

Bayesian QTL Mapping

Brian S. Yandell

University of Wisconsin-Madison

www.stat.wisc.edu/~yandell/statgen↑

UW-Madison, April 2008

April 2008

UW-Madison © Brian S. Yandell

1

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
5. model assessment
 - Bayes factors & model averaging
6. analysis of hyper data
7. software for Bayesian QTLs

April 2008

UW-Madison © Brian S. Yandell

2

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

April 2008

UW-Madison © Brian S. Yandell

3

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 = $MSE = (\text{bias})^2 + \text{variance}$

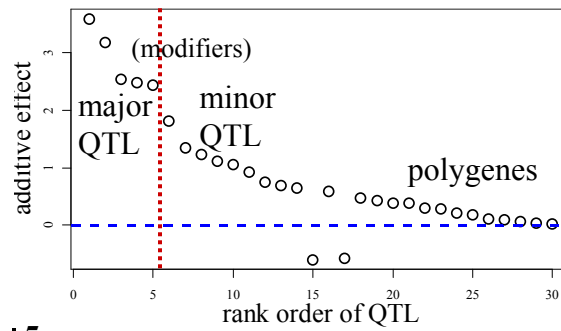
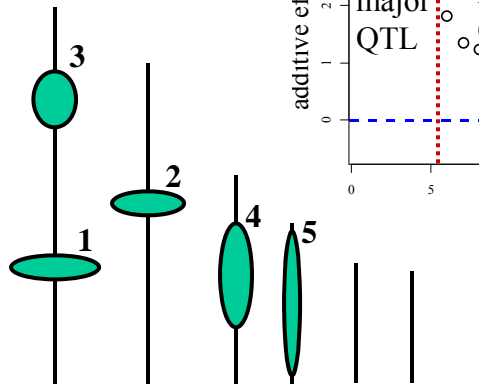
April 2008

UW-Madison © Brian S. Yandell

4

Pareto diagram of QTL effects

major QTL on linkage map



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

April 2008

UW-Madison © Brian S. Yandell

5

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

April 2008

UW-Madison © Brian S. Yandell

6

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

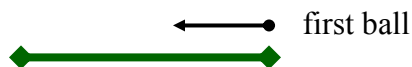
April 2008

UW-Madison © Brian S. Yandell

7

Who was Bayes?

- Reverend Thomas Bayes (1702-1761)
 - part-time mathematician
 - buried in Bunhill Cemetary, Moongate, London
 - famous paper in 1763 *Phil Trans Roy Soc London*
 - Barnard (1958 *Biometrika*), Press (1989) *Bayesian Statistics*
 - Stigler (1986) *History of Statistics*
 - Carlin Louis (1996); Gelman et al. (1995) books
 - Was Bayes the first with this idea? (Laplace)
- billiard balls on rectangular table
 - two balls tossed at random (uniform) on table
 - where is first ball if the second is to its **right** (**left**)?

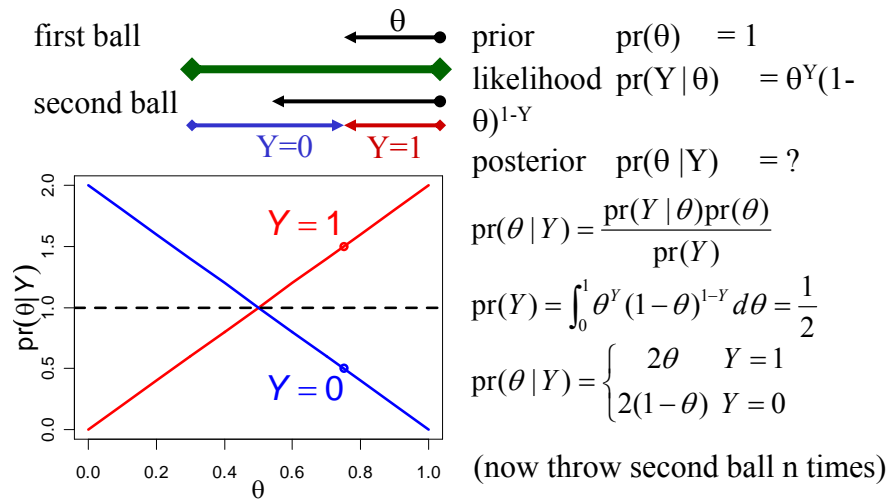


March 2011

UW-Madison © Brian S. Yandell

8

Where is the first ball?



March 2011

UW-Madison © Brian S. Yandell

9

What is Bayes Theorem?

- before and after observing data
 - prior: $\text{pr}(\theta) = \text{pr}(\text{parameters})$
 - posterior: $\text{pr}(\theta|Y) = \text{pr}(\text{parameters}|\text{data})$
- posterior = likelihood * prior / constant
 - usual likelihood of parameters given data
 - normalizing constant $\text{pr}(Y)$ depends only on data
 - constant often drops out of calculation

$$\text{pr}(\theta|Y) = \frac{\text{pr}(\theta, Y)}{\text{pr}(Y)} = \frac{\text{pr}(Y|\theta) \times \text{pr}(\theta)}{\text{pr}(Y)}$$

March 2011

UW-Madison © Brian S. Yandell

10

What is Probability?

Frequentist analysis

- chance over many trials
 - long run average
 - estimates
 - confidence intervals
 - long term frequency
 - hypothesis tests
 - p -values
- Type I error rate
 - reject null when true
 - chance of extreme result

Bayesian analysis

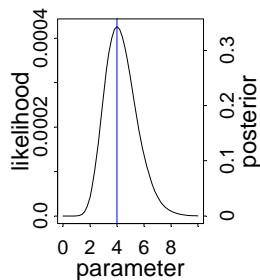
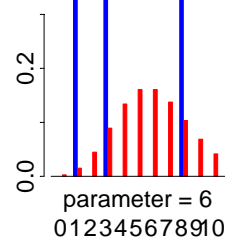
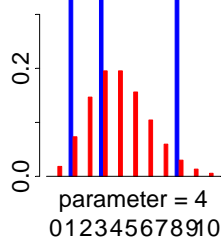
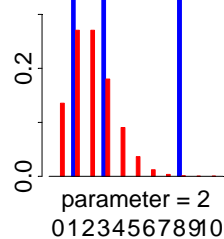
- uncertainty of true value
- prior
 - uncertainty before data
 - incorporate prior knowledge/experience
- posterior
 - uncertainty after analyzing current data
 - balance prior and data

March 2011

UW-Madison © Brian S. Yandell

11

Likelihood and Posterior Example



data : $Y = 1,3,8$

parameter : $\theta = ?$

$$\text{pr}(Y = y | \theta) = \frac{\theta^y e^{-\theta}}{y!}$$

(M. Newton, pers. comm.)

March 2011

UW-Madison © Brian S. Yandell

12

Frequentist or Bayesian?

- Frequentist approach
 - fixed parameters
 - range of values
 - maximize likelihood
 - ML estimates
 - find the peak
 - confidence regions
 - random region
 - invert a test
 - hypothesis testing
 - 2 nested models
- Bayesian approach
 - random parameters
 - distribution
 - posterior distribution
 - posterior mean
 - sample from dist
 - credible sets
 - fixed region given data
 - HPD regions
 - model selection/critique
 - Bayes factors

March 2011

UW-Madison © Brian S. Yandell

13

Frequentist or Bayesian?

- Frequentist approach
 - maximize over mixture of QT genotypes
 - locus profile likelihood
 - max over effects
 - HPD region for locus
 - natural for locus
 - 1-2 LOD drop
 - work to get effects
 - approximate shape of likelihood peak
- Bayesian approach
 - joint distribution over QT genotypes
 - sample distribution
 - joint effects & loci
 - HPD regions for
 - joint locus & effects
 - use density estimator

March 2011

UW-Madison © Brian S. Yandell

14

Choice of Bayesian priors

- elicited priors
 - higher weight for more probable parameter values
 - based on prior empirical knowledge
 - use previous study to inform current study
 - weather prediction, previous QTL studies on related organisms
- conjugate priors
 - convenient mathematical form
 - essential before computers, helpful now to simplify computation
 - large variances on priors reduces their influence on posterior
- non-informative priors
 - may have “no” information on unknown parameters
 - prior with all parameter values equally likely
 - may not sum to 1 (improper), which can complicate use
- **always** check sensitivity of posterior to choice of prior

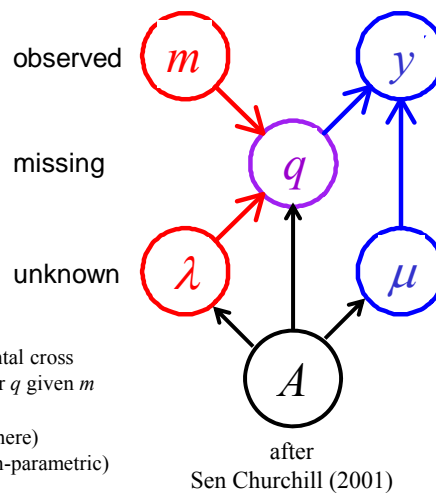
March 2011

UW-Madison © Brian S. Yandell

15

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)



April 2008

UW-Madison © Brian S. Yandell

16

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:

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

April 2008

UW-Madison © Brian S. Yandell

17

Bayes posterior vs. maximum likelihood (genetic architecture $A = \text{single QTL at } \lambda$)

- **LOD**: classical Log Odds
 - maximize likelihood over effects μ
 - R/qt1 scanone/scantwo: method = "em"
- **LPD**: Bayesian Log Posterior Density
 - average posterior over effects μ
 - R/qt1 scanone/scantwo: method = "imp"

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

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

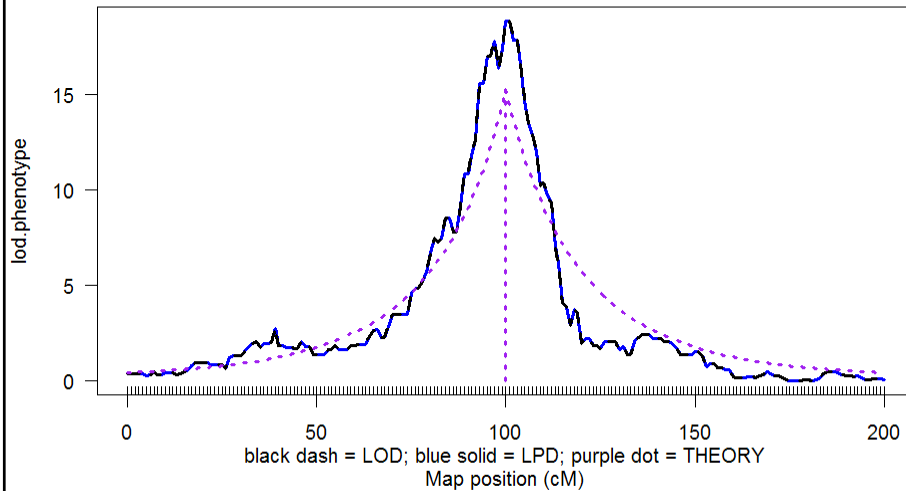
April 2008

UW-Madison © Brian S. Yandell

18

LOD & LPD: 1 QTL

n.ind = 100, 1 cM marker spacing



April 2008

UW-Madison © Brian S. Yandell

19

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?

