

# Computational Genomics

## Population Genetics: Quantitative Trait Locus (QTL) Mapping

Eric Xing

Lecture 4, January 25, 2007

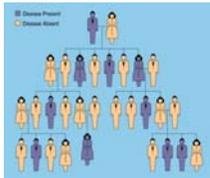
Reading: DTW book, Chap 13



## Phenotypical Traits

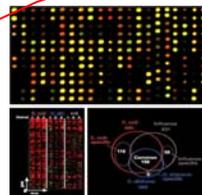
- Body measures:



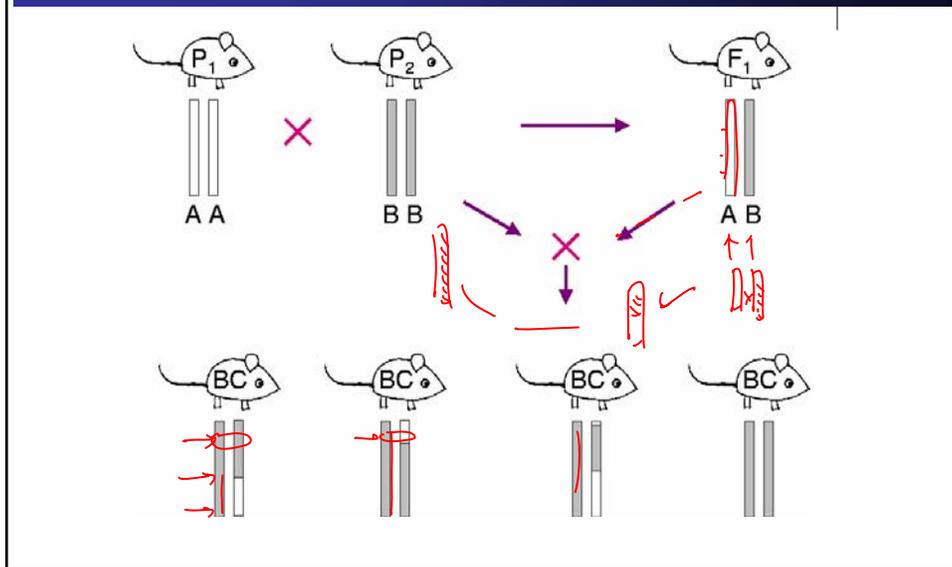
-  Disease susceptibility and drug response



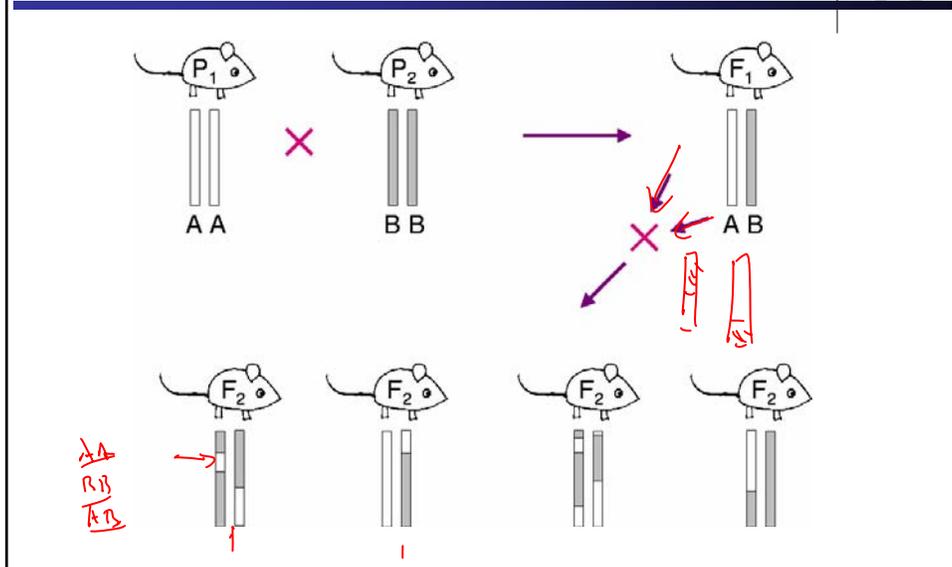
- Gene expression (microarray)



