

NCSU Summer Institute 2004  
QTL II  
Brian S. Yandell  
University of Wisconsin-Madison

- Model: selection for multiple QTL
- Pheno: extensions beyond normal data
- Bayes: interval mapping with prior info
- Traits: multiple phenotypes & microarrays

## contact information & resources

- email: [byandell@wisc.edu](mailto:byandell@wisc.edu)
- web: [www.stat.wisc.edu/~yandell/statgen](http://www.stat.wisc.edu/~yandell/statgen)
  - QTL & microarray resources
  - references, software, people
- thanks:
  - students: Chunfang “Amy” Jin, Fei Zou, Pat Gaffney, Jaya Satagopan
  - faculty/staff: Alan Attie, Hong Lan, Michael Newton, Christina Kendziorski, Tom Osborn, Jason Fine

## Model Selection for Multiple QTL

1. reality of multiple QTL 3-8
2. selecting a class of QTL models 9-15
3. comparing QTL models 16-24
  - QTL model selection criteria
  - issues of detecting epistasis
4. simulations and data studies 25-40
  - simulation with 8 QTL
  - plant BC, animal F2 studies
  - searching through QTL models

Model

NCSU QTL II: Yandell © 2004

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

Model

NCSU QTL II: Yandell © 2004

2

## 1. reality of multiple QTL

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

Model

NCSU QTL II: Yandell © 2004

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

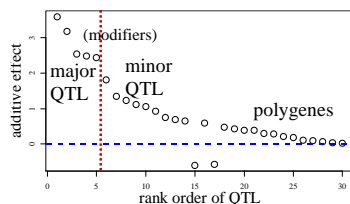
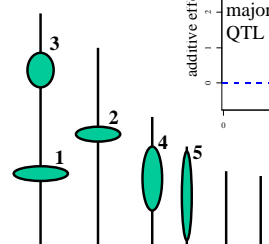
Model

NCSU QTL II: Yandell © 2004

4

## Pareto diagram of QTL effects

major QTL on linkage map



Model

NCSU QTL II: Yandell © 2004

5

## limits of multiple QTL?

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

Model

NCSU QTL II: Yandell © 2004

6

## QTL below detection level?

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

Model

NCSU QTL II: Yandell © 2004

7

## 2. selecting a class of QTL models

- phenotype distribution
  - normal (usual), binomial, Poisson, ...
  - exponential family, semi-parametric, nonparametric
- $\theta$  = gene action
  - additive (A) or general (A+D) effects
  - epistatic interactions (AA, AD, ..., or other types?)
- $\lambda$  = location of QTL
  - known locations?
  - widely spaced (no 2 in marker interval) or arbitrarily close?
- $m$  = number of QTL
  - single QTL?
  - multiple QTL: known or unknown number?

Model

NCSU QTL II: Yandell © 2004

8

## normal phenotype

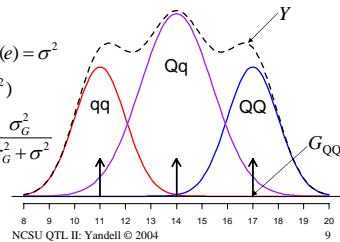
- trait = mean + genetic + environment
- genetic effect uncorrelated with environment
- $\text{pr}(\text{trait } Y \mid \text{genotype } Q, \text{effects } \theta)$

$$Y = G_Q + e$$

$$\text{var}(G_Q) = \sigma_G^2, \text{var}(e) = \sigma^2$$

$$\text{effects } \theta = (G_Q, \sigma^2)$$

$$\text{heritability } h^2 = \frac{\sigma_G^2}{\sigma_G^2 + \sigma^2}$$



Model

NCSU QTL II: Yandell © 2004

9

## two QTL with epistasis

- same phenotype model overview
 
$$Y = G_Q + e, \text{var}(e) = \sigma^2$$
- partition of genotypic value with epistasis
 
$$G_Q = \mu + \beta_1(Q) + \beta_2(Q) + \beta_{12}(Q)$$
- partition of genetic variance
 
$$\text{var}(G_Q) = \sigma_G^2 = \sigma_1^2 + \sigma_2^2 + \sigma_{12}^2$$

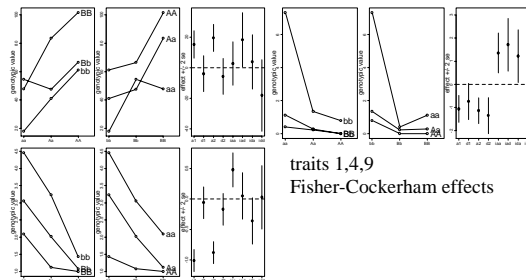
Model

NCSU QTL II: Yandell © 2004

10

## epistasis examples

(Doebley Stec Gustus 1995; Zeng pers. comm.)



Model

NCSU QTL II: Yandell © 2004

11

## multiple QTL with epistasis

- same overview model
 
$$Y = G_Q + e, \text{var}(e) = \sigma^2$$
- sum over multiple QTL in model  $M = \{1, 2, 12, \dots\}$ 

$$G_Q = \mu + \sum_{j \in M} \beta_j(Q)$$
- partition genetic variance in same manner
 
$$\text{var}(G_Q) = \sigma_G^2 = \sum_{j \in M} \sigma_j^2$$
- could restrict attention to 2-QTL interactions

Model

NCSU QTL II: Yandell © 2004

12

## model selection with epistasis

- additive by additive 2-QTL interaction
  - adds only 1 model degree of freedom (df) per pair
  - but could miss important kinds of interaction
- full epistasis adds many model df
  - 2 QTL in BC: 1 df (one interaction)
  - 2 QTL in F2: 4 df (AA, AD, DA, DD)
  - 3 QTL in F2: 20 df (3x4 d.f. 2-QTL, 8 d.f. 3-QTL)
- data-driven interactions (tree-structured)
  - contrasts comparing subsets of genotypes
  - double recessive or double dominant vs other genotypes
  - discriminant analysis based contrasts (Gibert and Le Roy 2003, 2004)
- some issues in model search
  - epistasis between significant QTL
    - check all possible pairs when QTL included?
    - allow higher order epistasis?
  - epistasis with non-significant QTL
    - whole genome paired with each significant QTL?
    - pairs of non-significant QTL?
  - Yi Xu (2000) *Genetics*; Yi, Xu, Allison (2003) *Genetics*; Yi (2004)