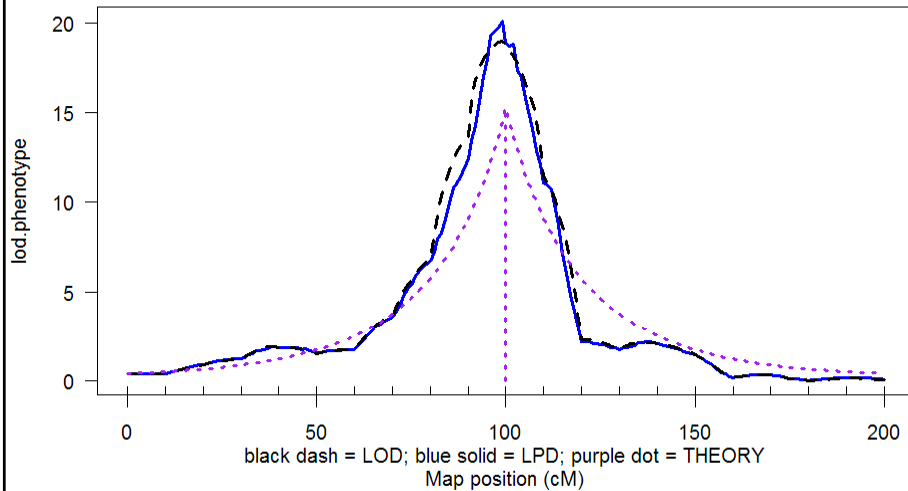
April 2008

UW-Madison © Brian S. Yandell

20

LOD & LPD: 1 QTL

n.ind = 100, 10 cM marker spacing



April 2008

UW-Madison © Brian S. Yandell

21

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

April 2008

UW-Madison © Brian S. Yandell

22

marginal LOD or LPD

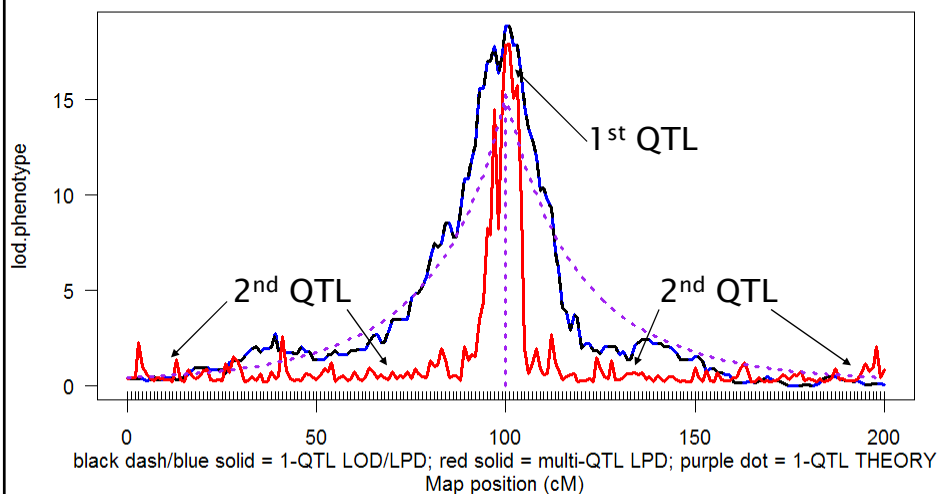
- What is contribution of a QTL adjusting for all others?
 - improvement in LPD due to QTL at locus λ
 - contribution due to main effects, epistasis, GxE?
- How does adjusted LPD *differ* from unadjusted LPD?
 - raised by removing variance due to unlinked QTL
 - raised or lowered due to bias of linked 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

April 2008

UW-Madison © Brian S. Yandell

23

LPD: 1 QTL vs. multi-QTL marginal contribution to LPD from QTL at λ

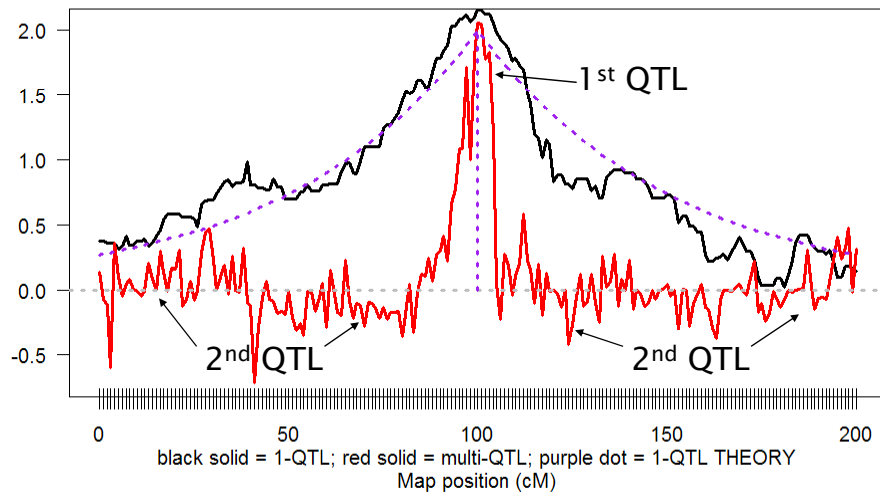


April 2008

UW-Madison © Brian S. Yandell

24

substitution effect: 1 QTL vs. multi-QTL
 single QTL effect vs. marginal effect from QTL at λ



April 2008

UW-Madison © Brian S. Yandell

25

3. Bayesian strategy for QTLs

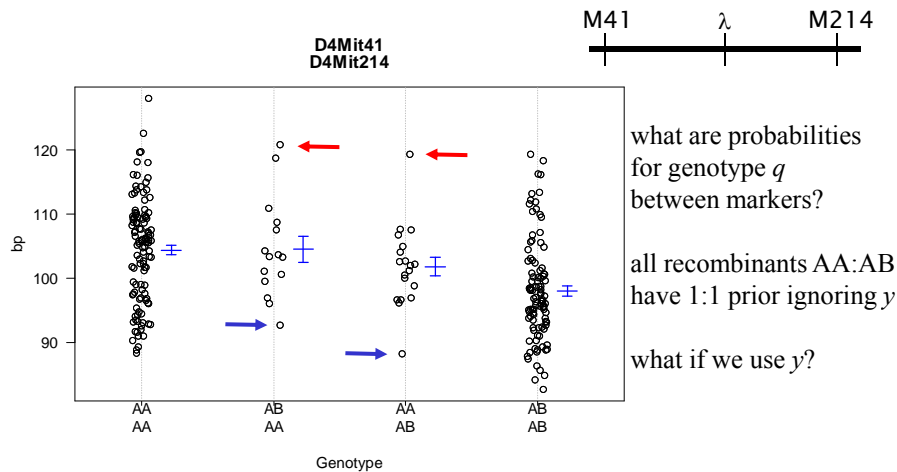
- augment data (y, m) with missing genotypes q
- build model for augmented data
 - genotypes (q) evaluated at loci (λ)
 - depends on flanking markers (m)
 - phenotypes (y) centered about effects (μ)
 - depends on missing genotypes (q)
 - λ and μ depend on genetic architecture (A)
 - How complicated is 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

April 2008

UW-Madison © Brian S. Yandell

26

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



April 2008

UW-Madison © Brian S. Yandell

27

posterior on QTL genotypes q

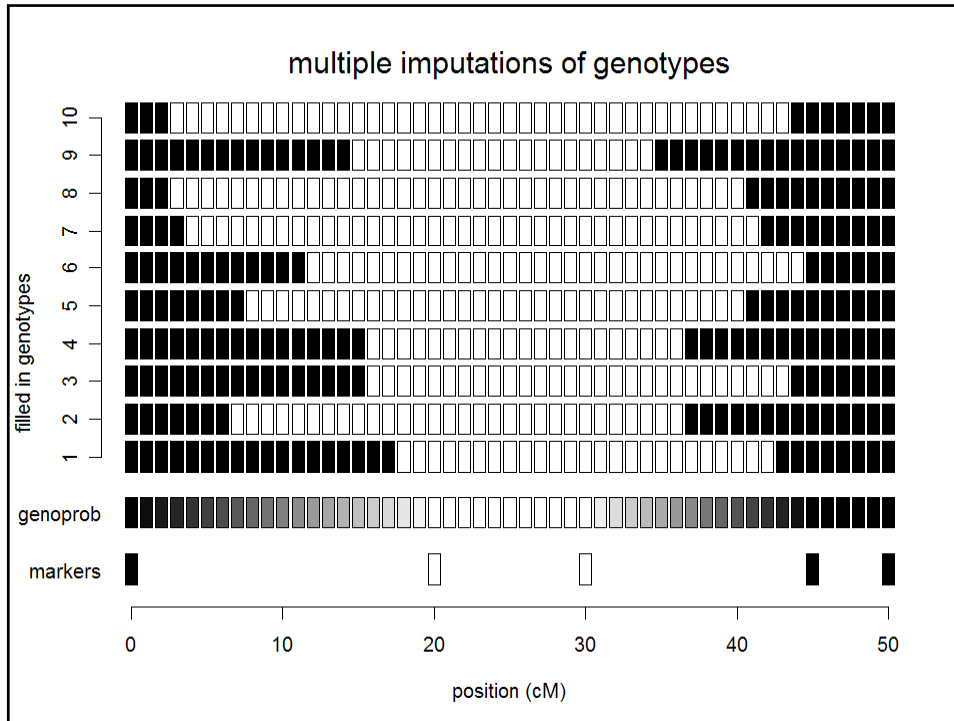
- full conditional of q given data, parameters
 - proportional to prior $\text{pr}(q | m, \lambda)$
 - weight toward q that agrees with flanking markers
 - proportional to likelihood $\text{pr}(y|q, \mu)$
 - weight toward q with similar phenotype values
 - posterior balances these two
- this *is* the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