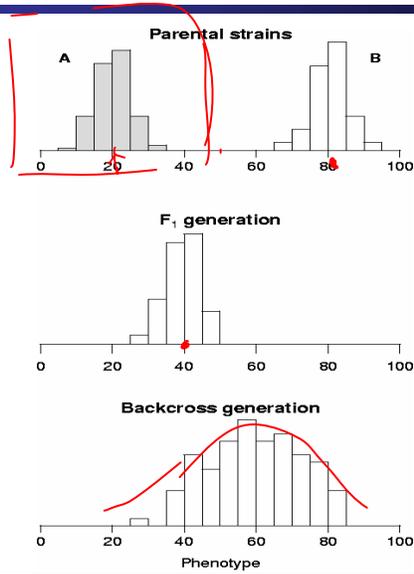
## Backcross experiment



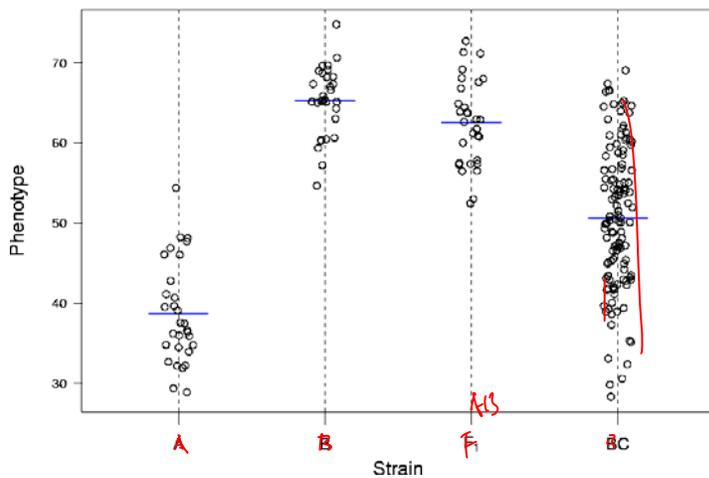
## F<sub>2</sub> intercross experiment



# Trait distributions: a classical view

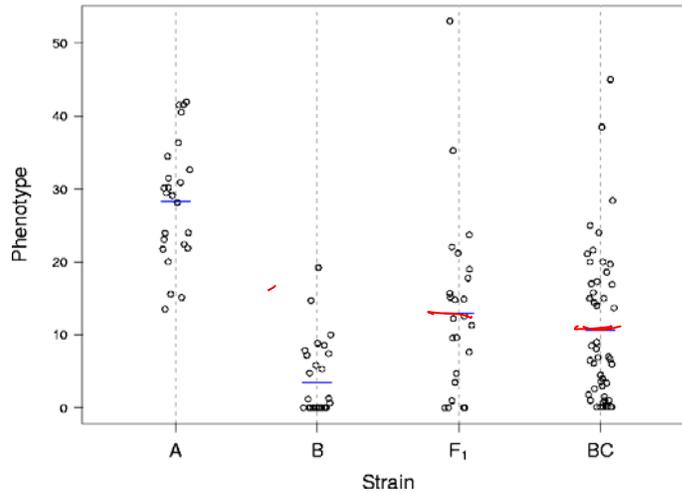


# Another representation of a trait distribution



Note the equivalent of dominance in our trait distributions.

## A second example



Note the approximate additivity in our trait distributions here.