Model

NCSU QTL II: Yandell © 2004

13

## 3. comparing QTL models

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

Model

NCSU QTL II: Yandell © 2004

14

## QTL likelihoods and parameters

- LOD or likelihood ratio compares model
  - $L(p) = \log$  likelihood for a particular model with  $p$  parameters
  - $\log(LR) = L(p_2) - L(p_1)$
  - $LOD = \log_{10}(LR) = \log(LR)/\log(10)$
- $p =$  number of model degrees of freedom
  - consider models with  $m$  QTL and all 2-QTL epistasis terms
  - BC:  $p = 1 + m + m(m-1)$
  - F2:  $p = 1 + 2m + 4m(m-1)$
- Bayesian information criterion balances complexity
  - $BIC(\delta) = -2 \log[L(p)] + \delta p \log(n)$
  - $n =$  number of individuals in study
  - $\delta =$  Broman's BIC adjustment

Model

NCSU QTL II: Yandell © 2004

15

## information criteria: likelihoods

- $L(p) =$  likelihood for model with  $p$  parameters
- common information criteria:
  - Akaike AIC =  $-2 \log[L(p)] + 2p$
  - Bayes/Schwartz BIC =  $-2 \log[L(p)] + p \log(n)$
  - BIC-delta  $BIC_\delta = -2 \log[L(p)] + \delta p \log(n)$
  - general form: IC =  $-2 \log[L(p)] + p D(n)$
- comparison of models
  - hypothesis testing: designed for one comparison
    - $2 \log[LR(p_1, p_2)] = L(p_2) - L(p_1)$
  - model selection: penalize complexity
    - $IC(p_1, p_2) = 2 \log[LR(p_1, p_2)] + (p_2 - p_1) D(n)$

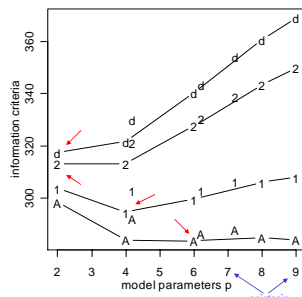
Model

NCSU QTL II: Yandell © 2004

16

## information criteria vs. model size

- WinQTL 2.0
- SCD data on F2
- A=AIC
- 1=BIC(1)
- 2=BIC(2)
- d=BIC( $\delta$ )
- models
  - 1,2,3,4 QTL
    - 2+5+9+2
  - epistasis
    - 2:2 AD



Model

NCSU QTL II: Yandell © 2004

17

## Bayes factors & BIC

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

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

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

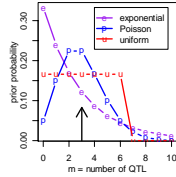
Model

NCSU QTL II: Yandell © 2004

18

## QTL Bayes factors

- $m$  = number of QTL
  - prior  $\text{pr}(m)$  chosen by user
  - posterior  $\text{pr}(m|Y,X)$ 
    - sampled marginal histogram
    - shape affected by prior  $\text{pr}(m)$
- pattern of QTL across genome
  - more complicated prior
  - posterior easily sampled



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

Model

NCSU QTL II: Yandell © 2004

19

## issues in computing Bayes factors

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

Model

NCSU QTL II: Yandell © 2004

20

## multiple QTL priors

- phenotype influenced by genotype & environment
  - $\text{pr}(Y|Q, \theta) \sim N(G_Q, \sigma^2)$ , or  $Y = G_Q + \text{environment}$
- partition genotype-specific mean into QTL effects
  - $G_Q = \text{mean} + \text{main effects} + \text{epistatic interactions}$
  - $G_Q = \mu + \beta(Q) = \mu + \sum_{j \in M} \beta_j(Q)$
- priors on mean and effects
  - $\mu \sim N(\mu_0, \kappa_0 \sigma^2)$  grand mean
  - $\beta(Q) \sim N(0, \kappa \sigma^2)$  model-independent genotypic effect
  - $\beta_j(Q) \sim N(0, \kappa_j \sigma^2 / |M|)$  effects down-weighted by size of  $M$
- determine hyper-parameters via Empirical Bayes

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

Model

NCSU QTL II: Yandell © 2004

21

## multiple QTL posteriors

- phenotype influenced by genotype & environment
  - $\text{pr}(Y|Q, \theta) \sim N(G_Q, \sigma^2)$ , or  $Y = \mu + G_Q + \text{environment}$
- relation of posterior mean to LS estimate

$$G_Q | Y, m \sim N(B_Q \hat{G}_Q, B_Q C_Q \sigma^2) \approx N(\hat{G}_Q, C_Q \sigma^2)$$

$$\text{LS estimate } \hat{G}_Q = \sum_i [\sum_{j \in M} \hat{\beta}_j(Q_i)] = \sum_i w_i Q_i$$

$$\text{variance } V(\hat{G}_Q) = \sum_i w_i^2 \sigma^2 = C_Q \sigma^2$$

$$\text{shrinkage } B_Q = \kappa / (\kappa + C_Q) \rightarrow 1$$

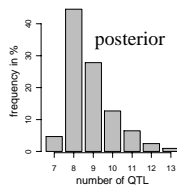
Model

NCSU QTL II: Yandell © 2004

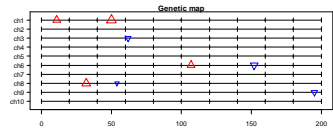
22

## 4. simulations and data studies

- simulated F2 intercross, 8 QTL
  - (Stephens, Fisch 1998)
  - $n=200$ , heritability = 50%
  - detected 3 QTL
- increase to detect all 8
  - $n=500$ , heritability to 97%



QTL	chr	loci	effect
1	1	11	-3
2	1	50	-5
3	3	62	+2
4	6	107	-3
5	6	152	+3
6	8	32	-4
7	8	54	+1
8	9	195	+2



