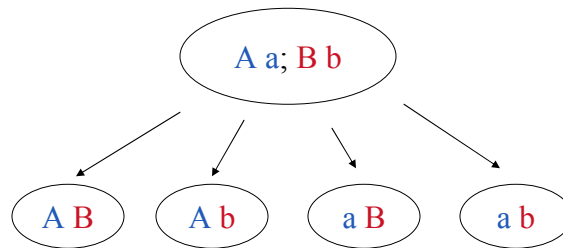
Mendel's two laws

- Modern genetics began with Mendel's experiments on garden peas. He studied seven contrasting pairs of characters, including:
 - The form of ripe seeds: round, wrinkled
 - The color of the seed albumen: yellow, green
 - The length of the stem: long, short
- **Mendel's first law:** Characters are controlled by pairs of genes which separate during the formation of the reproductive cells (meiosis)



Mendel's two laws

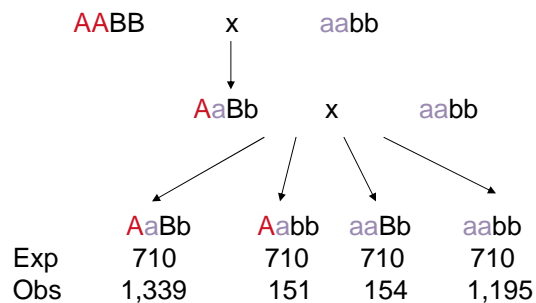
- **Mendel's second law:** When two or more pairs of gene segregate simultaneously, they do so independently.



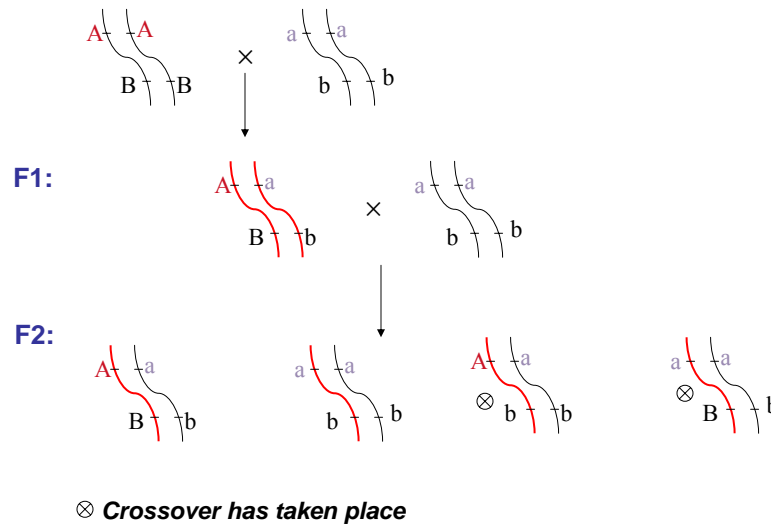
"Exceptions" to Mendel's Second Law

Morgan's fruitfly data (1909): 2,839 flies

Eye color A : red a : purple
Wing length B : normal b : vestigial



Morgan's explanation



Recombination

- *Parental types:* AaBb, aabb
- *Recombinants:* Aabb, aaBb
 - The proportion of recombinants between the two genes (or characters) is called the **recombination fraction** between these two genes.
- **Recombination fraction** It is usually denoted by r or θ . For Morgan's traits:

$$r = (151 + 154)/2839 = 0.107$$

If $r < 1/2$: two genes are said to be **linked**.

If $r = 1/2$: independent segregation (Mendel's second law).

Now we move on to (small) pedigrees.

One locus: founder probabilities



- **Founders** are individuals whose parents are not in the pedigree.
 - They may or may not be typed. Either way, we need to **assign probabilities** to their actual or possible genotypes.
 - This is usually done by assuming **Hardy-Weinberg equilibrium**. If the frequency of D is .01, $H-W$ says



$$\text{pr}(Dd) = 2 \times .01 \times .99$$

- Genotypes of *founder couples* are (usually) treated as **independent**.

