# 4 Interval Mapping for a Single QTL

- basic idea of interval mapping
- interval mapping by maximum likelihood
  - maximum likelihood using EM, MCMC
- Bayesian interval mapping
  - "natural" Bayesian priors
  - multiple imputation, MCMC
- bootstrapped variance estimates
- advantages & shortcomings of IM
- Haley-Knott regression approximation

# 4.1 basic idea of interval mapping

- study properties of likelihood at each possible QTL
  - treating QTL as missing data
  - assuming only a single QTL (for now)
- recall likelihood as mixture over unknown QTL
  - likelihood = product of sum of products
  - complicated to evaluate--requires iteration

$$L(\theta, \lambda \mid Y, X) = \mathrm{pr}(Y \mid X, \theta, \lambda)$$
$$= \mathrm{prod}_i \, \mathrm{pr}(Y_i \mid X_i, \theta, \lambda)$$
$$= \mathrm{prod}_i \, \mathrm{sum}_Q \, \mathrm{pr}(Q \mid X_i, \lambda) \mathrm{pr}(Y_i \mid Q, \theta)$$
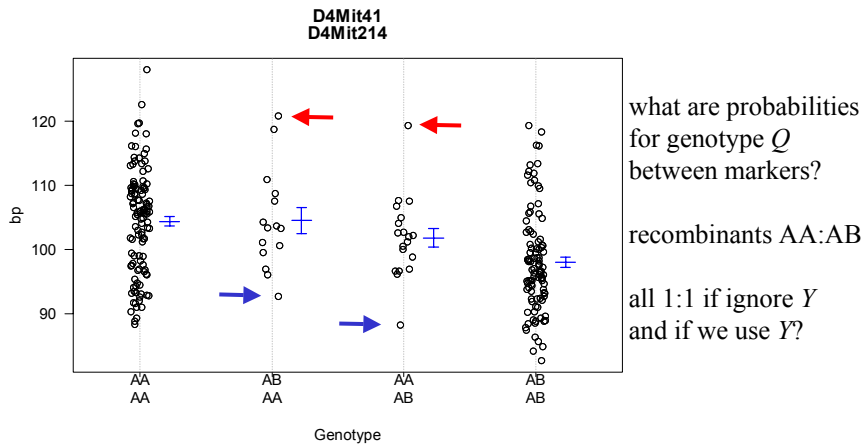
# uncertainty in QTL genotype *Q*

- how to improve guess on *Q* with data, parameters?
  - prior recombination: $\text{pr}(Q \mid X_i, \lambda)$
  - posterior recombination: $\text{pr}(Q \mid Y_i, X_i, \theta, \lambda)$
- main philosophies for assessing likelihood
  - maximum likelihood: study peak(s)
  - Bayesian analysis: study whole shape
- implementation methodologies
  - Expectation-Maximization (EM)
  - Markov chain Monte Carlo (MCMC)
  - multiple imputation
  - genetic algorithms, GEE, …

# posterior on QTL genotypes

- full conditional of *Q* given data, parameters
  - proportional to prior $\text{pr}(Q \mid X_i, \lambda)$
    - weight toward *Q* that agrees with flanking markers
  - proportional to likelihood $\text{pr}(Y_i \mid Q, \theta)$
    - weight toward *Q* so that group mean $G_Q \approx Y_i$
- phenotype and flanking markers may conflict
  - posterior recombination balances these two weights

$$\text{pr}(Q \mid Y_i, X_i, \theta, \lambda) = \frac{\text{pr}(Q \mid X_i, \lambda)\text{pr}(Y_i \mid Q, \theta)}{\text{pr}(Y_i \mid X_i, \theta, \lambda)}$$

# how does phenotype *Y* affect *Q*?

**D4Mit41**
**D4Mit214**



what are probabilities for genotype *Q* between markers?

recombinants AA:AB

all 1:1 if ignore *Y* and if we use *Y*?

---

# maximum likelihood (ML) idea

- pick QTL locus $\lambda$ (usually scan whole genome)
- find ML estimates of gene action $\theta$ given $\lambda$
- maximum likelihood at peak of likelihood
  - slope (derivative with respect to $\theta$) is zero
  - sometimes maximum is at a boundary (non-zero slope)
- slope is weighted average using posteriors for *Q*
  - cannot write estimate in "closed form"
  - need to know $\theta$ to estimate it!
  - iterate toward the maximum in some clever way

$$\frac{dL(\theta, \lambda \mid Y, X)}{d\theta} = \text{sum}_{i,Q} \, \text{pr}(Q \mid Y_i, X_i, \theta, \lambda) \frac{d \log(\text{pr}(Y_i \mid Q, \theta))}{d\theta}$$