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

April 2008

UW-Madison © Brian S. Yandell

28



Bayes for normal data

$Y = G + E$ posterior for single individual

environ $E \sim N(0, \sigma^2)$, σ^2 known

likelihood $\text{pr}(Y | G, \sigma^2) = N(Y | G, \sigma^2)$

prior $\text{pr}(G | \sigma^2, \mu, \kappa) = N(G | \mu, \sigma^2/\kappa)$

posterior $N(G | \mu + B_1(Y - \mu), B_1 \sigma^2)$

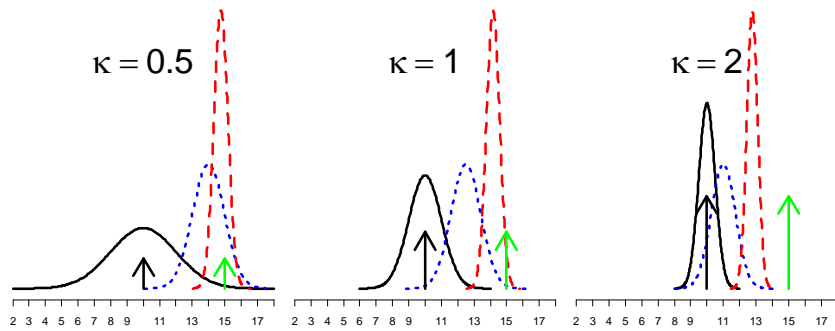
$Y_i = G + E_i$ posterior for sample of n individuals

shrinkage weights B_n go to 1

$$\text{pr}(G | Y, \sigma^2, \mu, \kappa) = N\left(G \mid \mu + B_n(\bar{Y}_\bullet - \mu), B_n \frac{\sigma^2}{n}\right)$$

$$\text{with } \bar{Y}_\bullet = \text{sum} \frac{Y_i}{n}, B_n = \frac{n}{\kappa + n} \rightarrow 1$$

effect of prior variance on posterior

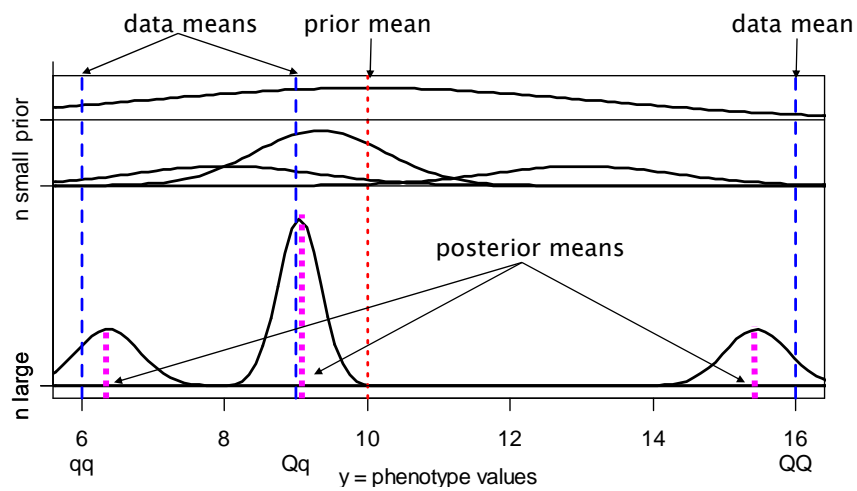


March 2011

UW-Madison © Brian S. Yandell

31

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



April 2008

UW-Madison © Brian S. Yandell

32

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

prior: $E(\mu_q) = \bar{y}_\bullet$ $V(\mu_q) = V(y)h_q^2$

posterior: $E(\mu_q | y) = \bar{y}_\bullet(1 - b_q) + \bar{y}_q b_q$ $V(\mu_q | y) = V(\bar{y}_q)b_q$

shrinkage: $b_q = 1 - \frac{V(\bar{y}_q)}{V(\bar{y}_q) + V(y)h_q^2} \approx 1$

April 2008

UW-Madison © Brian S. Yandell

33

Empirical Bayes: choosing hyper-parameters

How do we choose hyper-parameters μ, κ ?

Empirical Bayes: marginalize over prior

estimate μ, κ from marginal posterior

likelihood $\text{pr}(Y_i | Q_i, G, \sigma^2) = N(Y_i | G(Q_i), \sigma^2)$

prior $\text{pr}(G_Q | \sigma^2, \mu, \kappa) = N(G_Q | \mu, \sigma^2/\kappa)$

marginal $\text{pr}(Y_i | \sigma^2, \mu, \kappa) = N(Y_i | \mu, \sigma^2(\kappa + 1)/\kappa)$

estimates $\hat{\mu} = \bar{Y}_\bullet, s^2 = \text{sum}_i (Y_i - \bar{Y}_\bullet)^2 / n$

$\kappa \leq 1$ or $\kappa = \sigma^2/s^2$

EB posterior $\text{pr}(G_Q | Y) = N\left(G_Q \left| \bar{Y}_\bullet + \hat{B}_Q(\bar{Y}_Q - \bar{Y}_\bullet), \hat{B}_Q \frac{\sigma^2}{n_Q} \right.\right)$

March 2011

UW-Madison © Brian S. Yandell

34

What if variance σ^2 is unknown?

- recall that sample variance is proportional to chi-square
 - $\text{pr}(s^2 / \sigma^2) = \chi^2 (ns^2/\sigma^2 / n)$
 - or equivalently, $ns^2/\sigma^2 / \sigma^2 \sim \chi_n^2$
- conjugate prior is inverse chi-square
 - $\text{pr}(\sigma^2 / v, \tau^2) = \text{inv-}\chi^2 (\sigma^2 / v, \tau^2)$
 - or equivalently, $v\tau^2/\sigma^2 / v, \tau^2 \sim \chi_v^2$
 - empirical choice: $\tau^2 = s^2/3, v=6$
 - $E(\sigma^2 / v, \tau^2) = s^2/2, \text{Var}(\sigma^2 / v, \tau^2) = s^4/4$
- posterior given data
 - $\text{pr}(\sigma^2 / Y, v, \tau^2) = \text{inv-}\chi^2 (\sigma^2/v+n, (v\tau^2+ns^2)/(v+n))$

March 2011

UW-Madison © Brian S. Yandell

35

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 $H \leq 2^{n_{qtl}} (3^{n_{qtl}} \text{ for } F_2)$
- 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)

April 2008

UW-Madison © Brian S. Yandell

36

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

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

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

April 2008

UW-Madison © Brian S. Yandell

37

partition of multiple QTL effects

- partition genotype-specific mean into QTL effects

μ_q = mean + main effects + 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

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

April 2008

UW-Madison © Brian S. Yandell

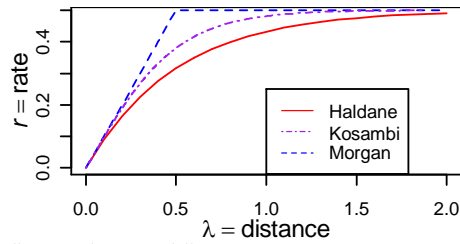
38

Recombination and Distance

- assume map and marker distances are known
- useful approximation for QTL linkage
 - Haldane map function: no crossover interference
 - independence implies crossover events are Poisson
- all computations consistent in approximation
 - rely on given map with known marker locations
 - 1-to-1 relation of distance to recombination
 - all map functions are approximate anyway

$$r = \frac{1}{2} (1 - e^{-2\lambda})$$

$$\lambda = -\frac{1}{2} \log(1 - 2r)$$



March 2011

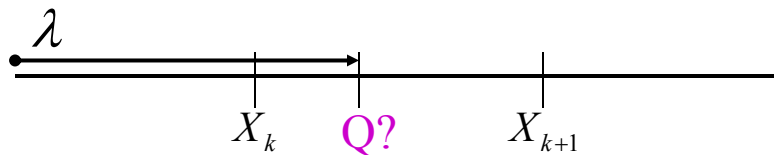
UW-Madison © Brian S. Yandell

39

recombination model $\text{pr}(Q|X, \lambda)$

- locus λ is distance along linkage map
 - identifies flanking marker region
- flanking markers provide good approximation
 - map assumed known from earlier study
 - inaccuracy slight using only flanking markers
 - extend to next flanking markers if missing data
 - could consider more complicated relationship
 - but little change in results

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



March 2011

UW-Madison © Brian S. Yandell

40

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
$$\text{pr}(\lambda | m, q) = \text{pr}(\lambda) \text{pr}(q | m, \lambda) / \text{constant}$$
 - constant determined by averaging
 - over all possible genotypes q
 - over all possible loci λ on entire map
- no easy way to write down posterior

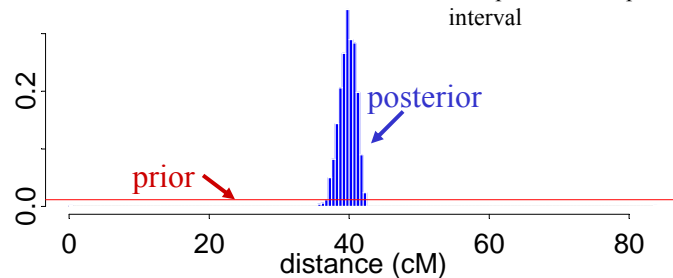
April 2008

UW-Madison © Brian S. Yandell

41

prior & posterior for QT locus

- prior information from other studies
 - concentrate on credible regions
 - use posterior of previous study as new prior
- no prior information on locus
 - uniform prior over genome
 - use framework map
 - choose interval proportional to length
 - then pick uniform position within interval



March 2011

UW-Madison © Brian S. Yandell

42

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.

April 2008

UW-Madison © Brian S. Yandell

43

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)

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

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

April 2008

UW-Madison © Brian S. Yandell

44

EM-MCMC duality

- EM approaches can be redone with MCMC
 - EM estimates & maximizes
 - MCMC draws random samples
 - simulated annealing: gradually cool towards peak
 - both can address same problem
- sometimes EM is hard (impossible) to use
- MCMC is tool of “last resort”
 - use exact methods if you can
 - try other approximate methods
 - be clever! (math, computing tricks)
 - very handy for hard problems in genetics

March 2011

UW-Madison © Brian S. Yandell

45

Why not Ordinary Monte Carlo?

