

# Bayesian Interval Mapping

1. what is goal of QTL study?	2-8
2. Bayesian QTL mapping	9-20
3. Markov chain sampling	21-27
4. sampling across architectures	28-34
5. epistatic interactions	35-42
6. comparing models	43-46

## 1. what is the goal of QTL study?

- uncover underlying biochemistry
  - identify how networks function, break down
  - find useful candidates for (medical) intervention
  - epistasis may play key role
  - statistical goal: maximize number of correctly identified QTL
- basic science/evolution
  - how is the genome organized?
  - identify units of natural selection
  - additive effects may be most important (Wright/Fisher debate)
  - statistical goal: maximize number of correctly identified QTL
- select “elite” individuals
  - predict phenotype (breeding value) using suite of characteristics (phenotypes) translated into a few QTL
  - statistical goal: minimize prediction error

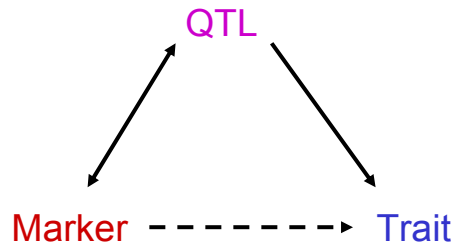
cross two inbred lines

→ linkage disequilibrium

→ associations

→ linked segregating QTL

(after Gary Churchill)



## pragmatics of multiple QTL

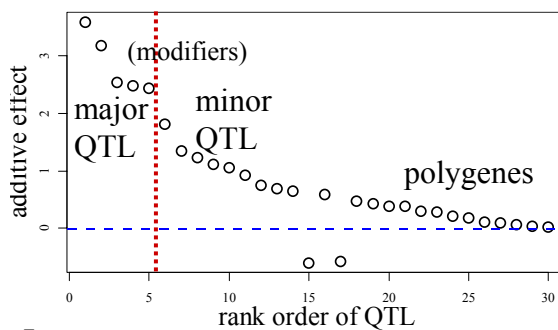
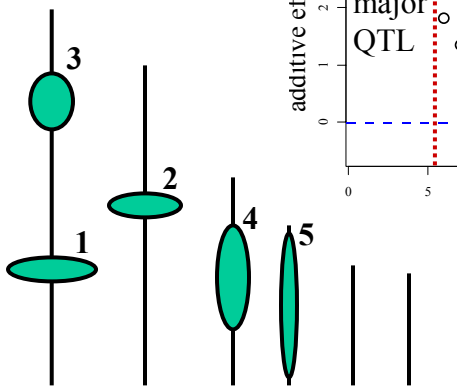
- evaluate some objective for model given data
  - classical likelihood
  - Bayesian posterior
- search over possible genetic architectures (models)
  - number and positions of loci
  - gene action: additive, dominance, epistasis
- estimate “features” of model
  - means, variances & covariances, confidence regions
  - marginal or conditional distributions
- art of model selection
  - how select “best” or “better” model(s)?
  - how to search over useful subset of possible models?

## advantages of multiple QTL approach

- improve statistical power, precision
  - increase number of QTL detected
  - better estimates of loci: less bias, smaller intervals
- improve inference of complex genetic architecture
  - patterns and individual elements of epistasis
  - appropriate estimates of means, variances, covariances
    - asymptotically unbiased, efficient
  - assess relative contributions of different QTL
- improve estimates of genotypic values
  - less bias (more accurate) and smaller variance (more precise)
  - mean squared error =  $MSE = (\text{bias})^2 + \text{variance}$

## Pareto diagram of QTL effects

major QTL on linkage map



## limits of multiple QTL?

- limits of statistical inference
  - power depends on sample size, heritability, environmental variation
  - “best” model balances fit to data and complexity (model size)
  - genetic linkage = correlated estimates of gene effects
- limits of biological utility
  - sampling: only see some patterns with many QTL
  - marker assisted selection (Bernardo 2001 *Crop Sci*)
    - 10 QTL ok, 50 QTL are too many
    - phenotype better predictor than genotype when too many QTL
    - increasing sample size may not give multiple QTL any advantage
  - hard to select many QTL simultaneously
    - $3^m$  possible genotypes to choose from