# Bayesian model posterior

- augment data (*Y,X*) with unknowns *Q*
  - study unknowns ($\theta,\lambda,Q$) given data (*Y,X*)
  - $Q \sim \mathrm{pr}(Q \mid Y_i, X_i, \theta, \lambda)$
- no longer need weighted average over *Q*
  - instead we average over *Q* to study parameters
  - $\mathrm{pr}(\theta,\lambda \mid Y,X) = \mathrm{sum}_Q \mathrm{pr}(\theta,\lambda,Q \mid Y,X)$
- study properties of posterior
  - need to specify priors for ($\theta,\lambda$)
  - denominator is very difficult to compute in practice
  - drawing samples from posterior in some clever way

$$\mathrm{pr}(\theta,\lambda,Q \mid Y,X) = \frac{\mathrm{pr}(Q \mid X,\lambda)\mathrm{pr}(Y \mid Q,\theta)\mathrm{pr}(\lambda \mid X)\mathrm{pr}(\theta)}{\mathrm{pr}(Y \mid X)}$$
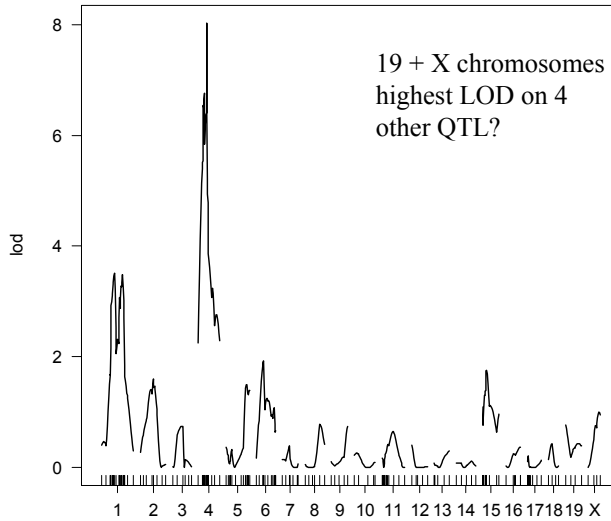
---

# 4.2 interval mapping by ML

- search whole genome for putative QTL
- "profile" likelihood across all possible $\lambda$
  - find ML estimate of $\theta$ given $\lambda$
  - ML estimate of ($\theta,\lambda$) at maximum over genome

$$L_0(\hat{\theta}_0 \mid Y) = \mathrm{prod}_i f(Y_i \mid \hat{\mu}, s^2)$$

$$L(\hat{\theta},\lambda \mid Y,X) = \mathrm{prod}_i \mathrm{sum}_Q \mathrm{pr}(Q \mid X,\lambda) f(Y_i \mid \hat{G}_Q, \hat{\sigma}^2_{pool})$$

$$LOD(\lambda) = \log_{10}\left( \frac{L(\hat{\theta},\lambda \mid Y,X)}{L_0(\hat{\theta}_0 \mid Y)} \right)$$
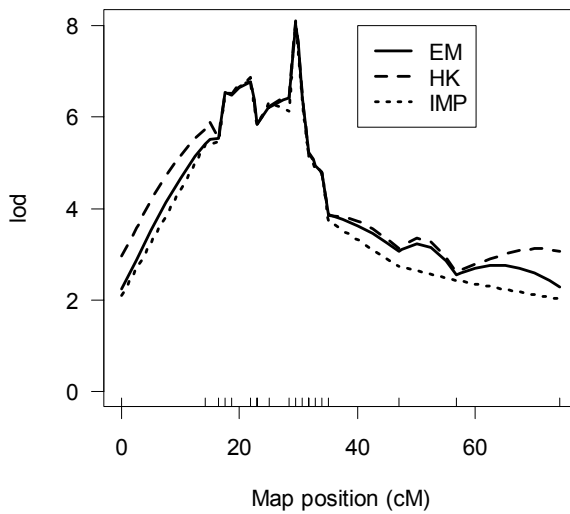
# LOD for hyper dataset



8

6

lod

4

2

0

19 + X chromosomes
highest LOD on 4
other QTL?

1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 X

# LOD($\lambda$) on chr 4 of hyper



8

6

lod

4

2

0

—— EM
– – – HK
······ IMP

0          20          40          60

Map position (cM)

EM "exact"
Haley-Knott regression
single imputation

all agree at the peak
and mostly at markers

note marker spacing

# EM method for interval mapping

- fix a possible QTL $\lambda$
- iterate between expectation & maximization
  - likelihood increases with each iteration
  - stop iterating when the change is "negligible"
- initial values
  - $P_{Qi} = \mathrm{pr}(Q \mid X_i, \lambda)$
    - recombination model in the absence of data
  - or use Haley-Knott regression estimates of $\theta$

# EM method for interval mapping

- E-step: estimate posterior recombination
  - $P_{Qi} = \mathrm{pr}(Q \mid Y_i, X_i, \theta, \lambda)$
  - estimate for every individual $i$, genotype $Q$
  - depends on effects $\theta$
- M-steps: maximize likelihood for $\theta$
  - may be many parameters
  - technical point: caution on parallel updates
  - solve system of equation: derivatives set to zero
  - depends on $P_{Qi}$

$$0 = \mathrm{sum}_{i,Q}\ P_{Qi}\ \frac{d\log(\mathrm{pr}(Y_i \mid Q,\theta))}{d\theta}$$

# 4.2.2 M-step for normal phenotype

- $Y = G_Q + e,\ e \sim N(0, \sigma^2)$
- $\text{pr}(Y \mid Q, \theta) = f(Y \mid G_Q,\ \sigma^2)$
- see notes in book for derivative details
- E-step estimates:

$$\hat{G}_Q = \text{sum}_i\ Y_i P_{Qi} / \text{sum}_i\ P_{Qi}$$

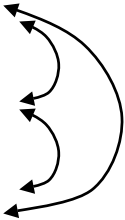$$\hat{\sigma}^2 = \text{sum}_{i,Q}\ \left(Y_i - \hat{G}_Q\right)^2 P_{Qi} / n$$

# 4.2.3 ML via MCMC

- basic idea of simulated annealing
  - start with non-informative priors on $(\theta, \lambda)$
  - sample from posterior (somehow…)
  - gradually shrink priors toward ML estimate
- slight difficulty
  - need to know $(\theta, \lambda)$ to sample from posterior
  - iteration leads to Markov chain
- point of this section
  - MCMC does not imply a Bayesian perspective!