Model

NCSU QTL II: Yandell © 2004

23

## loci pattern across genome

- notice which chromosomes have persistent loci
- best pattern found 42% of the time

### Chromosome

$m$	1	2	3	4	5	6	7	8	9	10	Count of 8000
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	2	0	2	1	0	218
9	2	0	1	0	0	3	0	2	1	0	218
9	2	0	1	0	0	2	0	2	2	0	198

Model

NCSU QTL II: Yandell © 2004

24

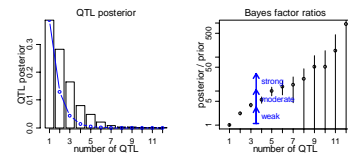
## *B. napus* 8-week vernalization whole genome study

- 108 plants from double haploid
  - similar genetics to backcross: follow 1 gamete
  - parents are Major (biennial) and Stellar (annual)
- 300 markers across genome
  - 19 chromosomes
  - average 6cM between markers
    - median 3.8cM, max 34cM
  - 83% markers genotyped
- phenotype is days to flowering
  - after 8 weeks of vernalization (cooling)
  - Stellar parent requires vernalization to flower
- Ferreira et al. (1994); Kole et al. (2001); Schranz et al. (2002)

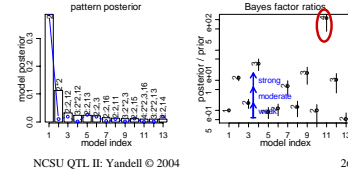
Model NCSU QTL II: Yandell © 2004 25

## Bayesian model assessment

row 1: # QTL  
row 2: pattern  
col 1: posterior  
col 2: Bayes factor  
note error bars on bf



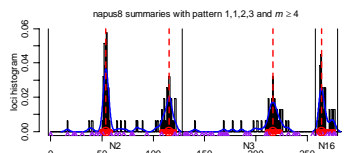
evidence suggests  
4-5 QTL  
N2(2-3),N3,N16



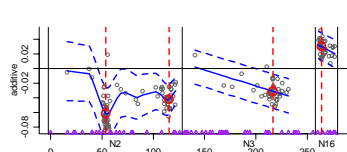
Model NCSU QTL II: Yandell © 2004 26

## Bayesian estimates of loci & effects

histogram of loci  
blue line is density  
red lines at estimates



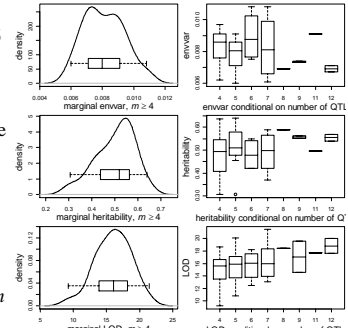
estimate additive effects  
(red circles)  
grey points sampled  
from posterior  
blue line is cubic spline  
dashed line for 2 SD



Model NCSU QTL II: Yandell © 2004 27

## Bayesian model diagnostics

pattern: N2(2),N3,N16  
col 1: density  
col 2: boxplots by *m*



environmental variance  
 $\sigma^2 = .008$ ,  $\sigma = .09$   
heritability  
 $h^2 = 52\%$   
LOD = 16  
(highly significant)