- independent samples of joint distribution
- chaining (or peeling) of effects
$$\text{pr}(\theta|Y,Q) = \text{pr}(G_Q | Y, Q, \sigma^2) \text{pr}(\sigma^2 | Y, Q)$$
- possible analytically here given genotypes Q
- Monte Carlo: draw N samples from posterior
 - sample variance σ^2
 - sample genetic values G_Q given variance σ^2
- but we know markers X , not genotypes Q !
 - would have messy average over possible Q
 - $\text{pr}(\theta|Y,X) = \text{sum}_Q \text{pr}(\theta|Y,Q) \text{pr}(Q|Y,X)$

March 2011

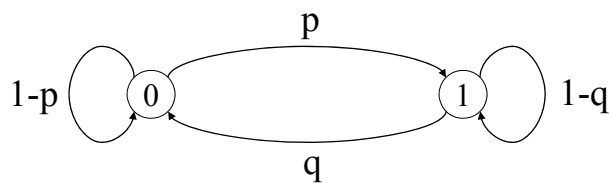
UW-Madison © Brian S. Yandell

46

What is a Markov chain?

- future given present is independent of past
- update chain based on current value
 - can make chain arbitrarily complicated
 - chain converges to stable pattern $\pi()$ we wish to study

$$\text{pr}(1) = p/(p + q)$$

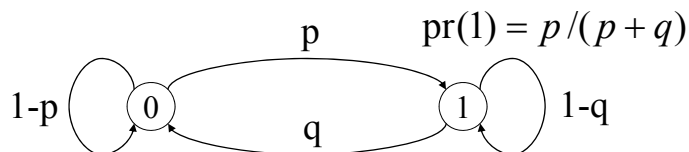


March 2011

UW-Madison © Brian S. Yandell

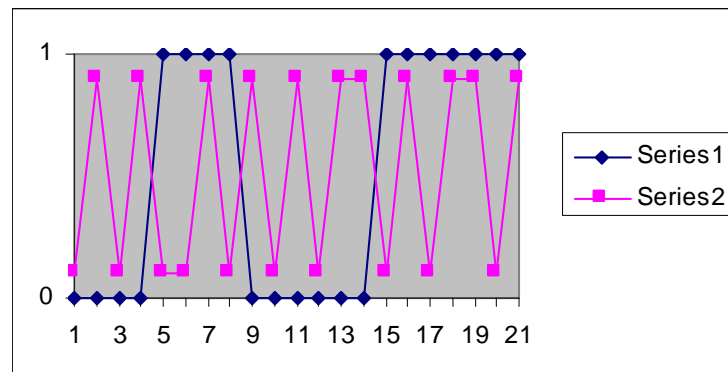
47

Markov chain idea



Mitch's

other
pubs



March 2011

UW-Madison © Brian S. Yandell

48

Markov chain Monte Carlo

- can study arbitrarily complex models
 - need only specify how parameters affect each other
 - can reduce to specifying full conditionals
- construct Markov chain with “right” model
 - joint posterior of unknowns as limiting “stable” distribution
 - update unknowns given data and all other unknowns
 - sample from full conditionals
 - cycle at random through all parameters
 - next step depends only on current values
- nice Markov chains have nice properties
 - sample summaries make sense
 - consider almost as random sample from distribution
 - ergodic theorem and all that stuff

March 2011

UW-Madison © Brian S. Yandell

49

MCMC sampling of (λ, q, μ)

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

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

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

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



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

April 2008

UW-Madison © Brian S. Yandell

50

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

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

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

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

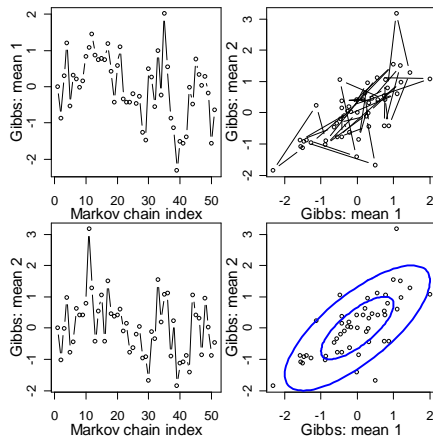
April 2008

UW-Madison © Brian S. Yandell

51

Gibbs sampler samples: $\rho = 0.6$

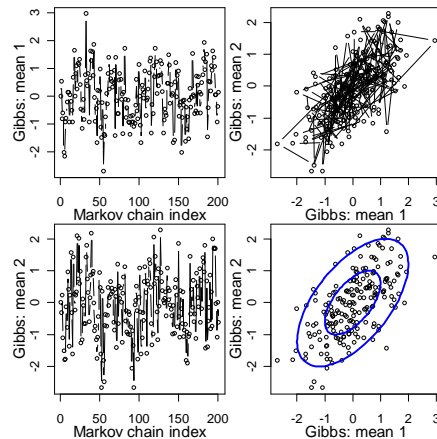
$N = 50$ samples



April 2008

UW-Madison © Brian S. Yandell

$N = 200$ samples

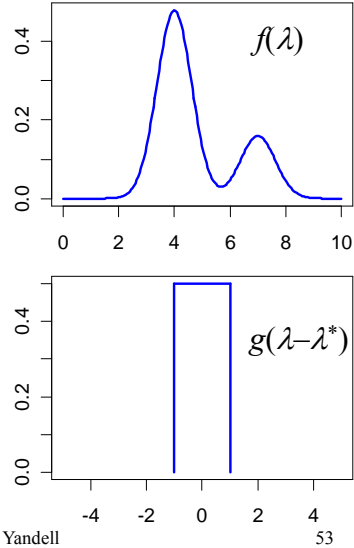


52

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

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

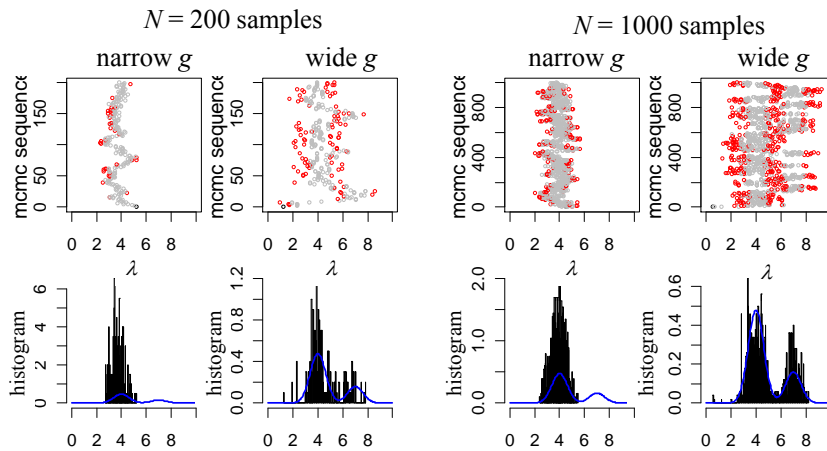


April 2008

UW-Madison © Brian S. Yandell

53

Metropolis-Hastings samples

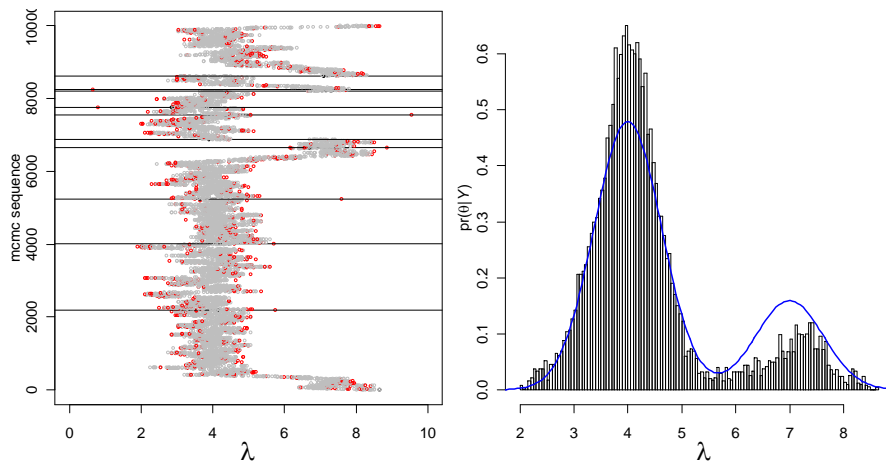


April 2008

UW-Madison © Brian S. Yandell

54

MCMC realization



added twist: occasionally propose from whole domain

April 2008

UW-Madison © Brian S. Yandell

55

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

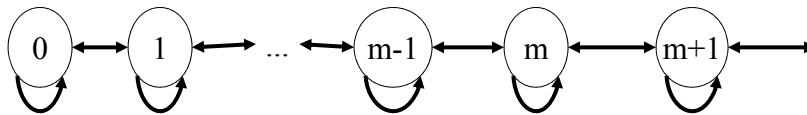
April 2008

UW-Madison © Brian S. Yandell

56

Markov chain for number m

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

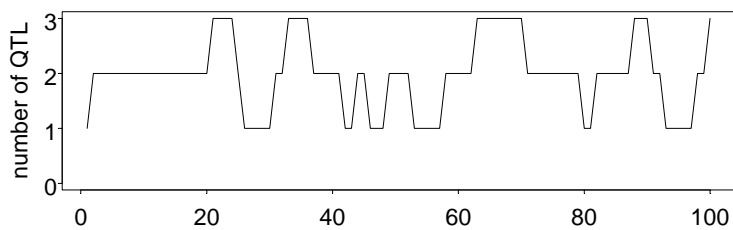
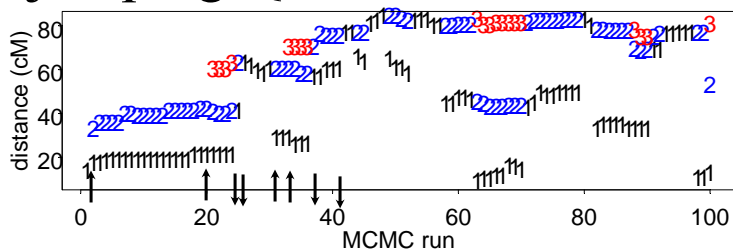


March 2011

UW-Madison © Brian S. Yandell

57

jumping QTL number and loci



March 2011

UW-Madison © Brian S. Yandell

58

Whole Genome 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
 - $|\gamma| = \Sigma\gamma =$ size of genetic architecture
 - $\lambda =$ loci determined implicitly by γ
 - $\beta =$ genotypic effects (main and epistatic)
 - $\mu =$ reference