# 4.3 Bayesian interval mapping

- sample missing genotypes $Q$
- decouple effects $\theta$ from QTL $\lambda$
- but $Q$ depends on $(\theta, \lambda)$ and vice versa
- also need to specify priors

$$\lambda \sim \frac{\text{pr}(Q \mid X, \lambda)\text{pr}(\lambda \mid X)}{\text{pr}(Q \mid X)}$$

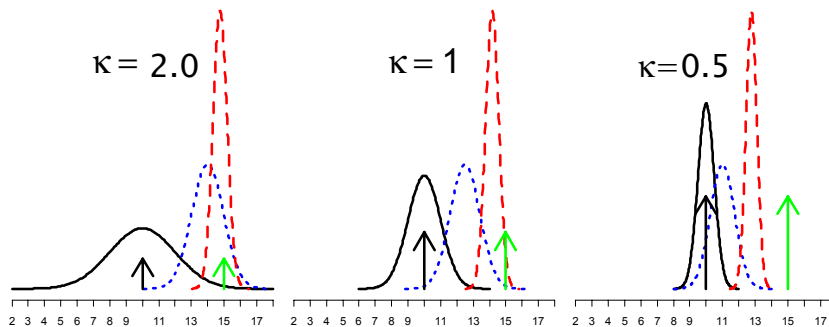$$Q \sim \text{pr}(Q \mid Y_i, X_i, \theta, \lambda)$$

$$\theta \sim \frac{\text{pr}(Y \mid Q, \theta)\text{pr}(\theta)}{\text{pr}(Y \mid Q)}$$

# 4.3.1 Bayesian priors for QTL

- locus $\lambda$ may be uniform over genome
  - $\text{pr}(\lambda \mid X) = 1 / $ length of genome
- missing genotypes $Q$
  - $\text{pr}(Q \mid X, \lambda)$
  - recombination model is formally a prior
- effects $\theta = (G, \sigma^2)$, $G = (G_{QQ}, G_{Qq}, G_{qq})$
  - conjugate priors for normal phenotype
  - $G_Q \sim N(\mu, \kappa\sigma^2)$
  - $\sigma^2 \sim$ inverse-$\chi^2(\nu, \tau^2)$, or $\nu\tau^2 / \sigma^2 \sim \chi^2$

# effect of prior variance on posterior



$\kappa = 2.0$      $\kappa = 1$      $\kappa = 0.5$

normal prior, posterior for $n = 1$, posterior for $n = 5$ , true mean
(solid black)  (dotted blue)     (dashed red)    (green arrow)

---

# details of phenotype priors

- priors depend on "hyper-parameters"
- $G_Q \sim N(\mu,\ \kappa\sigma^2)$
  - center around phenotype grand mean
  - $\kappa\sigma^2 \approx \sigma_G^2 =$ genetic variance
  - $\kappa \approx \sigma_G^2 / \sigma^2 = h^2 / (1 - h^2)$
  - $h^2 = \sigma_G^2 / (\sigma_G^2 + \sigma^2) =$ heritability
- $\sigma^2 \sim$ inverse-$\chi^2(\nu, \tau^2)$, or $\nu\tau^2 / \sigma^2 \sim \chi^2$
  - $\tau^2 \approx s^2 =$ total sample variance
  - $\nu =$ prior degrees of freedom = small integer

# Bayes for normal data

$Y = G + E$  posterior for single individual

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

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

prior  $\mathrm{pr}(G \mid \sigma^2, \mu, \kappa) = N(G \mid \mu, \kappa \sigma^2)$

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

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

          shrinkage weights $B_n$ go to 1

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

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

---

# posterior by QT genetic value

$Y = G_Q + E$  genetic  $Q = qq, Qq, QQ$

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

                parameters  $\theta = (G, \sigma^2)$

likelihood  $\mathrm{pr}(Y \mid Q, G, \sigma^2) = N(Y \mid G_Q, \sigma^2)$

prior  $\mathrm{pr}(G_Q \mid \sigma^2, \mu, \kappa) = N(G_Q \mid \mu, \kappa \sigma^2)$

posterior:

$$\mathrm{pr}(G_Q \mid Y, Q, \sigma^2, \mu, \kappa) = N\left( G_Q \,\middle|\, \mu + B_Q(\bar{Y}_Q - \mu), B_Q \frac{\sigma^2}{n_Q} \right)$$

$$n_Q = \mathrm{count}\{Q_i = Q\}, \bar{Y}_Q = \sum_{\{i:Q_i=Q\}} \frac{Y_i}{n_Q}, B_Q = \frac{\kappa n_Q}{\kappa n_Q + 1} \to 1$$

## Empirical Bayes: choosing hyper-parameters

How do we choose hyper-parameters $\mu, \kappa$?