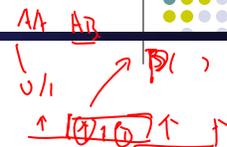
## QTL mapping

- Data

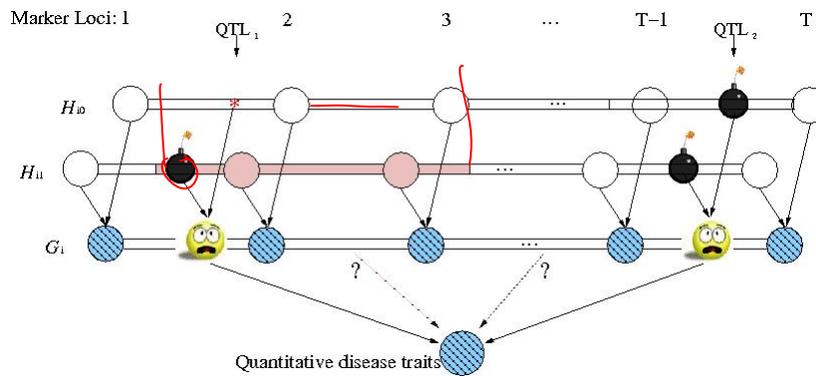
- Phenotypes:  $y_i$  = trait value for mouse  $i$
- Genotype:  $x_{ij}$  = 1/0 (i.e., A/H) of mouse  $i$  at marker  $j$  (backcross); need three states for intercross
- Genetic map: Locations of markers

- Goals

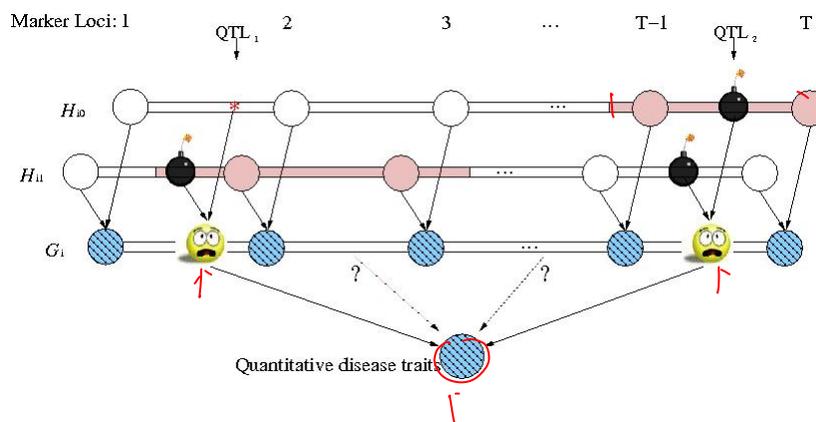
- Identify the (or at least one) genomic region, called quantitative trait locus = QTL, that contributes to variation in the trait
- Form confidence intervals for the QTL location
- Estimate QTL effects



## QTL mapping (BC)



## QTL mapping (F2)



## Models: Recombination



- We assume no chromatid or crossover interference.
- ⇒ points of exchange (crossovers) along chromosomes are distributed as a Poisson process, rate 1 in genetic distance
- ⇒ the marker genotypes  $\{x_{ij}\}$  form a Markov chain along the chromosome for a backcross; what do they form in an  $F_2$  intercross?

$$\frac{1}{d} \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

## Models: Genotype $\rightarrow$ Phenotype



- Let  $y =$  phenotype,  $\sim N(\mu, \sigma^2)$   
 $g =$  whole genome genotype  $\mu(\text{---})$   
 $\sigma^2(\text{---})$
- Imagine a small number of QTL with genotypes  $g_1, \dots, g_p$  ( $2^p$  or  $3^p$  distinct genotypes for BC, IC resp, why?).

We assume

$$E(y|g) = \mu(\overset{\downarrow}{g_1}, \dots, \overset{\downarrow}{g_p}), \quad \text{var}(y|g) = \sigma^2(g_1, \dots, g_p)$$

## Models: Genotype → Phenotype



- **Homoscedacity** (constant variance)

$$\sigma^2(g_1, \dots, g_p) = \sigma^2 \text{ (constant)}$$

- **Normality** of residual variation

$$y|g \sim N(\mu_g, \sigma^2)$$

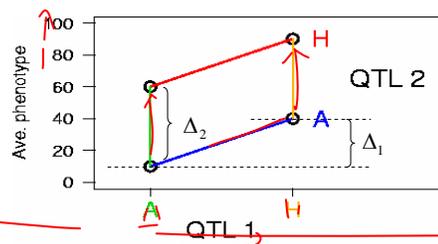
- **Additivity:**

$$\mu(g_1, \dots, g_p) = \mu + \sum \Delta_j g_j \quad (g_j = 0/1 \text{ for BC})$$