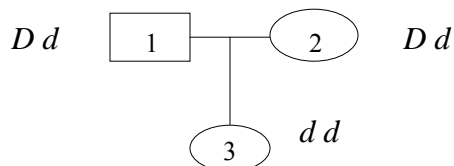


$$\text{pr}(\text{pop } Dd, \text{mom } dd) = (2 \times .01 \times .99) \times (.99)^2$$

One locus: transmission probabilities



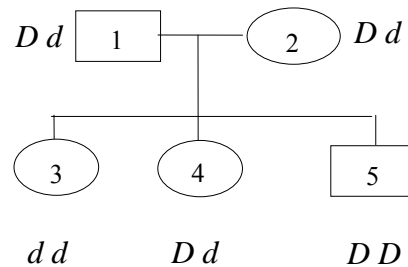
- Children get their genes from their parents' genes, independently, according to **Mendel's laws**;



$$\begin{aligned} \text{pr}(\text{kid } 3 \text{ } dd \mid \text{pop } 1 \text{ } Dd \text{ \& mom } 2 \text{ } Dd) \\ = 1/2 \times 1/2 \end{aligned}$$

- The inheritances are independent for different children.

One locus: transmission probabilities - II



$$\begin{aligned} & \text{pr}(3 \text{ } dd \text{ \& } 4 \text{ } Dd \text{ \& } 5 \text{ } DD \mid 1 \text{ } Dd \text{ \& } 2 \text{ } Dd) \\ &= (1/2 \times 1/2) \times (2 \times 1/2 \times 1/2) \times (1/2 \times 1/2). \end{aligned}$$

- The factor 2 comes from summing over the two mutually exclusive and equiprobable ways 4 can get a D and a d .

One locus: penetrance probabilities



- Independent Penetrance Model:
 - Pedigree analyses usually suppose that, given the genotype at all loci, and in some cases age and sex, the chance of having a particular phenotype depends only on genotype at one locus, and is independent of all other factors: genotypes at other loci, environment, genotypes and phenotypes of relatives, etc.

- Complete penetrance:

DD



$$\text{pr}(\text{affected} \mid DD) = 1$$

- Incomplete penetrance:

DD



$$\text{pr}(\text{affected} \mid DD) = .8$$

One locus: penetrance - II



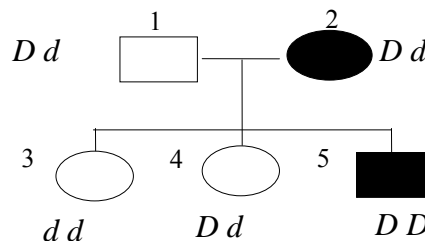
- Age and sex-dependent penetrance:



$DD (45)$

$$\text{pr}(\text{affected} \mid DD, \text{male}, 45 \text{ y.o.}) = .6$$

One locus: putting it all together



- Assume
 - Penetrances: $\text{pr}(\text{affected} \mid dd) = .1$, $\text{pr}(\text{affected} \mid Dd) = .3$, $\text{pr}(\text{affected} \mid DD) = .8$,
 - and that allele D has frequency .01.
 - In general, shaded means affected, blank means unaffected.
- The probability of this pedigree is the product:

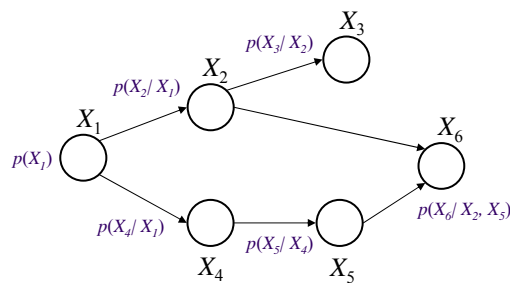
$$(2 \times .01 \times .99 \times .7) \times (2 \times .01 \times .99 \times .3) \times (1/2 \times 1/2 \times .9) \times (2 \times 1/2 \times 1/2 \times .7) \times (1/2 \times 1/2 \times .8)$$

One locus: putting it all together - II



- To write the likelihood of a pedigree:
 - we begin by multiplying founder gene frequencies,
 - followed by founder penetrances.
 - next we multiply transmission probabilities,
 - followed by penetrance probabilities of offspring, using their independence given parental genotypes.
 - If there are missing or incomplete data, we must sum over all mutually exclusive possibilities compatible with the observed data.
- Two algorithms:
 - The general strategy of beginning with founders, then non-founders, and multiplying and summing as appropriate, has been codified in what is known as the **Elston-Stewart algorithm** for calculating probabilities over pedigrees. It is one of the two widely used approaches.
 - The other is termed the **Lander-Green algorithm** and takes a quite different approach.
 - Both are hidden Markov models, both have compute time/space limitations with multiple individuals/loci (see next) , and extending them beyond their current limits is the ongoing outstanding problem.