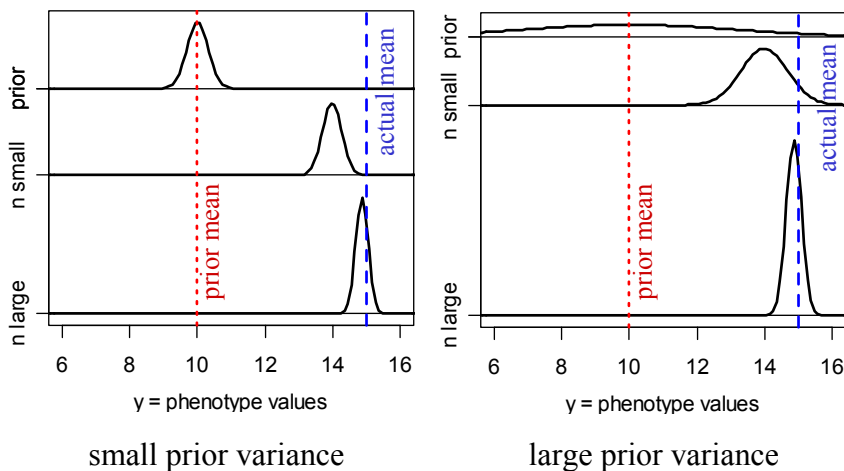
## QTL below detection level?

- problem of selection bias
  - QTL of modest effect only detected sometimes
  - their effects are biased upwards when detected
- probability that QTL detected
  - avoids sharp in/out dichotomy
  - avoid pitfalls of one “best” model
  - examine “better” models with more probable QTL
- build  $m$  = number of QTL detected into QTL model
  - directly allow uncertainty in genetic architecture
  - model selection over genetic architecture

## 2. Bayesian QTL mapping

- Reverend Thomas Bayes (1702-1761)
  - part-time mathematician
  - buried in Bunhill Cemetary, Moongate, London
  - famous paper in 1763 *Phil Trans Roy Soc London*
  - was Bayes the first with this idea? (Laplace?)
- basic idea (from Bayes' original example)
  - two billiard balls tossed at random (uniform) on table
  - where is first ball if the second is to its **left**?
    - prior: anywhere on the table
    - posterior: more likely toward **right** end of table

## Bayes posterior for normal data

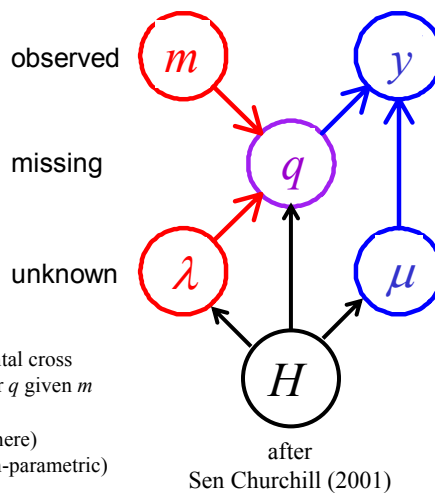


# Bayes posterior for normal data

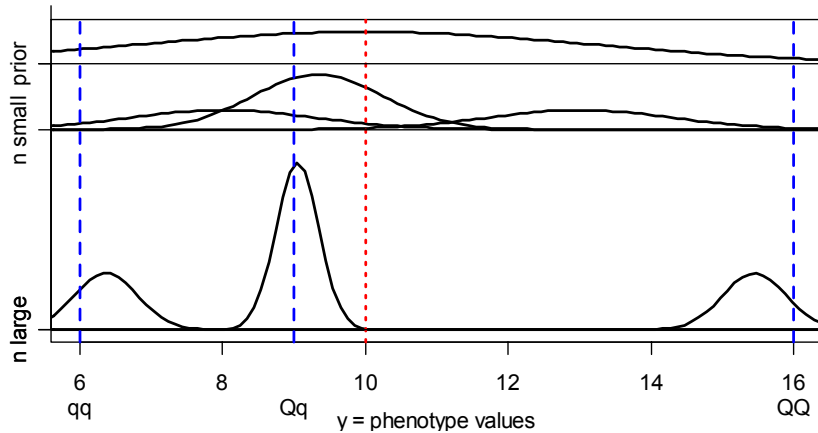
model	$y_i = \mu + e_i$
environment	$e \sim N(0, \sigma^2), \sigma^2 \text{ known}$
likelihood	$y \sim N(\mu, \sigma^2)$
prior	$\mu \sim N(\mu_0, \kappa\sigma^2), \kappa \text{ known}$
posterior:	mean tends to sample mean
single individual	$\mu \sim N(\mu_0 + b_1(y_1 - \mu_0), b_1\sigma^2)$
sample of $n$ individuals	$\mu \sim N(b_n \bar{y}_\bullet + (1 - b_n)\mu_0, b_n\sigma^2 / n)$ with $\bar{y}_\bullet = \sum_{i=1, \dots, n} y_i / n$
fudge factor (shrinks to 1)	$b_n = \frac{\kappa n}{\kappa n + 1} \rightarrow 1$