but note change with *m*

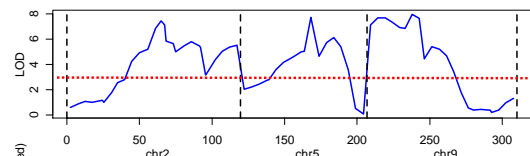
Model NCSU QTL II: Yandell © 2004 28

## studying diabetes in an F2

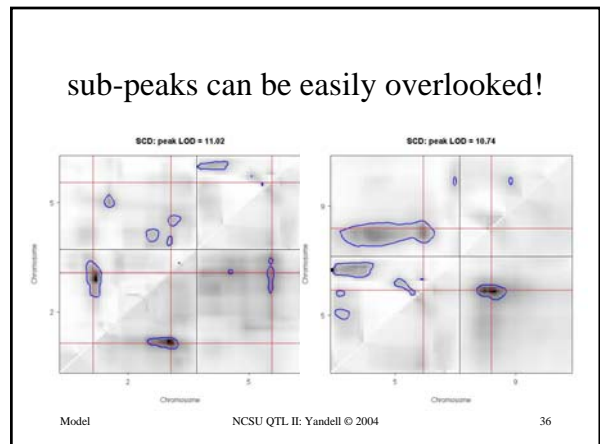
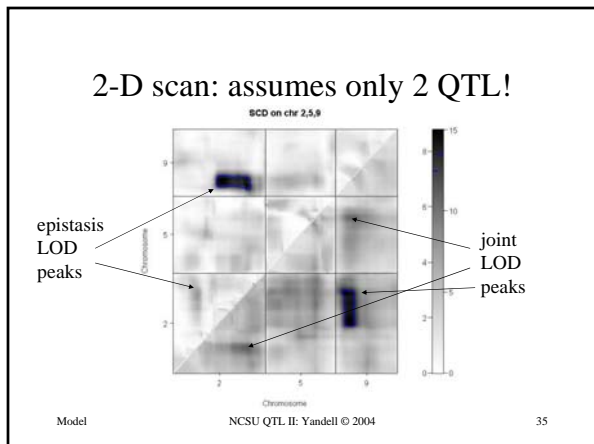
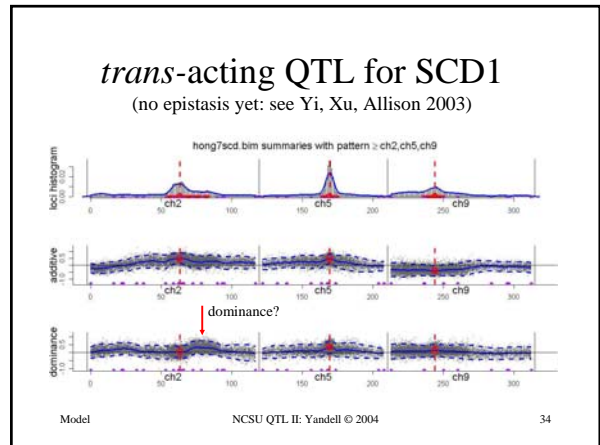
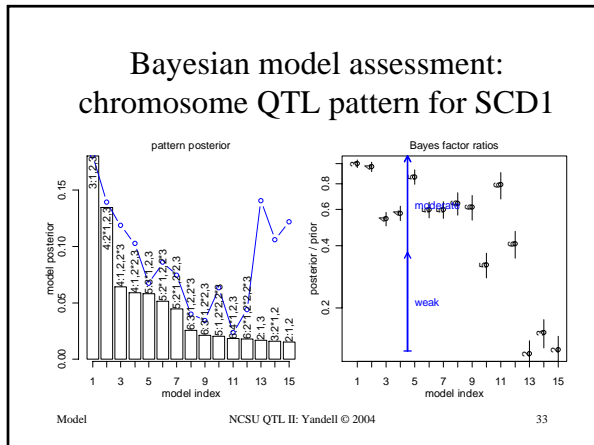
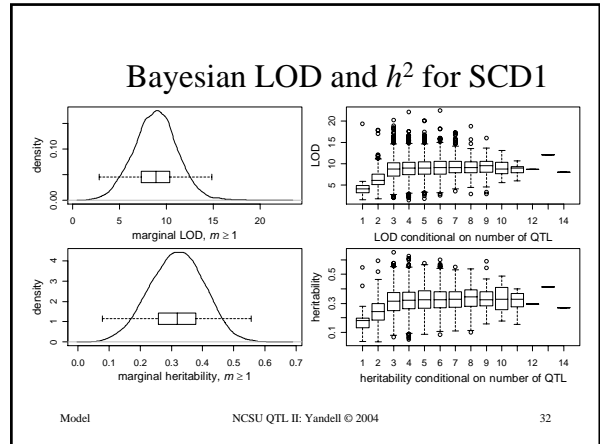
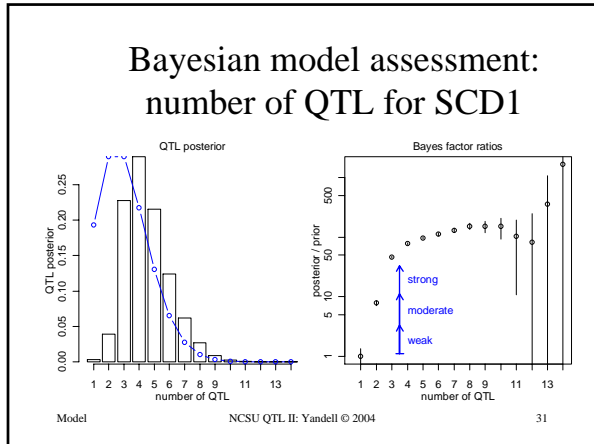
- segregating cross of inbred lines
  - B6.ob x BTBR.ob → F1 → F2
  - selected mice with ob/ob alleles at leptin gene (chr 6)
  - measured and mapped body weight, insulin, glucose at various ages (Stoehr et al. 2000 Diabetes)
  - sacrificed at 14 weeks, tissues preserved
- gene expression data
  - Affymetrix microarrays on parental strains, F1
    - key tissues: adipose, liver, muscle,  $\beta$ -cells
    - novel discoveries of differential expression (Nadler et al. 2000 PNAS; Lan et al. 2002 in review; Ntambi et al. 2002 PNAS)
  - RT-PCR on 108 F2 mice liver tissues
    - 15 genes, selected as important in diabetes pathways
    - SCD1, PEPCK, ACO, FAS, GPAT, PPARgamma, PPARalpha, G6Pase, PDL, ...

Model NCSU QTL II: Yandell © 2004 29

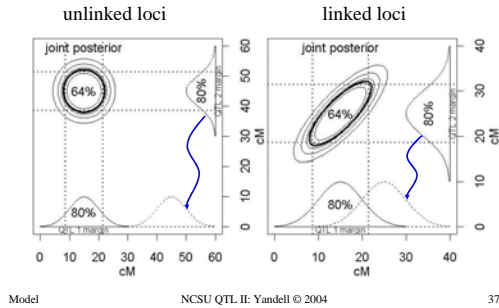
## Multiple Interval Mapping SCD1: multiple QTL plus epistasis!



Model NCSU QTL II: Yandell © 2004 30



## 1-D and 2-D marginals $\text{pr}(\text{QTL at } \lambda \mid Y, X, m)$



## false detection rates and thresholds

- multiple comparisons: test QTL across genome
  - size =  $\text{pr}(\text{LOD}(\lambda) > \text{threshold} \mid \text{no QTL at } \lambda)$
  - threshold guards against a single false detection
    - very conservative on genome-wide basis
  - difficult to extend to multiple QTL
- positive false discovery rate (Storey 2001)
  - $\text{pFDR} = \text{pr}(\text{no QTL at } \lambda \mid \text{LOD}(\lambda) > \text{threshold})$
  - Bayesian posterior HPD region based on threshold
    - $A = \{\lambda \mid \text{LOD}(\lambda) > \text{threshold}\} \approx \{\lambda \mid \text{pr}(\lambda \mid Y, X, m) \text{ large}\}$
  - extends naturally to multiple QTL

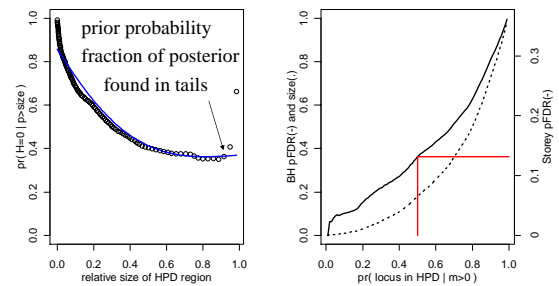
Model NCSU QTL II: Yandell © 2004 38

## pFDR and QTL posterior

- positive false detection rate
  - $\text{pFDR} = \text{pr}(\text{no QTL at } \lambda \mid Y, X, \lambda \text{ in } A)$
  - $\text{pFDR} = \frac{\text{pr}(H=0) * \text{size}}{\text{pr}(m=0) * \text{size} + \text{pr}(m>0) * \text{power}}$
  - power = posterior =  $\text{pr}(\text{QTL in } A \mid Y, X, m > 0)$
  - size = (length of  $A$ ) / (length of genome)
- extends to other model comparisons
  - $m = 1$  vs.  $m = 2$  or more QTL
  - pattern = ch1, ch2, ch3 vs. pattern > 2\*ch1, ch2, ch3

