

Bayesian QTL Mapping

Brian S. Yandell

University of Wisconsin-Madison

www.stat.wisc.edu/~yandell/statgen↑

UW-Madison, April 2010

outline

1. What is the goal of QTL study?
2. Bayesian vs. classical QTL study
3. Bayesian strategy for QTLs
4. model search using MCMC
 - Gibbs sampler and Metropolis-Hastings
 - reversible jump MCMC
5. model assessment
 - Bayes factors & model averaging
6. analysis of hyper data
7. software for Bayesian QTLs

1. what is the goal of QTL study?

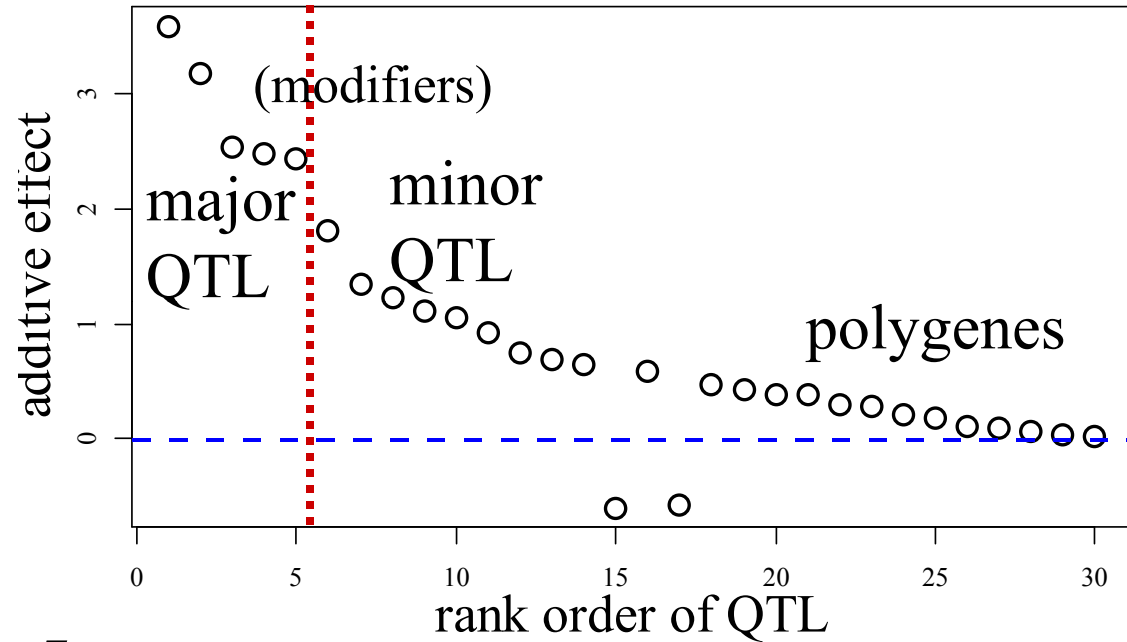
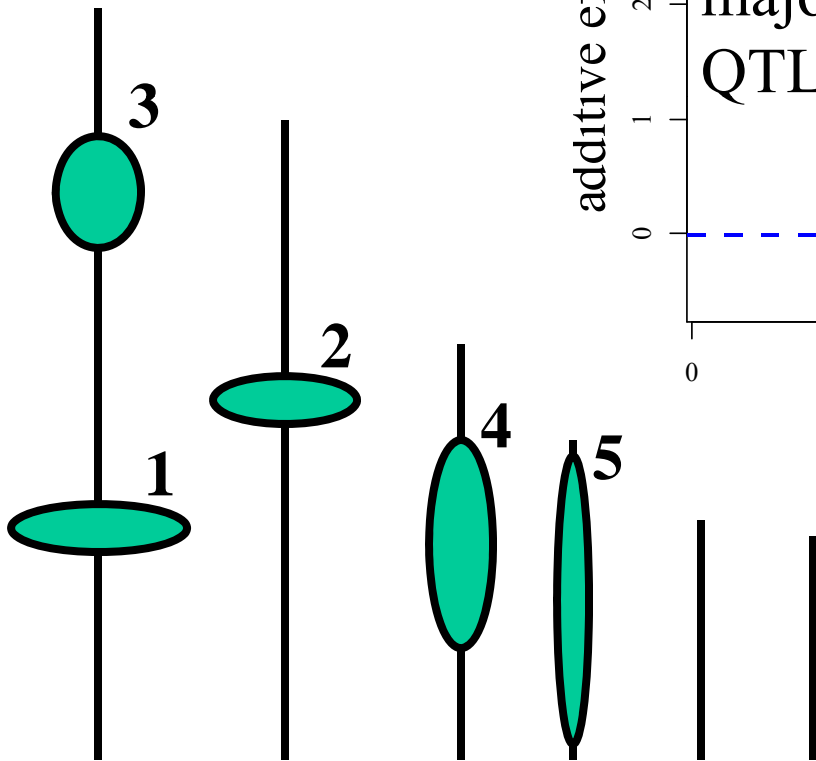
- uncover underlying biochemistry
 - identify how networks function, break down
 - find useful candidates for (medical) intervention
 - epistasis may play key role
 - statistical goal: maximize number of correctly identified QTL
- basic science/evolution
 - how is the genome organized?
 - identify units 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

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



Intuitive idea of ellipses:
Horizontal = significance
Vertical = support interval

check QTL in context of genetic architecture

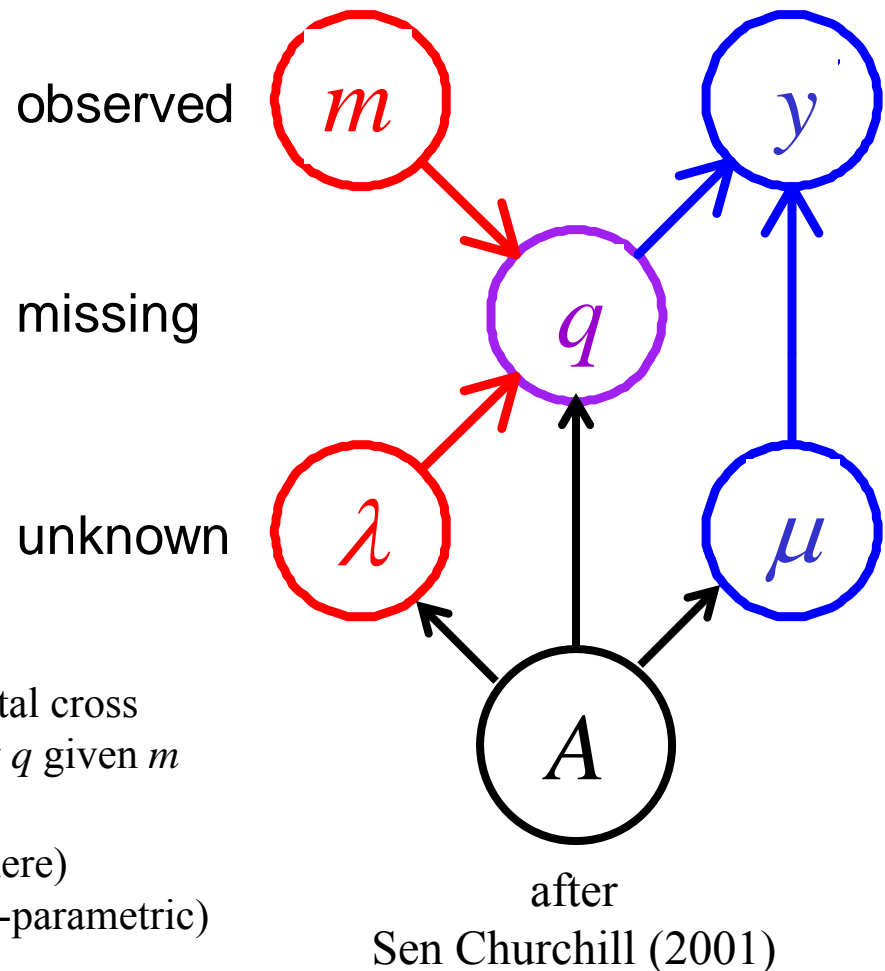
- scan for each QTL adjusting for all others
 - adjust for linked and unlinked QTL
 - adjust for linked QTL: reduce bias
 - adjust for unlinked QTL: reduce variance
 - adjust for environment/covariates
- examine entire genetic architecture
 - number and location of QTL, epistasis, GxE
 - model selection for best genetic architecture

2. Bayesian vs. classical QTL study

- classical study
 - *maximize* over unknown effects
 - *test* for detection of QTL at loci
 - model selection in stepwise fashion
- Bayesian study
 - *average* over unknown effects
 - *estimate* chance of detecting QTL
 - sample all possible models
- both approaches
 - average over missing QTL genotypes
 - scan over possible loci

QTL model selection: key players

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



likelihood and posterior

$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{constant}} : \text{Bayes' rule}$$

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

likelihood mixes over missing QTL genotypes :

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

Bayes posterior vs. maximum likelihood

(genetic architecture $A = \text{single QTL at } \lambda$)

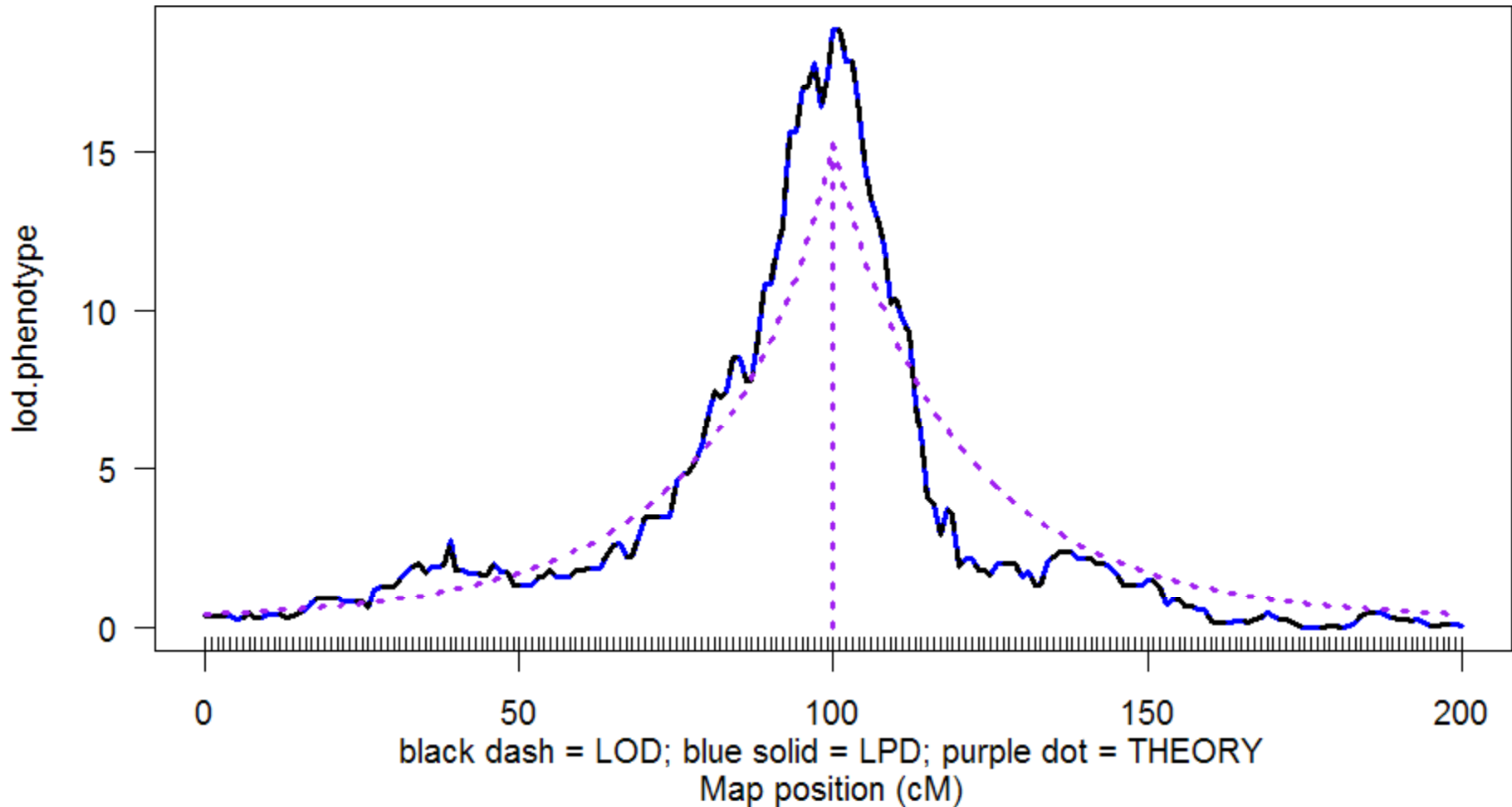
- LOD: classical Log ODds
 - maximize likelihood over effects μ
 - R/qtl scanone/scantwo: method = "em"
- *LPD*: Bayesian *Log Posterior Density*
 - average posterior over effects μ
 - R/qtl scanone/scantwo: method = "imp"

$$\text{LOD}(\lambda) = \log_{10} \left(\max_{\mu} \text{pr}(y | m, \mu, \lambda) \right)$$

$$\text{LPD}(\lambda) = \log_{10} \left(\text{pr}(\lambda | m) \sum_{\mu} \text{pr}(y | m, \mu, \lambda) \text{pr}(\mu) \right)$$

LOD & LPD: 1 QTL

n.ind = 100, 1 cM marker spacing

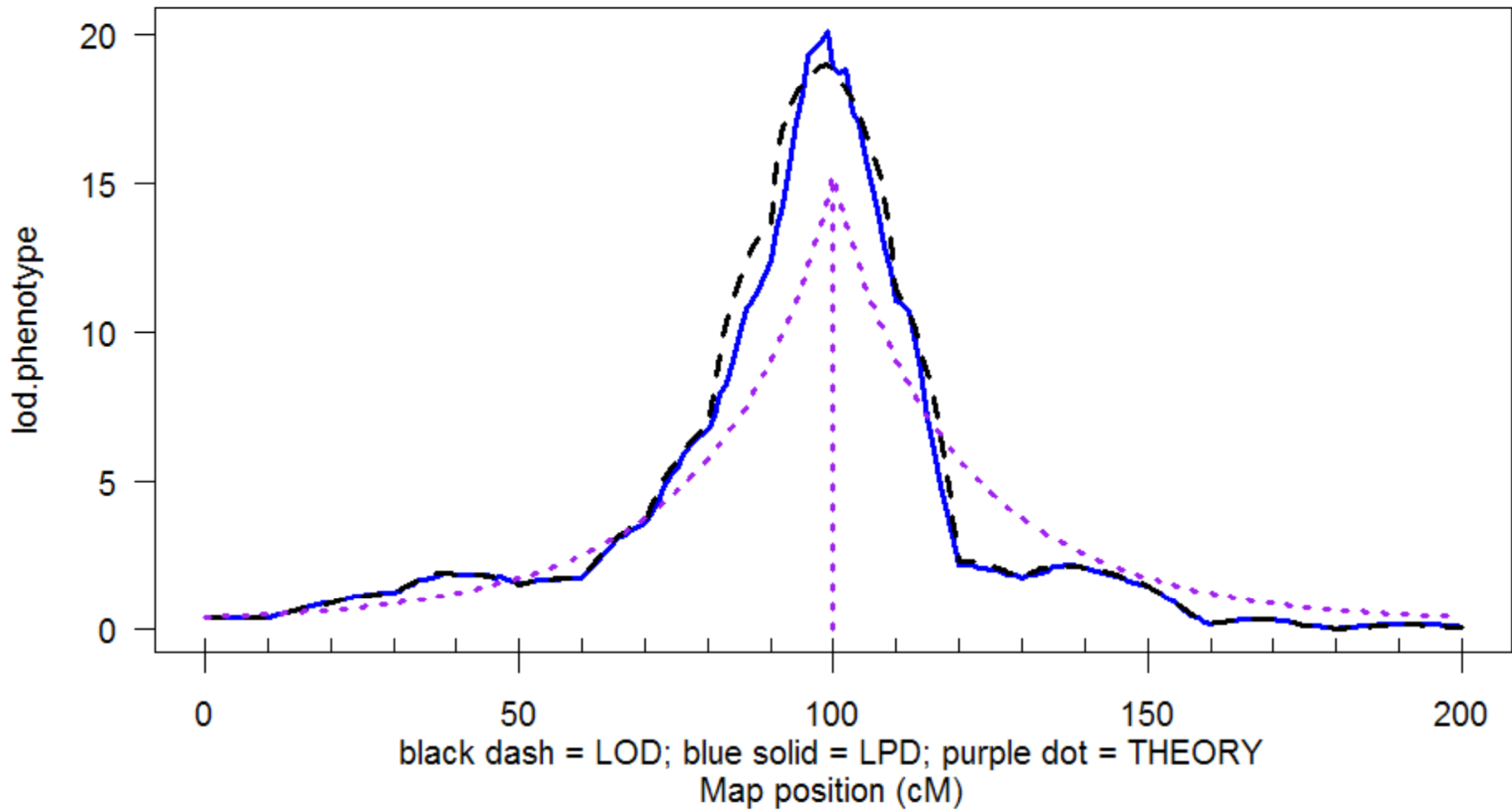


Simplified likelihood surface 2-D for BC locus and effect

- locus λ and effect $\Delta = \mu_2 - \mu_1$
- profile likelihood along ridge
 - maximize likelihood at each λ for Δ
 - symmetric in Δ around MLE given λ
- weighted average of posterior
 - average likelihood at each λ with weight $p(\Delta)$
 - how does prior $p(\Delta)$ affect symmetry?