# Bayesian QTL: key players

- observed measurements
  - $y$  = phenotypic trait
  - $m$  = markers & linkage map
  - $i$  = individual index (1, ...,  $n$ )
- missing data
  - missing marker data
  - $q$  = QT genotypes
    - alleles QQ, Qq, or qq at locus
- unknown quantities
  - $\lambda$  = QT locus (or loci)
  - $\mu$  = phenotype model parameters
  - $H$  = QTL model/genetic architecture
- $\text{pr}(q|m, \lambda, H)$  genotype model
  - grounded by linkage map, experimental cross
  - recombination yields multinomial for  $q$  given  $m$
- $\text{pr}(y|q, \mu, H)$  phenotype model
  - distribution shape (assumed normal here)
  - unknown parameters  $\mu$  (could be non-parametric)



## pr(y|q,μ) phenotype model



13

## Bayes posterior QTL means

posterior centered on sample genotypic mean  
but shrunken slightly toward overall mean

prior:  $\mu_q \sim N(\bar{y}_\bullet, \kappa\sigma^2)$

posterior:  $\mu_q \sim N(b_q \bar{y}_q + (1 - b_q) \bar{y}_\bullet, b_q \sigma^2 / n_q)$

$$n_q = \text{count}\{q_i = q\}, \bar{y}_q = \frac{\text{sum}_{\{q_i=q\}} y_i}{n_q}$$

fudge factor:  $b_q = \frac{\kappa n_q}{\kappa n_q + 1} \rightarrow 1$

QTL 2: Bayes

Seattle SISG: Yandell © 2006

14

## partition of multiple QTL effects

- partition genotype-specific mean into QTL effects

$$\mu_q = \text{mean} + \text{main effects} + \text{epistatic interactions}$$

$$\mu_q = \mu + \beta_q = \mu + \sum_{j \text{ in } H} \beta_{qj}$$

- priors on mean and effects

$$\mu \sim N(\mu_0, \kappa_0 \sigma^2) \quad \text{grand mean}$$

$$\beta_q \sim N(0, \kappa_1 \sigma^2) \quad \text{model-independent genotypic effect}$$

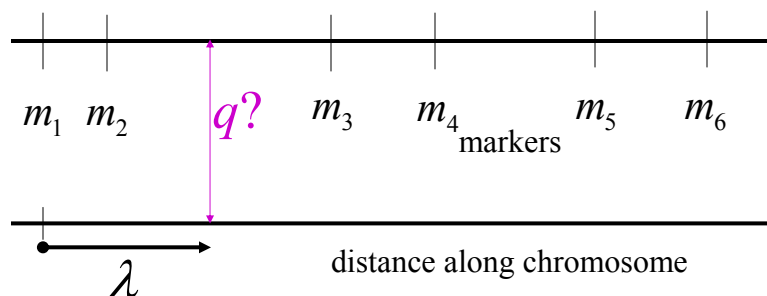
$$\beta_{qj} \sim N(0, \kappa_1 \sigma^2 / |H|) \quad \text{effects down-weighted by size of } H$$

- determine hyper-parameters via empirical Bayes

$$\mu_0 \approx \bar{Y}_\bullet \quad \text{and} \quad \kappa_1 \approx \frac{h^2}{1-h^2} = \frac{\sigma_G^2}{\sigma^2}$$

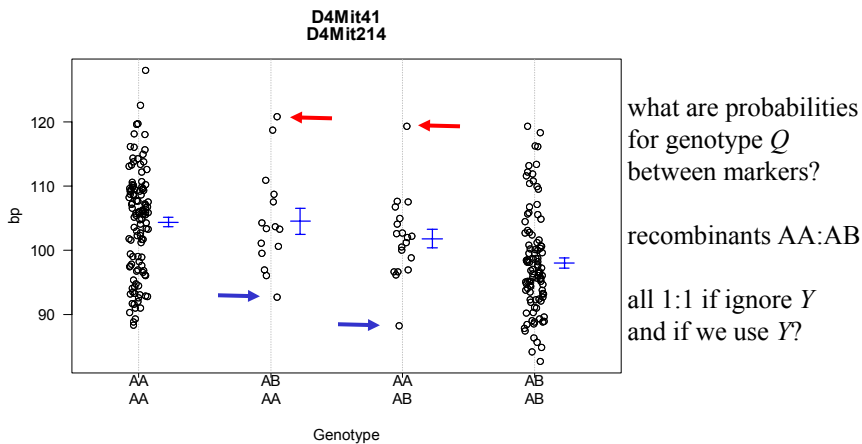
## $\text{pr}(q|m, \lambda)$ recombination model

$$\text{pr}(q|m, \lambda) = \text{pr}(\text{geno} | \text{map}, \text{locus}) \approx \text{pr}(\text{geno} | \text{flanking markers}, \text{locus})$$