Model NCSU QTL II: Yandell © 2004 39

## pFDR for SCD1 analysis





## Extending the Phenotype Model

1. limitations of parametric models 2-9
  - diagnostic tools for QTL analysis
  - QTL mapping with other parametric "families"
  - quick fixes via data transformations
2. semi-parametric approaches 10-24
3. non-parametric approaches 25-31
- bottom line for normal phenotype model
  - may work well to pick up loci
  - may be poor at estimating effects if data not normal

Pheno

NCSU QTL II: Yandell © 2005

1

## 1. limitations of parametric models

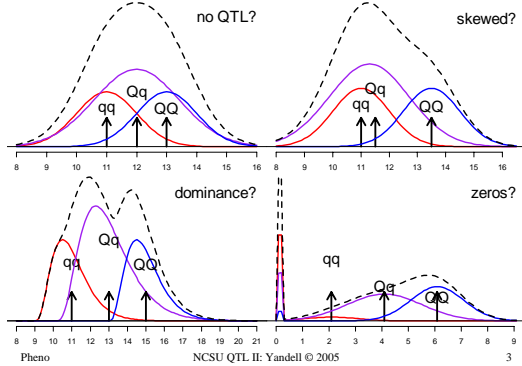
- measurements not normal
  - categorical traits: counts (*e.g.* number of tumors)
    - use methods specific for counts
    - binomial, Poisson, negative binomial
  - traits measured over time and/or space
    - survival time (*e.g.* days to flowering)
    - developmental process; signal transduction between cells
    - TP Speed (pers. comm.); Ma, Casella, Wu (2002)
- false positives due to miss-specified model
  - how to check model assumptions?
- want more robust estimates of effects
  - parametric: only center (mean), spread (SD)
  - shape of distribution may be important

Pheno

NCSU QTL II: Yandell © 2005

2

## what if data are far away from ideal?



Pheno

NCSU QTL II: Yandell © 2005

3

## diagnostic tools for QTL (Hackett 1997)

- illustrated with BC, adapt regression diagnostics
- normality & equal variance (fig. 1)
  - plot fitted values vs. residuals--football shaped?
  - normal scores plot of residuals--straight line?
- number of QTL: likelihood profile (fig. 2)
  - flat shoulders near LOD peak: evidence for 1 vs. 2 QTL
- genetic effects
  - effect estimate near QTL should be  $(1-2r)a$
  - plot effect vs. location

Pheno

NCSU QTL II: Yandell © 2005

4

## marker density & sample size: 2 QTL

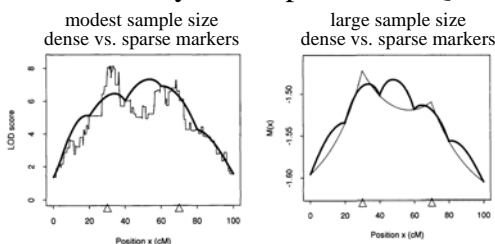


FIGURE 1.—The two-QTL true model with a QTL at 30 cM and a second QTL of somewhat smaller effect at 70 cM (true locations indicated by  $\Delta$ ). A normal single-QTL model is assumed and the LOD score for 100 simulated individuals is given for dense markers (thin curve) and markers at 30-cM intervals (bold curve).

Wright Kong (1997 Genetics)

Pheno

NCSU QTL II: Yandell © 2005

5

## robust locus estimate for non-normal phenotype

large sample size & dense marker map:  
no need for normality

but what happens for modest sample sizes?

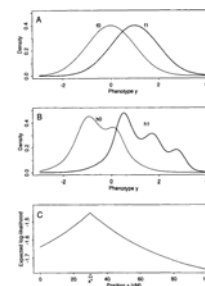


FIGURE 2.—Misspecification of the phenotype model. (A) The assumed distribution  $f$  and  $g$ . (B) The true distributions  $f$ ,  $g$ , and  $h$ . (C) The expected likelihood across the chromosome when the markers are dense. Despite the misspecification, the function is maximized at exactly the true location  $x^* = 30$  cM (indicated by  $\Delta$ ).

Wright Kong (1997 Genetics)

Pheno

NCSU QTL II: Yandell © 2005

6

## What shape is your histogram?

- histogram conditional on known QT genotype
  - $\text{pr}(Y|qq, \theta)$  model shape with genotype qq
  - $\text{pr}(Y|Qq, \theta)$  model shape with genotype Qq
  - $\text{pr}(Y|QQ, \theta)$  model shape with genotype QQ
- is the QTL at a given locus  $\lambda$ ?
  - no QTL  $\text{pr}(Y|qq, \theta) = \text{pr}(Y|Qq, \theta) = \text{pr}(Y|QQ, \theta)$
  - QTL present mixture if genotype unknown
- mixture across possible genotypes
  - sum over  $Q = qq, Qq, QQ$
  - $\text{pr}(Y|X, \lambda, \theta) = \sum_Q \text{pr}(Q|X, \lambda) \text{pr}(Y|Q, \theta)$

Pheno

NCSU QTL II: Yandell © 2005

7

## interval mapping likelihood

- likelihood: basis for scanning the genome
  - product over  $i = 1, \dots, n$  individuals
  - $L(\theta, \lambda|Y) = \text{product}_i \text{pr}(Y_i|X_i, \lambda)$
  - $= \text{product}_i \sum_Q \text{pr}(Q|X_i, \lambda) \text{pr}(Y_i|Q, \theta)$
- problem: unknown phenotype model
  - parametric  $\text{pr}(Y|Q, \theta) = f(Y | \mu, G_Q, \sigma^2)$
  - semi-parametric  $\text{pr}(Y|Q, \theta) = f(Y) \exp(Y\beta_Q)$
  - non-parametric  $\text{pr}(Y|Q, \theta) = F_Q(Y)$