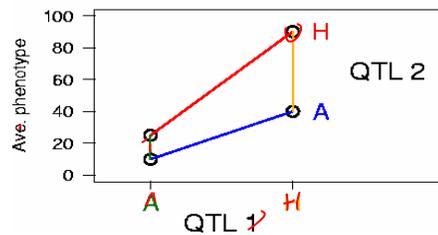
- **Epistasis:** Any deviations from additivity.

$$\mu(g_1, \dots, g_p) = \mu + \sum \Delta_j g_j + \sum \omega_{ij} g_i g_j$$

## Additivity, or non-additivity (BC)



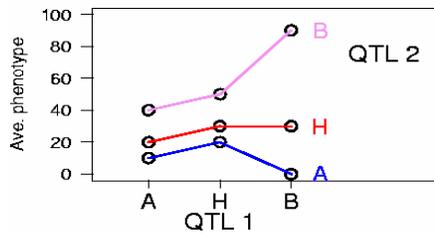
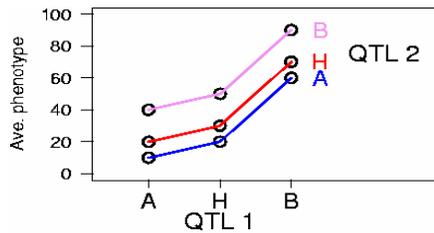
The effect of QTL 1 is the same, irrespective of the genotype of QTL 2, and vice versa.



Epistatic QTLs

$$\Delta_i \sim p(\cdot | g_j)$$

## Additivity or non-additivity: F2



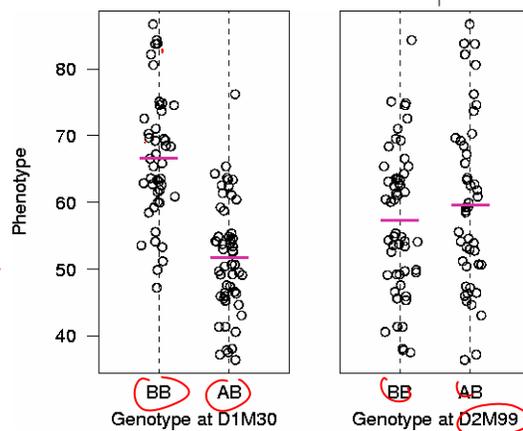
## The simplest method: ANOVA

SP?



- Split subjects into groups according to genotype at a single marker
- Do a t-test/ANOVA
- Repeat for each marker

t-test/ANOVA will tell whether there is sufficient evidence to say that measurements from one condition (i.e., genotype) differ significantly from another



- LOD score =  $\log_{10}$  likelihood ratio, comparing single-QTL model to the “no QTL anywhere” model.

# ANOVA at marker loci



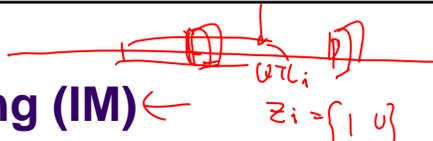
## Advantages

- Simple
- Easily incorporate covariates (sex, env, treatment ...)
- Easily extended to more complex models

## Disadvantages

- Must exclude individuals with missing genotype data
- Imperfect information about QTL location
- Suffers in low density scans
- Only considers one QTL at a time