## how does phenotype $Y$ improve posterior for genotype $Q$ ?



## posterior on QTL genotypes

- full conditional for  $q$  depends data for individual  $i$ 
  - proportional to prior  $\text{pr}(q | m_i, \lambda)$ 
    - weight toward  $q$  that agrees with flanking markers
  - proportional to likelihood  $\text{pr}(y_i | q, \mu)$ 
    - weight toward  $q$  so that group mean  $\mu_q \approx y_i$
- phenotype and prior recombination may conflict
  - posterior recombination balances these two weights
  - this is “E step” in EM for classical QTL analysis

$$\text{pr}(q | y_i, m_i, \mu, \lambda) = \frac{\text{pr}(q | m_i, \lambda) \text{pr}(y_i | q, \mu)}{\text{pr}(y_i | m_i, \mu, \lambda)}$$

## Bayesian model posterior

- augment data  $(y, m)$  with unknowns  $q$
- study unknowns  $(\mu, \lambda, q)$  given data  $(y, m)$ 
  - properties of posterior  $\text{pr}(\mu, \lambda, q | y, m)$
- sample from posterior in some clever way
  - multiple imputation or MCMC

$$\text{pr}(q, \mu, \lambda | y, m) = \frac{\text{pr}(y | q, \mu) \text{pr}(q | m, \lambda) \text{pr}(\mu) \text{pr}(\lambda | m)}{\text{pr}(y | m)}$$

$$\text{pr}(\mu, \lambda | y, m) = \text{sum}_q \text{pr}(q, \mu, \lambda | y, m)$$

## Bayesian priors for QTL

- missing genotypes  $q$ 
  - $\text{pr}(q | m, \lambda)$
  - recombination model is formally a prior
- effects  $(\mu, \sigma^2)$ 
  - prior =  $\text{pr}(\mu_q | \sigma^2) \text{pr}(\sigma^2)$
  - use conjugate priors for normal phenotype
    - $\text{pr}(\mu_q | \sigma^2) = \text{normal}$
    - $\text{pr}(\sigma^2) = \text{inverse chi-square}$
- each locus  $\lambda$  may be uniform over genome
  - $\text{pr}(\lambda | m) = 1 / \text{length of genome}$
- combined prior
  - $\text{pr}(q, \mu, \lambda | m) = \text{pr}(q | m, \lambda) \text{pr}(\mu) \text{pr}(\lambda | m)$

### 3. Markov chain sampling of architectures

- construct Markov chain around posterior
  - want posterior as stable distribution of Markov chain
  - in practice, the chain tends toward stable distribution
    - initial values may have low posterior probability
    - burn-in period to get chain mixing well
- hard to sample  $(q, \mu, \lambda, H)$  from joint posterior
  - update  $(q, \mu, \lambda)$  from full conditionals for model  $H$
  - update genetic architecture  $H$

$$(q, \mu, \lambda, H) \sim \text{pr}(q, \mu, \lambda, H | y, m)$$

$$(q, \mu, \lambda, H)_1 \rightarrow (q, \mu, \lambda, H)_2 \rightarrow \dots \rightarrow (q, \mu, \lambda, H)_N$$

### MCMC sampling of $(\lambda, q, \mu)$

- Gibbs sampler
  - genotypes  $q$
  - effects  $\mu$
  - *not* loci  $\lambda$

$$q \sim \text{pr}(q | y_i, m_i, \mu, \lambda)$$

$$\mu \sim \frac{\text{pr}(y | q, \mu) \text{pr}(\mu)}{\text{pr}(y | q)}$$

$$\lambda \sim \frac{\text{pr}(q | m, \lambda) \text{pr}(\lambda | m)}{\text{pr}(q | m)}$$



- Metropolis-Hastings sampler
  - extension of Gibbs sampler
  - does not require normalization
    - $\text{pr}(q | m) = \sum_{\lambda} \text{pr}(q | m, \lambda) \text{pr}(\lambda)$

## full conditional for locus

- cannot easily sample from locus full conditional
$$\begin{aligned}\text{pr}(\lambda | y, m, \mu, q) &= \text{pr}(\lambda | m, q) \\ &= \text{pr}(q | m, \lambda) \text{pr}(\lambda) / \text{constant}\end{aligned}$$
- constant is very difficult to compute explicitly
  - must average over all possible loci  $\lambda$  over genome
  - must do this for every possible genotype  $q$
- Gibbs sampler will not work in general
  - but can use method based on ratios of probabilities
  - Metropolis-Hastings is extension of Gibbs sampler

## Gibbs sampler idea

- toy problem
  - want to study two correlated effects
  - could sample directly from their bivariate distribution
- instead use Gibbs sampler:
  - sample each effect from its full conditional given the other
  - pick order of sampling at random
  - repeat many times