Empirical Bayes:  marginalize over prior

estimate $\mu, \kappa$ from marginal posterior

| | |
|---|---|
| likelihood | $\text{pr}(Y_i \mid Q_i, G, \sigma^2) = \text{N}(Y_i \mid G(Q_i), \sigma^2)$ |
| prior | $\text{pr}(G_Q \mid \sigma^2, \mu, \kappa) = \text{N}(G_Q \mid \mu, \kappa\sigma^2)$ |
| marginal | $\text{pr}(Y_i \mid \sigma^2, \mu, \kappa) = \text{N}(Y_i \mid \mu, (\kappa+1)\sigma^2)$ |
| estimates | $\hat{\mu} = \overline{Y}_\bullet, s^2 = \text{sum}_i (Y_i - \overline{Y}_\bullet)^2 / n$ |
| | $\hat{\sigma}^2 = s^2 / (\kappa+1) \approx s^2 / (1 - h^2)$ |
| EB posterior | $\text{pr}(G_Q \mid Y) = \text{N}\left( G_Q \middle| \overline{Y}_\bullet + B_Q(\overline{Y}_Q - \overline{Y}_\bullet), B_Q \dfrac{\hat{\sigma}^2}{n_Q} \right)$ |

---

# What if variance $\sigma^2$ is unknown?

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

# joint effects posterior details

$Y_i = G(Q_i) + E_i$      genetic      $Q_i = qq, Qq, QQ$

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

parameters      $\theta = (G, \sigma^2)$

likelihood      $pr(Y_i \mid Q_i, G, \sigma^2) = N(Y_i \mid G(Q_i), \sigma^2)$

prior      $pr(G_Q \mid \sigma^2, \mu, \kappa) = N(G_Q \mid \mu, \sigma^2/\kappa)$

$pr(\sigma^2 \mid \nu, \tau^2) = \text{inv-}\chi^2 (\sigma^2 \mid \nu, \tau^2)$

posterior:      $pr(G_Q \mid Y, Q, \sigma^2, \mu, \kappa) = N\left( G_Q \middle| \overline{Y}_. + B_Q(\overline{Y}_Q - \overline{Y}_.), B_Q \dfrac{\sigma^2}{n_Q} \right)$

$$pr(\sigma^2 \mid Y, Q, G_Q, \nu, \tau^2) = \text{inv-}\chi^2 \left( \sigma^2 \middle| \nu + n, \dfrac{\nu\tau^2 + n s_Q^2}{\nu + n} \right)$$

with $B_Q = \dfrac{n_Q}{\kappa + n_Q}, s_Q^2 = \text{sum}_i \left( Y_i - G(Q_i) \right)^2 / n$

---

# 4.3.2 Bayesian multiple imputation

- basic idea
  - impute multiple copies of missing genotypes $Q$
    - sample $Q \sim pr(Q \mid X, \lambda)$
    - weighted to appear as draws from posterior
  - average out gene effects $\theta$
  - study posterior for putative QTL $\lambda$
- most effective for multiple QTL
  - use single QTL to introduce idea
  - consider all loci as possible QTL
    - sample on grid $\Lambda$ of `pseudomarkers' (every 2cM)
    - similar to interval map scan of whole genome

# importance sampling idea

- draw samples from one distribution
  - $Q_1, Q_2, Q_3, ..., Q_n \sim f(Q)$
- weight them appropriately by $\omega(Q)$
- sample summaries from distribution $g(Q)$
  - $g(Q) = f(Q)\omega(Q)$ / constant
  - mean for $f$ = $\text{sum}_i Q_i$ / $n$
  - mean for $g$ = $\text{sum}_i Q_i \omega(Q_i)$ / $\text{sum}_i \omega(Q_i)$

---

# example: mean copies of Q

| genotype | qq | Qq | QQ | sum |
|---|---|---|---|---|
| Q copies | 0 | 1 | 2 | |
| true $g$ | 0.25 | 0.5 | 0.25 | 1.0 |
| draw $f$ | 1/3 | 1/3 | 1/3 | 1.0 |
| weight $\omega$ | 1 | 2 | 1 | |
| $f\times\omega$ | 1/3 | 2/3 | 1/3 | 4/3 |
| importance sampling | | | | $g = f\times 0.75\omega$ |
| sample | 30 | 30 | 40 | 100 |
| mean $f$ | 0×30 | 1×30 | 2×40 | 110/100=1.1 |
| mean $g$ | 0×1×30 | 1×2×30 | 2×1×40 | 140/130=1.08 |

# what are appropriate weights?

- ideally draw genotype from posterior
  - want sample $Q \sim g(Q) = \text{sum}_\theta \, \text{pr}(Q \mid Y,X,\theta,\lambda)\text{pr}(\theta)$
  - but have sample $Q \sim f(Q) = \text{pr}(Q \mid X,\lambda)$
- appropriate weights
  - $\omega(Q,\lambda \mid Y,X) = \text{pr}(\lambda \mid X) \, \text{sum}_\theta \, \text{pr}(Y \mid Q,\theta)\text{pr}(\theta)$
- estimate marginal posterior for QTL $\lambda$
  - draw $N$ samples from prior at each QTL $\lambda$