# Interval mapping (IM)

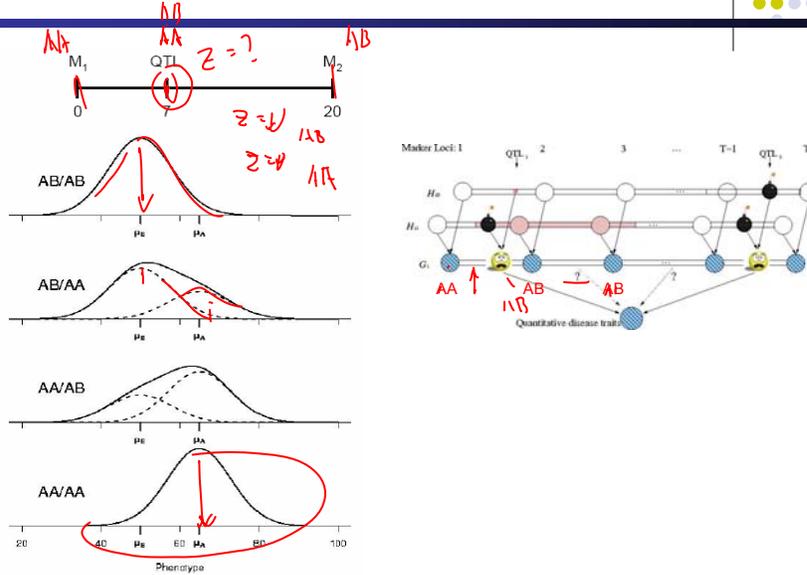


- Consider any one position in the genome as the location for a putative QTL
- For a particular mouse, let  $z = 1/0$  if (unobserved) genotype at QTL is AB/AA
- Calculate  $\Pr(z = 1 \mid \text{marker data of an interval bracketing the QTL})$ 
  - Assume no meiotic interference
  - Need only consider flanking typed markers
  - May allow for the presence of genotyping errors
- Given genotype at the QTL, phenotype is distributed as

$$y_i | z_i \sim \text{Normal}(\mu_{z_i}, \sigma^2)$$

- Given marker data, phenotype follows a *mixture* of normal distributions

## IM: the mixture model



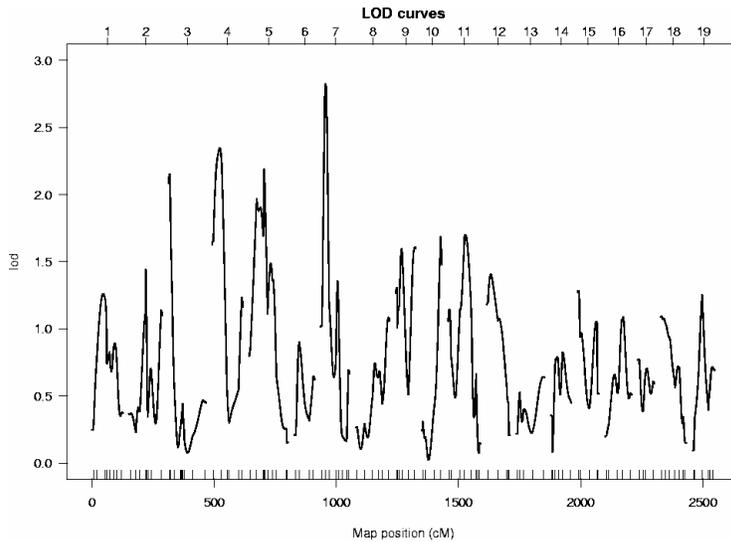
## IM: estimation and LOD scores

- Use a version of the EM algorithm to obtain estimates of  $\mu_{AA}$ ,  $\mu_{AB}$ , and  $\sigma$  (an *iterative* algorithm)
- Calculate the LOD score

$$\text{LOD} = \log_{10} \left\{ \frac{P(\text{data} | \hat{\mu}_{AA}, \hat{\mu}_{AB})}{P(\text{data} | \text{no QTL})} \right\}$$

- Repeat for all other genomic positions (in practice, at 0.5 cM steps along genome)