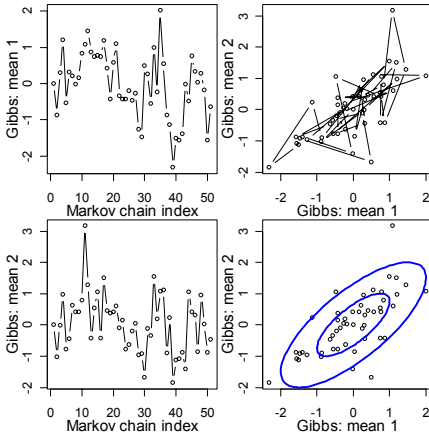
$$\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \sim N \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$$

$$\mu_1 \sim N(\rho\mu_2, 1 - \rho^2)$$

$$\mu_2 \sim N(\rho\mu_1, 1 - \rho^2)$$

# Gibbs sampler samples: $\rho = 0.6$

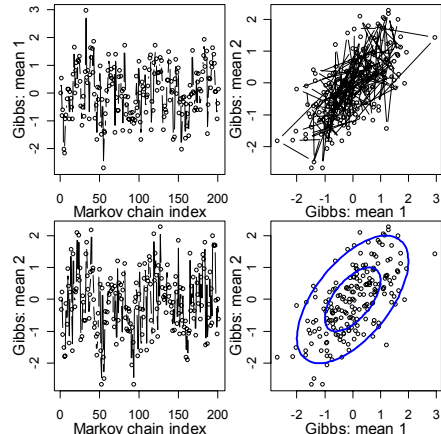
$N = 50$  samples



QTL 2: Bayes

Seattle SISG: Yandell © 2006

$N = 200$  samples

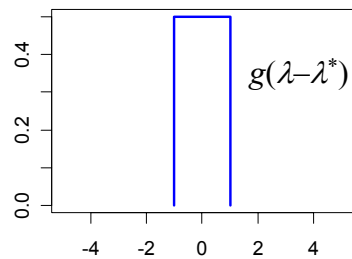
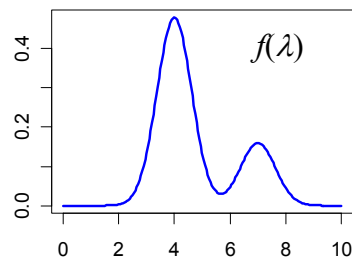


25

# Metropolis-Hastings idea

- want to study distribution  $f(\lambda)$ 
  - take Monte Carlo samples
    - unless too complicated
  - take samples using ratios of  $f$
- Metropolis-Hastings samples:
  - propose new value  $\lambda^*$ 
    - near (?) current value  $\lambda$
    - from some distribution  $g$
  - accept new value with prob  $a$ 
    - Gibbs sampler:  $a = 1$  always

$$a = \min\left(1, \frac{f(\lambda^*)g(\lambda - \lambda^*)}{f(\lambda)g(\lambda^* - \lambda)}\right)$$