    $Q_1, Q_2, Q_3, ..., Q_N \sim \text{pr}(Q \mid X,\lambda)$

    $\text{pr}(\lambda \mid Y,X) = \text{sum}_Q \, \omega(Q,\lambda \mid Y,X) \, \text{pr}(Q \mid X,\lambda) \,/\, \text{constant}$

    $\qquad\qquad \approx \text{sum}_j \, \omega(Q_j,\lambda \mid Y,X) \,/\, \text{constant}$
  - constant is summed over all $\lambda$, but not actually needed

# relating weights to posterior

- posterior is simply averaged over $\theta$
- weights comprise terms except $\text{pr}(Q \mid X,\lambda)$
- estimating weights: see Sen & Churchill

$$\text{pr}(\lambda,Q \mid Y,X) = \text{sum}_\theta \, \text{pr}(\theta,\lambda,Q \mid Y,X)$$

$$= \frac{\text{pr}(Q \mid X,\lambda)\text{pr}(\lambda \mid X) \, \text{sum}_\theta \, \text{pr}(Y \mid Q,\theta)\text{pr}(\theta)}{\text{pr}(Y \mid X)}$$

$$= \text{pr}(Q \mid X,\lambda)\omega(Q,\lambda \mid Y,X) \,/\, \text{pr}(Y \mid X)$$

# estimating effects via imputation

- multiple imputation averages over effects
- difficult to study posterior of effects directly
- can estimate usual summaries

$$E(\theta \mid Y, X) = \text{sum}_Q \, E(\theta \mid Y, Q) \, \text{pr}(Q \mid Y, X)$$
$$= \text{sum}_{Q,\lambda} \, E(\theta \mid Y, Q) \, \text{pr}(Q \mid X, \lambda) \omega(Q, \lambda \mid Y, X) / \text{pr}(Y \mid X)$$
$$\approx \text{sum}_{\lambda, j} \, E(\theta \mid Y, Q_j) \omega(Q_j, \lambda \mid Y, X) / \text{constant}$$

---

# 4.3.3 Bayesian MCMC

- Markov chain Monte Carlo
  - Monte Carlo samples along a Markov chain
- What is a Markov chain?
- What is MCMC?
  - Sampling from full conditionals
  - Gibbs sampler, Metropolis-Hastings

# 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

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

# Markov chain idea
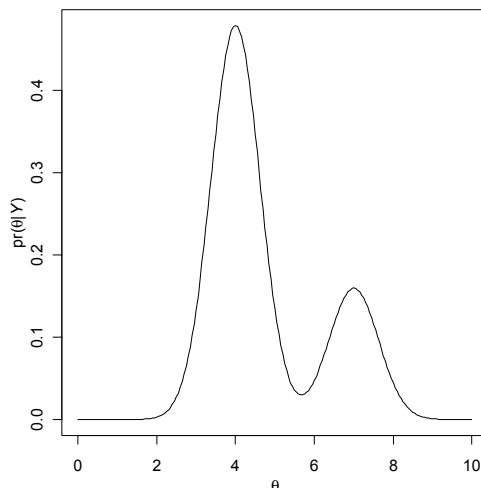
$$\text{pr}(1) = p/(p+q)$$

# 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

# Markov chain Monte Carlo idea

have posterior $\mathrm{pr}(\theta|Y)$
want to draw samples

propose $\theta \sim \mathrm{pr}(\theta|Y)$
(ideal: Gibbs sample)

propose new $\theta$ "nearby"
accept if more probable
toss coin if less probable
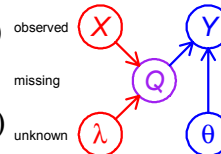based on relative heights
(Metropolis-Hastings)

# MCMC realization



added twist: occasionally propose from whole domain

---

# marginal posteriors

- joint posterior
  - $pr(\lambda, Q, \theta | Y, X) = pr(\theta) pr(\lambda) pr(Q|X, \lambda) pr(Y|Q, \theta)$ /constant
- genetic effects
  - $pr(\theta | Y, X) = \text{sum}_Q pr(\theta | Y, Q) pr(Q | Y, X)$ observed
- QTL locus
  - $pr(\lambda | Y, X) = \text{sum}_Q pr(\lambda | X, Q) pr(Q | Y, X)$ unknown
- QTL genotypes more complicated
  - $pr(Q | Y, X) = \text{sum}_{\lambda, \theta} pr(Q | Y, X, \lambda, \theta) pr(\lambda, \theta | Y, X)$
  - impossible to separate $\lambda$ and $\theta$ in sum

missing

# Why not Ordinary Monte Carlo?

- independent samples of joint distribution
- chaining (or peeling) of effects

$$\text{pr}(\theta|Y,Q)=\text{pr}(G_Q \mid Y,Q,\sigma^2)\,\text{pr}(\sigma^2 \mid Y,Q)$$

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

# MCMC Idea for QTLs

- 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
- update components from full conditionals
  - update effects $\theta$ given genotypes & traits
  - update locus $\lambda$ given genotypes & marker map
  - update genotypes $Q$ given traits, marker map, locus & effects

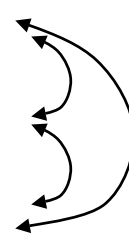$$(\lambda,Q,\theta) \sim \text{pr}(\lambda,Q,\theta \mid Y,X)$$
$$(\lambda,Q,\theta)_1 \rightarrow (\lambda,Q,\theta)_2 \rightarrow \cdots \rightarrow (\lambda,Q,\theta)_N$$