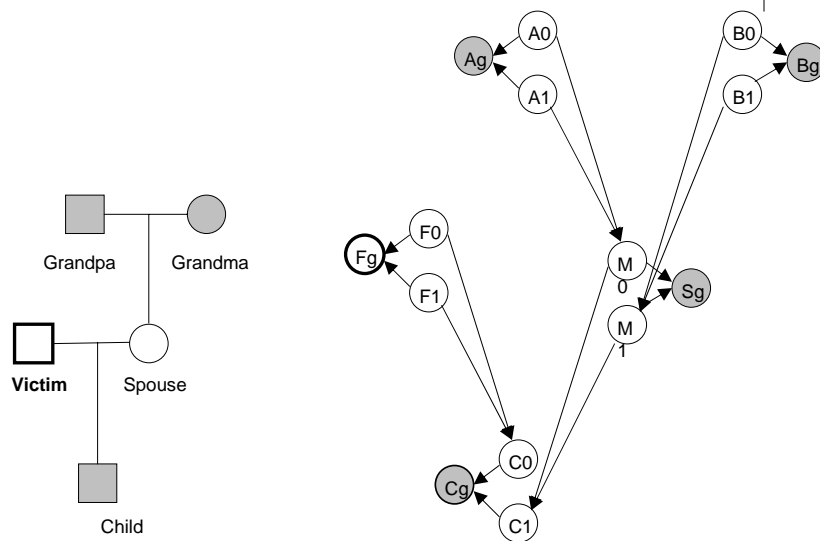
Probabilistic Graphical Models



- The joint distribution on (X_1, X_2, \dots, X_N) factors according to the “parent-of” relations defined by the edges E :

$$p(X_1, X_2, X_3, X_4, X_5, X_6) = p(X_1) p(X_2/X_1) p(X_3/X_2) p(X_4/X_1) p(X_5/X_4) p(X_6/X_2, X_5)$$

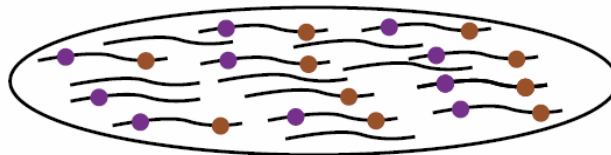
Pedigree as Graphical Models: the allele network



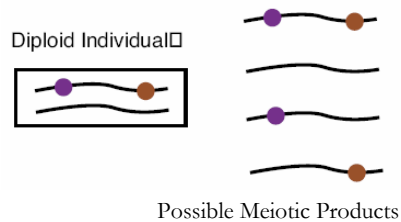
Linkage Disequilibrium



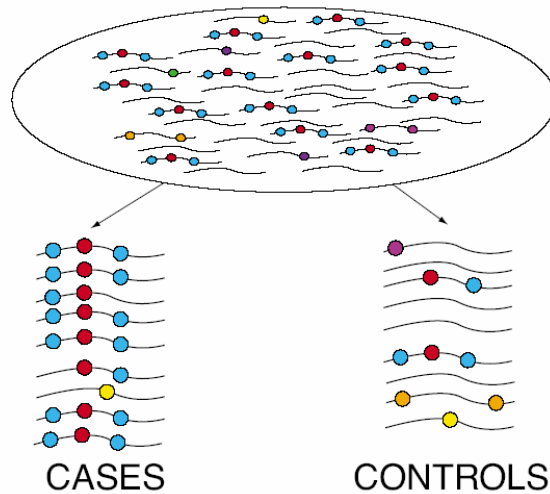
- LD is the non-random association of alleles at different sites



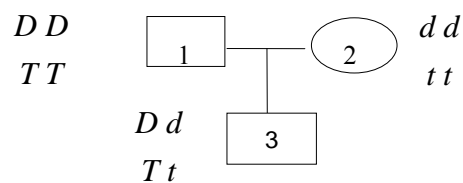
- Genetic recombination breaks down LD



Linkage Disequilibrium in Gene Mapping



Two loci: linkage and recombination



- Son 3 produces sperm with $D-T$, $D-t$, $d-T$ or $d-t$ in proportions:

	T	t	
D	$(1-\theta)/2$	$\theta/2$	1/2
d	$\theta/2$	$(1-\theta)/2$	1/2
	1/2	1/2	

no recomb.