Pheno

NCSU QTL II: Yandell © 2005

8

## useful models & transformations

- binary trait (yes/no, hi/lo, ...)
- map directly as another marker
- categorical: break into binary traits?
- mixed binary/continuous: condition on  $Y > 0$ ?
- known model for biological mechanism
  - counts Poisson
  - fractions binomial
  - clustered negative binomial
- transform to stabilize variance
  - counts  $\sqrt{Y} = \text{sqrt}(Y)$
  - concentration  $\log(Y)$  or  $\log(Y+c)$
  - fractions  $\arcsin(\sqrt{Y})$
- transform to symmetry (approx. normal)
  - fraction  $\log(Y/(1-Y))$  or  $\log((Y+c)/(1+c-Y))$
- empirical transform based on histogram
  - watch out: hard to do well even without mixture
  - probably better to map untransformed, then examine residuals

Pheno

NCSU QTL II: Yandell © 2005

9

## 2. semi-parametric QTL

- phenotype model  $\text{pr}(Y|Q, \theta) = f(Y)\exp(Y\beta_Q)$ 
  - unknown parameters  $\theta = (f, \beta)$ 
    - $f(Y)$  is a (unknown) density if there is no QTL
    - $\beta = (\beta_{qq}, \beta_{Qq}, \beta_{QQ})$
    - $\exp(Y\beta_Q)$  'tilts'  $f$  based on genotype  $Q$  and phenotype  $Y$
- test for QTL at locus  $\lambda$ 
  - $\beta_Q = 0$  for all  $Q$ , or  $\text{pr}(Y|Q, \theta) = f(Y)$
- includes many standard phenotype models
  - normal  $\text{pr}(Y|Q, \theta) = N(G_Q, \sigma^2)$
  - Poisson  $\text{pr}(Y|Q, \theta) = \text{Poisson}(G_Q)$
  - exponential, binomial, ..., but not negative binomial

Pheno

NCSU QTL II: Yandell © 2005

10

## QTL for binomial data

- approximate methods: marker regression
  - Zeng (1993,1994); Visscher et al. (1996); McIntyre et al. (2001)
- interval mapping, CIM
  - Xu Atchley (1996); Yi Xu (2000)
  - $Y \sim \text{binomial}(1, \pi)$ ,  $\pi$  depends on genotype  $Q$
  - $\text{pr}(Y|Q) = (\pi_Q)^Y (1 - \pi_Q)^{(1-Y)}$
  - substitute this phenotype model in EM iteration
- or just map it as another marker!
  - but may have complex

Pheno

NCSU QTL II: Yandell © 2005

11

## EM algorithm for binomial QTL

- E-step: posterior probability of genotype  $Q$ 

$$\text{pr}(Q | Y_i, X_i, \lambda, \pi_Q) = \frac{\text{pr}(Q | X_i, \lambda) (\pi_Q)^{Y_i} (1 - \pi_Q)^{(1-Y_i)}}{\text{sum}_Q \text{ of numerator}}$$
- M-step: MLE of binomial probability  $\pi_Q$ 

$$\pi_Q = \frac{\text{sum}_i Y_i \text{pr}(Q | Y_i, X_i, \lambda, \pi_Q)}{\text{sum}_i \text{pr}(Q | Y_i, X_i, \lambda, \pi_Q)}$$

Pheno

NCSU QTL II: Yandell © 2005

12

## threshold or latent variable idea

- "real", unobserved phenotype  $Z$  is continuous
- observed phenotype  $Y$  is ordinal value
  - no/yes; poor/fair/good/excellent
  - $\text{pr}(Y = j) = \text{pr}(\tau_{j-1} < Z \leq \tau_j)$
  - $\text{pr}(Y \leq j) = \text{pr}(Z \leq \tau_j)$
- use logistic regression idea (Hackett Weller 1995)
  - substitute new phenotype model in to EM algorithm
  - or use Bayesian posterior approach
  - extended to multiple QTL (papers in press)

$$\text{pr}(Y \leq j | Q) = \text{pr}(Z \leq \tau_j | Q) = [1 + \exp(\mu + G_Q - \tau_j)]^{-1}$$

Pheno

NCSU QTL II: Yandell © 2005

13

## quantitative & qualitative traits

- Broman (2003): spike in phenotype
  - large fraction of phenotype has one value
  - map binary trait (is/is not that value)
  - map continuous trait given not that value
- multiple traits
  - Williams et al. (1999)
    - multiple binary & normal traits
    - variance component analysis
  - Corander Sillanpaa (2002)
    - multiple discrete & continuous traits
    - latent (unobserved) variables

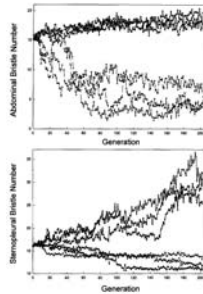
Pheno

NCSU QTL II: Yandell © 2005

14

## other parametric approaches

- Poisson counts
  - Mackay Fry (1996)
    - trait = bristle number
  - Shepel et al (1998)
    - trait = tumor count
- negative binomial
  - Lan *et al.* (2001)
    - number of tumors
- exponential
  - Jansen (1992)



Mackay Fry (1996 *Genetics*)