LOD & LPD: 1 QTL

n.ind = 100, 10 cM marker spacing



likelihood and posterior

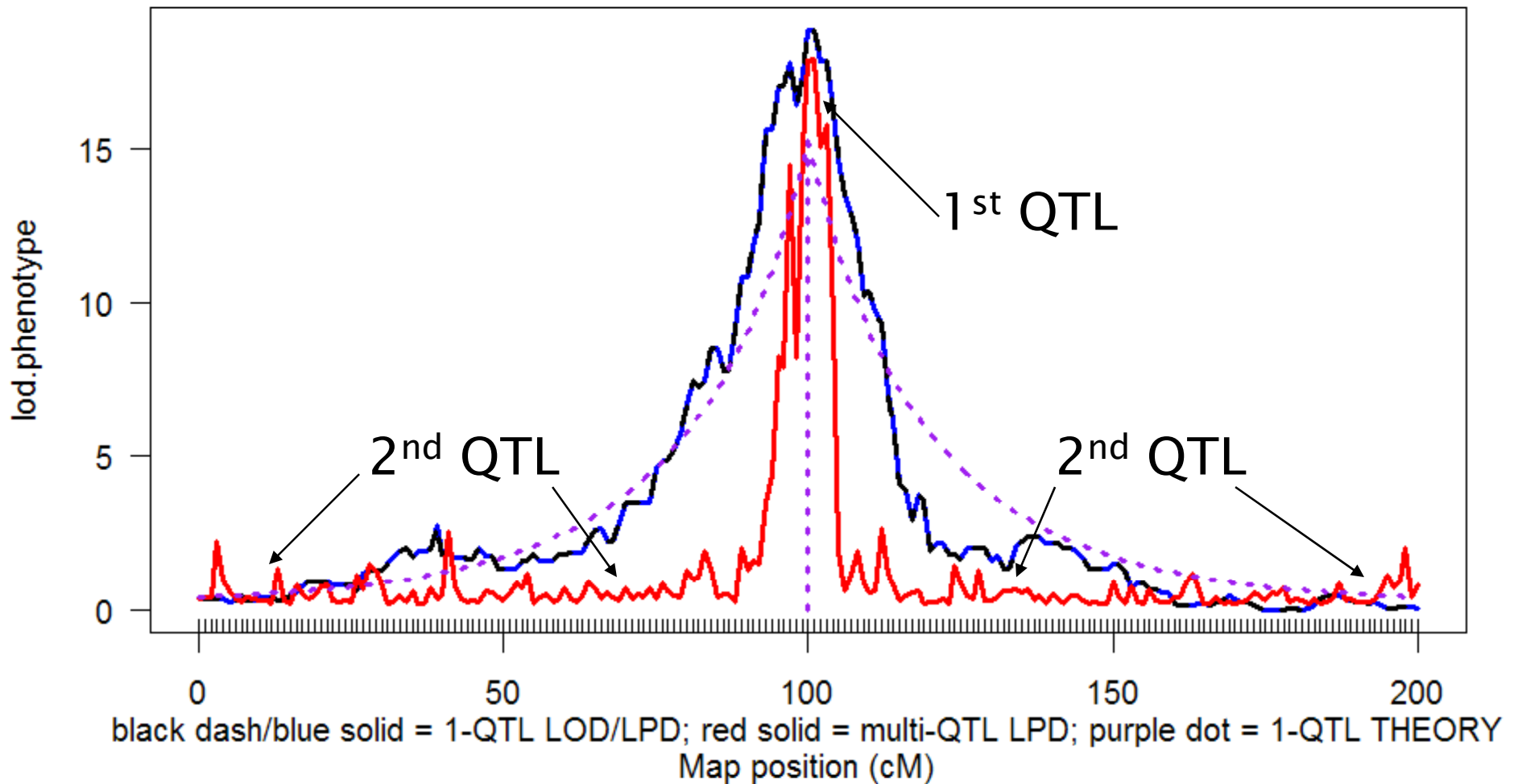
- likelihood relates “known” data (y, m, q) to unknown values of interest (μ, λ, A)
 - $\text{pr}(y, q | m, \mu, \lambda, A) = \text{pr}(y | q, \mu, A) \text{pr}(q | m, \lambda, A)$
 - mix over unknown genotypes (q)
- posterior turns likelihood into a distribution
 - weight likelihood by priors
 - rescale to sum to 1.0
 - posterior = likelihood * prior / constant

marginal LOD or LPD

- What is contribution of a QTL adjusting for all others?
 - improvement in LPD due to QTL at locus λ
 - contribution due to main effects, epistasis, GxE?
- How does adjusted LPD *differ* from unadjusted LPD?
 - raised by removing variance due to unlinked QTL
 - raised or lowered due to bias of linked QTL
 - analogous to Type III adjusted ANOVA tests
- can ask these same questions using classical LOD
 - see Broman's newer tools for multiple QTL inference

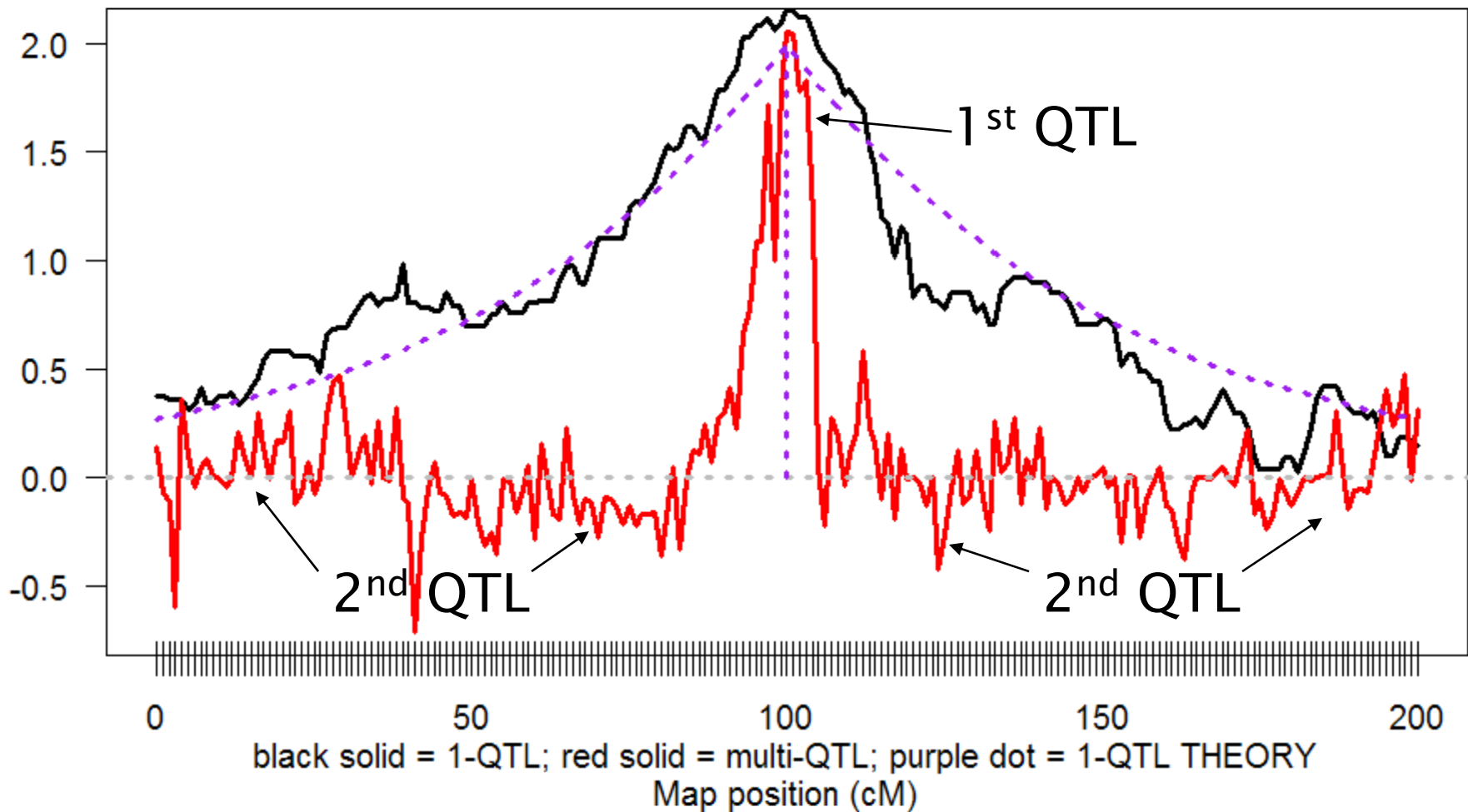
LPD: 1 QTL vs. multi-QTL

marginal contribution to LPD from QTL at λ



substitution effect: 1 QTL vs. multi-QTL

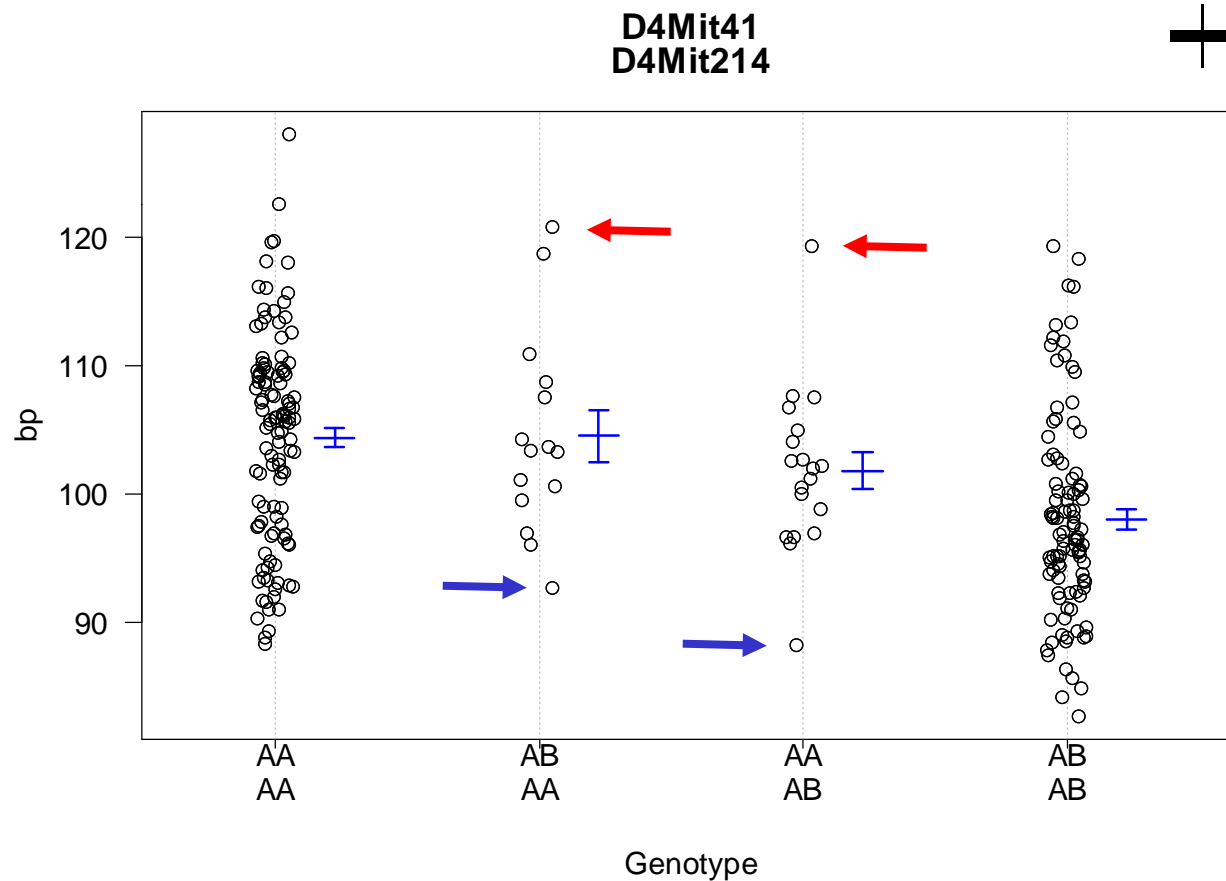
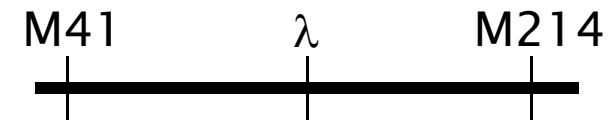
single QTL effect vs. marginal effect from QTL at λ



3. Bayesian strategy for QTLs

- augment data (y, m) with missing genotypes q
- build model for augmented data
 - genotypes (q) evaluated at loci (λ)
 - depends on flanking markers (m)
 - phenotypes (y) centered about effects (μ)
 - depends on missing genotypes (q)
 - λ and μ depend on genetic architecture (A)
 - How complicated is model? number of QTL, epistasis, etc.
- sample from model in some clever way
- infer most probable genetic architecture
 - estimate loci, their main effects and epistasis
 - study properties of estimates

do phenotypes help to guess genotypes? posterior on QTL genotypes q



what are probabilities
for genotype q
between markers?

all recombinants AA:AB
have 1:1 prior ignoring y

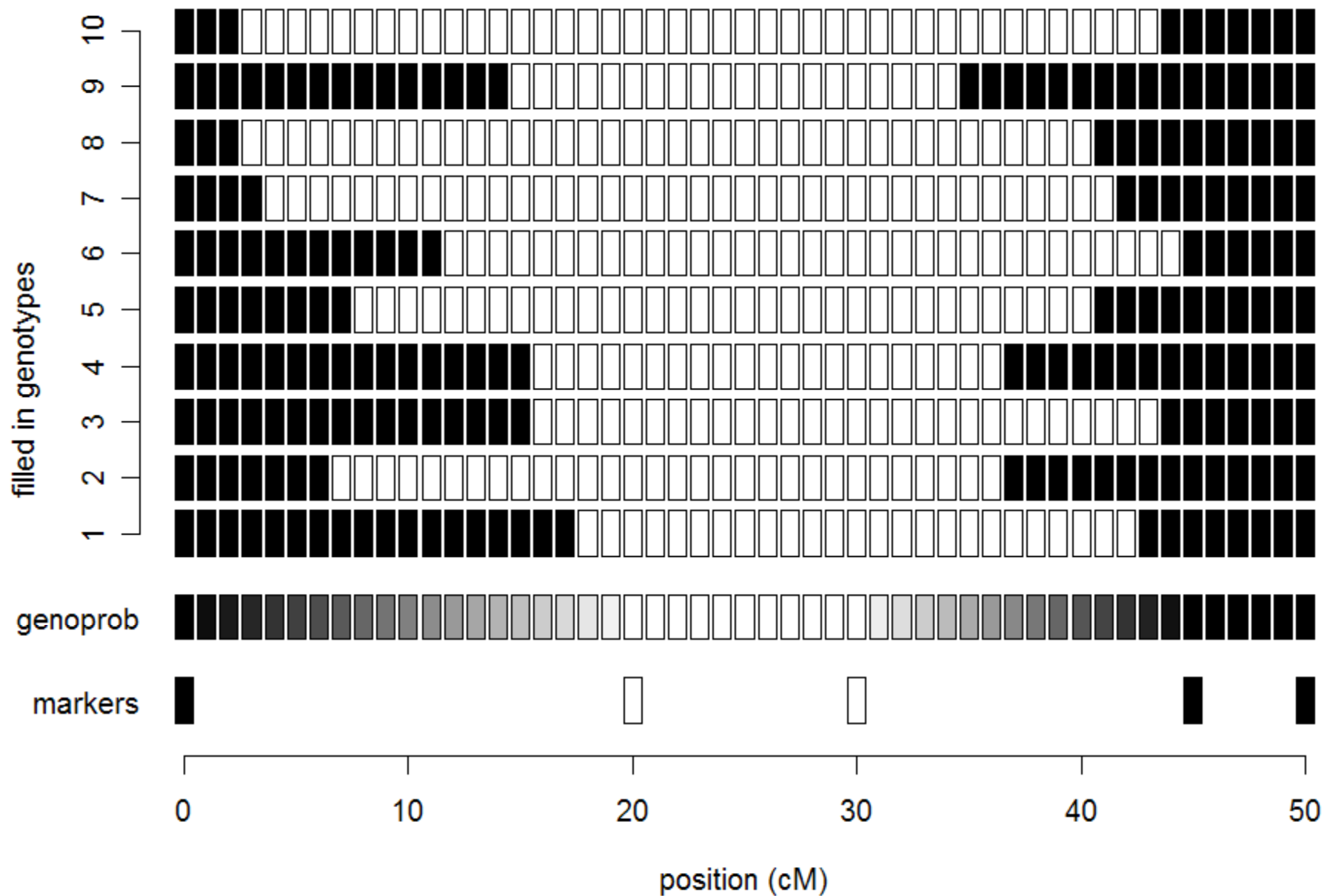
what if we use y ?

posterior on QTL genotypes q

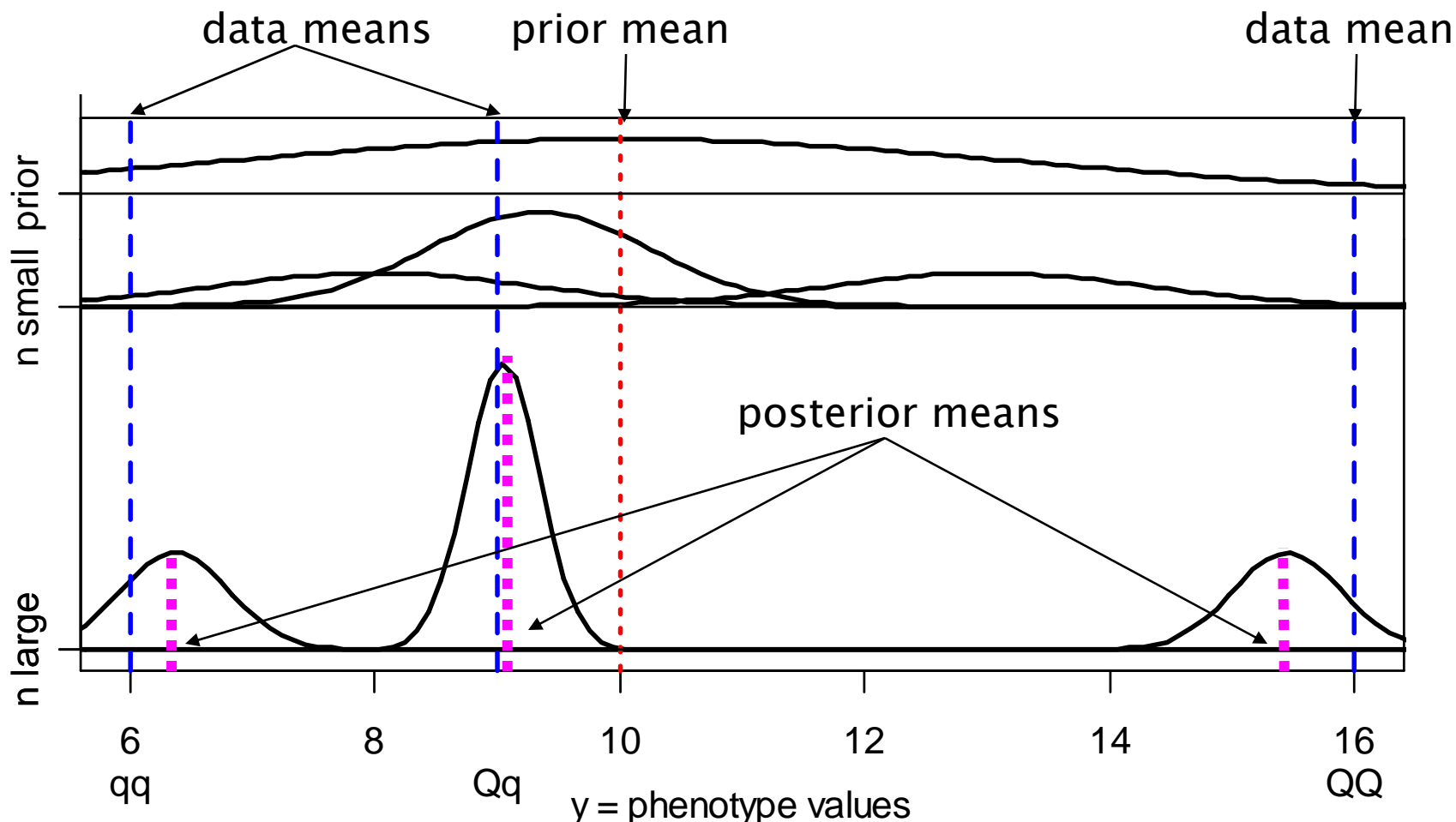
- full conditional of q given data, parameters
 - proportional to prior $\text{pr}(q | m, \lambda)$
 - weight toward q that agrees with flanking markers
 - proportional to likelihood $\text{pr}(y/q, \mu)$
 - weight toward q with similar phenotype values
 - posterior balances these two
- this *is* the E-step of EM computations

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

multiple imputations of genotypes



where are the genotypic means?
 (phenotype mean for genotype q is μ_q)



prior & posteriors: genotypic means μ_q

- prior for genotypic means
 - centered at grand mean
 - variance related to heritability of effect
 - hyper-prior on variance (details omitted)
- posterior
 - shrink genotypic means toward grand mean
 - shrink variance of genotypic mean

$$\begin{array}{lll} \text{prior:} & E(\mu_q) = \bar{y} & V(\mu_q) = V(y)h_q^2 \\ \text{posterior:} & E(\mu_q | y) = \bar{y} \cdot (1 - b_q) + \bar{y}_q b_q & V(\mu_q | y) = V(\bar{y}_q) b_q \\ \text{shrinkage:} & b_q = 1 - \frac{V(\bar{y}_q)}{V(\bar{y}_q) + V(y)h_q^2} \approx 1 & \end{array}$$

multiple QTL phenotype model

- phenotype affected by genotype & environment

$$E(y/q) = \mu_q = \beta_0 + \sum_{\{j \text{ in } H\}} \beta_j(q)$$

number of terms in QTL model $H \leq 2^{n_{qtl}}$ ($3^{n_{qtl}}$ for F_2)

- partition genotypic mean into QTL effects

$$\mu_q = \beta_0 + \beta_1(q_1) + \beta_2(q_2) + \beta_{12}(q_1, q_2)$$

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

- partition prior and posterior (details omitted)

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

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

$$h_q^2 = \frac{\sigma_q^2}{\sigma_q^2 + \sigma^2} = h_1^2 + h_2^2 + h_{12}^2$$

partition of multiple QTL effects

- partition genotype-specific mean into QTL effects

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

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

- priors on mean and effects

$\mu \sim N(\mu_0, \kappa_0 \sigma^2)$ grand mean

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

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

- determine hyper-parameters via empirical Bayes

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

Where are the loci λ on the genome?