April 2008

UW-Madison © Brian S. Yandell

59

Methods of Model Search

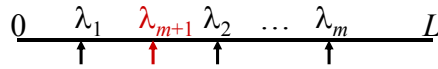
- Reversible jump (transdimensional) MCMC
 - sample possible loci (λ determines possible γ)
 - collapse to model containing just those QTL
 - bookkeeping when model dimension changes
- Composite model with indicators
 - include all terms in model: β and γ
 - sample possible architecture (γ determines 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)

April 2008

UW-Madison © Brian S. Yandell

60

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

April 2008

UW-Madison © Brian S. Yandell

61

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

$$\begin{array}{l} \curvearrowright nqtl = 1 : Y = \beta_0 + \beta_1(q_1) + e \\ \curvearrowright nqtl = 2 : Y = \beta_0 + \beta_1(q_1) + \beta_2(q_2) + e \end{array}$$

April 2008

UW-Madison © Brian S. Yandell

62

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)

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

April 2008

UW-Madison © Brian S. Yandell

63

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

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

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

April 2008

UW-Madison © Brian S. Yandell

64

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)

April 2008

UW-Madison © Brian S. Yandell

65

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

March 2011

UW-Madison © Brian S. Yandell

66

reversible jump choices

action step: draw one of three choices

(m = number of QTL in model)

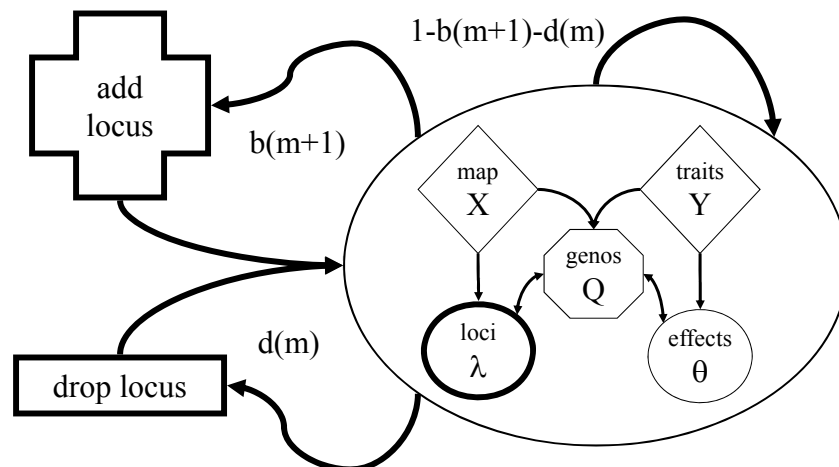
- update step with probability $1-b(m+1)-d(m)$
 - update current model
 - loci, effects, genotypes as before
- add a locus with probability $b(m+1)$
 - propose a new locus
 - innovate effect and genotypes at new locus
 - decide whether to accept the “birth” of new locus
- drop a locus with probability $d(m)$
 - pick one of existing loci to drop
 - decide whether to accept the “death” of locus

March 2011

UW-Madison © Brian S. Yandell

67

RJ-MCMC updates

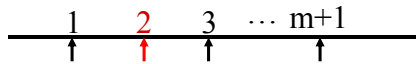


March 2011

UW-Madison © Brian S. Yandell

68

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

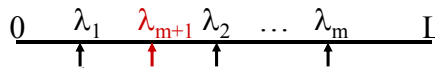
$$q_d(r; m + 1) = \frac{1}{m + 1}$$

March 2011

UW-Madison © Brian S. Yandell

69

propose to add a locus



- propose a new locus
 - uniform chance over genome
 - 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 ...

$$q_b(\lambda) = 1/L$$

March 2011

UW-Madison © Brian S. Yandell

70

acceptance of reversible jump

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

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

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

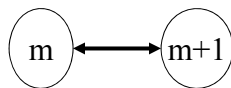
March 2011

UW-Madison © Brian S. Yandell

71

acceptance of reversible jump

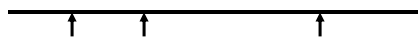
- move probabilities



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

- birth & death proposals

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



- Jacobian between models

- fudge factor
- see stepwise regression example

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

March 2011

UW-Madison © Brian S. Yandell

72

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

March 2011

UW-Madison © Brian S. Yandell

73

model selection in regression

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

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

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

March 2011

UW-Madison © Brian S. Yandell

74

slope estimate for 1 regressor

recall least squares estimate of slope

note relation of slope to correlation

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

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

March 2011

UW-Madison © Brian S. Yandell

75

2 correlated regressors

slopes adjusted for other regressors

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

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

March 2011

UW-Madison © Brian S. Yandell

76

Gibbs Sampler for Model 1

- mean $\mu \sim \phi\left(\eta + B_n(\bar{Y} - \eta), B_n \frac{\sigma^2}{n}\right), B_n = \frac{n}{n + \kappa}$
- slope $a \sim \phi\left(B_n \frac{\sum_{i=1}^n (Q_{i1} - \bar{Q}_1)(Y_i - \bar{Y})}{ns_1^2}, B_n \frac{\sigma^2}{ns_1^2}\right)$
- variance $\sigma^2 \sim \text{inv-}\chi^2\left(v + n, \frac{v\tau^2 + \sum_{i=1}^n (Y_i - \bar{Y} - a(Q_{i1} - \bar{Q}_1))^2}{v + n}\right)$

March 2011

UW-Madison © Brian S. Yandell

77

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 $\sigma^2 \sim \text{inv-}\chi^2\left(v + n, \frac{v\tau^2 + \sum_{i=1}^n \left(Y_i - \bar{Y} - \sum_{k=1}^2 a_k(Q_{ik} - \bar{Q}_k)\right)^2}{v + n}\right)$

March 2011

UW-Madison © Brian S. Yandell

78

updates from 2->1

- drop 2nd regressor
- adjust other regressor

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

$$a_2 \rightarrow 0$$

March 2011

UW-Madison © Brian S. Yandell

79

updates from 1->2

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

$$z \sim \phi(0,1), \quad J = \frac{\sigma}{s_{21}\sqrt{n}}$$

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

$$a_1 \rightarrow a - a_2 c_{21} = a - z \times c_{21} J - \hat{a}_2 c_{21}$$

March 2011

UW-Madison © Brian S. Yandell

80

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 |
| $1 \rightarrow 2$ | $(\mu, a, \sigma^2; z)$ | $z \sim \phi(0, 1)$ | $\begin{cases} a_2 \rightarrow \hat{a}_2 + z \times J \\ a_1 \rightarrow a - a_2 c_{21} \end{cases}$ |
| $2 \rightarrow 1$ | $(\mu, a_1, a_2, \sigma^2)$ | | $\begin{cases} a \rightarrow a_1 + a_2 c_{21} \\ z \rightarrow 0 \end{cases}$ |

March 2011

UW-Madison © Brian S. Yandell

81

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)

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

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

March 2011

UW-Madison © Brian S. Yandell

82

Jacobian & the calculus

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

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

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

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

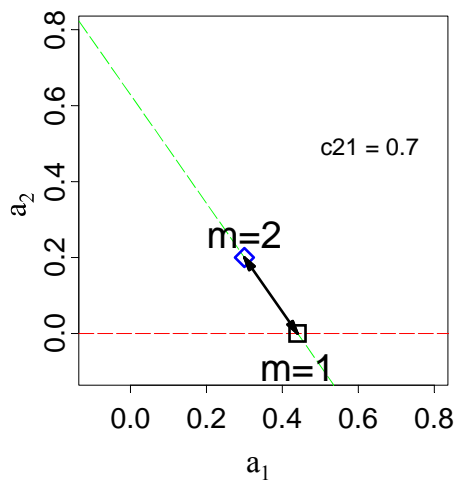
March 2011

UW-Madison © Brian S. Yandell

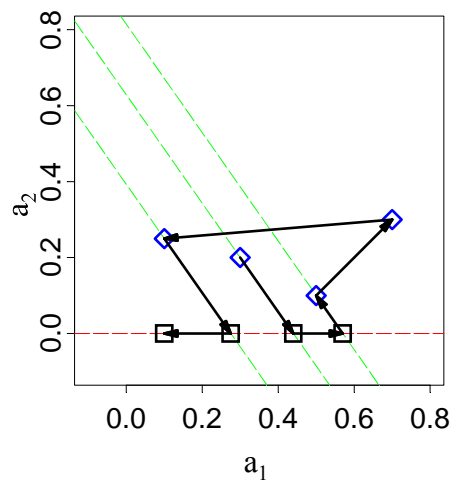
83

geometry of reversible jump

Move Between Models



Reversible Jump Sequence



March 2011

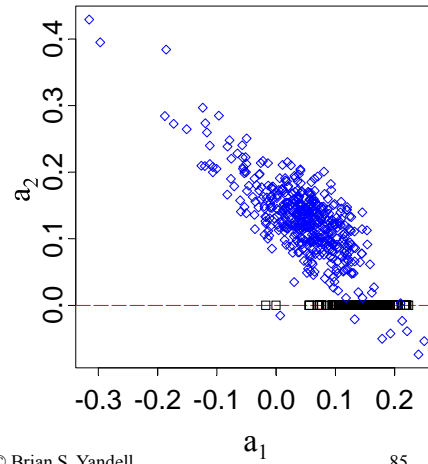
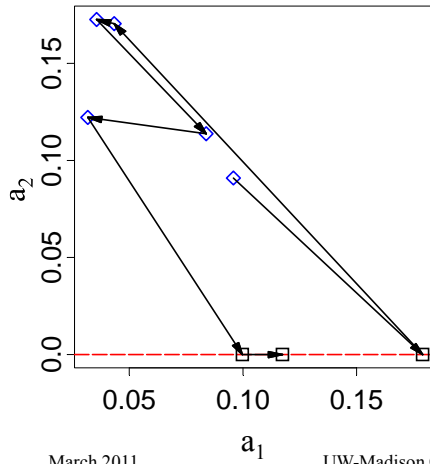
UW-Madison © Brian S. Yandell