Pheno

NCSU QTL II: Yandell © 2005

15

## semi-parametric empirical likelihood

- phenotype model  $\text{pr}(Y|Q, \theta) = f(Y) \exp(Y\beta_Q)$ 
  - "point mass" at each measured phenotype  $Y_i$
  - subject to distribution constraints for each  $Q$ :  
 $1 = \sum_i f(Y_i) \exp(Y_i \beta_Q)$
- non-parametric empirical likelihood (Owen 1988)
 
$$L(\theta, \lambda | Y, X) = \text{product}_i [\text{sum}_Q \text{pr}(Q|X_i, \lambda) f(Y_i) \exp(Y_i \beta_Q)]$$

$$= \text{product}_i f(Y_i) [\text{sum}_Q \text{pr}(Q|X_i, \lambda) \exp(Y_i \beta_Q)]$$

$$= \text{product}_i f(Y_i) w_i$$
  - weights  $w_i = w(Y_i|X_i, \beta, \lambda)$  rely only on flanking markers
    - 4 possible values for BC, 9 for F2, etc.
- profile likelihood:  $L(\lambda | Y, X) = \max_{\theta} L(\theta, \lambda | Y, X)$

Pheno

NCSU QTL II: Yandell © 2005

16

## semi-parametric formal tests

- partial empirical LOD
  - Zou, Fine, Yandell (2002 *Biometrika*)
- conditional empirical LOD
  - Zou, Fine (2003 *Biometrika*); Jin, Fine, Yandell (2004)
- has same formal behavior as parametric LOD
  - single locus test: approximately  $\chi^2$  with 1 d.f.
  - genome-wide scan: can use same critical values
  - permutation test: possible with some work
- can estimate cumulative distributions
  - nice properties (converge to Gaussian processes)

Pheno

NCSU QTL II: Yandell © 2005

17

## partial empirical likelihood

$$\log(L(\theta, \lambda | Y, X)) = \sum_i \log(f(Y_i)) + \log(w_i)$$

now profile with respect to  $\beta, \lambda$

$$\log(L(\beta, \lambda | Y, X)) = \sum_i \log(f_i) + \log(w_i) + \sum_Q \alpha_Q (1 - \sum_i f_i \exp(Y_i \beta_Q))$$

partial likelihood: set Lagrange multipliers  $\alpha_Q$  to 0  
force  $f$  to be a distribution that sums to 1

point mass density estimates

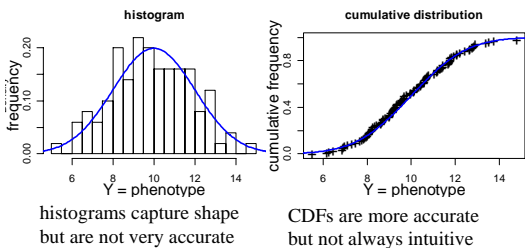
$$f_i = (\sum_i w_i)^{-1} \text{ with } w_i = \sum_Q \exp(Y_i \beta_Q) \text{pr}(Q | X_i, \lambda)$$

Pheno

NCSU QTL II: Yandell © 2005

18

## histograms and CDFs



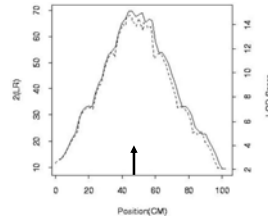
Pheno

NCSU QTL II: Yandell © 2005

19

## rat study of breast cancer Lan *et al.* (2001 *Genetics*)

- rat backcross
  - two inbred strains
    - Wistar-Furth susceptible
    - Wistar-Kyoto resistant
  - backcross to WF
  - 383 females
  - chromosome 5, 58 markers
- search for resistance genes
- $Y = \#$  mammary carcinomas
- where is the QTL?

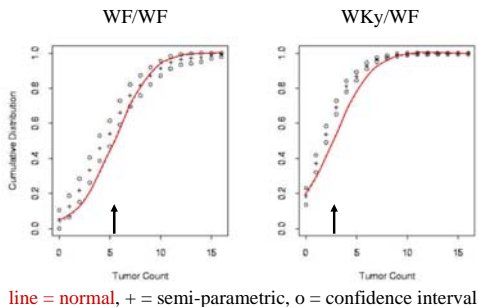


Pheno

NCSU QTL II: Yandell © 2005

20

## what shape histograms by genotype?



Pheno

NCSU QTL II: Yandell © 2005

21

## conditional empirical LOD

- partial empirical LOD has problems
  - tests for F2 depends on unknown weights
  - difficult to generalize to multiple QTL
- conditional empirical likelihood unbiased
  - examine genotypes given phenotypes
  - does not depend on  $f(Y)$
  - $\text{pr}(X_i)$  depends only on mating design
  - unbiased for selective genotyping (Jin *et al.* 2004)

$$\text{pr}(X_i|Y_i, \theta, \lambda, Q) = \exp(Y_i \beta_Q) \text{pr}(Q|X_i|\lambda) \text{pr}(X_i) / \text{constant}$$

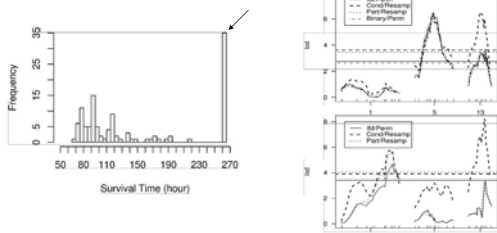
Pheno

NCSU QTL II: Yandell © 2005

22

## spike data example

Boyartchuk *et al.* (2001); Broman (2003)  
133 markers, 20 chromosomes  
116 female mice  
*Listeria monocytogenes* infection



Pheno

NCSU QTL II: Yandell © 2005

23

## new resampling threshold method

- EM locally approximates LOD by quadratic form
- use local covariance of  $\beta$  estimates to further approximate
  - relies on  $n$  independent standard normal variates  $Z = (Z_1, \dots, Z_n)$
  - one set of variates  $Z$  for the entire genome!
- repeatedly resample independent standard normal variates  $Z$ 
  - no need to recompute maximum likelihood on new samples
  - intermediate EM calculations used directly
- evaluate threshold as with usual permutation test
  - extends naturally to multiple QTL
- results shown in previous figure

$$\begin{aligned} \text{LOD}(\lambda) &\approx n \hat{\beta}^T(\lambda) S(\lambda) \hat{\beta}(\lambda) \approx Z^T C^T(\lambda) S(\lambda) C(\lambda) Z \\ \text{cov}(\sqrt{n} \hat{\beta}(\lambda)) &= -C^T(\lambda) C(\lambda) \\ \sqrt{n} \hat{\beta}(\lambda) &\approx C(\lambda) Z, \text{ with } Z \sim N(0, I) \end{aligned}$$