Two loci: linkage and recombination - II



- Son produces sperm with DT , Dt , dT or dt in proportions:

	T	t	
D	$(1-\theta)/2$	$\theta/2$	1/2
d	$\theta/2$	$(1-\theta)/2$	1/2
	1/2	1/2	

$\theta = 1/2$: independent assortment (*cf* Mendel) **unlinked** loci

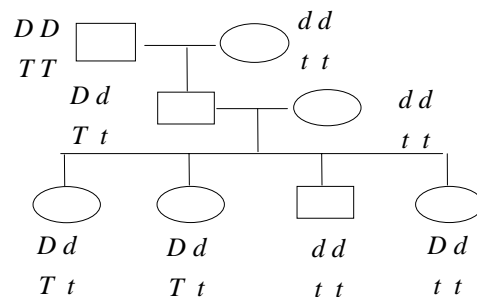
$\theta < 1/2$: **linked** loci

$\theta \approx 0$: **tightly linked** loci

Note: $\theta > 1/2$ is never observed

If the loci are linked, then $D-T$ and $d-t$ are *parental*, and $D-t$ and $d-T$ are *recombinant* haplotypes.

Two loci: estimation of recombination fractions



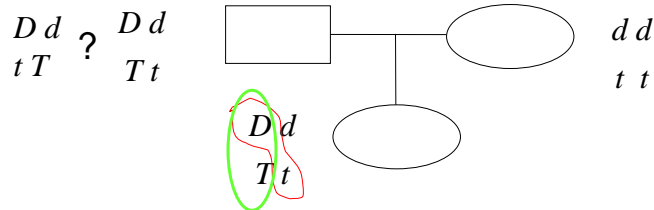
Recombination only discernible in the father. Here $\hat{\theta} = 1/4$ (why?)

This is called the **phase-known double backcross pedigree**.

Two loci: phase



- Suppose we have data on two linked loci as follows:



- Was the daughter's $D-T$ from her father a parental or recombinant combination?
 - This is the problem of **phase**: did father get $D-T$ from one parent and $d-t$ from the other? If so, then the daughter's **paternally derived haplotype** is **parental**.
 - If father got $D-t$ from one parent and $d-T$ from the other, these would be parental, and daughter's paternally derived haplotype would be recombinant.

Two loci: dealing with phase



- Phase is usually regarded as unknown genetic information, specifically, in parental origin of alleles at heterozygous loci.
- Sometimes it can be inferred with certainty from genotype data on parents.
- Often it can be inferred with high probability from genotype data on several children.
- In general genotype data on relatives helps, but does not necessarily determine phase.
- In practice, probabilities must be calculated under all phases compatible with the observed data, and added together. The need to do so is the main reason linkage analysis is computationally intensive, especially with multilocus analyses.

Two loci: founder probabilities



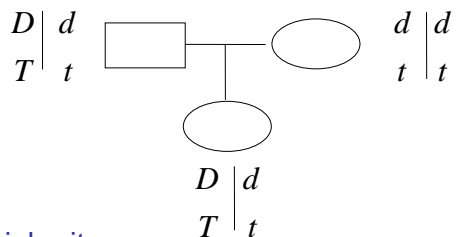
- Two-locus founder probabilities are typically calculated assuming **linkage equilibrium**, i.e. independence of genotypes across loci.
- If D and d have frequencies .01 and .99 at one locus, and T and t have frequencies .25 and .75 at a second, linked locus, this assumption means that DT , Dt , dT and dt have frequencies .01 x .25, .01 x .75, .99 x .25 and .99 x .75 respectively. Together with Hardy-Weinberg, this implies that

$$\begin{array}{c} Dd \\ Tt \end{array} \quad \square$$

$$\begin{aligned} \text{pr}(DdTt) &= (2 \times .01 \times .99) \times (2 \times .25 \times .75) \\ &= 2 \times (.01 \times .25) \times (.99 \times .75) + 2 \times (.01 \times .75) \times (.99 \times .25). \end{aligned}$$

- This last expression adds haplotype pair probabilities.

Two loci: transmission probabilities



- Haplotype inheritance:
 - Initially, this must be done with haplotypes, so that account can be taken of recombination.
 - Then terms like that below are summed over possible phases.
 - Here only the father can exhibit recombination: mother is **uninformative**.
- $$\begin{aligned} &\text{pr}(\text{kid } DT/dt \mid \text{pop } DT/dt \text{ \& mom } dt/dt) \\ &= \text{pr}(\text{kid } DT \mid \text{pop } DT/dt) \times \text{pr}(\text{kid } dt \mid \text{mom } dt/dt) \\ &= (1-0)/2 \times 1. \end{aligned}$$