- prior over genome for QTL positions
 - flat prior = no prior idea of loci
 - or use prior studies to give more weight to some regions
- posterior depends on QTL genotypes q
$$\text{pr}(\lambda \mid m, q) = \text{pr}(\lambda) \text{pr}(q \mid m, \lambda) / \text{constant}$$
 - constant determined by averaging
 - over all possible genotypes q
 - over all possible loci λ on entire map
- no easy way to write down posterior

model fit with multiple imputation

(Sen and Churchill 2001)

- pick a genetic architecture
 - 1, 2, or more QTL
- fill in missing genotypes at ‘pseudomarkers’
 - use prior recombination model
- use clever weighting (importance sampling)
- compute LPD, effect estimates, etc.

What is the genetic architecture A ?

- components of genetic architecture
 - how many QTL?
 - where are loci (λ)? how large are effects (μ)?
 - which pairs of QTL are epistatic?
- use priors to weight posterior
 - toward guess from previous analysis
 - improve efficiency of sampling from posterior
 - increase samples from architectures of interest

4. QTL Model Search using MCMC

- 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
- sample QTL model components from full conditionals
 - sample locus λ given q, A (using Metropolis-Hastings step)
 - sample genotypes q given λ, μ, y, A (using Gibbs sampler)
 - sample effects μ given q, y, A (using Gibbs sampler)
 - sample QTL model A given λ, μ, y, q (using Gibbs or M-H)

$$(\lambda, q, \mu, A) \sim \text{pr}(\lambda, q, \mu, A \mid y, m)$$

$$(\lambda, q, \mu, A)_1 \rightarrow (\lambda, q, \mu, A)_2 \rightarrow \cdots \rightarrow (\lambda, q, \mu, A)_N$$

MCMC sampling of (λ, q, μ)

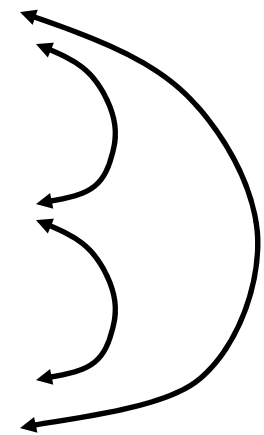
- Gibbs sampler

- genotypes q
- effects μ
- *not* loci λ

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

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

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



- Metropolis-Hastings sampler

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

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

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

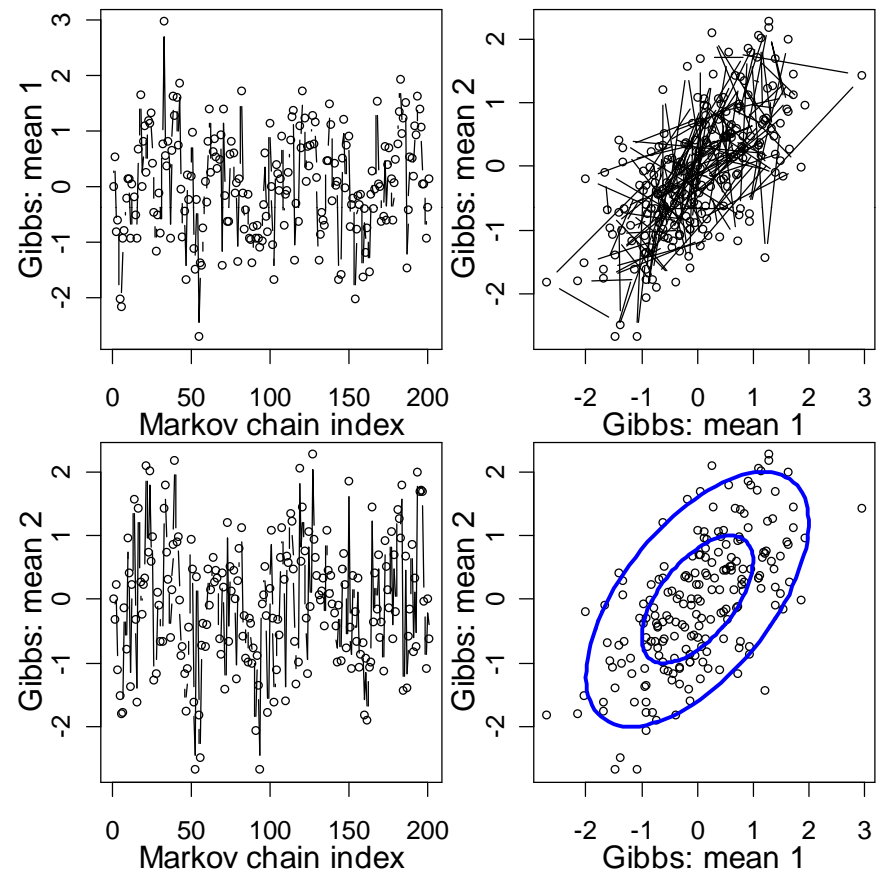
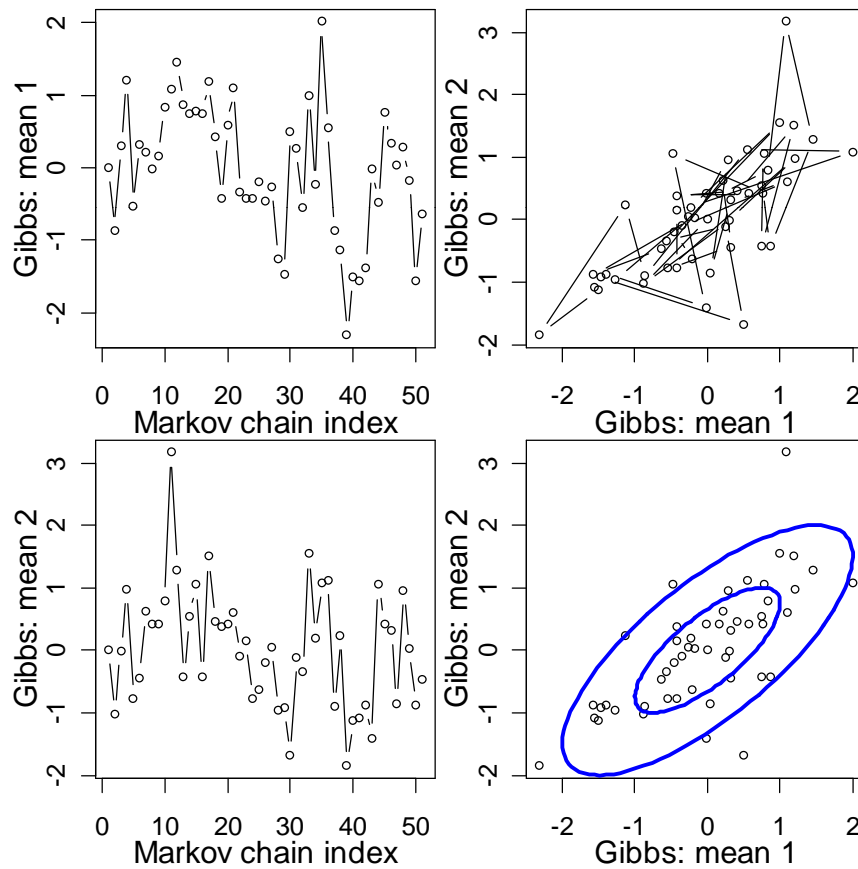
$$\mu_1 \sim N(\rho\mu_2, 1 - \rho^2)$$

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

Gibbs sampler samples: $\rho = 0.6$

$N = 50$ samples

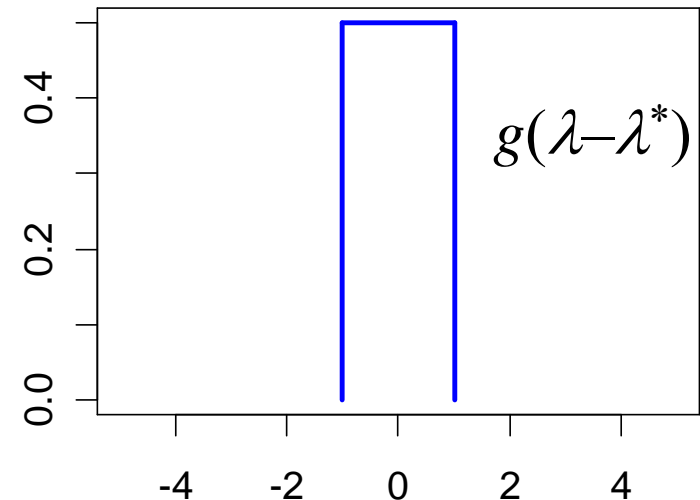
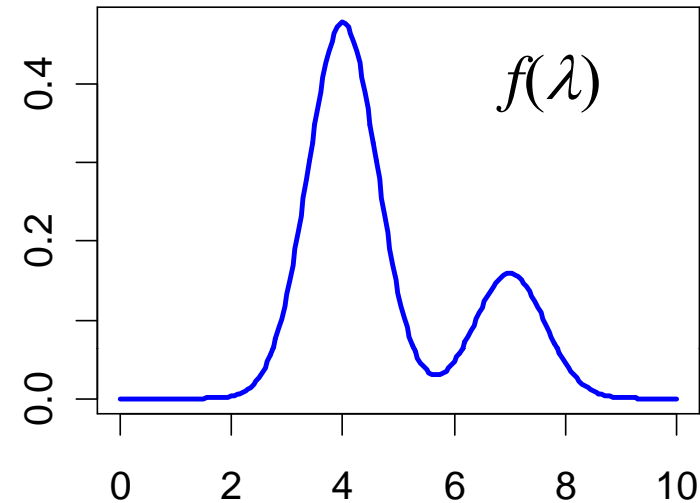
$N = 200$ samples



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 λ^*
 - near (?) current value λ
 - from some distribution g
 - accept new value with prob a
 - Gibbs sampler: $a = 1$ always

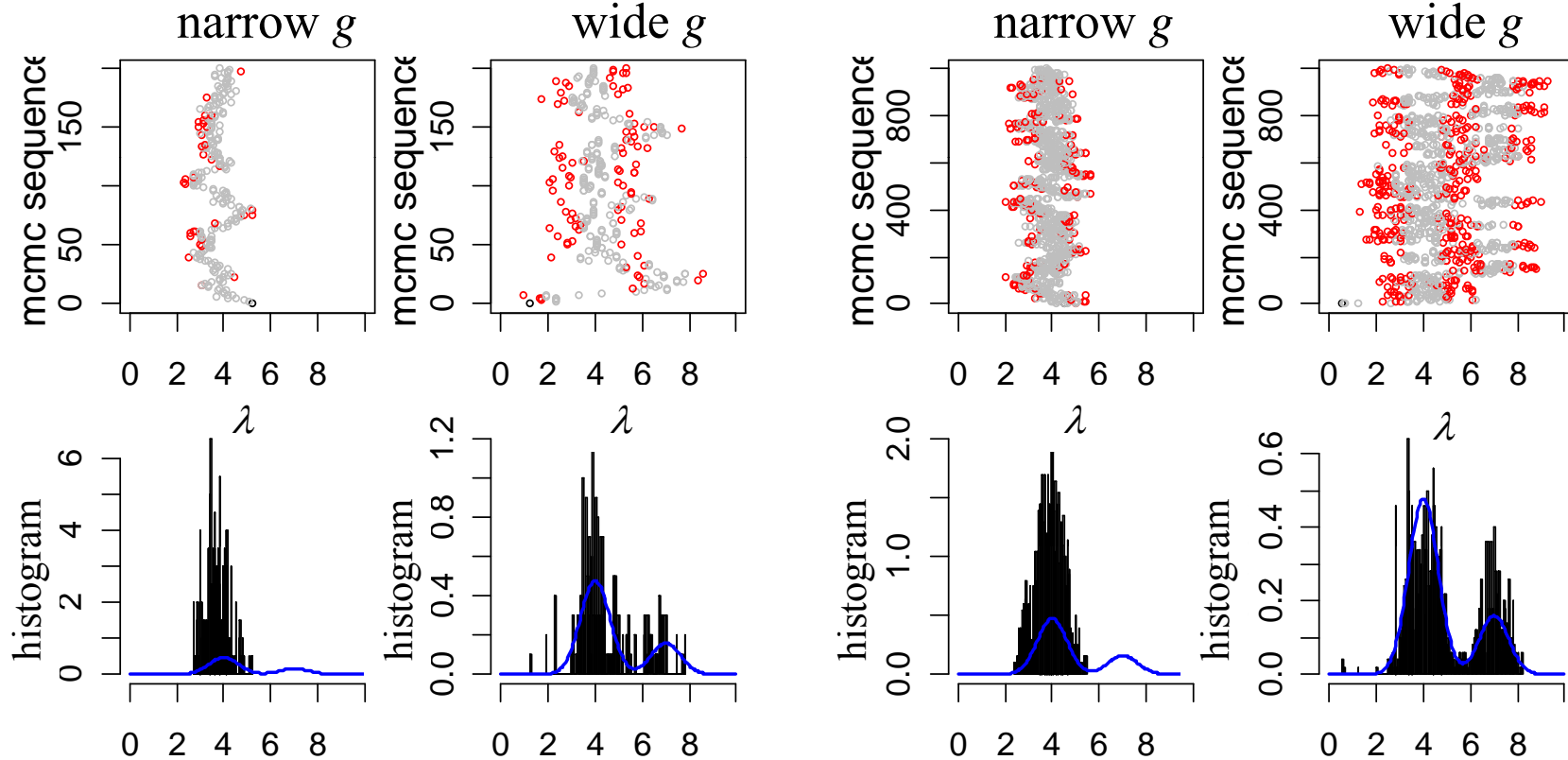
$$a = \min\left(1, \frac{f(\lambda^*)g(\lambda^* - \lambda)}{f(\lambda)g(\lambda - \lambda^*)}\right)$$



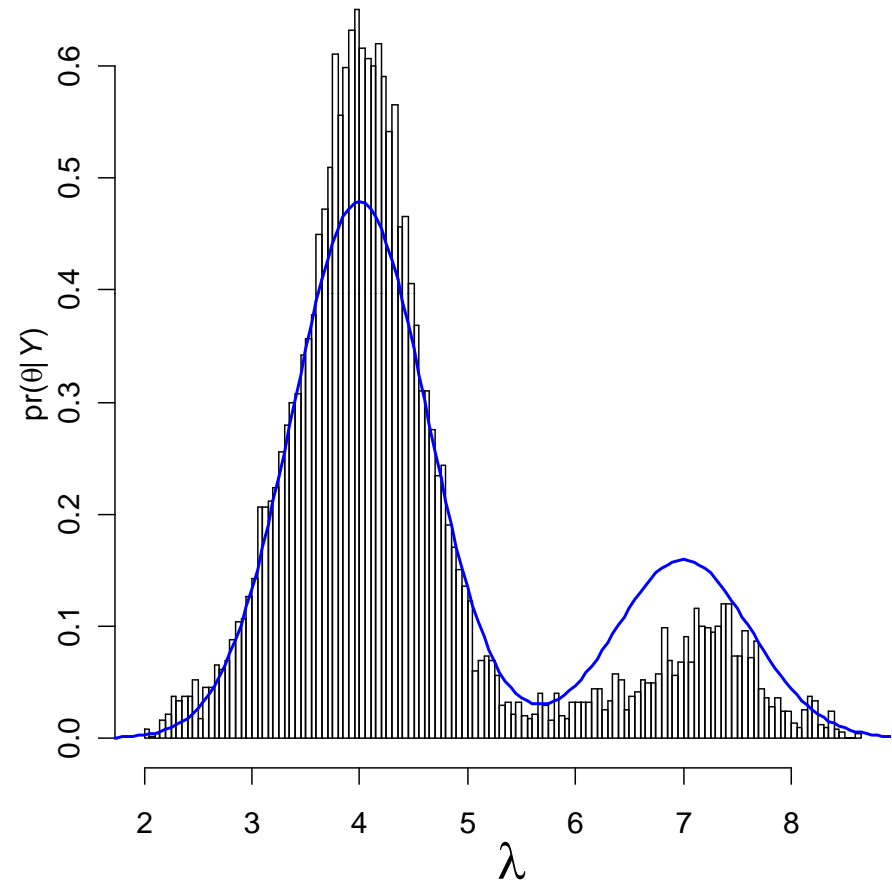
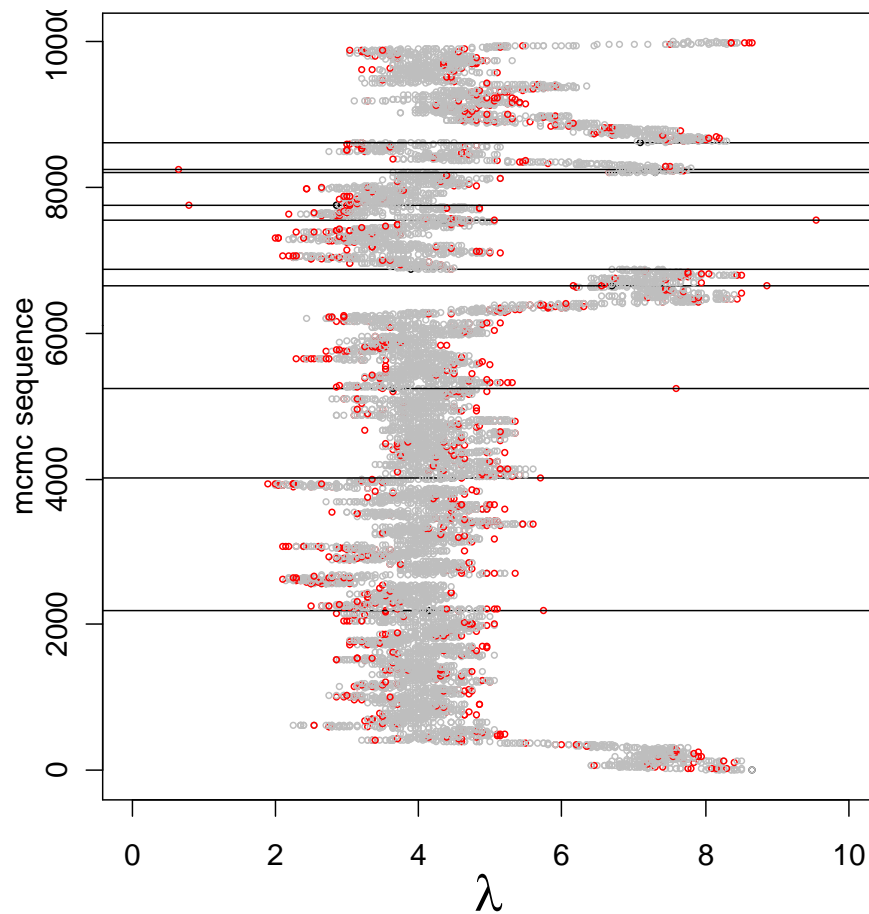
Metropolis-Hastings samples