# sample from full conditionals

- hard to sample from joint posterior
- update each unknown given all others
- examine posterior: keep terms with unknown
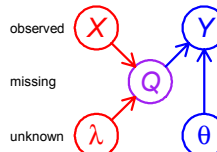- normalizing denominator make a distribution

$$\lambda \sim \frac{\mathrm{pr}(Q\,|\,X,\lambda)\mathrm{pr}(\lambda\,|\,X)}{\mathrm{pr}(Q\,|\,X)}$$

$$Q \sim \mathrm{pr}(Q\,|\,Y_i,X_i,\theta,\lambda)$$

$$\theta \sim \frac{\mathrm{pr}(Y\,|\,Q,\theta)\mathrm{pr}(\theta)}{\mathrm{pr}(Y\,|\,Q)}$$

---

# sample from full conditionals for model with *m* QTL

observed $X$    $Y$

missing $Q$

unknown $\lambda$   $\theta$

- hard to sample from joint posterior
  - $\mathrm{pr}(\lambda,Q,\theta\,|\,Y,X) = \mathrm{pr}(\theta)\mathrm{pr}(\lambda)\mathrm{pr}(Q|X,\lambda)\mathrm{pr}(Y|Q,\theta)$ /constant
- easy to sample parameters from full conditionals
  - full conditional for genetic effects
    - $\mathrm{pr}(\theta|Y,X,\lambda,Q) = \mathrm{pr}(\theta|Y,Q) = \mathrm{pr}(\theta)\,\mathrm{pr}(Y|Q,\theta)$ /constant
  - full conditional for QTL locus
    - $\mathrm{pr}(\lambda\,|Y,X,\theta,Q) = \mathrm{pr}(\lambda\,|X,Q) = \mathrm{pr}(\lambda)\,\mathrm{pr}(Q|X,\lambda)$ /constant
  - full conditional for QTL genotypes
    - $\mathrm{pr}(Q|Y,X,\lambda,\theta) = \mathrm{pr}(Q|X,\lambda)\,\mathrm{pr}(Y|Q,\theta)$ /constant

# Gibbs sampler idea

- want to study two correlated normals
- could sample directly from bivariate normal
- Gibbs sampler:
  - sample each from its full conditional
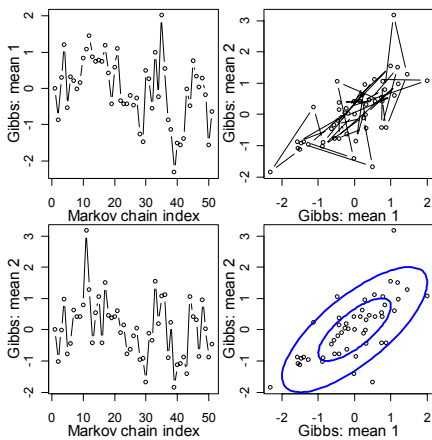  - pick order of sampling at random
  - repeat $N$ times

$$\binom{\theta_1}{\theta_2} \bigg| \mu, \rho \sim N\left(\binom{\mu_1}{\mu_2}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

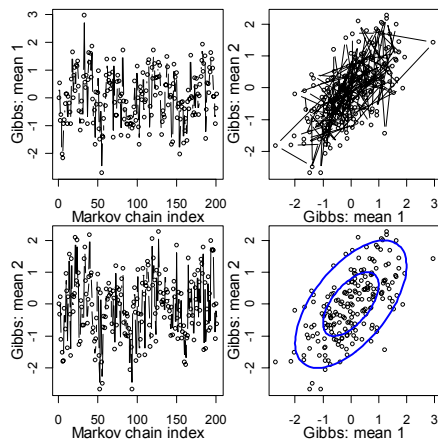$$\theta_1 \mid \theta_2, \mu, \rho \sim N\left(\mu_1 + \rho(\theta_2 - \mu_2), 1 - \rho^2\right)$$

$$\theta_2 \mid \theta_1, \mu, \rho \sim N\left(\mu_2 + \rho(\theta_1 - \mu_1), 1 - \rho^2\right)$$
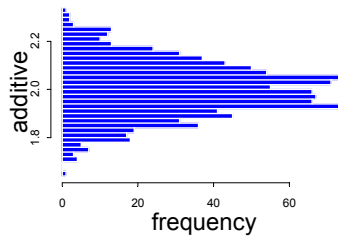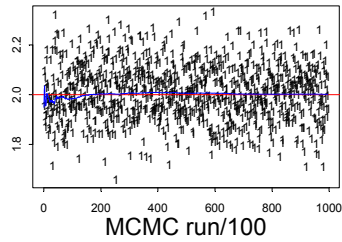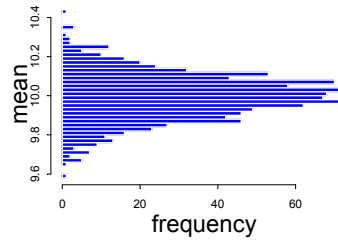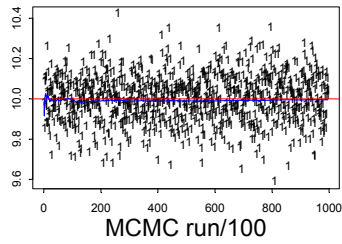
# Gibbs sampler samples: $\rho = 0.6$

