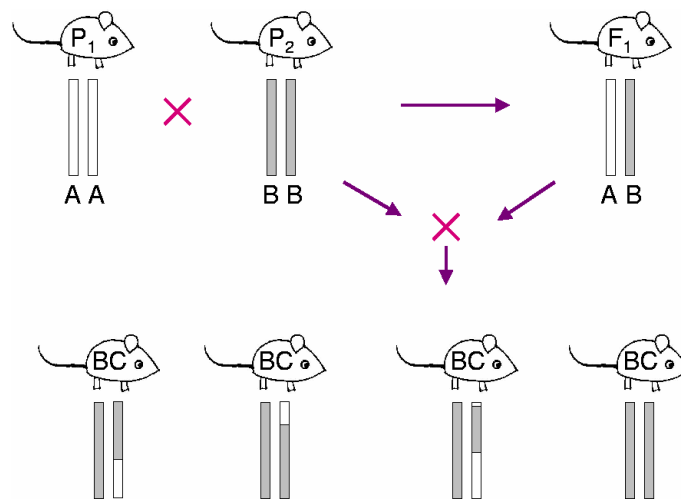


# Quantitative Trait Locus (QTL) Mapping

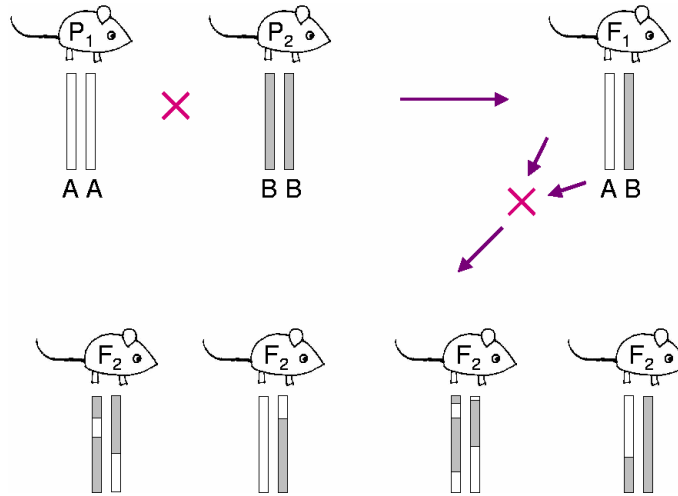


10-810, CMB lecture 9---Eric Xing

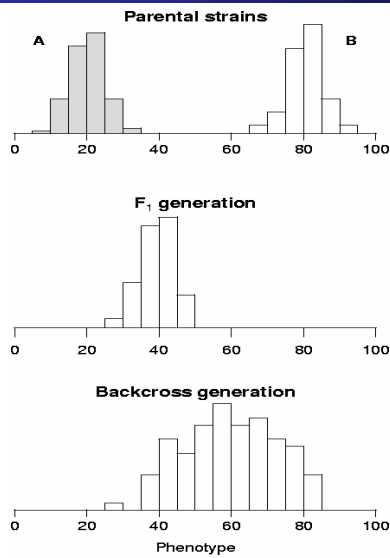
## Backcross experiment



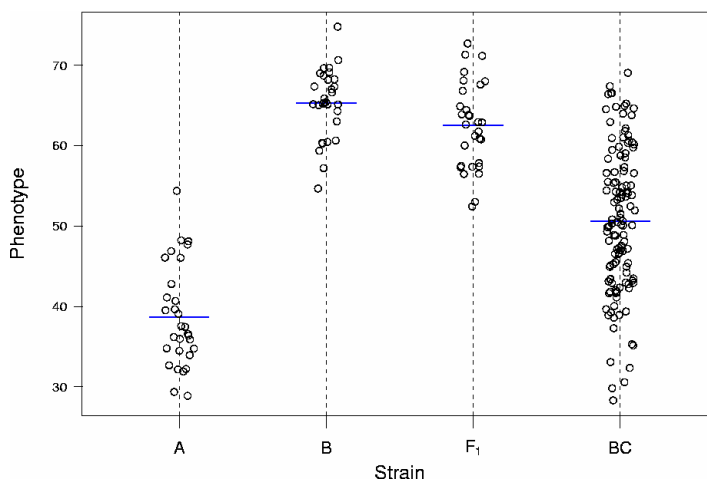
## F<sub>2</sub> intercross experiment



## Trait distributions: a classical view

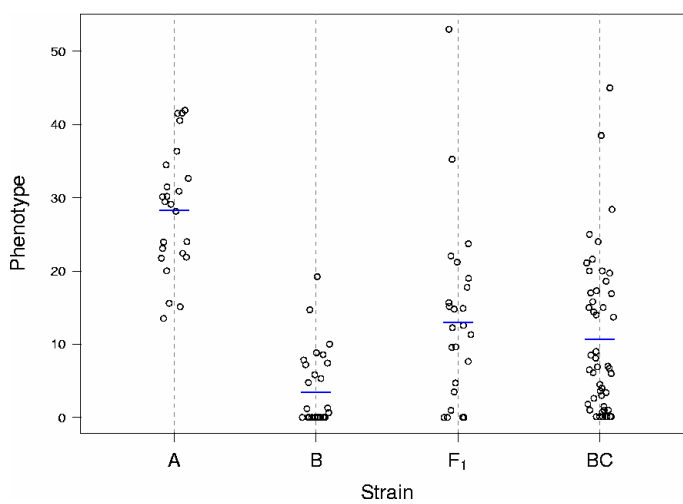