$N = 200$ samples

$N = 1000$ samples



MCMC realization



added twist: occasionally propose from whole domain

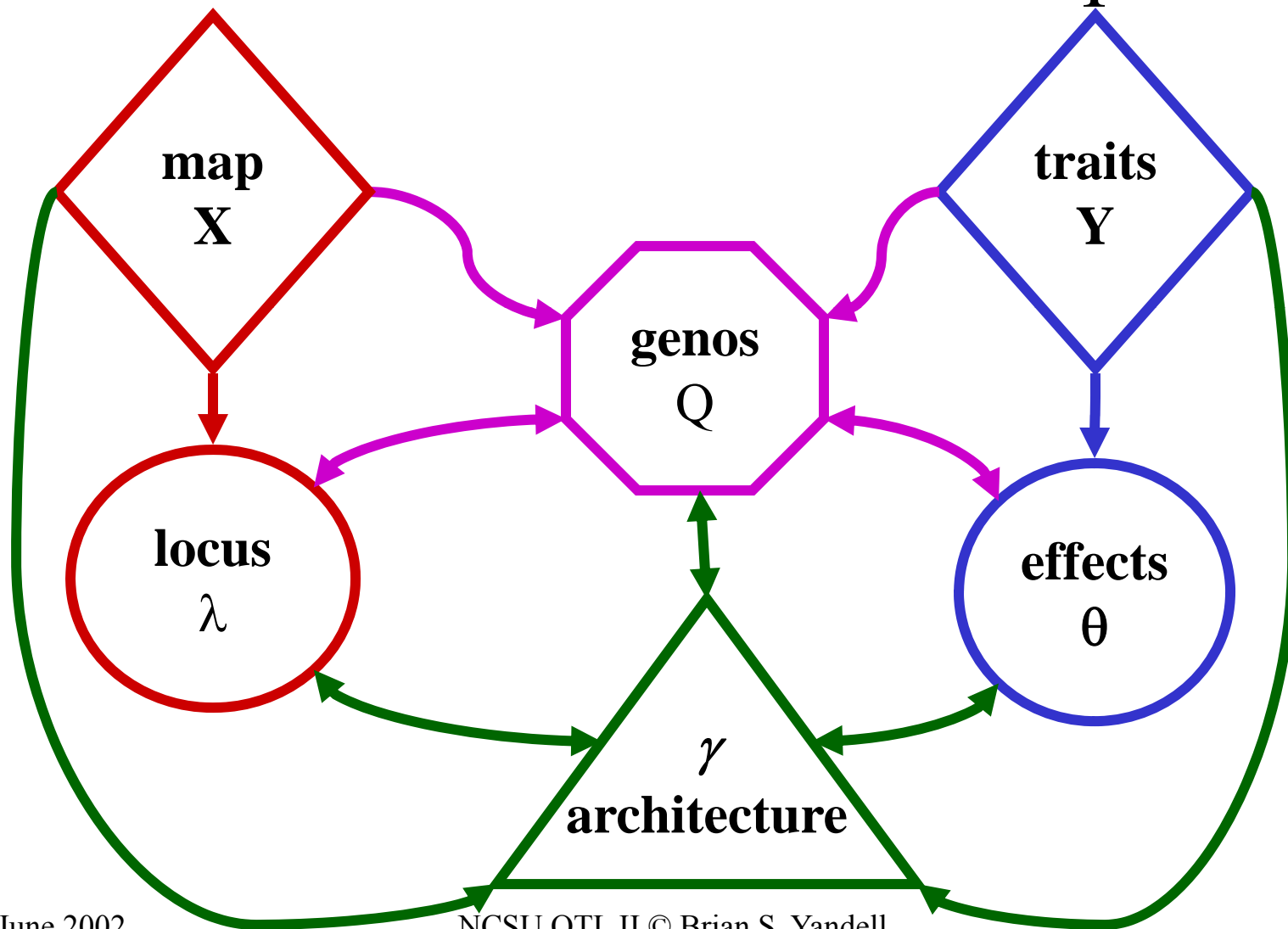
Multiple QTL Phenotype Model

- $E(y) = \mu + \beta(q) = \mu + X\Gamma\beta$
 - $y = n$ phenotypes
 - $X = n \times L$ design matrix
 - in theory covers whole genome of size L cM
 - X determined by genotypes and model space
 - only need terms associated with $q = n \times n_{\text{QTL}}$ genotypes at QTL
 - $\Gamma = \text{diag}(\gamma) =$ genetic architecture
 - $\gamma = 0, 1$ indicators for QTLs or pairs of QTLs
 - $|\gamma| = \sum \gamma =$ size of genetic architecture
 - $\lambda =$ loci determined implicitly by γ
 - $\beta =$ genotypic effects (main and epistatic)
 - $\mu =$ reference

methods of model search

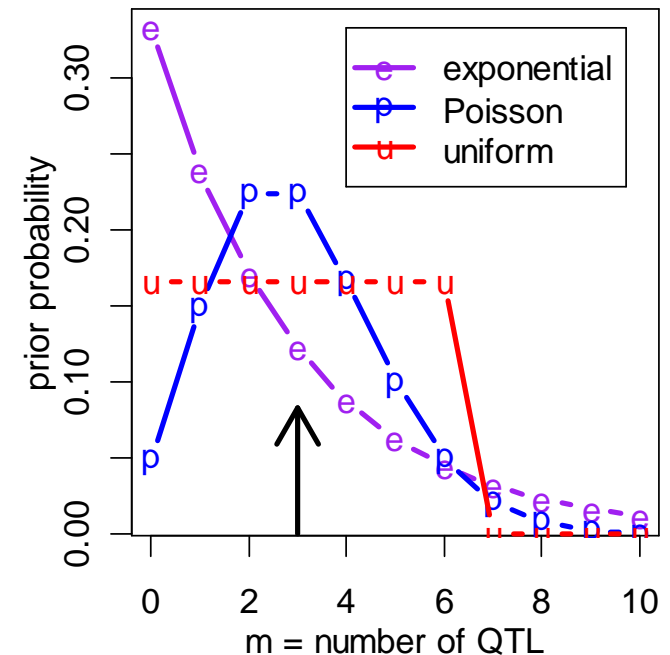
- Reversible jump (transdimensional) MCMC
 - sample possible loci (λ determines possible γ)
 - collapse to model containing just those QTL
 - bookkeeping when model dimension changes
- Composite model with indicators
 - include all terms in model: β and γ
 - sample possible architecture (γ determines λ)
 - can use LASSO-type prior for model selection
- Shrinkage model
 - set $\gamma = 1$ (include all loci)
 - allow variances of β to differ (shrink coefficients to zero)

RJ-MCMC full conditional updates



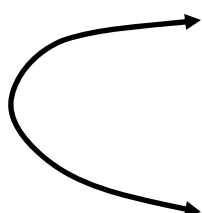
index model by number of QTL

- model changes with number of QTL
 - analogous to stepwise regression if Q known
 - use reversible jump MCMC to change number
 - book keeping to compare models
 - change of variables between models
- what prior on number of QTL?
 - uniform over some range
 - Poisson with prior mean
 - exponential with prior mean



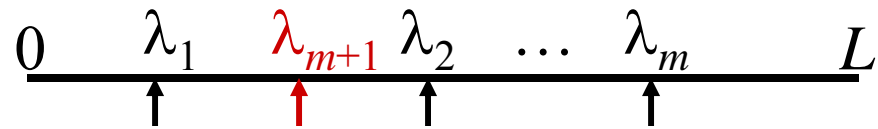
reversible jump MCMC

- consider known genotypes q at 2 known loci λ
 - models with 1 or 2 QTL
- M-H step between 1-QTL and 2-QTL models
 - model changes dimension (via careful bookkeeping)
 - consider mixture over QTL models H


$$nqtl = 1 : Y = \beta_0 + \beta_1(q_1) + e$$

$$nqtl = 2 : Y = \beta_0 + \beta_1(q_1) + \beta_2(q_2) + e$$

sampling across QTL models A

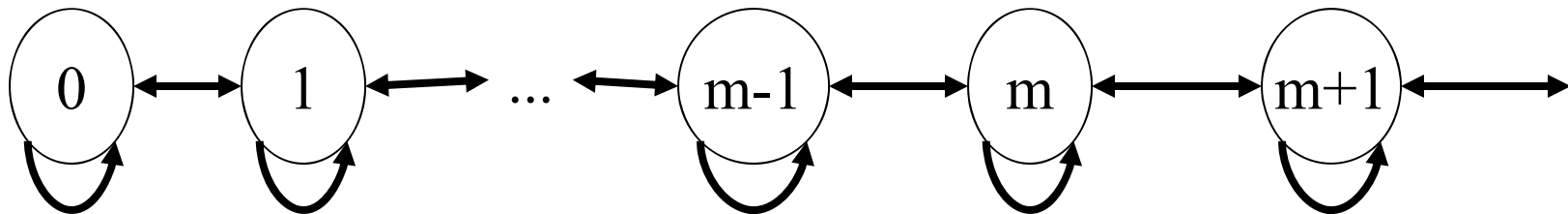
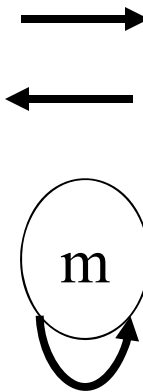


action steps: draw one of three choices

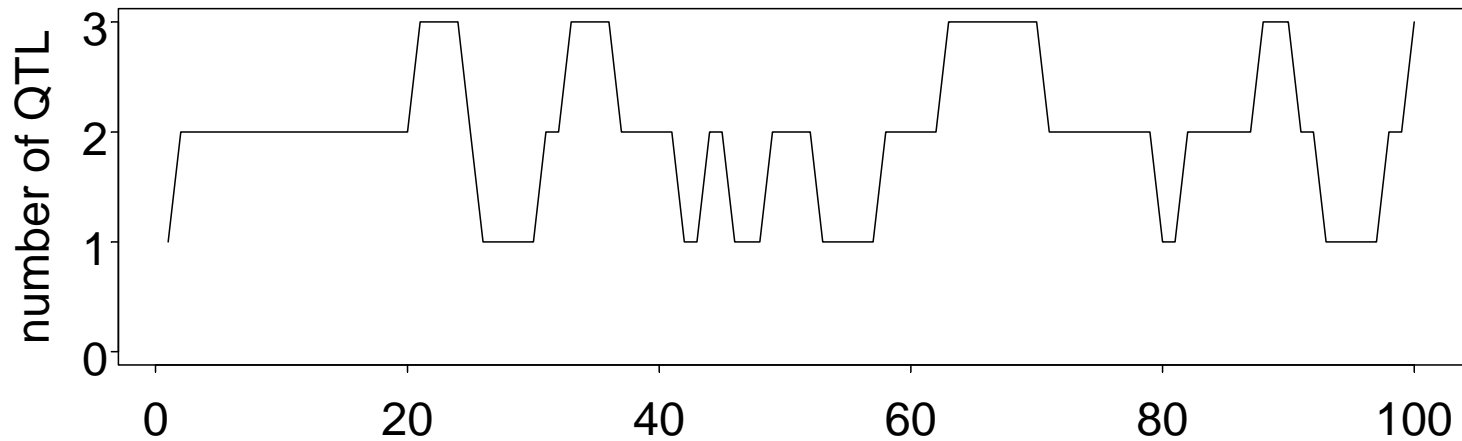
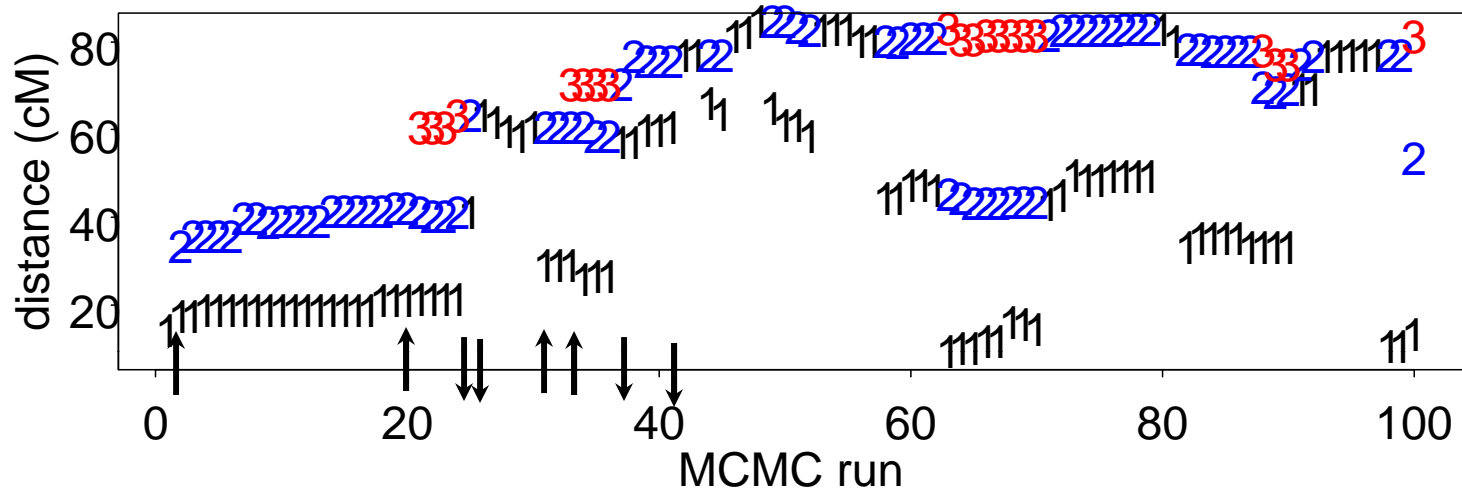
- update QTL model A with probability $1-b(A)-d(A)$
 - update current model using full conditionals
 - sample QTL loci, effects, and genotypes
- add a locus with probability $b(A)$
 - propose a new locus along genome
 - innovate new genotypes at locus and phenotype effect
 - decide whether to accept the “birth” of new locus
- drop a locus with probability $d(A)$
 - propose dropping one of existing loci
 - decide whether to accept the “death” of locus

Markov chain for number m

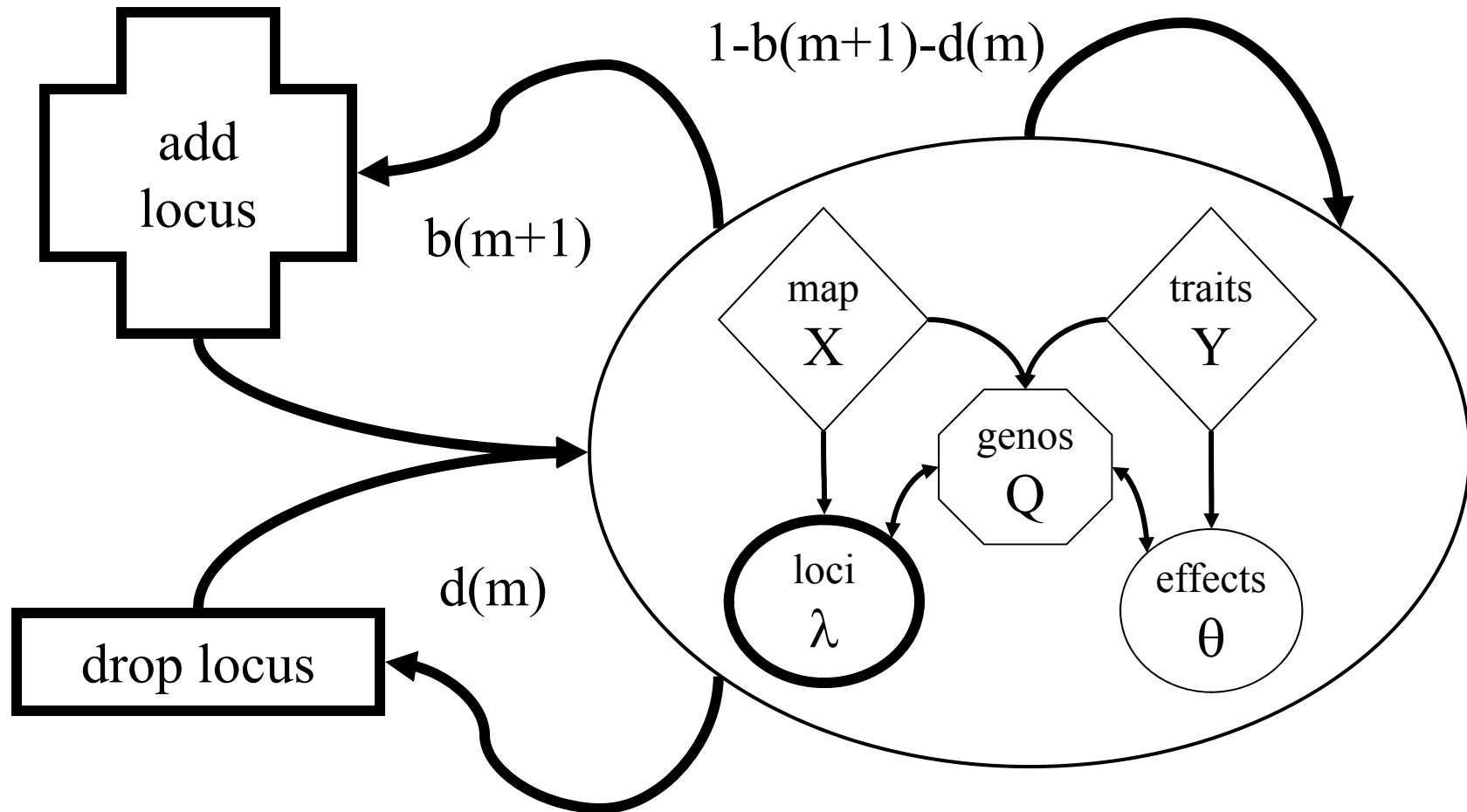
- add a new locus
- drop a locus
- update current model



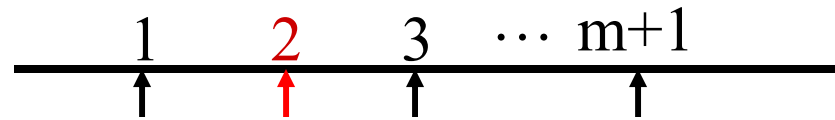
jumping QTL number and loci



RJ-MCMC updates



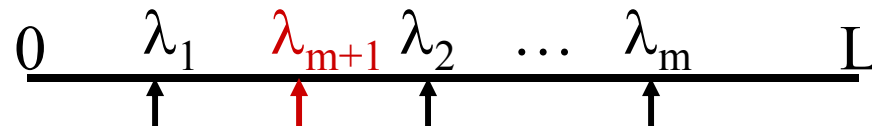
propose to drop a locus



- choose an existing locus
 - equal weight for all loci ?
 - more weight to loci with small effects?
- “drop” effect & genotypes at old locus
 - adjust effects at other loci for collinearity
 - this is reverse jump of Green (1995)
- check acceptance ...
 - do not drop locus, effects & genotypes
 - until move is accepted

$$q_d(r; m + 1) = \frac{1}{m + 1}$$

propose to add a locus



- propose a new locus
 - uniform chance over genome $q_b(\lambda) = 1/L$
 - actually need to be more careful (R van de Ven, pers. comm.)
 - choose interval between loci already in model (include $0, L$)
 - probability proportional to interval length $(\lambda_2 - \lambda_1)/L$
 - uniform chance within this interval $1/(\lambda_2 - \lambda_1)$
 - need genotypes at locus & model effect
- innovate effect & genotypes at new locus
 - draw genotypes based on recombination (prior)
 - no dependence on trait model yet
 - draw effect as in Green's reversible jump
 - adjust for collinearity: modify other parameters accordingly
- check acceptance ...