84

QT additive reversible jump

a short sequence

first 1000 with $m < 3$



March 2011

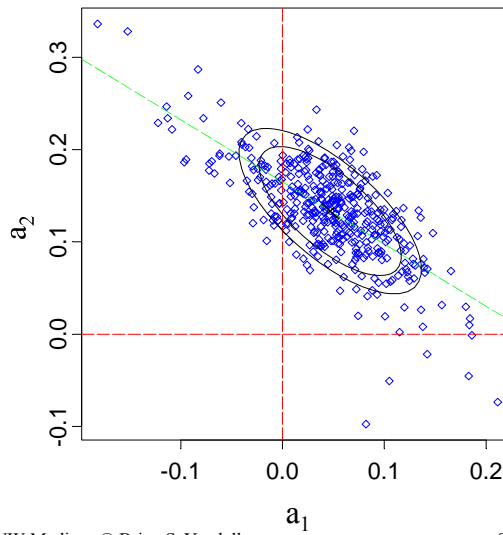
UW-Madison © Brian S. Yandell

85

credible set for additive

90% & 95% sets
based on normal

regression line
corresponds to
slope of updates

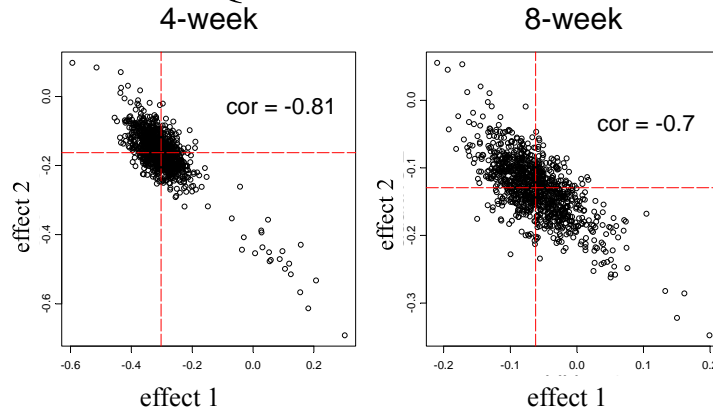


March 2011

UW-Madison © Brian S. Yandell

86

collinear QTL = correlated effects



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

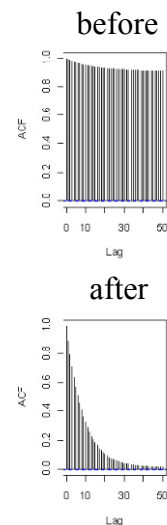
April 2008

UW-Madison © Brian S. Yandell

87

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



March 2011

UW-Madison © Brian S. Yandell

88

8 QTL simulation (Stevens Fisch 1998)

- $n = 200, h^2 = .5$
– SF detected 3 QTL
- Bayesian IM

| n | h^2 | detect |
|-----|-------|--------|
| 200 | .5 | 2 |
| 200 | .8 | 4 |
| 500 | .9 | 7 |
| 500 | .97 | 8 |

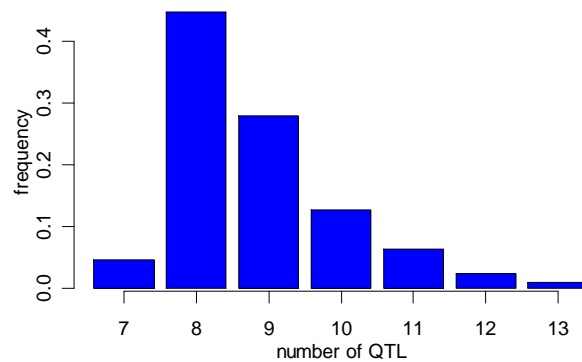
| QTL No. j | Location, λ_j | | | Additive effect α_j | Dominance Effect δ_j |
|----------------|-------------------------|-------------------------|---------------|-------------------------------|--------------------------------|
| | Chrom. λ_j^c | Marker λ_j^m | Position (cM) | | |
| 1 | 1 | 1 | 11 | -3 | 0 |
| 2 | 1 | 3 | 10 | -5 | 0 |
| 3 | 3 | 4 | 2 | 2 | 0 |
| 4 | 6 | 6 | 7 | -3 | 0 |
| 5 | 6 | 8 | 12 | 3 | 0 |
| 6 | 8 | 2 | 12 | -4 | 0 |
| 7 | 8 | 3 | 14 | 1 | 0 |
| 8 | 9 | 10 | 15 | 2 | 0 |

March 2011

UW-Madison © Brian S. Yandell

89

posterior number of QTL



geometric prior with mean 0.5
seems to have no influence on posterior here

March 2011

UW-Madison © Brian S. Yandell

90

posterior genetic architecture

| | Chromosome count vector | | | | | | | | | | |
|---|-------------------------|---|---|---|---|---|---|---|---|----|-------|
| m | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | Count |
| 8 | 2 | 0 | 1 | 0 | 0 | 2 | 0 | 2 | 1 | 0 | 3371 |
| 9 | 3 | 0 | 1 | 0 | 0 | 2 | 0 | 2 | 1 | 0 | 751 |
| 7 | 2 | 0 | 1 | 0 | 0 | 2 | 0 | 1 | 1 | 0 | 377 |
| 9 | 2 | 0 | 1 | 0 | 0 | 3 | 0 | 2 | 1 | 0 | 218 |
| 9 | 2 | 0 | 1 | 0 | 0 | 2 | 0 | 2 | 2 | 0 | 198 |

March 2011

UW-Madison © Brian S. Yandell

91

Bayesian model averaging

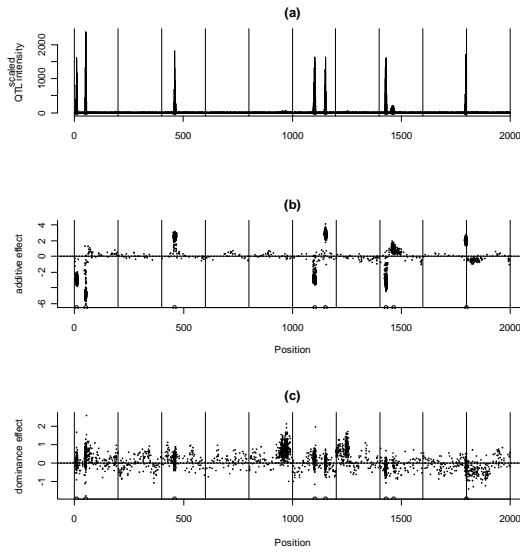
- average summaries over multiple architectures
- avoid selection of “best” model
- focus on “better” models
- examples in data talk later

April 2008

UW-Madison © Brian S. Yandell

92

model averaging for 8 QTL

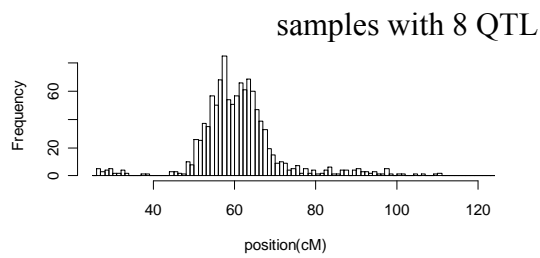
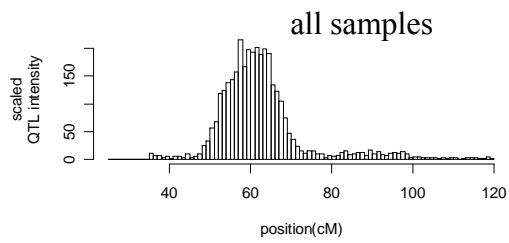


March 2011

UW-MADISON © DEBBI S. TANIGUCHI

93

model averaging: focus on chr. 8



March 2011

94

5. Model Assessment

- balance model fit against model complexity

| | | |
|----------------|-------------------|------------------|
| model fit | smaller model | bigger model |
| prediction | miss key features | fits better |
| interpretation | may be biased | no bias |
| parameters | easier | more complicated |
| | 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 $\text{pr}(\text{data} | \text{model})$

April 2008

UW-Madison © Brian S. Yandell

95

Bayes factors

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

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

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

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

April 2008

UW-Madison © Brian S. Yandell

96

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

April 2008

UW-Madison © Brian S. Yandell

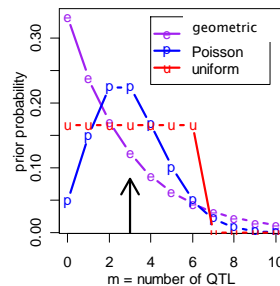
97

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

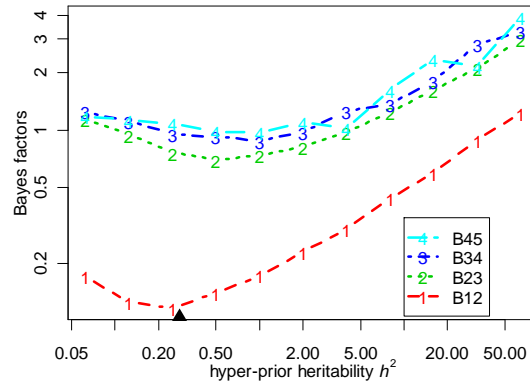


April 2008

UW-Madison © Brian S. Yandell

98

BF sensitivity to fixed prior for effects



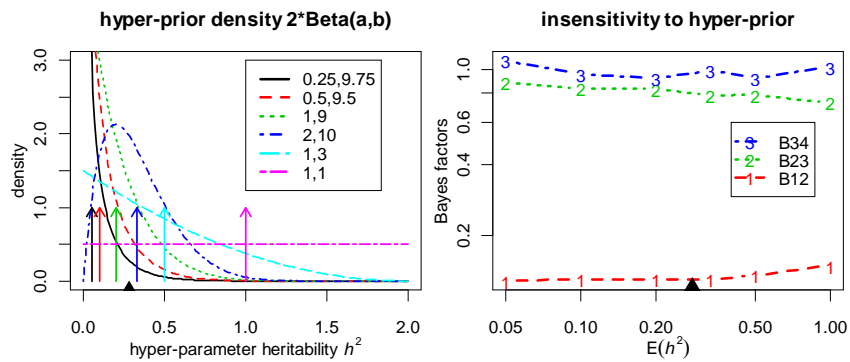
$$\beta_{qj} \sim N(0, \sigma_G^2 / m), \sigma_G^2 = h^2 \sigma_{\text{total}}^2, h^2 \text{ fixed}$$

April 2008

UW-Madison © Brian S. Yandell

99

BF insensitivity to random effects prior



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

April 2008

UW-Madison © Brian S. Yandell

100

How sensitive is posterior to choice of prior?