## Another representation of a trait distribution



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  $i$

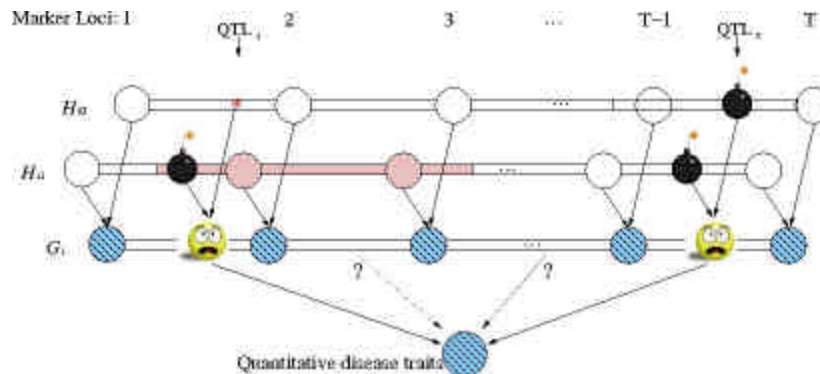
Genotype:  $x_{ij} = 1/0$  (i.e., A/H) of mouse  $i$  at marker  $j$  (backcross);  
need two dummy variables for intercross

Genetic map: Locations of markers

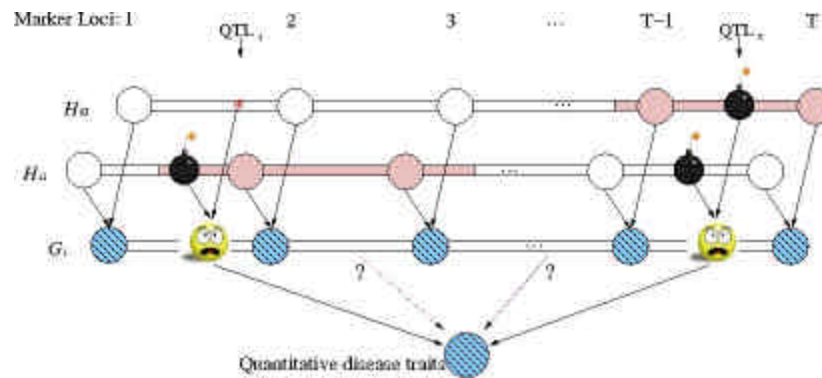
### Goals

- 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?