acceptance of reversible jump

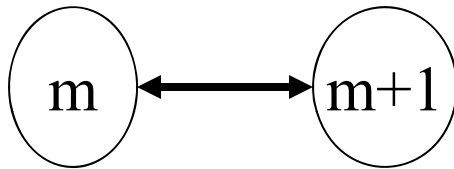
- accept birth of new locus with probability
 $\min(1, A)$
- accept death of old locus with probability
 $\min(1, 1/A)$

$$A = \frac{\text{pr}(\theta_{m+1}, m+1 | Y, X)}{\text{pr}(\theta_m, m | Y, X)} \times \frac{d(m+1)}{b(m)} \frac{q_b(\lambda_{m+1})}{q_d(r; m+1)} \frac{1}{J}$$

$$\theta_m = (Q, \theta, \lambda, m)$$

acceptance of reversible jump

- move probabilities



$$\frac{d(m+1)}{b(m)}$$

- birth & death proposals



$$\frac{q_b(\lambda_{m+1})}{q_d(r; m+1)}$$

- Jacobian between models
 - fudge factor
 - see stepwise regression example

$$J = \frac{\sigma}{s_{r|others} \sqrt{n}}$$

reversible jump details

- reversible jump MCMC details
 - can update model with m QTL
 - have basic idea of jumping models
 - now: careful bookkeeping between models
- RJ-MCMC & Bayes factors
 - Bayes factors from RJ-MCMC chain
 - components of Bayes factors

reversible jump idea

- expand idea of MCMC to compare models
- adjust for parameters in different models
 - augment smaller model with innovations
 - constraints on larger model
- calculus “change of variables” is key
 - add or drop parameter(s)
 - carefully compute the Jacobian
- consider stepwise regression
 - Mallick (1995) & Green (1995)
 - efficient calculation with Hausholder decomposition

model selection in regression

- known regressors (e.g. markers)
 - models with 1 or 2 regressors
- jump between models
 - centering regressors simplifies calculations

$$m = 1 : Y_i = \mu + a(Q_{i1} - \bar{Q}_1) + e_i$$

$$m = 2 : Y_i = \mu + a_1(Q_{i1} - \bar{Q}_1) + a_2(Q_{i2} - \bar{Q}_2) + e_i$$

slope estimate for 1 regressor

recall least squares estimate of slope

note relation of slope to correlation

$$\hat{a} = \frac{r_{1y} s_y}{s_1}, \quad r_{1y} = \frac{\sum_{i=1}^n (Q_{i1} - \bar{Q}_1)(Y_i - \bar{Y}) / n}{s_1 s_y}$$

$$s_1^2 = \sum_{i=1}^n (Q_{i1} - \bar{Q}_1)^2 / n, \quad s_y^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 / n$$

2 correlated regressors

slopes adjusted for other regressors

$$\hat{a}_1 = \frac{(r_{1y} - r_{12}r_{2y})s_y}{s_1} = \hat{a} - \frac{r_{2y}s_y}{s_2}c_{21}, \quad c_{21} = \frac{r_{12}s_2}{s_1}$$

$$\hat{a}_2 = \frac{(r_{2y} - r_{12}r_{1y})s_y}{s_2}, \quad s_{2.1}^2 = \frac{\sum_{i=1}^n (Q_{i2} - \bar{Q}_2 - c_{21}(Q_{i1} - \bar{Q}_1))^2}{n}$$

Gibbs Sampler for Model 1

- mean $\mu \sim \phi\left(\eta + B_n(\bar{Y}_\cdot - \eta), B_n \frac{\sigma^2}{n}\right), B_n = \frac{n}{n + \kappa}$
- slope $a \sim \phi\left(B_n \frac{\sum_{i=1}^n (Q_{i1} - \bar{Q}_1)(Y_i - \bar{Y}_\cdot)}{nS_1^2}, B_n \frac{\sigma^2}{nS_1^2}\right)$
- variance $\sigma^2 \sim \text{inv} - \chi^2\left(v + n, \frac{v\tau^2 + \sum_{i=1}^n (Y_i - \bar{Y}_\cdot - a(Q_{i1} - \bar{Q}_1))^2}{v + n}\right)$

Gibbs Sampler for Model 2

- mean $\mu \sim \phi \left(\eta + B_n (\bar{Y} - \eta), B_n \frac{\sigma^2}{n} \right)$
- slopes $a_2 \sim \phi \left(B_n \frac{\sum_{i=1}^n (Q_{i2} - \bar{Q}_2)(Y_i - \bar{Y} - a_1(Q_{i1} - \bar{Q}_1))}{nS_{2.1}^2}, B_n \frac{\sigma^2}{nS_{2.1}^2} \right)$
- variance $\sigma^2 \sim \text{inv-}\chi^2 \left(v + n, \frac{v\tau^2 + \sum_{i=1}^n \left(Y_i - \bar{Y} - \sum_{k=1}^2 a_k (Q_{ik} - \bar{Q}_k) \right)^2}{v + n} \right)$

updates from 2->1

- drop 2nd regressor
- adjust other regressor

$$a \rightarrow a_1 + a_2 c_{21}$$

$$a_2 \rightarrow 0$$

updates from 1->2

- add 2nd slope, adjusting for collinearity
- adjust other slope & variance

$$z \sim \phi(0,1), \quad J = \frac{\sigma}{s_{2.1}\sqrt{n}}$$
$$a_2 \rightarrow \hat{a}_2 + z \times J, \quad \hat{a}_2 = \frac{\sum_{i=1}^n (Q_{i2} - \bar{Q}_2)(Y_i - \hat{\mu} - \hat{a}_1(Q_{i1} - \bar{Q}_1))}{ns_{2.1}^2}$$
$$a_1 \rightarrow a - a_2 c_{21} = a - z \times c_{21} J - \hat{a}_2 c_{21}$$

model selection in regression

- known regressors (e.g. markers)
 - models with 1 or 2 regressors
- jump between models
 - augment with new innovation z

	m	parameters	innovations	transformations
$1 \rightarrow 2$		$(\mu, a, \sigma^2; z)$	$z \sim \phi(0, 1)$	$\left\{ \begin{array}{l} a_2 \rightarrow \hat{a}_2 + z \times J \\ a_1 \rightarrow a - a_2 c_{21} \end{array} \right\}$
$2 \rightarrow 1$		$(\mu, a_1, a_2, \sigma^2)$		$\left\{ \begin{array}{l} a \rightarrow a_1 + a_2 c_{21} \\ z \rightarrow 0 \end{array} \right\}$

change of variables

- change variables from model 1 to model 2
- calculus issues for integration
 - need to formally account for change of variables
 - infinitesimal steps in integration (db)
 - involves partial derivatives (next page)

$$\begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \begin{bmatrix} 1 & -c_{21}J & -c_{21} \\ 0 & J & 1 \end{bmatrix} \times \begin{pmatrix} a \\ z \\ \hat{a}_2 \end{pmatrix} = g(a; z | Y, Q_1, Q_2)$$

$$\int \pi(a_1, a_2 | Y, Q_1, Q_2) da_1 da_2 = \int \pi(a; z | Y, Q_1, Q_2) J da dz$$

Jacobian & the calculus

- Jacobian sorts out change of variables
 - careful: easy to mess up here!

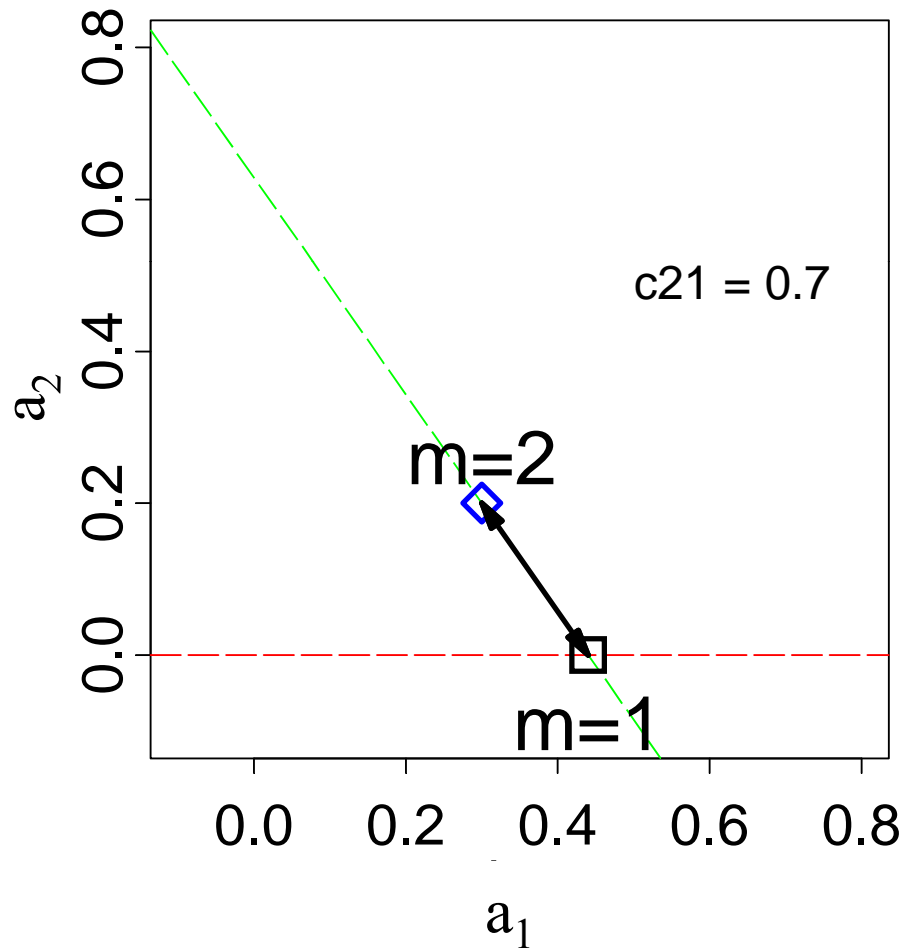
$$g(a; z) = (a_1, a_2), \quad \frac{\partial g(a; z)}{\partial a \partial z} = \begin{bmatrix} 1 & -c_{21}J \\ 0 & J \end{bmatrix}$$

$$\left| \det \left(\begin{bmatrix} 1 & -c_{21}J \\ 0 & J \end{bmatrix} \right) \right| = |1 \times J - 0 \times (-c_{21}J)| = J$$

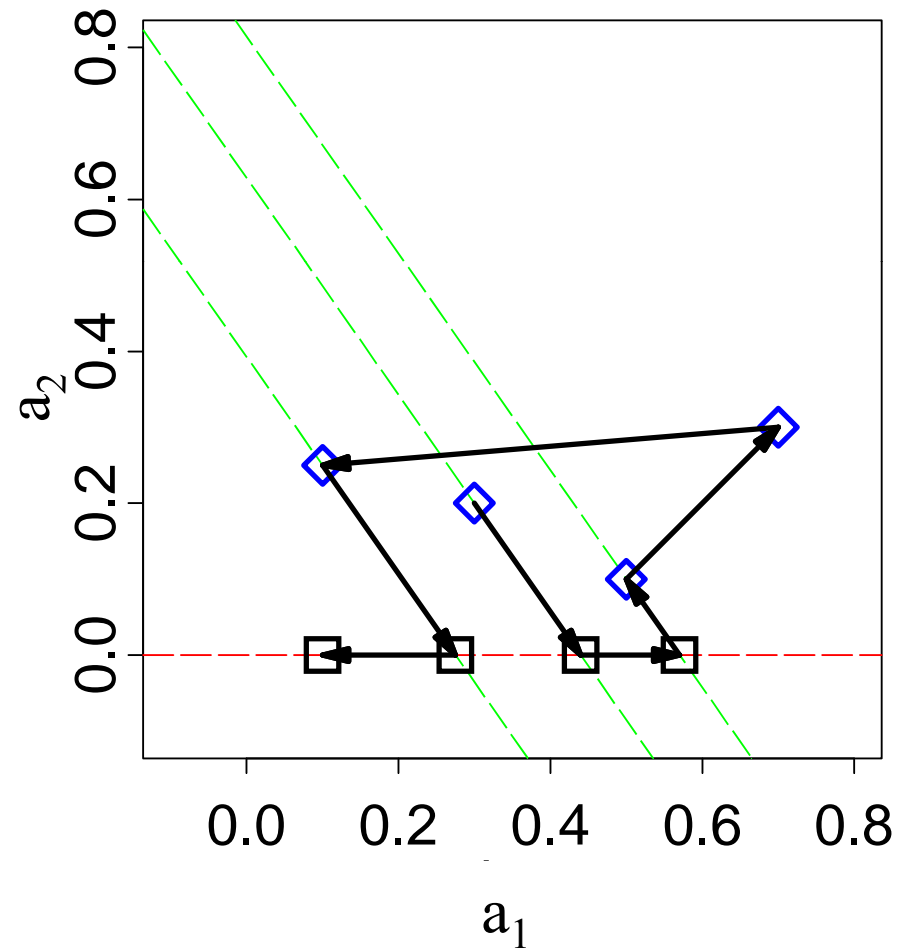
$$da_1 da_2 = \left| \det \left(\frac{\partial g(\mu, a, \sigma^2; z)}{\partial a \partial z} \right) \right| da_1 da_2 = J da dz$$

geometry of reversible jump

Move Between Models



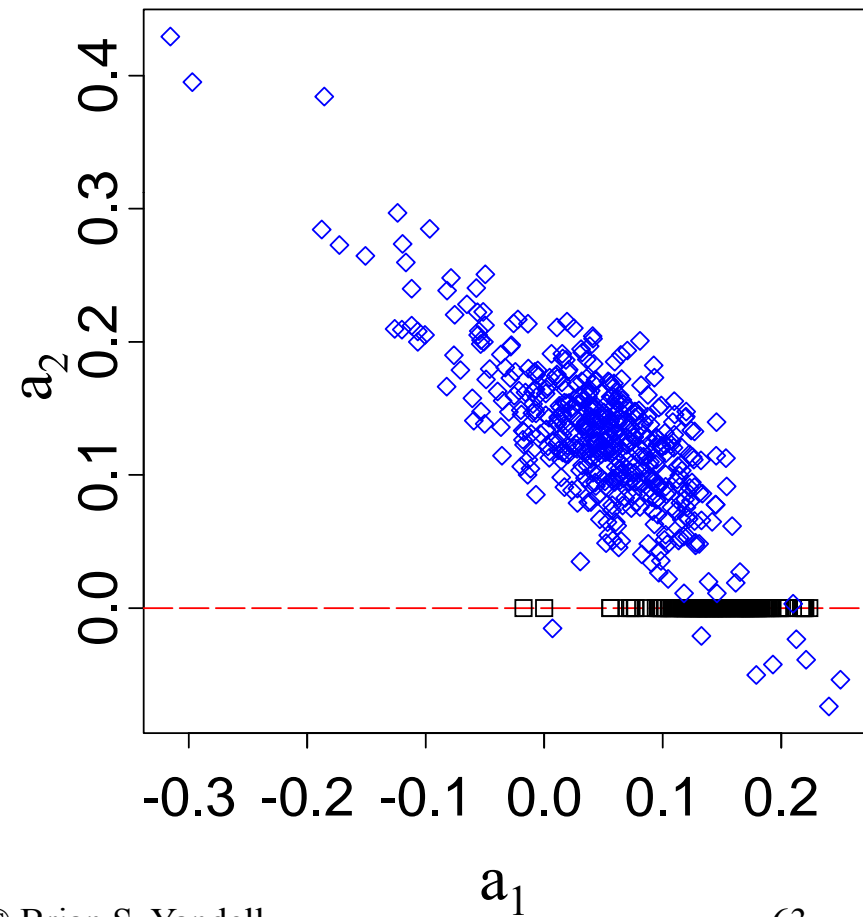
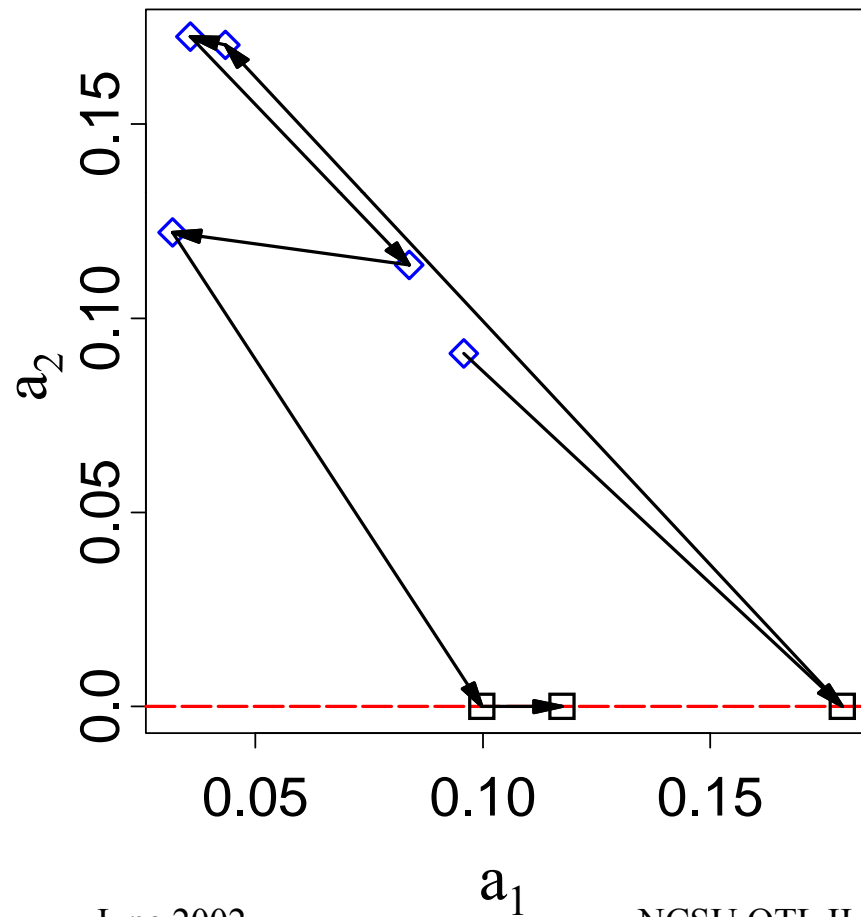
Reversible Jump Sequence



QT additive reversible jump

a short sequence

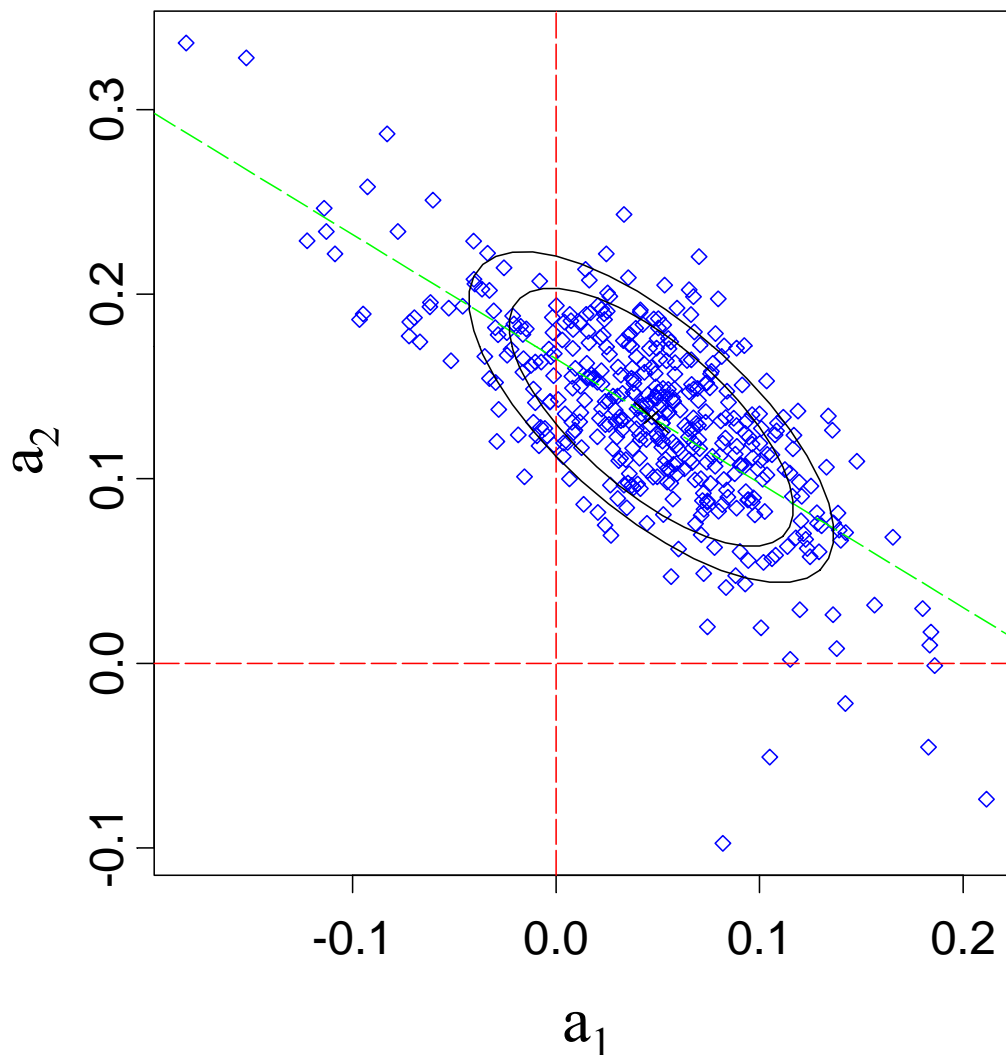
first 1000 with $m < 3$



credible set for additive

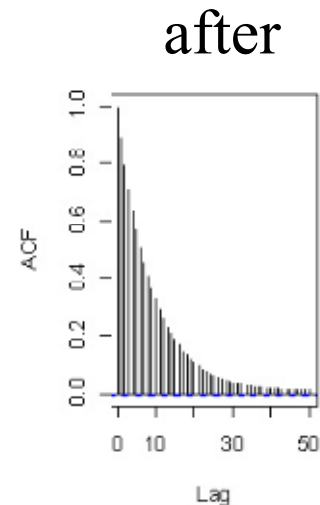
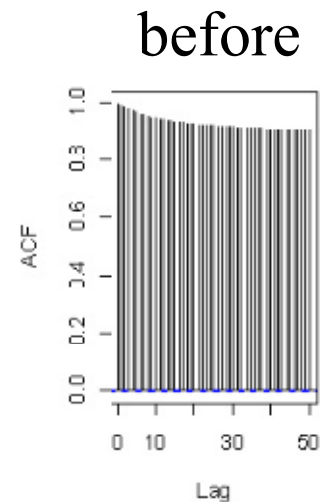
90% & 95% sets
based on normal

regression line
corresponds to
slope of updates



multivariate updating of effects

- more computations when $m > 2$
- avoid matrix inverse
 - Cholesky decomposition of matrix
- simultaneous updates
 - effects at all loci
- accept new locus based on
 - sampled new genos at locus
 - sampled new effects at all loci
- also long-range positions updates