- simulations with 0, 1 or 2 QTL
 - strong effects (additive = 2, variance = 1)
 - linked loci 36cM apart
- differences with number of QTL
 - clear differences by actual number
 - works well with 100,000, better with 1M
- effect of Poisson prior mean
 - larger prior mean shifts posterior up
 - but prior does not take over

March 2011

UW-Madison © Brian S. Yandell

101

simulation study: prior

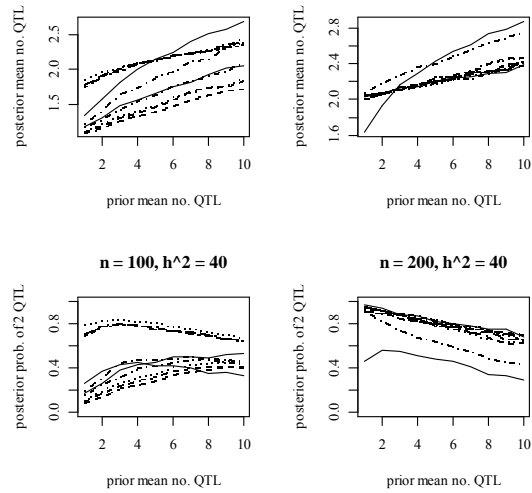
- 2 QTL at 15, 65cM
- $n = 100, 200$; $h^2 = 40\%$
- vary prior mean from 1 to 10 QTL
 - Poisson prior
- 10 independent simulations
- examine posterior mean, probability