## Models: Genotype® Phenotype



Let  $y$  = phenotype,  
 $g$  = whole genome genotype

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

$$E(y|g) = m(g_1, \dots, g_p), \quad \text{var}(y|g) = s^2(g_1, \dots, g_p)$$

## Models: Genotype® Phenotype, ctd



**Homoscedacity** (constant variance)

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

**Normality** of residual variation

$$y|g \sim N(m_g, s^2)$$

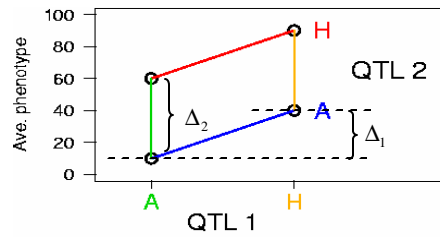
**Additivity:**

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

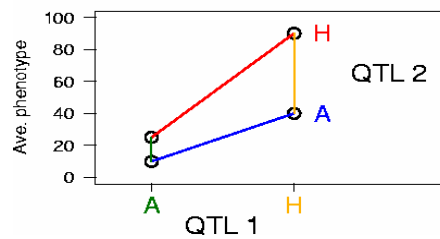
**Epistasis**: Any deviations from additivity.

$$m(g_1, \dots, g_p) = m + \sum D_j g_j + \sum w_{ij} g_i g_j$$

## Additivity, or non-additivity (BC)

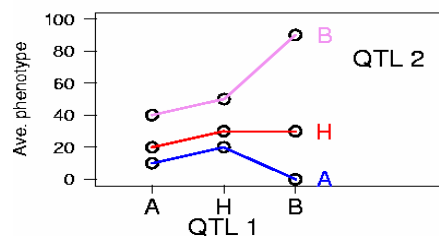
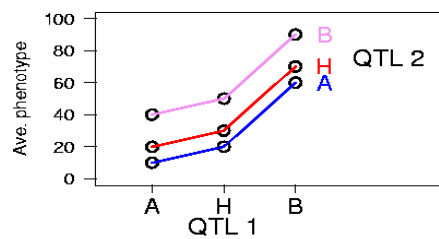


The effect of QTL 1 is the same, irrespective of the genotype of QTL 2, and vice versa.



Epistatic QTLs  
 $\Delta_i \sim p(\mid g_j)$

## Additivity or non-additivity: F<sub>2</sub>

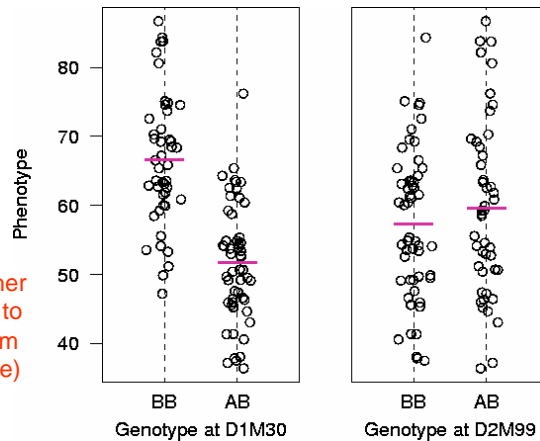


## The simplest method: ANOVA



- Split subjects into groups according to genotype at a 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

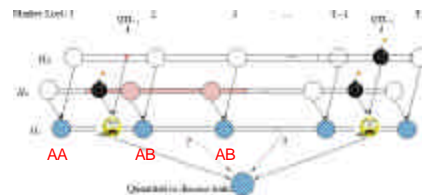
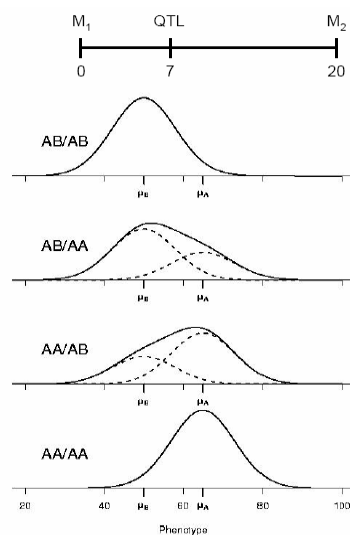