Gibbs sampler with loci indicators

- partition genome into intervals
 - at most one QTL per interval
 - interval = 1 cM in length
 - assume QTL in middle of interval
- use loci to indicate presence/absence of QTL in each interval
 - $\gamma = 1$ if QTL in interval
 - $\gamma = 0$ if no QTL
- Gibbs sampler on loci indicators
 - see work of Nengjun Yi (and earlier work of Ina Hoeschele)

$$Y = \beta_0 + \gamma_1 \beta_1(q_1) + \gamma_2 \beta_2(q_1) + e$$

Bayesian shrinkage estimation

- soft loci indicators
 - strength of evidence for λ_j depends on variance of β_j
 - similar to $\gamma > 0$ on grey scale
- include all possible loci in model
 - pseudo-markers at 1cM intervals
- Wang et al. (2005 *Genetics*)
 - Shizhong Xu group at U CA Riverside

$$Y = \beta_0 + \beta_1(q_1) + \beta_2(q_1) + \dots + e$$

$$\beta_j(q_j) \sim N(0, \sigma_j^2), \sigma_j^2 \sim \text{inverse - chisquare}$$

epistatic interactions

- model space issues
 - Fisher-Cockerham partition vs. tree-structured?
 - 2-QTL interactions only?
 - general interactions among multiple QTL?
 - retain model hierarchy (include main QTL)?
- 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 et al. (2005, 2007)

5. Model Assessment

- balance model fit against model complexity

	smaller model	bigger model
model fit	miss key features	fits better
prediction	may be biased	no bias
interpretation	easier	more complicated
parameters	low variance	high variance

- information criteria: penalize likelihood by model size
 - compare $IC = -2 \log L(\text{model} | \text{data}) + \text{penalty}(\text{model size})$
- Bayes factors: balance posterior by prior choice
 - compare $\text{pr}(\text{data} | \text{model})$

Bayes factors

- ratio of model likelihoods
 - ratio of posterior to prior odds for architectures
 - average over unknown effects (μ) and loci (λ)

$$BF = \frac{\text{pr}(\text{data} \mid \text{model } A_1)}{\text{pr}(\text{data} \mid \text{model } A_2)}$$

- roughly equivalent to BIC
 - BIC maximizes over unknowns
 - BF averages over unknowns

$$2 \log_{10}(BF) = 2LOD + (\text{change in model size}) \log_{10}(n)$$

issues in computing Bayes factors

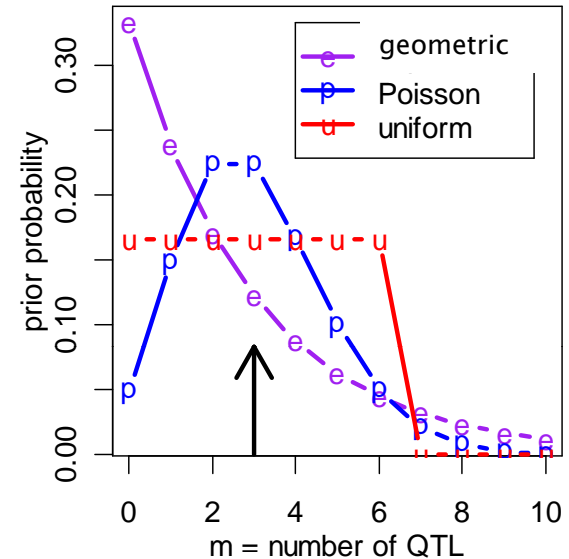
- *BF* insensitive to shape of prior on A
 - geometric, Poisson, uniform
 - precision improves when prior mimics posterior
- *BF* sensitivity to prior variance on effects θ
 - prior variance should reflect data variability
 - resolved by using hyper-priors
 - automatic algorithm; no need for user tuning
- easy to compute Bayes factors from samples
 - apply Bayes' rule and solve for $\text{pr}(y / m, A)$
 - $\text{pr}(A / y, m) = \text{pr}(y / m, A) \text{pr}(A | m) / \text{constant}$
 - $\text{pr}(\text{data}|\text{model}) = \text{constant} * \text{pr}(\text{model}|\text{data}) / \text{pr}(\text{model})$
 - posterior $\text{pr}(A / y, m)$ is marginal histogram

Bayes factors and genetic model A

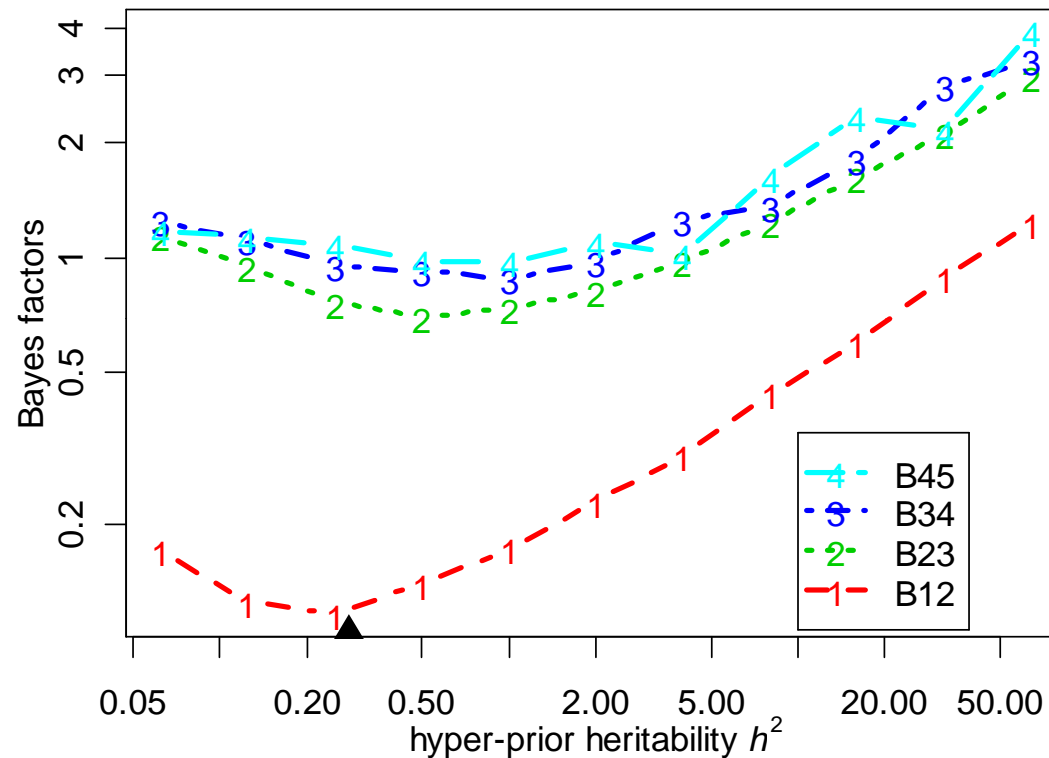
- $|A|$ = number of QTL
 - prior $\text{pr}(A)$ chosen by user
 - posterior $\text{pr}(A/y, m)$
 - sampled marginal histogram
 - shape affected by prior $\text{pr}(A)$

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

- pattern of QTL across genome
- gene action and epistasis

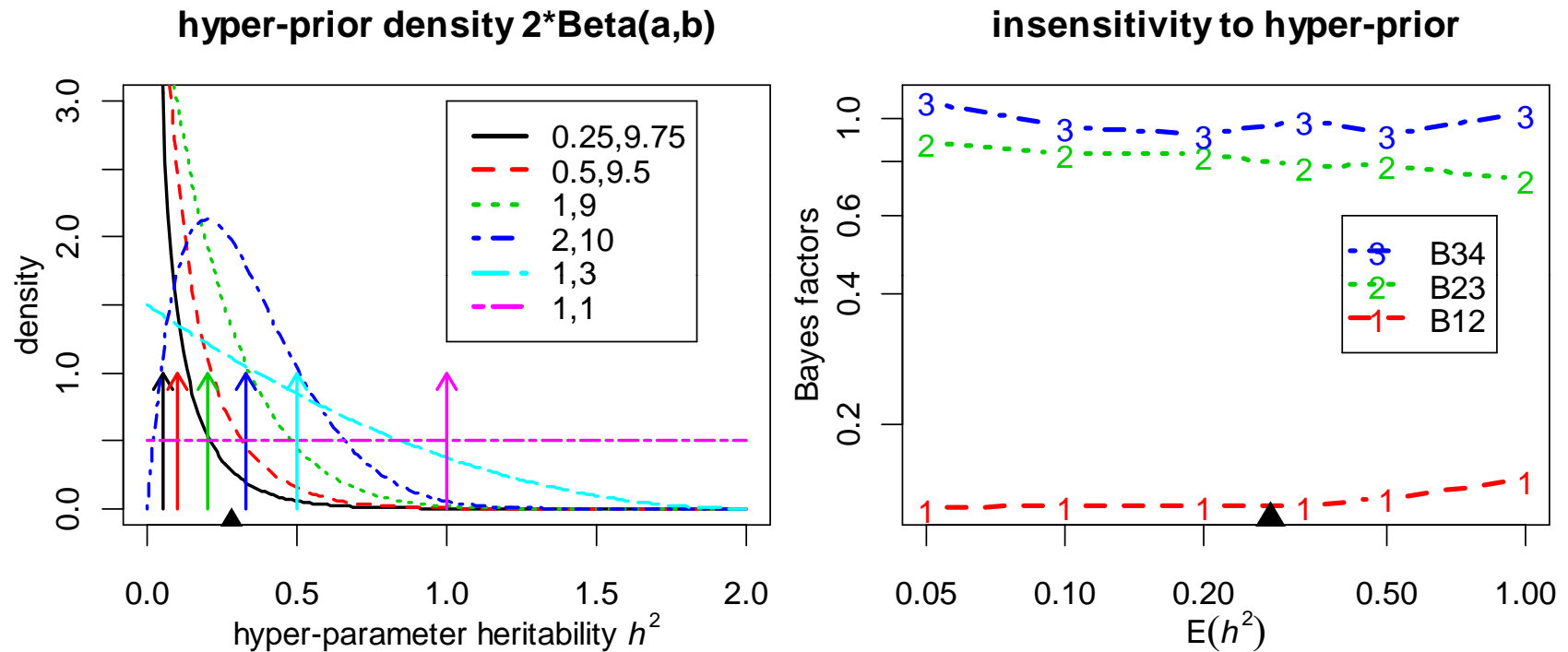


BF sensitivity to fixed prior for effects



$$\beta_{qj} \sim \mathbf{N}(0, \sigma_G^2 / m), \sigma_G^2 = h^2 \sigma_{\text{total}}^2, h^2 \text{ fixed}$$

BF insensitivity to random effects prior



$$\beta_{qj} \sim N(0, \sigma_G^2 / m), \sigma_G^2 = h^2 \sigma_{\text{total}}^2, \frac{1}{2} h^2 \sim \text{Beta}(a, b)$$

marginal BF scan by QTL

- compare models with and without QTL at λ
 - average over all possible models
 - estimate as ratio of samples with/without QTL
- scan over genome for peaks
 - $2\log(\text{BF})$ seems to have similar properties to LPD

$$BF_{\lambda} = \frac{\text{pr}(y \mid m, \text{model with } \lambda)}{\text{pr}(y \mid m, \text{model without } \lambda)}$$

Bayesian model averaging

- average summaries over multiple architectures
- avoid selection of “best” model
- focus on “better” models
- examples in data talk later

6. analysis of hyper data

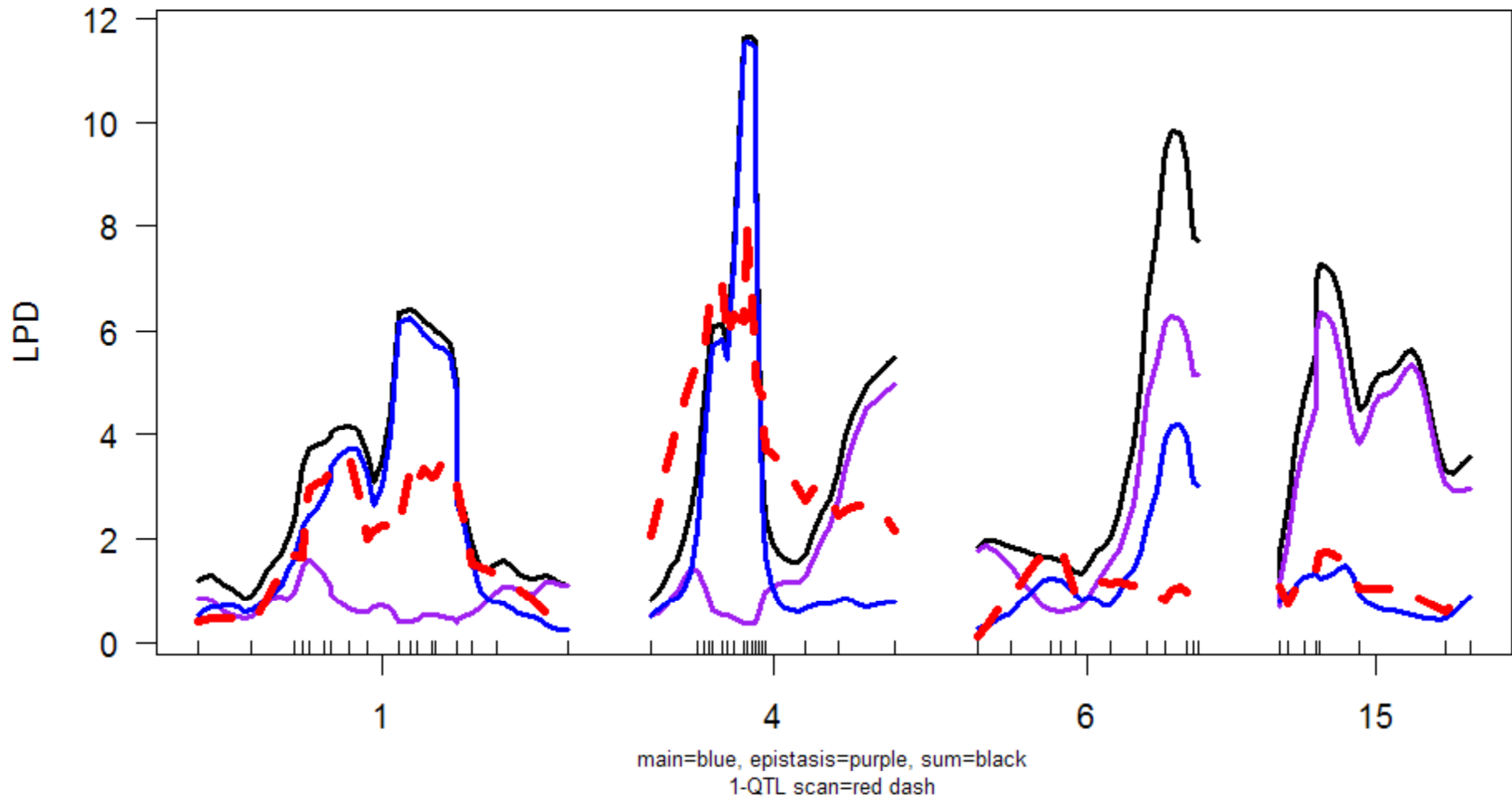
- marginal scans of genome
 - detect significant loci
 - infer main and epistatic QTL, GxE
- infer most probable genetic architecture
 - number of QTL
 - chromosome pattern of QTL with epistasis
- diagnostic summaries
 - heritability, unexplained variation

marginal scans of genome

- LPD and $2\log(\text{BF})$ “tests” for each locus
- estimates of QTL effects at each locus
- separately infer main effects and epistasis
 - main effect for each locus (blue)
 - epistasis for loci paired with another (purple)
 - identify epistatic QTL in 1-D scan
 - infer pairing in 2-D scan