Two Loci: Penetrance



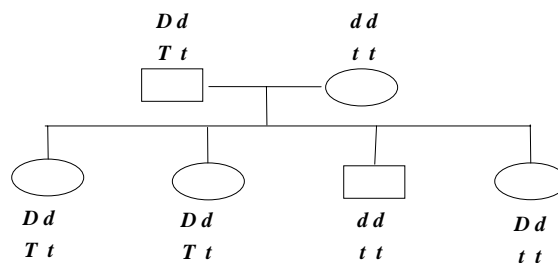
- In all standard linkage programs, different parts of phenotype are conditionally independent given all genotypes, and two-loci penetrances split into products of one-locus penetrances.
- Assuming the penetrances for DD, Dd and dd given earlier, and that T,t are two alleles at a co-dominant marker locus.

$$\begin{aligned}
 & \Pr(\text{affected} \ \& \ Tt \mid DD, Tt) \\
 &= \Pr(\text{affected} \mid DD, Tt) \times \Pr(Tt \mid DD, Tt) \\
 &= 0.8 \times 1
 \end{aligned}$$

Two loci: phase unknown double backcross



- We assume below pop is as likely to be DT/dt as Dt/dT .



$$\begin{aligned}
 & \Pr(\text{all data} \mid \theta) \\
 &= \Pr(\text{parents' data} \mid \theta) \times \Pr(\text{kids' data} \mid \text{parents' data}, \theta) \\
 &= \Pr(\text{parents' data}) \times \{[(1-\theta)/2]^3 \times \theta/2 + [(\theta/2)^3 \times (1-\theta)/2]\}
 \end{aligned}$$

This is then maximised in θ , in this case numerically. Here $\hat{\theta} = 0.25$

Log (base 10) odds or LOD scores



- Suppose $\text{pr}(\text{data} \mid \theta)$ is the likelihood function of a recombination fraction θ generated by some 'data', and $\text{pr}(\text{data} \mid 1/2)$ is the same likelihood when $\theta = 1/2$.

- Statistical theory tells us that the ratio

$$L = \text{pr}(\text{data} \mid \theta^*) / \text{pr}(\text{data} \mid 1/2)$$

provides a basis for deciding whether $\theta = \theta^*$ rather than $\theta = 1/2$.

- This can equally well be done with $\text{Log}_{10}L$, i.e.

$$\text{LOD}(\theta^*) = \text{Log}_{10}\{\text{pr}(\text{data} \mid \theta^*) / \text{pr}(\text{data} \mid 1/2)\}$$

measures the relative strength of the data for $\theta = \theta^*$ rather than $\theta = 1/2$. Usually we write θ , not θ^* and calculate the function **LOD**(θ).

Facts about/interpretation of LOD scores



1. Positive LOD scores suggests stronger support for θ^* than for $1/2$, negative LOD scores the reverse.
2. Higher LOD scores means stronger support, lower means the reverse.
3. LODs are additive across independent pedigrees, and under certain circumstances can be calculated sequentially.
4. For a single two-point linkage analysis, the threshold $\text{LOD} \approx 3$ has become the de facto standard for "establishing linkage", i.e. rejecting the null hypothesis of no linkage.
5. When more than one locus or model is examined, the remark in 4 must be modified, sometimes dramatically.

Assumptions underpinning most 2-point human linkage analyses



- **Founder Frequencies:** Hardy-Weinberg, random mating at each locus. Linkage equilibrium across loci, **known** allele frequencies; founders independent.
- **Transmission:** Mendelian segregation, no mutation.
- **Penetrance:** single locus, no room for dependence on relatives' phenotypes or environment. **Known** (including phenocopy rate).
- **Implicit:** phenotype and genotype data **correct**, marker order and location correct
- **Comment:** Some analyses are *robust*, others can be *very sensitive* to violations of some of these assumptions. Non-standard linkage analyses can be developed.

Beyond two-point human linkage analysis



- The real challenge is multipoint linkage analysis, but going there would take more time than we have today.
- Next in importance is dealing with two-locus penetrances.



Acknowledgements

Melanie Bahlo, WEHI
Hongyu Zhao, Yale
Karl Broman, Johns Hopkins
Nusrat Rabbee, UCB



References

www.netspace.org/MendelWeb

HLK Whitehouse: **Towards an Understanding of the Mechanism of Heredity**, 3rd ed. Arnold 1973

Kenneth Lange: **Mathematical and statistical methods for genetic analysis**, Springer 1997

Elizabeth A Thompson: **Statistical inference from genetic data on pedigrees**, CBMS, IMS, 2000.

Jurg Ott : **Analysis of human genetic linkage**, 3rd edn
Johns Hopkins University Press 1999

JD Terwilliger & J Ott : **Handbook of human genetic linkage**, Johns Hopkins University Press 1994