## LOD score curves

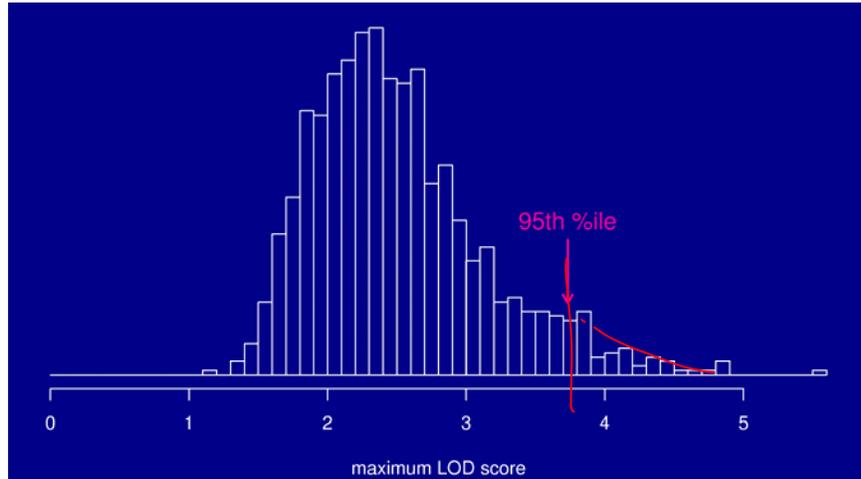


## LOD thresholds



- To account for the genome-wide search, compare the observed LOD scores to the distribution of the maximum LOD score, genome-wide, that would be obtained if there were no QTL anywhere.
- **LOD threshold** = 95th %ile of the distribution of genome-wide maxLOD, when there are no QTL anywhere
- **Derivations:**
  - Analytical calculations (Lander & Botstein, 1989)
  - Simulations
  - Permutation tests (Churchill & Doerge, 1994).

# Permutation distribution for trait4



# Interval mapping

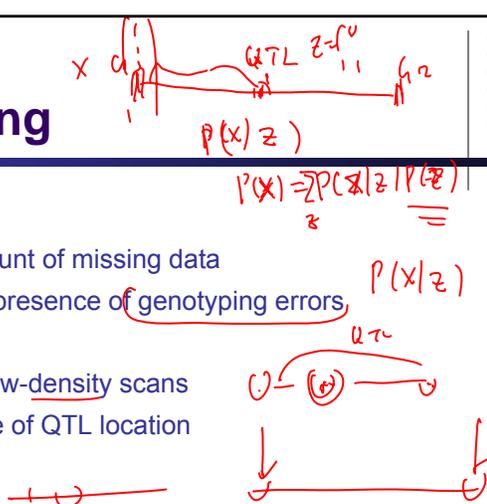


## Advantages

- Make proper account of missing data
- Can allow for the presence of genotyping errors
- Pretty pictures
- Higher power in low-density scans
- Improved estimate of QTL location

## Disadvantages

- Greater computational effort
- Requires specialized software
- More difficult to include covariates?
- Only considers one QTL at a time



## Multiple QTL methods



### Why consider multiple QTL at once?

- To separate linked QTL. If two QTL are close together on the same chromosome, our one-at-a-time strategy may have problems finding either (e.g. if they work in opposite directions, or interact). Our LOD scores won't make sense either.
- To permit the investigation of interactions. It may be that interactions greatly strengthen our ability to find QTL, though this is not clear.
- To reduce residual variation. If QTL exist at loci other than the one we are currently considering, they should be in our model. For if they are not, they will be in the error, and hence reduce our ability to detect the current one. See below.

## The problem



- $n$  backcross subjects;  $M$  markers in all, with at most a handful expected to be near QTL



$x_{ij}$  = genotype (0/1) of mouse  $i$  at marker  $j$

$y_i$  = phenotype (trait value) of mouse  $i$

$$Y_i = \mu + \sum_{j=1}^M \Delta_j x_{ij} + \varepsilon_j \quad \text{Which } \Delta_j \neq 0?$$

*(Handwritten red annotations: a circle around the summation term, and a red arrow pointing to the equation)*

⇒ Variable selection in linear models (regression)

*(Handwritten red notation:  $f(\rightarrow)$ )*

## Finding QTL as model selection



### Select class of models

- Additive models
- Additive plus pairwise interactions
- Regression trees

### Compare models ( $\gamma$ )

- $BIC_{\delta}(\gamma) = \log RSS(\gamma) + \gamma(\delta \log n/n)$
- Sequential permutation tests

### Search model space

- Forward selection (FS)
- Backward elimination (BE)
- FS followed by BE
- MCMC

### Assess performance

- Maximize no QTL found;
- control false positive rate

## Acknowledgements



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

## References



[www.netspace.org/MendelWeb](http://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