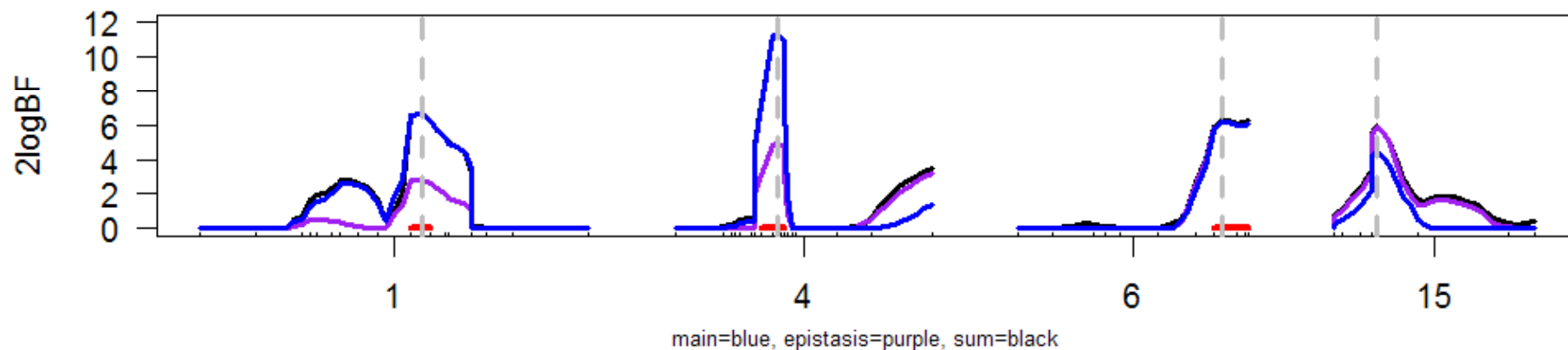
hyper data: scanone

LPD of bp for main+epistasis+sum

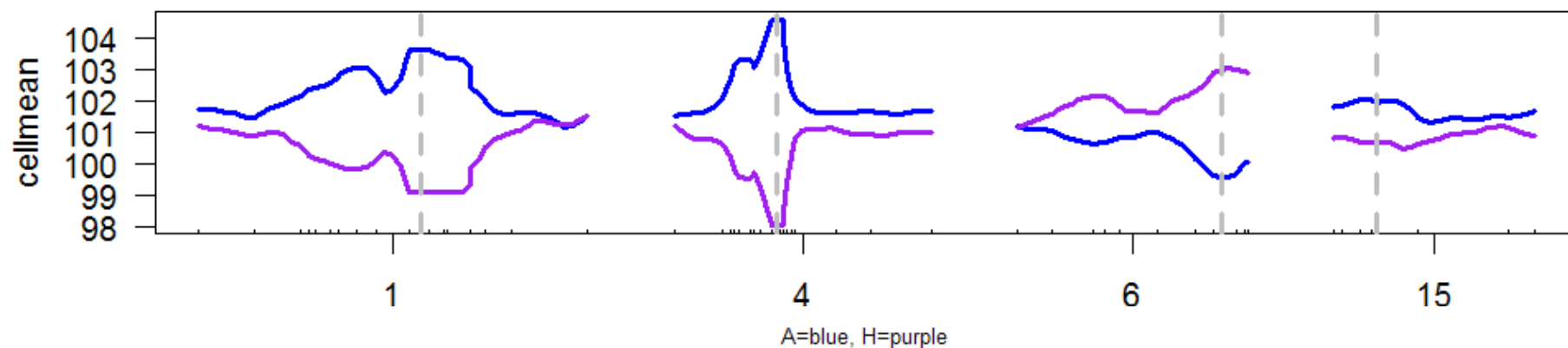


2log(BF) scan with 50% HPD region

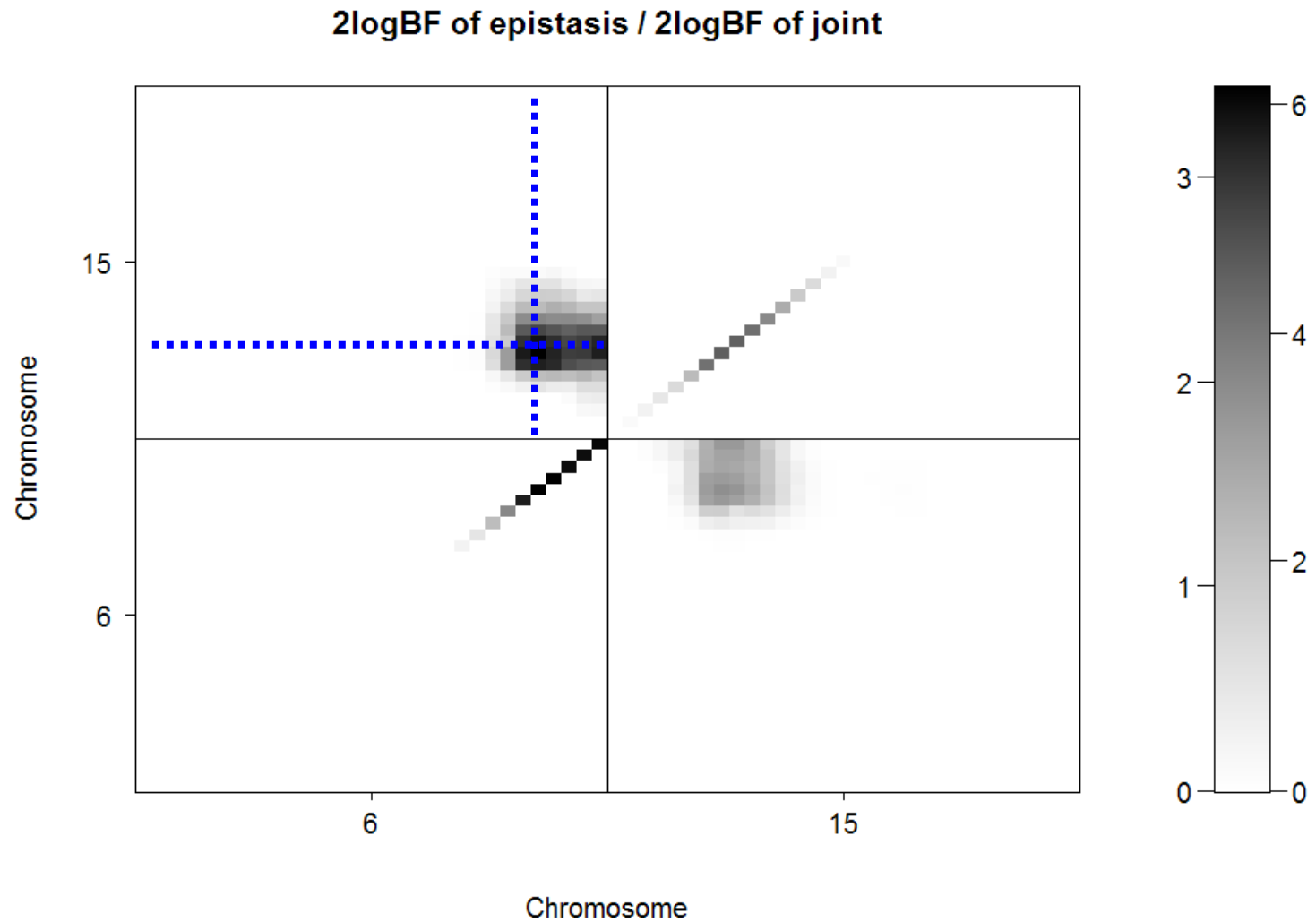
2logBF of bp for main+epistasis+sum



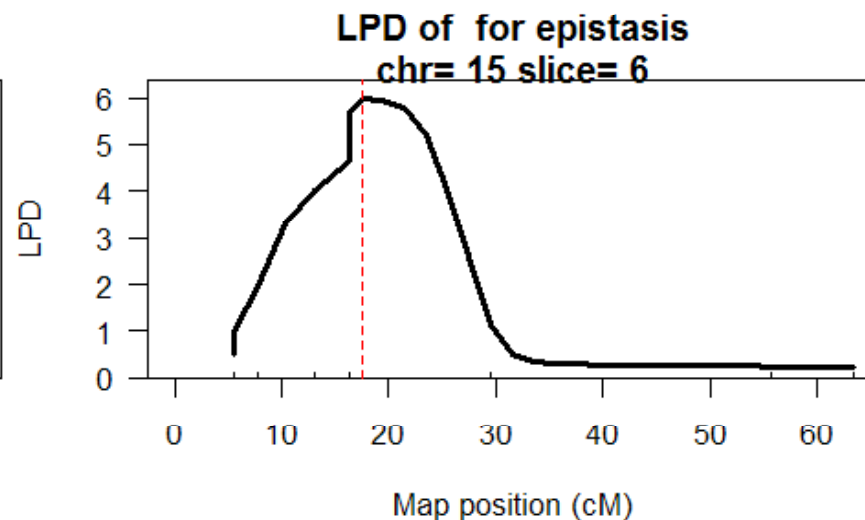
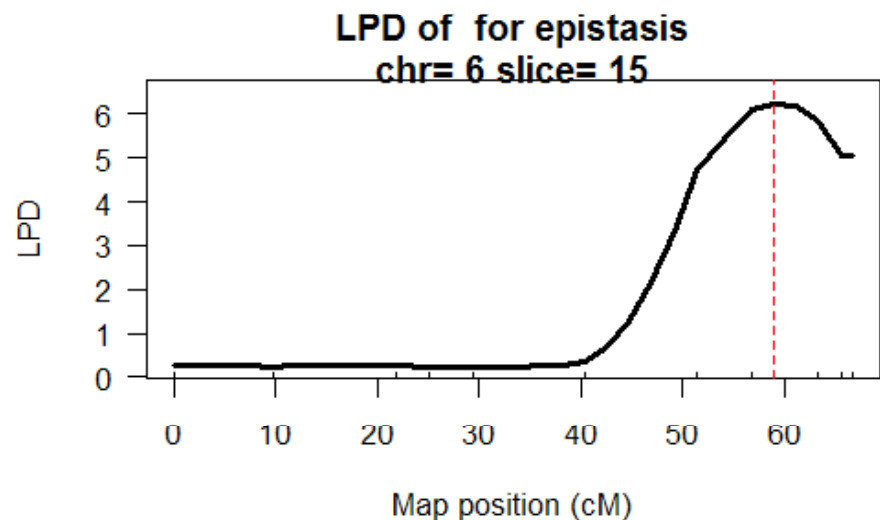
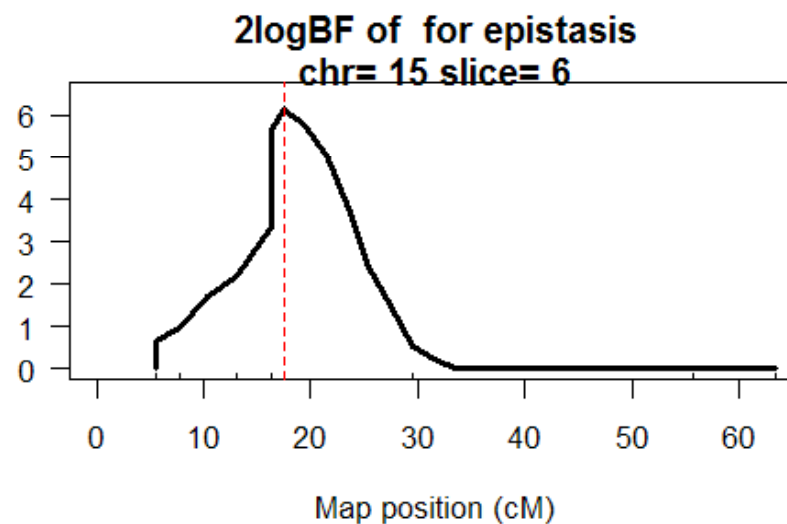
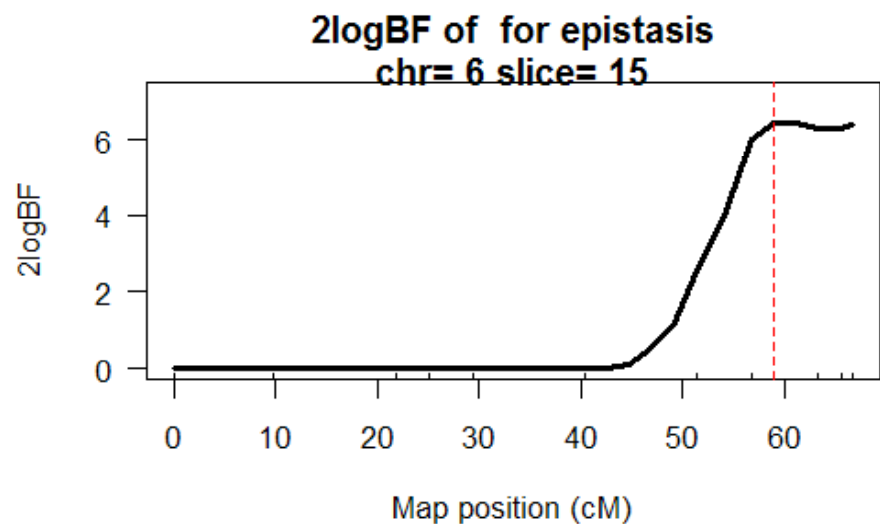
cellmean of bp for A+H



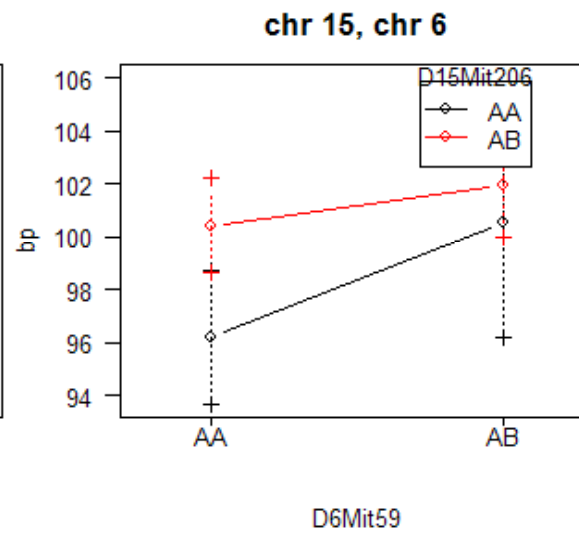
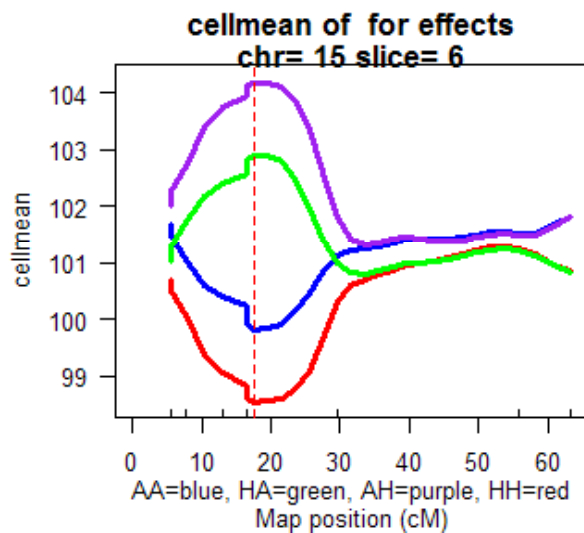
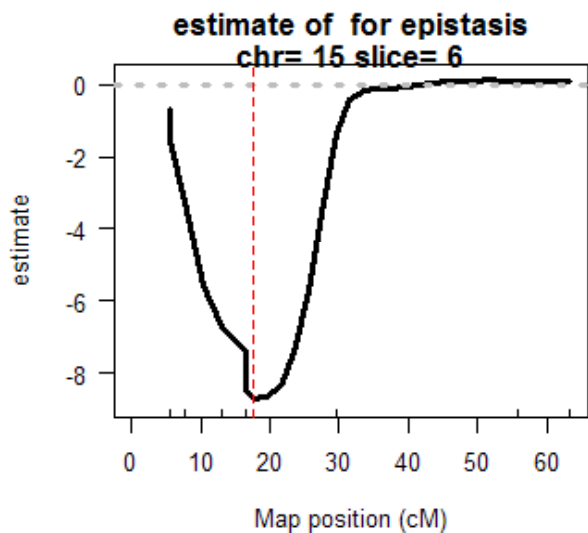
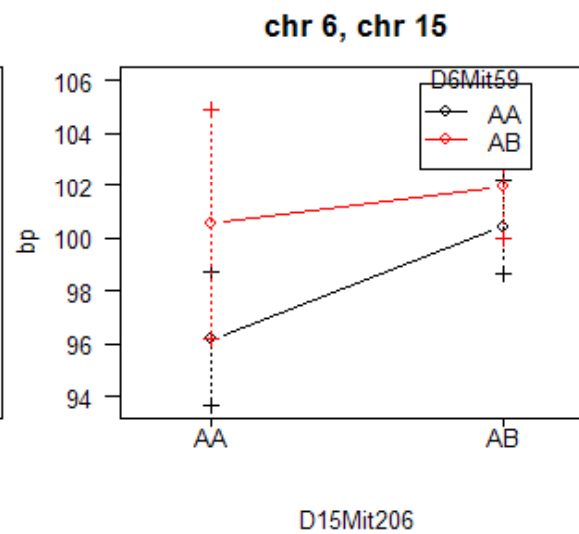
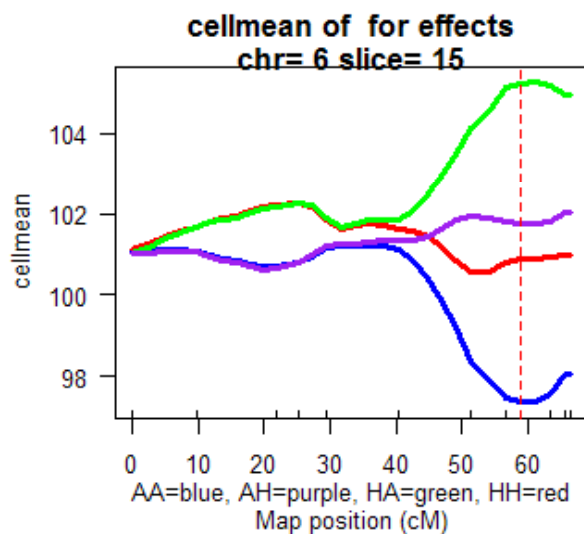
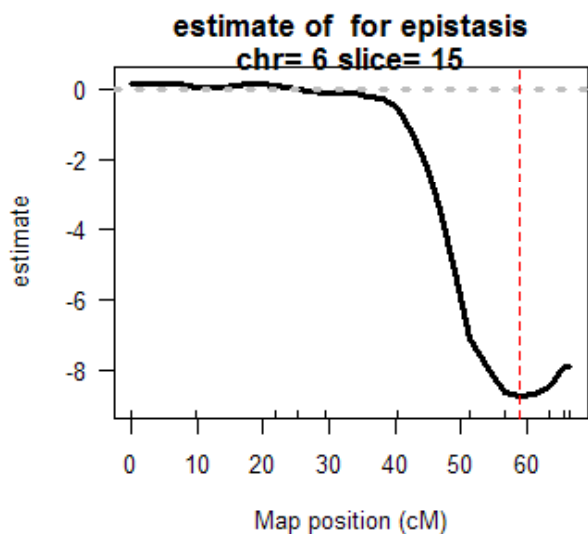
2-D plot of $2\log\text{BF}$: chr 6 & 15



1-D Slices of 2-D scans: chr 6 & 15



1-D Slices of 2-D scans: chr 6 & 15

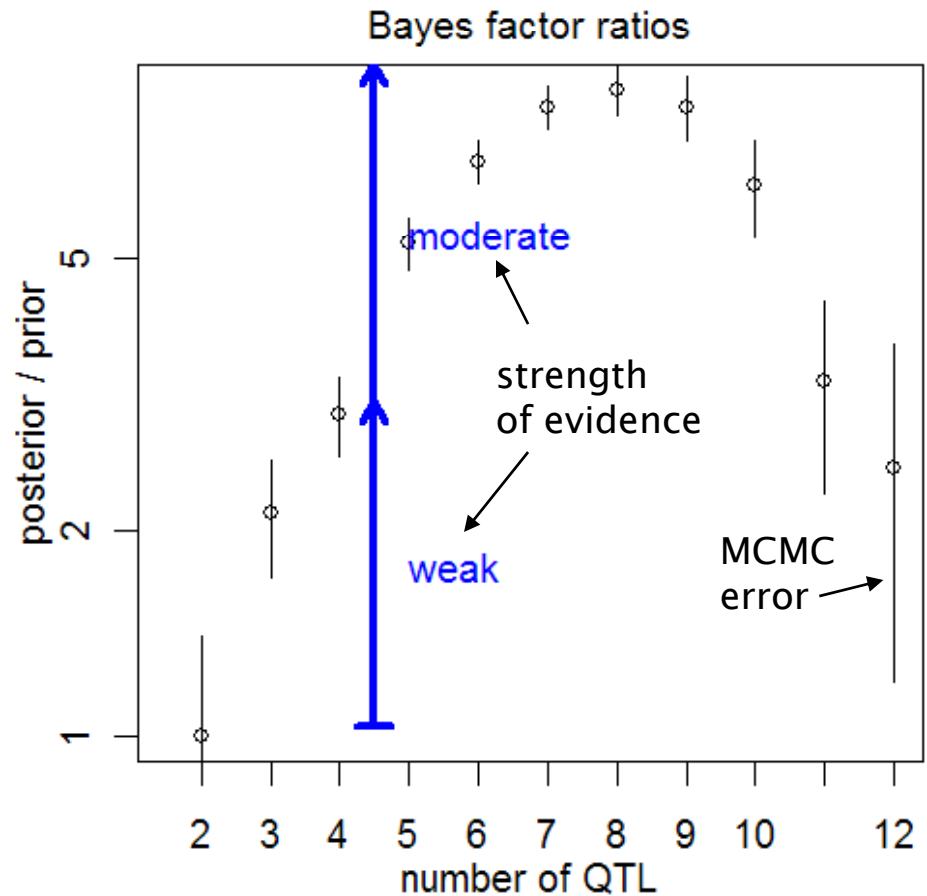
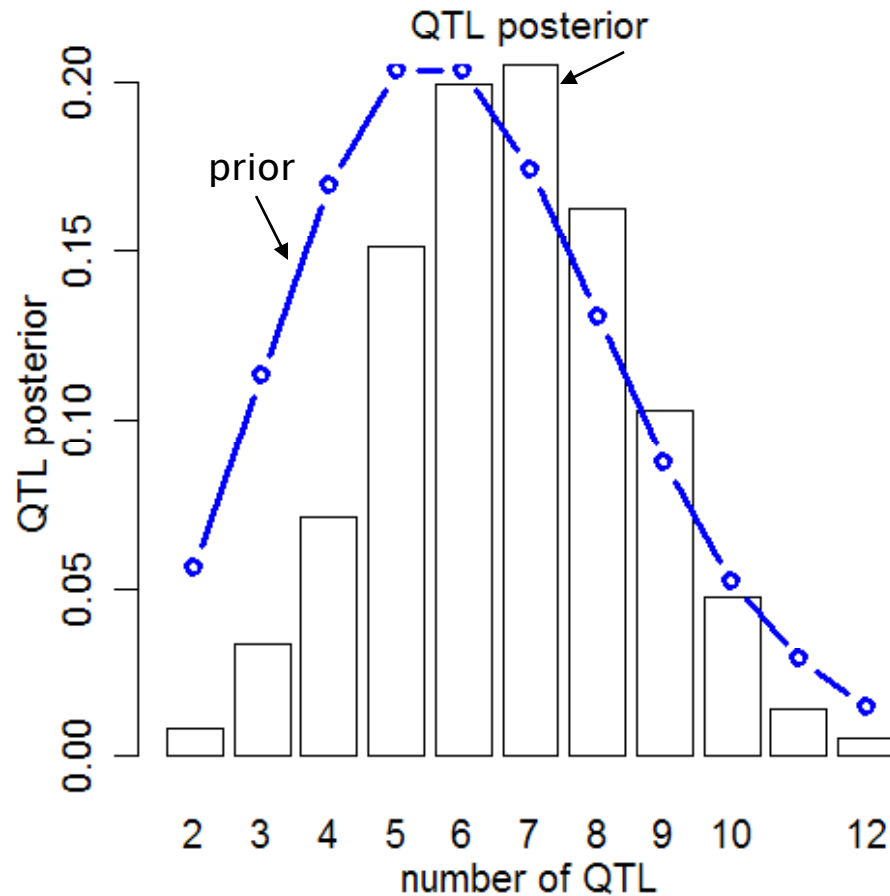


What is best genetic architecture?