# Gibbs Sampler: effects & genotypes

- for given locus $\lambda$, can sample effects $\theta$ and genotypes $Q$
  - effects parameter vector $\theta = (G, \sigma^2)$ with $G = (G_{qq}, G_{Qq}, G_{QQ})$
  - missing genotype vector $Q = (Q_1, Q_2, \ldots, Q_n)$
- Gibbs sampler: update one at a time via full conditionals
  - randomly select order of unknowns
  - update each given current values of all others, locus $\lambda$ and data $(Y, X)$
    - sample variance $\sigma^2$ given $Y, Q$ and genetic values $G$
    - sample genotype $Q_i$ given markers $X_i$ and locus $\lambda$
  - can do block updates if more efficient
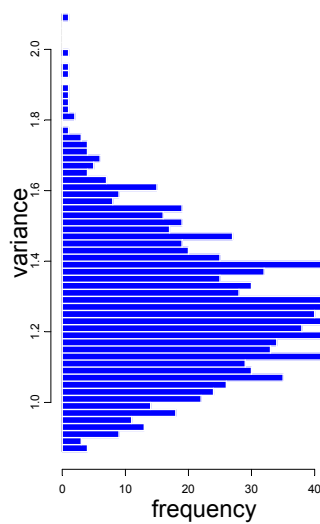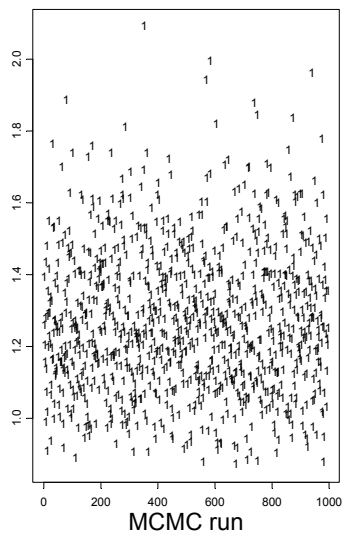    - sample all genetic values $G$ given $Y, Q$ and variance $\sigma^2$

# phenotype model: alternate form

- genetic value $G_Q$ in "cell means" form easy
- but often useful to model effects directly
  - sort out additive and dominance effects
  - useful for reduced models with multiple QTL
    - QTL main effects and interactions (pairwise, 3-way, etc.)
- we only consider additive effects here
  - $G_{qq} = \mu - a$, $G_{Qq} = \mu$, $G_{QQ} = \mu + a$
- recoding for regression model
  - $Q_i = -1$ for genotype qq
  - $Q_i = 0$ for genotype Qq
  - $Q_i = 1$ for genotype QQ
  - $G(Q_i) = \mu + a Q_i$

# MCMC run of mean & additive

# MCMC run for variance

# missing marker data

- sample missing marker data a la QT genotypes
- full conditional for missing markers depends on
  - flanking markers
  - possible flanking QTL
- can explicitly decompose by individual $i$
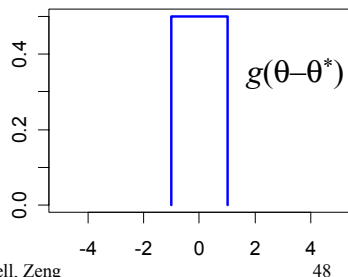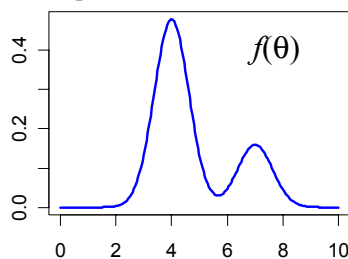  - binomial (or trinomial) probability

$$X_{ik} = \text{aa}, \text{Aa or AA}$$
$$\text{pr}(X_{ik} \mid Y_i, X_i, Q_i, \theta, \lambda) = \text{pr}(X_{ik} \mid X_i, Q_i, \lambda)$$

---

# Metropolis-Hastings idea

- want to study distribution $f(\theta)$
- take Monte Carlo samples
  - unless too complicated
- Metropolis-Hastings samples:
  - current sample value $\theta$
  - propose new value $\theta^*$
    - from some distribution $g(\theta, \theta^*)$
    - Gibbs sampler: $g(\theta, \theta^*) = f(\theta^*)$
  - accept new value with prob $A$
    - Gibbs sampler: $A = 1$

$$A = \min\left(1, \frac{f(\theta^*)g(\theta^*, \theta)}{f(\theta)g(\theta, \theta^*)}\right)$$

# Metropolis-Hastings samples



$N = 200$ samples

narrow $g$      wide $g$

$N = 1000$ samples

narrow $g$      wide $g$

Broman, Churchill, Yandell, Zeng

---

# full conditional for locus

- cannot easily sample from locus full conditional

  $pr(\lambda \,|Y,X,\theta,Q)$      $= pr(\lambda \,|X,Q)$

                          $= pr(\lambda) \, pr(Q|X,\lambda) \,/\text{constant}$

- cannot explicitly determine full conditional
  - difficult to normalize
  - need to average over all possible genotypes over entire map