March 2011

UW-Madison © Brian S. Yandell

102

posterior m depends on prior

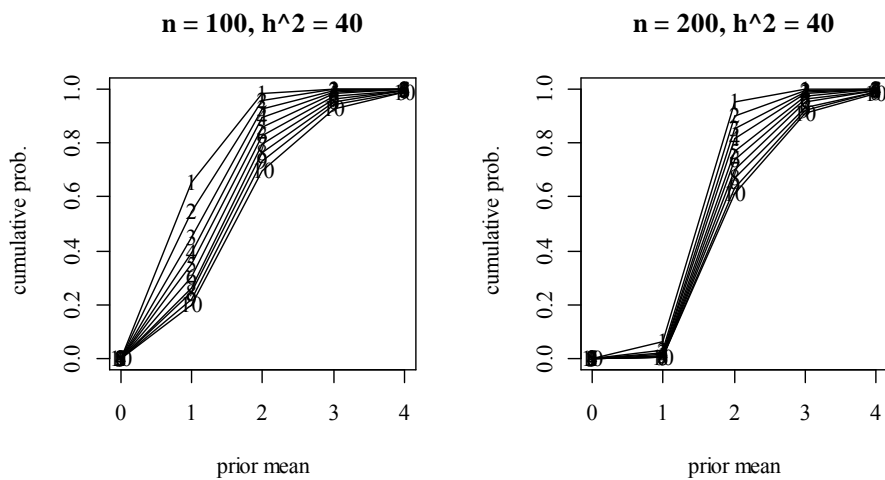


March 2011

UW-Madison © Brian S. Yandell

103

cumulative posterior as prior mean changes

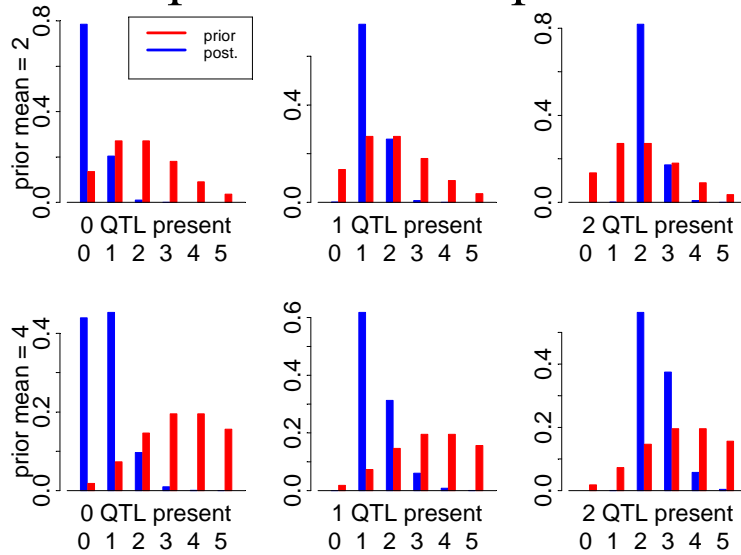


March 2011

UW-Madison © Brian S. Yandell

104

effect of prior mean on posterior m

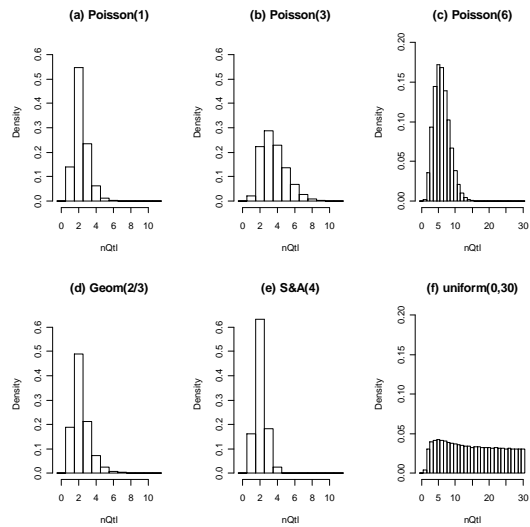


March 2011

UW-Madison © Brian S. Yandell

105

effect of prior shape on posterior



March 2011

UW-Madison © Brian S. Yandell

106

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

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

April 2008

UW-Madison © Brian S. Yandell

107

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

April 2008

UW-Madison © Brian S. Yandell

108

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

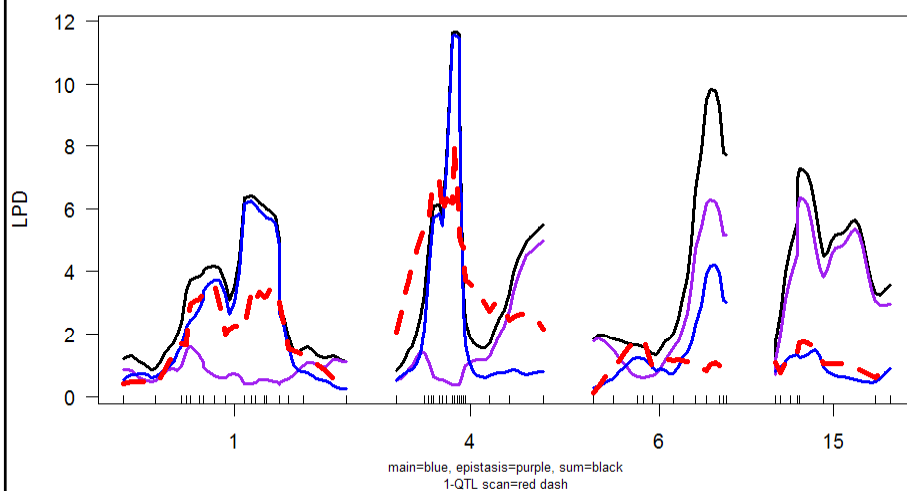
April 2008

UW-Madison © Brian S. Yandell

109

hyper data: scanone

LPD of bp for main+epistasis+sum

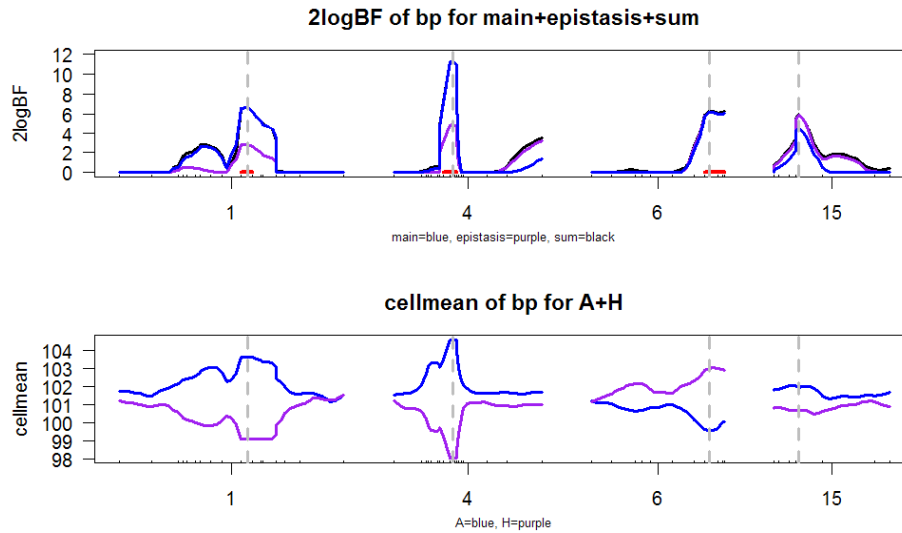


April 2008

UW-Madison © Brian S. Yandell

110

2log(BF) scan with 50% HPD region

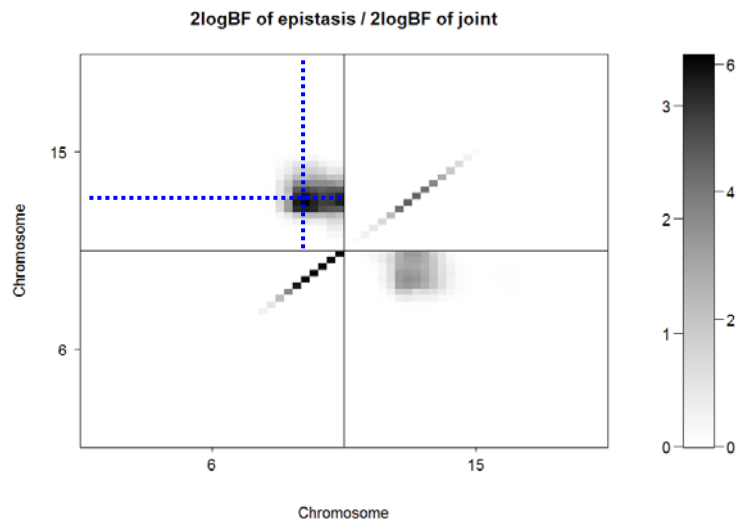


April 2008

UW-Madison © Brian S. Yandell

111

2-D plot of 2logBF: chr 6 & 15

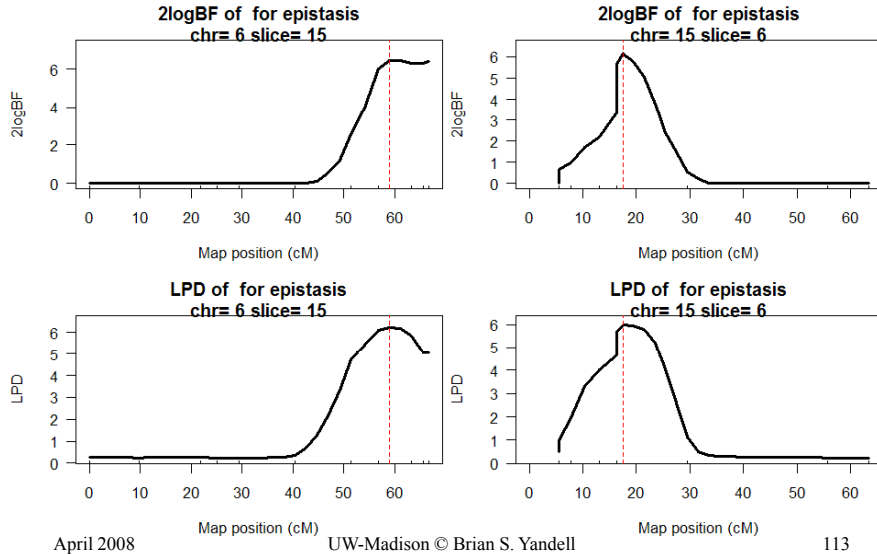


April 2008

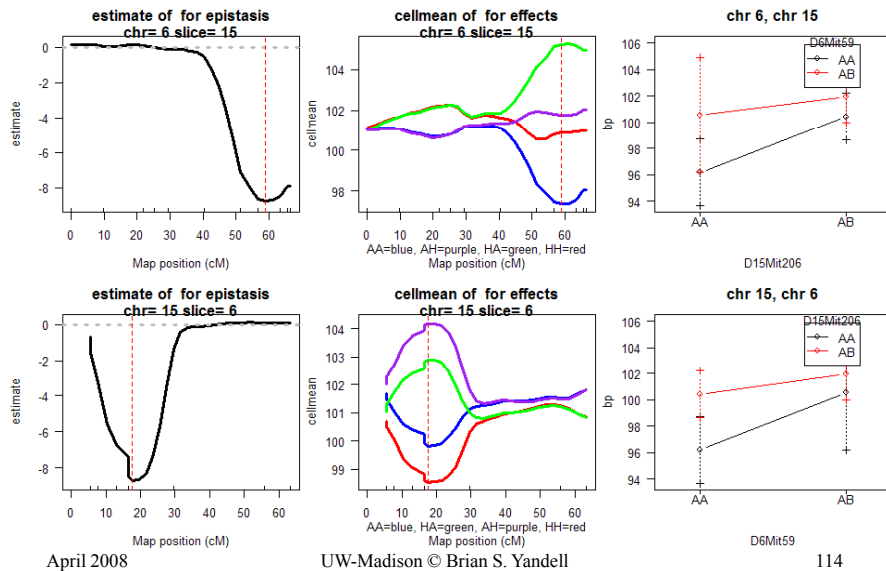
UW-Madison © Brian S. Yandell

112

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?

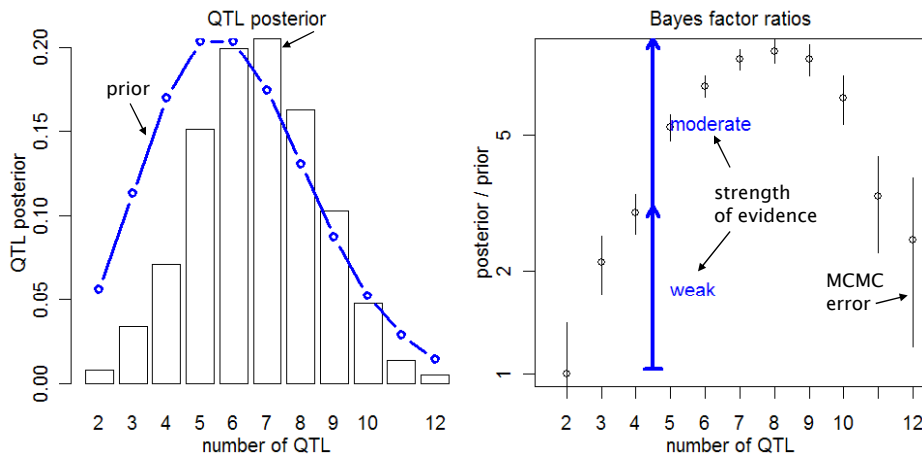
- How many QTL?
- What is pattern across chromosomes?
- examine posterior relative to prior
 - prior determined ahead of time
 - posterior estimated by histogram/bar chart
 - Bayes factor ratio = $\text{pr}(\text{model}|\text{data}) / \text{pr}(\text{model})$

April 2008

UW-Madison © Brian S. Yandell

115

How many QTL? posterior, prior, Bayes factor ratios



April 2008

UW-Madison © Brian S. Yandell

116

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 |

April 2008

UW-Madison © Brian S. Yandell

117

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

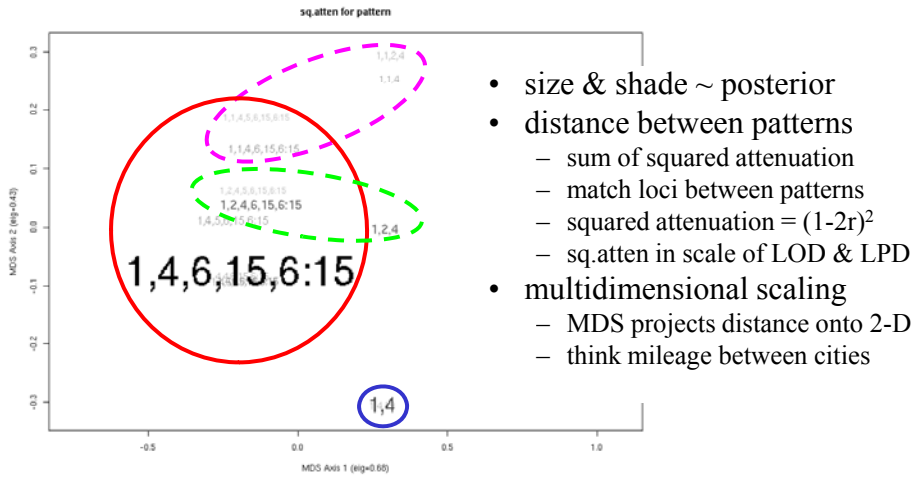
| | chrom | locus | locus.LCL | locus.UCL | n.qtl |
|-----|-------|-------|-----------|-----------|-----------|
| 247 | 1 | 69.9 | 24.44875 | 95.7985 | 0.8026667 |
| 245 | 4 | 29.5 | 14.20000 | 74.3000 | 0.8800000 |
| 248 | 6 | 59.0 | 13.83333 | 66.7000 | 0.7096667 |
| 246 | 15 | 19.5 | 13.10000 | 55.7000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

April 2008

UW-Madison © Brian S. Yandell

118

how close are other patterns?



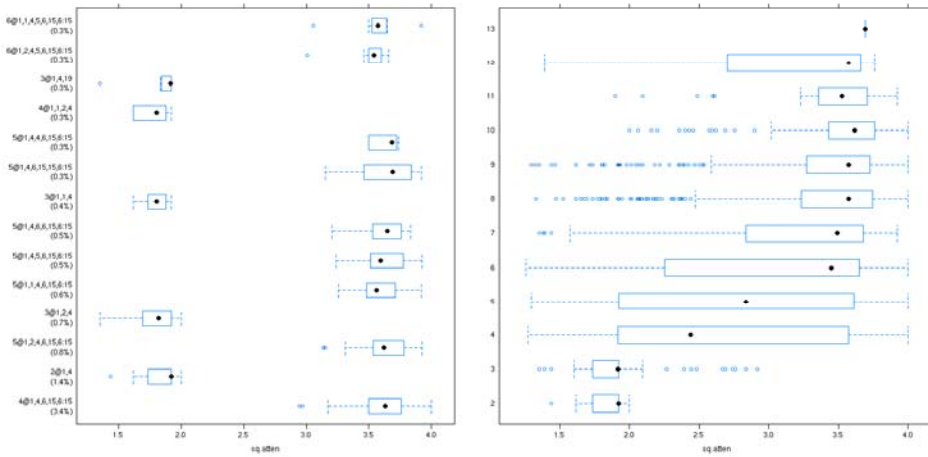
- size & shade ~ posterior
- distance between patterns
 - sum of squared attenuation
 - match loci between patterns
 - squared attenuation = $(1-2r)^2$
 - sq.atten in scale of LOD & LPD
- multidimensional scaling
 - MDS projects distance onto 2-D
 - think mileage between cities

April 2008

UW-Madison © Brian S. Yandell

119

how close are other patterns?

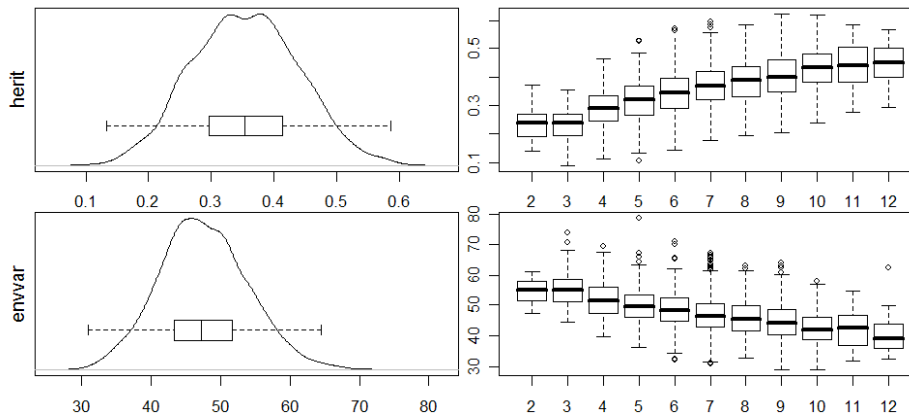


April 2008

UW-Madison © Brian S. Yandell

120

diagnostic summaries



April 2008

UW-Madison © Brian S. Yandell

121

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

April 2008

UW-Madison © Brian S. Yandell

122

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)

April 2008

UW-Madison © Brian S. Yandell

123

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

April 2008

UW-Madison © Brian S. Yandell

124

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