- How many QTL?
- What is pattern across chromosomes?
- examine posterior relative to prior
 - prior determined ahead of time
 - posterior estimated by histogram/bar chart
 - Bayes factor ratio = $\text{pr}(\text{model}|\text{data}) / \text{pr}(\text{model})$

How many QTL?

posterior, prior, Bayes factor ratios



most probable patterns

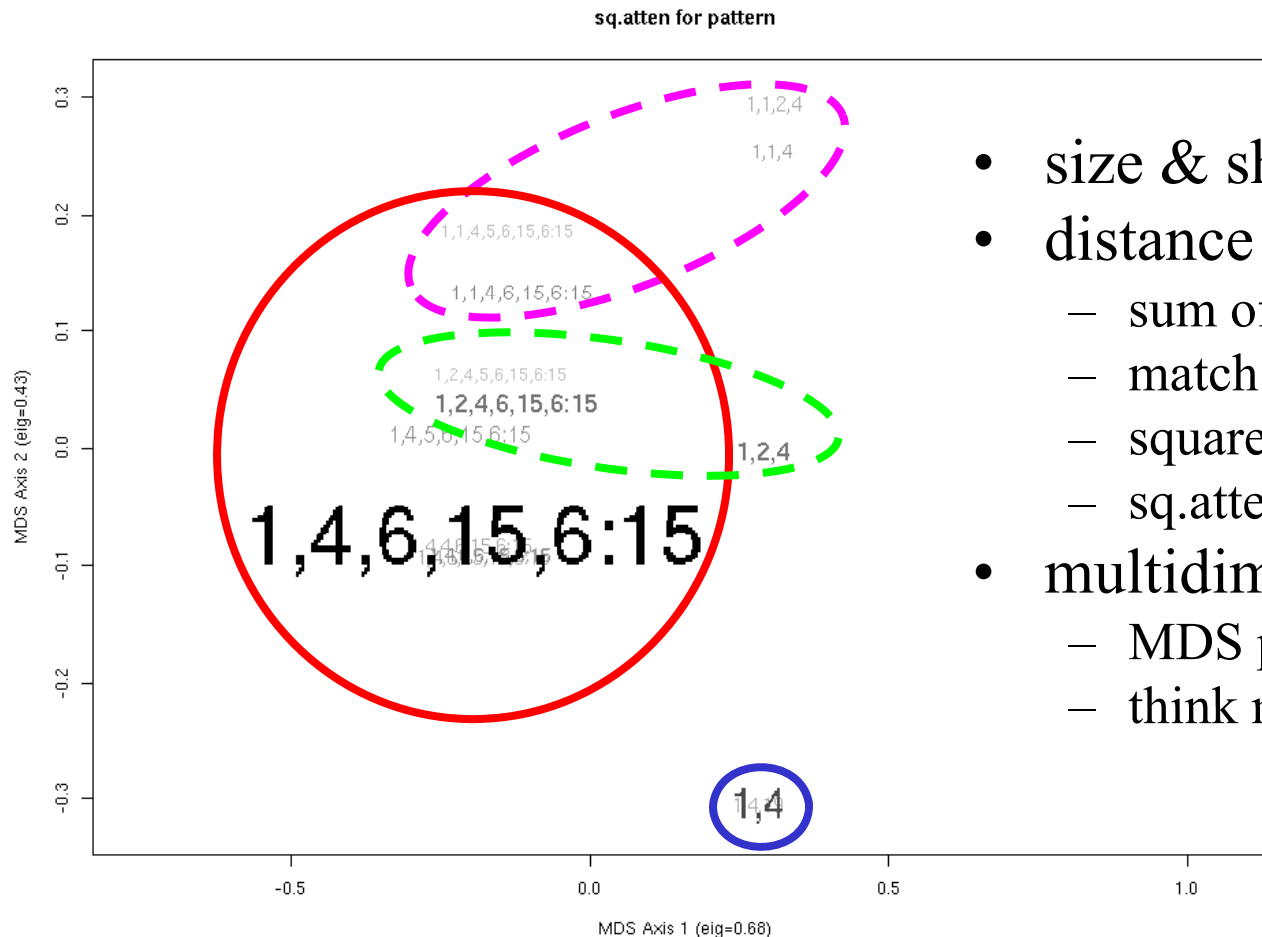
	nqtl	posterior	prior	bf	bfse
1,4,6,15,6:15	5	0.03400	2.71e-05	24.30	2.360
1,4,6,6,15,6:15	6	0.00467	5.22e-06	17.40	4.630
1,1,4,6,15,6:15	6	0.00600	9.05e-06	12.80	3.020
1,1,4,5,6,15,6:15	7	0.00267	4.11e-06	12.60	4.450
1,4,6,15,15,6:15	6	0.00300	4.96e-06	11.70	3.910
1,4,4,6,15,6:15	6	0.00300	5.81e-06	10.00	3.330
1,2,4,6,15,6:15	6	0.00767	1.54e-05	9.66	2.010
1,4,5,6,15,6:15	6	0.00500	1.28e-05	7.56	1.950
1,2,4,5,6,15,6:15	7	0.00267	6.98e-06	7.41	2.620
1,4	2	0.01430	1.51e-04	1.84	0.279
1,1,2,4	4	0.00300	3.66e-05	1.59	0.529
1,2,4	3	0.00733	1.03e-04	1.38	0.294
1,1,4	3	0.00400	6.05e-05	1.28	0.370
1,4,19	3	0.00300	5.82e-05	1.00	0.333

what is best estimate of QTL?

- find most probable pattern
 - 1,4,6,15,6:15 has posterior of 3.4%
- estimate locus across all nested patterns
 - Exact pattern seen ~100/3000 samples
 - Nested pattern seen ~2000/3000 samples
- estimate 95% confidence interval using quantiles

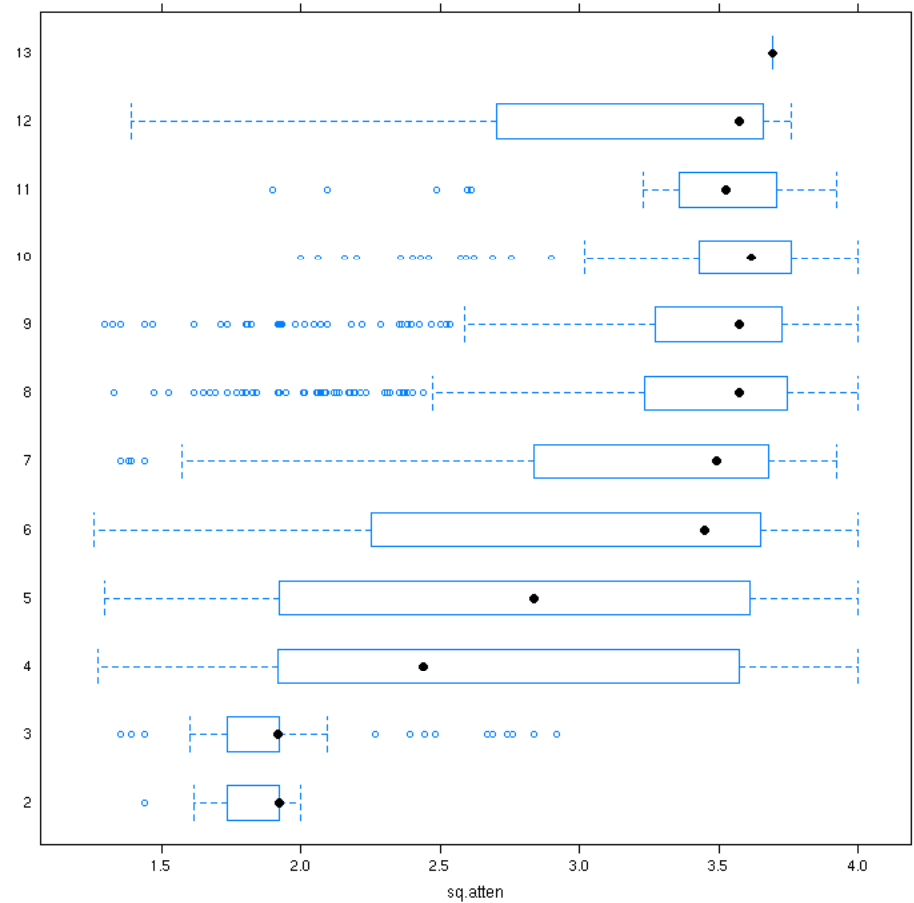
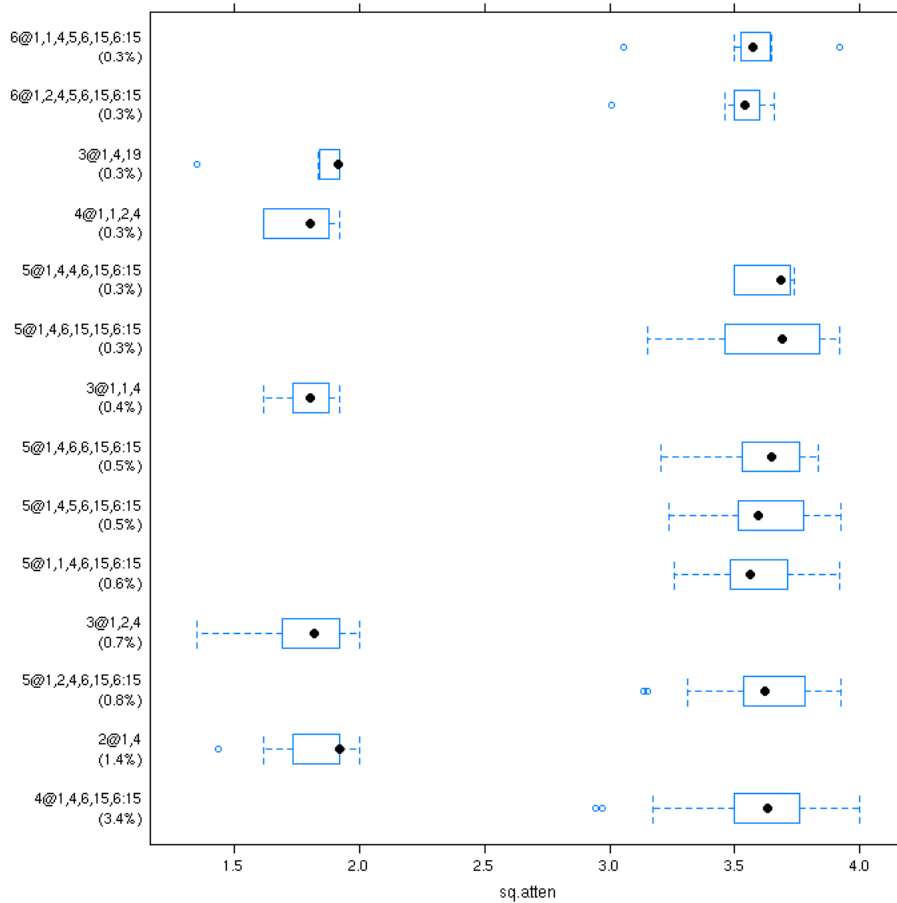
	chrom	locus	locus.LCL	locus.UCL	n.qtl
247	1	69.9	24.44875	95.7985	0.8026667
245	4	29.5	14.20000	74.3000	0.8800000
248	6	59.0	13.83333	66.7000	0.7096667
246	15	19.5	13.10000	55.7000	0.8450000

how close are other patterns?

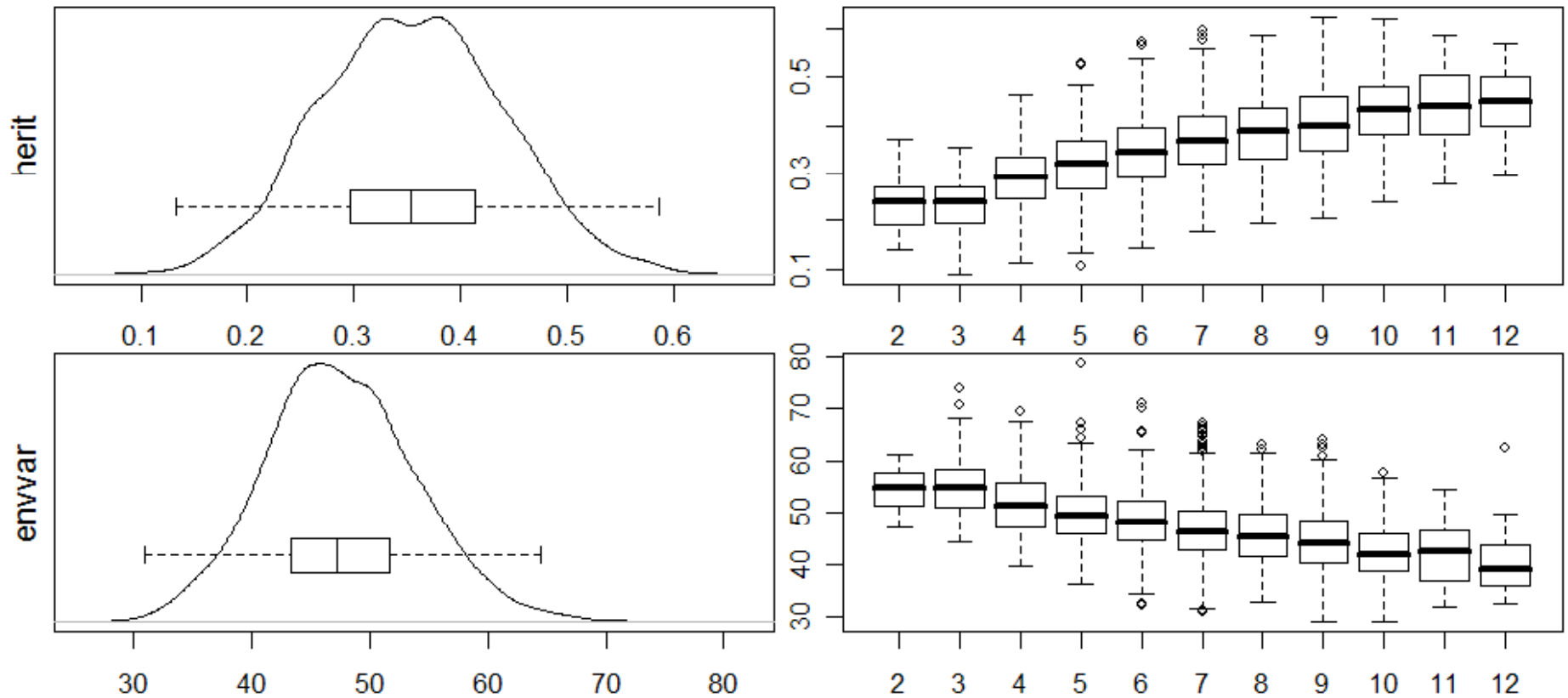


- size & shade ~ posterior
- distance between patterns
 - sum of squared attenuation
 - match loci between patterns
 - squared attenuation = $(1-2r)^2$
 - sq.atten in scale of LOD & LPD
- multidimensional scaling
 - MDS projects distance onto 2-D
 - think mileage between cities

how close are other patterns?



diagnostic summaries



7. Software for Bayesian QTLs

R/qtlbim

- publication
 - CRAN release Fall 2006
 - Yandell et al. (2007 *Bioinformatics*)
- properties
 - cross-compatible with R/qtl
 - epistasis, fixed & random covariates, GxE
 - extensive graphics

R/qtlbim: software history

- Bayesian module within WinQTLCart
 - WinQTLCart output can be processed using R/bim
- Software history
 - initially designed (Satagopan Yandell 1996)
 - major revision and extension (Gaffney 2001)
 - R/bim to CRAN (Wu, Gaffney, Jin, Yandell 2003)
 - R/qtlbim total rewrite (Yandell et al. 2007)

other Bayesian software for QTLs

- R/bim*: Bayesian Interval Mapping
 - Satagopan Yandell (1996; Gaffney 2001) CRAN
 - no epistasis; reversible jump MCMC algorithm
 - version available within WinQTLCart (statgen.ncsu.edu/qtlcart)
- R/qtl*
 - Broman et al. (2003 Bioinformatics) CRAN
 - multiple imputation algorithm for 1, 2 QTL scans & limited mult-QTL fits
- Bayesian QTL / Multimapper
 - Sillanpää Arjas (1998 Genetics) www.rni.helsinki.fi/~mjs
 - no epistasis; introduced posterior intensity for QTLs
- (no released code)
 - Stephens & Fisch (1998 Biometrics)
 - no epistasis
- R/bqtl
 - C Berry (1998 TR) CRAN
 - no epistasis, Haley Knott approximation

* Jackson Labs (Hao Wu, Randy von Smith) provided crucial technical support

many thanks

Karl Broman

Jackson Labs

Gary Churchill

Hao Wu

Randy von Smith

U AL Birmingham

David Allison

Nengjun Yi

Tapan Mehta

Samprit Banerjee

Ram Venkataraman

Daniel Shriner

USDA Hatch, NIH/NIDDK (Attie), NIH/R01 (Yi, Broman)

Tom Osborn

David Butruille

Marcio Ferrera

Josh Udahl

Pablo Quijada

Alan Attie

Jonathan Stoehr

Hong Lan

Susie Clee

Jessica Byers

Mark Keller

Michael Newton

Hyuna Yang

Daniel Sorensen

Daniel Gianola

Liang Li

my students

Jaya Satagopan

Fei Zou

Patrick Gaffney

Chunfang Jin

Elias Chaibub

W Whipple Neely

Jee Young Moon