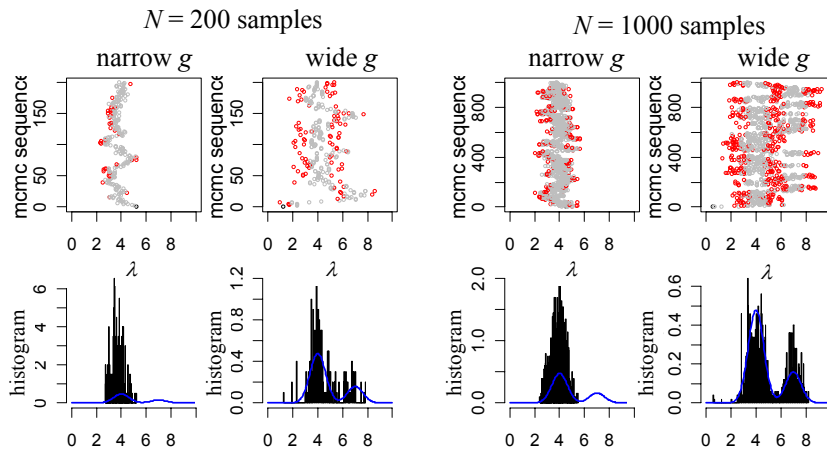


QTL 2: Bayes

Seattle SISG: Yandell © 2006

26

# Metropolis-Hastings samples



## 4. sampling across architectures

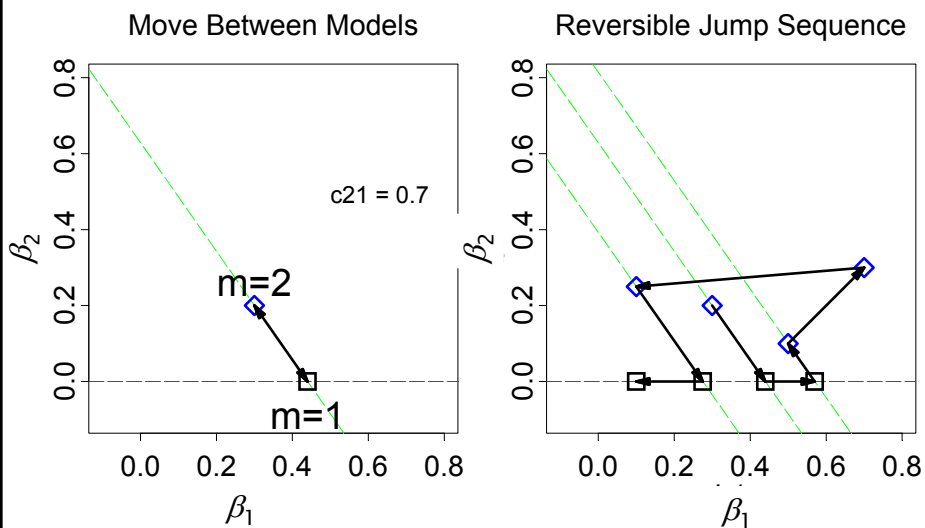
- search across genetic architectures  $M$  of various sizes
  - allow change in number of QTL
  - allow change in types of epistatic interactions
- methods for search
  - reversible jump MCMC
  - Gibbs sampler with loci indicators
- complexity of epistasis
  - Fisher-Cockerham effects model
  - general multi-QTL interaction & limits of inference

## model selection in regression

- consider known genotypes  $q$  at 2 known loci  $\lambda$ 
  - models with 1 or 2 QTL
- jump between 1-QTL and 2-QTL models
- adjust parameters when model changes
  - $\beta_{q1}$  estimate changes between models 1 and 2
  - due to collinearity of QTL genotypes

$$\begin{array}{l} \curvearrowright m = 1 : \mu_q = \mu + \beta_{q1} \\ \curvearrowright m = 2 : \mu_q = \mu + \beta_{q1} + \beta_{q2} \end{array}$$

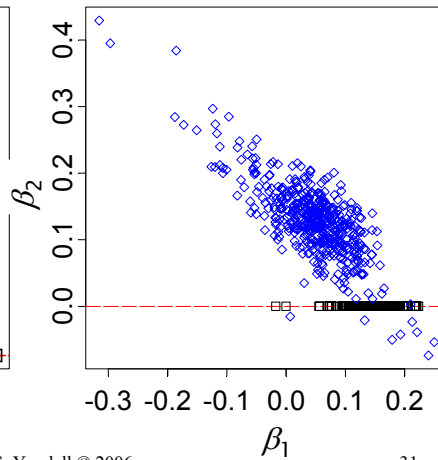
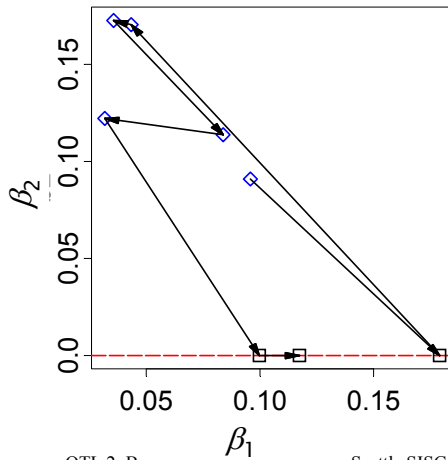
## geometry of reversible jump



## geometry allowing $q$ and $\lambda$ to change

a short sequence

first 1000 with  $m < 3$



QTL 2: Bayes

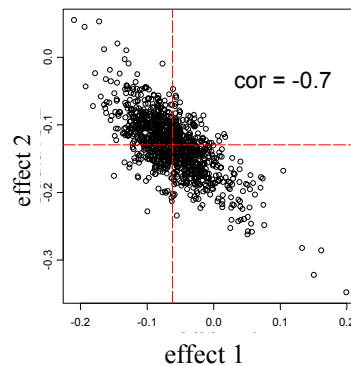
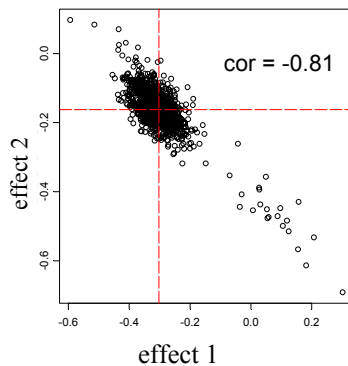
Seattle SISG: Yandell © 2006

31

## collinear QTL = correlated effects

4-week

8-week



- linked QTL = collinear genotypes
  - correlated estimates of effects (negative if in coupling phase)
  - sum of linked effects usually fairly constant