## Interval mapping (IM)



- 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  $\Pr(z = 1 \mid \text{marker data})$ 
  - 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  $N(\mu + ?z, s^2)$
- Given marker data, phenotype follows a *mixture* of normal distributions

## IM: the mixture model



## IM: estimation and LOD scores

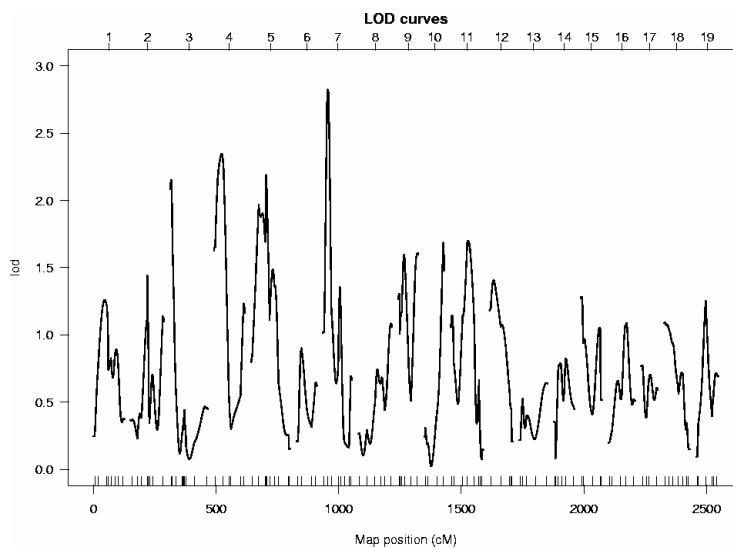


- Use a version of the EM algorithm to obtain estimates of  $\mu_{AA}$ ,  $\mu_{AB}$ , and  $s$  (an *iterative* algorithm)
- Calculate the LOD score

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

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

## LOD score curves



## LOD thresholds



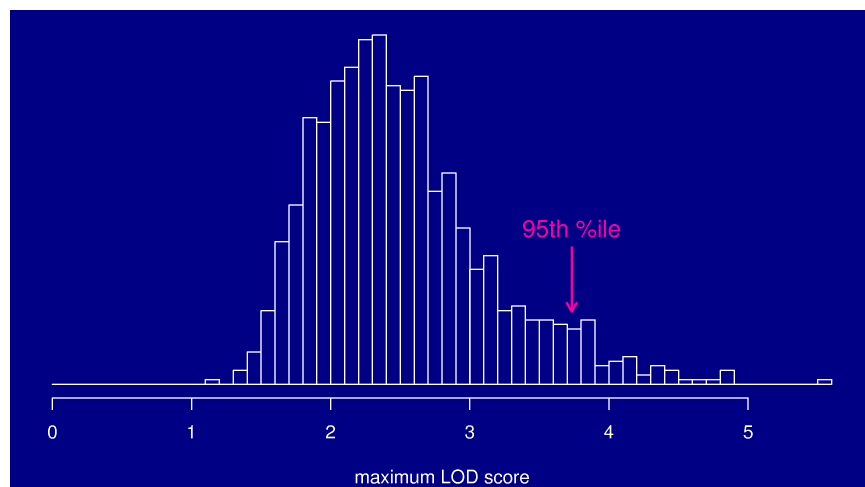
To account for the genome-wide search, compare the observed LOD scores to the distribution of the maximum LOD score, genome-wide, that would be obtained if there were no QTL anywhere.

**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  $i$

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

⇒ Variable selection in linear models (regression)

## Finding QTL as model selection



### Select class of models

- Additive models
- Additive plus pairwise interactions
- Regression trees

### Search model space

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

### Compare models ( $g$ )

- $BIC_d(g) = \log \text{RSS}(g) + g(d \log n/n)$
- Sequential permutation tests

### Assess performance

- Maximize no QTL found;
- control false positive rate



---

## Acknowledgements

Karl Broman, Johns Hopkins  
Nusrat Rabbee, UCB