- Gibbs sampler will not work
  - but can use method based on ratios of probabilities…
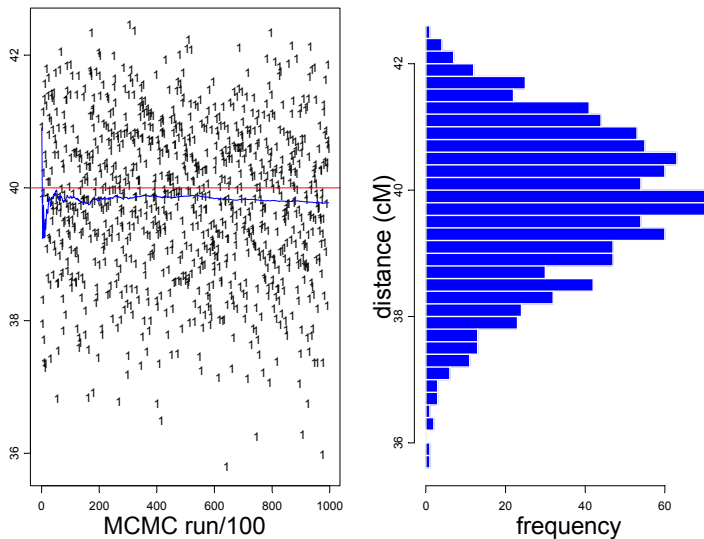
Broman, Churchill, Yandell, Zeng

# Metropolis-Hastings Step

- pick new locus based upon current locus
  - propose new locus from distribution $q(\ )$
    - pick value near current one?
    - pick uniformly across genome?
  - accept new locus with probability $a()$
- Gibbs sampler is special case of M-H
  - always accept new proposal
- acceptance insures right stable distribution
  - accept new proposal with probability $A$
  - otherwise stick with current value

$$A(\lambda_{old}, \lambda_{new}) = \min\left(1, \frac{\pi(\lambda_{new} \mid \mathbf{x}^*)q(\lambda_{new}, \lambda_{old})}{\pi(\lambda_{old} \mid \mathbf{x}^*)q(\lambda_{old}, \lambda_{new})}\right)$$

# MCMC Run for 1 locus at 40cM

# Care & Use of MCMC

- sample chain for long run (100,000-1,000,000)
  - longer for more complicated likelihoods
  - use diagnostic plots to assess "mixing"
- standard error of estimates
  - use histogram of posterior
  - compute variance of posterior--just another summary
- studying the Markov chain
  - Monte Carlo error of series (Geyer 1992)
    - time series estimate based on lagged auto-covariances
  - convergence diagnostics for "proper mixing"

---

# 4.4 bootstrapped variance estimates

- (re)sample $(Y_i, X_i)$ with replacement
  - create bootstrap sample "new" data of size $n$
  - estimate loci $\lambda$ and effects $\theta$
- repeat this $N$ times
- construct summaries of these
  - mean, variance, median, percentile
- construct 95% confidence intervals for $\lambda$ and $\theta$
  - (2.5%ile, 97.5%ile)
  - order estimates, pick number $.025N$ and $.975N$

# 4.5 advantages & shortcomings of IM

- advantages over single marker analysis
  - can infer position and effect of QTL
  - estimated locations & effects almost unbiased
    - if only one segregating QTL per chromosome
  - requires fewer individuals for detection of QTL

# shortcomings of IM

- not an interval test
  - cannot say whether or not QTL is in an interval
  - not independent of effects of QTL outside interval
- can give false positives due to linkage
  - high LOD score due to nearby QTL
  - less of a problem for unlinked QTL
- can detect "ghost QTL"
  - higher peak between two linked QTL
  - estimated position and effect are biased

# 4.6 Haley-Knott Regression Approximation

- likelihood mixes over missing genotypes
  - normal data $\rightarrow$ mixture of normals
- approximate mixture by one normal
  - just estimate mean and variance
- advantages
  - works well for closely spaced markers
  - mean is correct
  - can exploit flanking markers for missing data
  - calculations are easy and fast (PLABQTL)
- disadvantages
  - variance depends on marker genotypes and spacings
  - approximation errors accumulate for multiple QTL

# Haley-Knott regression idea

- replace missing genotypes $Q$ by expected values
  - $P_i = E(Q \mid X_i, \lambda) = \text{sum}_Q \ \text{pr}(Q \mid X_i, \lambda) \ Q$
- fit regression model (e.g. additive gene action)
  - $Y_i = \mu + \alpha P_i + e_i$ , $i = 1, \ldots, n$
- assume constant variance
- correct mean
  - $E(Y_i \mid X_i, \theta, \lambda) = P_i$
- wrong variance
  - $V(Y_i \mid X_i, \theta, \lambda) = \sigma^2 \ \text{sum}_Q \ [\text{pr}(Q \mid X_i, \lambda)]^2$

# Haley-Knott and EM

- both use expected value of genotypes $Q$
  - HK: $P_i = E(Q \mid X_i, \lambda)$ = prior expectation
  - EM: $P_i = E(Q \mid Y_i, X_i, \theta, \lambda)$ = prior expectation
- both solve regression problems for effects
- difference is in iteration
  - HK is first step iteration
  - EM iterates E-step and M-steps to convergence