QTL 2: Bayes

Seattle SISG: Yandell © 2006

32



# reversible jump MCMC idea



- Metropolis-Hastings updates: draw one of three choices
  - update  $m$ -QTL model with probability  $1-b(m+1)-d(m)$ 
    - update current model using full conditionals
    - sample  $m$  QTL loci, effects, and genotypes
  - add a locus with probability  $b(m+1)$ 
    - propose a new locus and innovate new genotypes & genotypic effect
    - decide whether to accept the “birth” of new locus
  - drop a locus with probability  $d(m)$ 
    - propose dropping one of existing loci
    - decide whether to accept the “death” of locus
- Satagopan Yandell (1996, 1998); Sillanpaa Arjas (1998); Stevens Fisch (1998)
  - these build on RJ-MCMC idea of Green (1995); Richardson Green (1997)

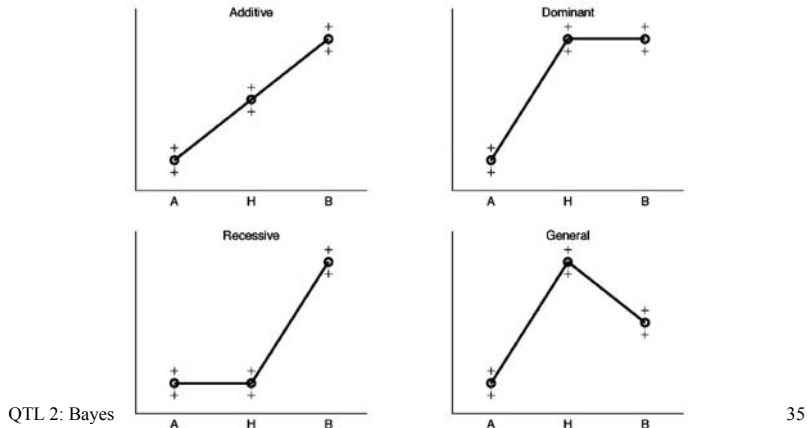
# Gibbs sampler with loci indicators

- consider only QTL at pseudomarkers
  - every 1-2 cM
  - modest approximation with little bias
- use loci indicators in each pseudomarker
  - $\delta = 1$  if QTL present
  - $\delta = 0$  if no QTL present
- Gibbs sampler on loci indicators  $\delta$ 
  - relatively easy to incorporate epistasis
  - Yi, Yandell, Churchill, Allison, Eisen, Pomp (2005 *Genetics*)
    - (see earlier work of Nengjun Yi and Ina Hoeschele)

$$\mu_q = \mu + \delta_1 \beta_{q1} + \delta_2 \beta_{q2}$$

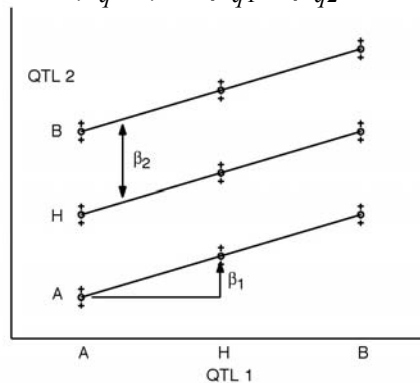
## 5. Gene Action and Epistasis

additive, dominant, recessive, general effects  
of a single QTL (Gary Churchill)



## additive effects of two QTL (Gary Churchill)

$$\mu_q = \mu + \beta_{q1} + \beta_{q2}$$



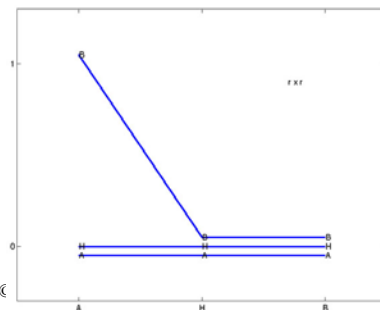
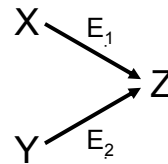
# Epistasis (Gary Churchill)

The allelic state at one locus can mask or uncover the effects of allelic variation at another.

- W. Bateson, 1907.

## epistasis in parallel pathways (GAC)

- Z keeps trait value low
- neither  $E_1$  nor  $E_2$  is rate limiting
- loss of function alleles are segregating from parent A at  $E_1$  and from parent B at  $E_2$



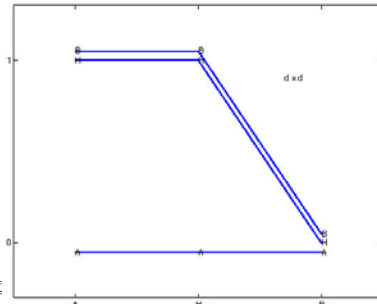
## epistasis in a serial pathway (GAC)

- Z keeps trait value high



- neither  $E_1$  nor  $E_2$  is rate limiting

- loss of function alleles are segregating from parent B at  $E_1$  and from parent A at  $E_2$



QTL 2: Bayes

Seattle SISG: Yandell ©

## QTL with epistasis

- same phenotype model overview

$$y = \mu_q + e, \text{var}(e) = \sigma^2$$

- partition of genotypic value with epistasis

$$\mu_q = \mu + \beta_{q1} + \beta_{q2} + \beta_{q12}$$

- partition of genetic variance & heritability

$$\text{var}(\mu_q) = \sigma_G^2 = \sigma_1^2 + \sigma_2^2 + \sigma_{12}^2$$

$$h^2 = \frac{\sigma_G^2}{\sigma_G^2 + \sigma^2} = h_1^2 + h_2^2 + h_{12}^2$$

QTL 2: Bayes

Seattle SISG: Yandell © 2006

40

# epistatic interactions

- model space issues
  - 2-QTL interactions only?
    - or general interactions among multiple QTL?
  - partition of effects
    - Fisher-Cockerham or tree-structured or ?
- model search issues
  - epistasis between significant QTL
    - check all possible pairs when QTL included?
    - allow higher order epistasis?
  - epistasis with non-significant QTL
    - whole genome paired with each significant QTL?
    - pairs of non-significant QTL?
- Yi Xu (2000) *Genetics*; Yi, Xu, Allison (2003) *Genetics*; Yi *et al.* (2005) *Genetics*

# limits of epistatic inference

- power to detect effects
  - epistatic model size grows exponentially
    - $|H| = 3^{nqtl}$  for general interactions
  - power depends on ratio of  $n$  to model size
    - want  $n / |H|$  to be fairly large (say  $> 5$ )
    - $n = 100$ ,  $nqtl = 3$ ,  $n / |H| \approx 4$
- empty cells mess up adjusted (Type 3) tests
  - missing  $q_1Q_2 / q_1Q_2$  or  $q_1Q_2q_3 / q_1Q_2q_3$  genotype
  - null hypotheses not what you would expect
  - can confound main effects and interactions
  - can bias AA, AD, DA, DD partition

## 6. comparing QTL models

- balance model fit with model "complexity"
  - want maximum likelihood
  - without too complicated a model
- information criteria quantifies the balance
  - Bayes information criteria (BIC) for likelihood
  - Bayes factors for Bayesian approach

## Bayes factors & BIC

$$B_{12} = \frac{\text{pr}(\text{model}_1 | Y) / \text{pr}(\text{model}_2 | Y)}{\text{pr}(\text{model}_1) / \text{pr}(\text{model}_2)} = \frac{\text{pr}(Y | \text{model}_1)}{\text{pr}(Y | \text{model}_2)}$$

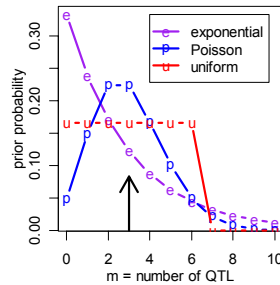
- what is a Bayes factor?
  - ratio of posterior odds to prior odds
  - ratio of model likelihoods
- BF is equivalent to  $LR$  statistic when
  - comparing two nested models
  - simple hypotheses (e.g. 1 vs 2 QTL)
- BF is equivalent to Bayes Information Criteria (BIC)
  - for general comparison of any models
  - want Bayes factor to be substantially larger than 1 (say 10 or more)

$$-2 \log(B_{12}) = -2 \log(LR) - (p_2 - p_1) \log(n)$$

## Bayes factors and genetic model $H$

- $H$  = number of QTL
  - prior  $\text{pr}(H)$  chosen by user
  - posterior  $\text{pr}(H|y,m)$ 
    - sampled marginal histogram
    - shape affected by prior  $\text{pr}(H)$

$$BF_{H,H+1} = \frac{\text{pr}(H|y,m)/\text{pr}(H)}{\text{pr}(H+1|y,m)/\text{pr}(H+1)}$$



- pattern of QTL across genome
- gene action and epistasis

## issues in computing Bayes factors

- $BF$  insensitive to shape of prior on  $nqtl$ 
  - geometric, Poisson, uniform
  - precision improves when prior mimics posterior
- $BF$  sensitivity to prior variance on effects  $\theta$ 
  - prior variance should reflect data variability
  - resolved by using hyper-priors
    - automatic algorithm; no need for user tuning
- easy to compute Bayes factors from samples
  - sample posterior using MCMC
  - posterior  $\text{pr}(nqtl|y,m)$  is marginal histogram