Pheno

NCSU QTL II: Yandell © 2005

24

### 3. non-parametric methods

- phenotype model  $\text{pr}(Y|Q, \theta) = F_Q(Y)$ 
  - $\theta = F = (F_{qq}, F_{Qq}, F_{Qq})$  arbitrary distribution functions
- interval mapping Wilcoxon rank-sum test
  - replaced  $Y$  by  $\text{rank}(Y)$ 
    - (Kruglyak Lander 1995; Poole Drinkwater 1996; Broman 2003)
  - claimed no estimator of QTL effects
- non-parametric shift estimator
  - semi-parametric shift (Hodges-Lehmann)
    - Zou (2001) thesis, Zou, Yandell, Fine (2002 in review)
  - non-parametric cumulative distribution
    - Fine, Zou, Yandell (2001 in review)
- stochastic ordering (Hoff et al. 2002)

Pheno

NCSU QTL II: Yandell © 2005

25

### rank-sum QTL methods

- phenotype model  $\text{pr}(Y|Q, \theta) = F_Q(Y)$
- replace  $Y$  by  $\text{rank}(Y)$  and perform IM
  - extension of Wilcoxon rank-sum test
  - fully non-parametric (Kruglyak Lander 1995; Poole Drinkwater 1996)
- Hodges-Lehmann estimator of shift  $\beta$ 
  - most efficient if  $\text{pr}(Y|Q, \theta) = F(Y+Q\beta)$
  - find  $\beta$  that matches medians
    - problem: genotypes  $Q$  unknown
    - resolution: Haley-Knott (1992) regression scan
  - works well in practice, but theory is elusive
    - Zou, Yandell Fine (*Genetics*, in review)

Pheno

NCSU QTL II: Yandell © 2005

26

### non-parametric QTL CDFs

- estimate non-parametric phenotype model
  - cumulative distributions  $F_Q(y) = \text{pr}(Y \leq y | Q)$
  - can use to check parametric model validity
- basic idea:
 
$$\text{pr}(Y \leq y | X, \lambda) = \sum_Q \text{pr}(Q|X, \lambda) F_Q(y)$$
  - depends on  $X$  only through flanking markers
  - few possible flanking marker genotypes
    - 4 for BC, 9 for F2, etc.

Pheno

NCSU QTL II: Yandell © 2005

27

### finding non-parametric QTL CDFs

- cumulative distribution  $F_Q(y) = \text{pr}(Y \leq y | Q)$
- $F = \{F_Q, \text{all possible QT genotypes } Q\}$ 
  - BC with 1 QTL:  $F = \{F_{QQ}, F_{Qq}\}$
- find  $F$  to minimize over all phenotypes  $y$ 

$$\sum_i [I(Y_i \leq y) - \sum_Q \text{pr}(Q|X, \lambda) F_Q(y)]^2$$
- looks complicated, but simple to implement

Pheno

NCSU QTL II: Yandell © 2005

28

### non-parametric CDF properties

- readily extended to censored data
  - time to flowering for non-vernalized plants
  - Fine, Zou, Yandell (2004 *Biometrics J*)
- nice large sample properties
  - estimates of  $F(y) = \{F_Q(y)\}$  jointly normal
  - point-wise, experiment-wise confidence bands
- more robust to heavy tails and outliers
- can use to assess parametric assumptions

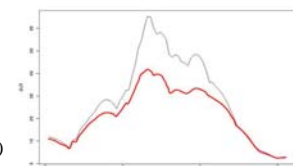
Pheno

NCSU QTL II: Yandell © 2005

29

### what QTL influence flowering time? no vernalization: censored survival

- *Brassica napus*
  - Major female
    - needs vernalization
  - Stellar male
    - insensitive
  - 99 double haploids
- $Y = \log(\text{days to flower})$ 
  - over 50% Major at QTL never flowered
  - log not fully effective



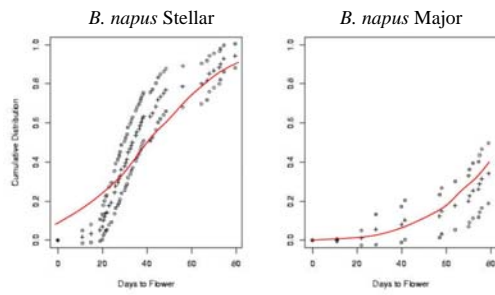
grey = normal, red = non-parametric

Pheno

NCSU QTL II: Yandell © 2005

30

## what shape is flowering distribution?



line = normal, + = non-parametric, o = confidence interval

## Bayesian Interval Mapping

1. Who was Bayes? 2-6
  - What is Bayes theorem?
2. Bayesian inference for QTL 7-14
3. Markov chain sampling 15-29
  - for fixed number of QTL  $m$
4. Sampling across architectures 30-40
  - handling epistasis

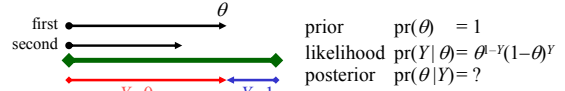
Bayes

NCSU QTL II: Yandell © 2005

1

## 1. who was Bayes? what is Bayes theorem?

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



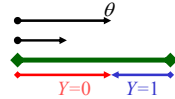
Bayes

NCSU QTL II: Yandell © 2005

2

## what is Bayes theorem?

- where is first ball if the second is to its left (right)?
- prior: probability of parameter before observing data
  - $\text{pr}(\theta) = \text{pr}(\text{parameter})$
  - equal chance of being anywhere on the table
- posterior: probability of parameter after observing data
  - $\text{pr}(\theta|Y) = \text{pr}(\text{parameter} | \text{data})$
  - more likely to left if first ball is toward the right end of table
- likelihood: probability of data given parameters
  - $\text{pr}(Y|\theta) = \text{pr}(\text{data} | \text{parameter})$
  - basis for classical statistical inference
- Bayes theorem
  - posterior = likelihood \* prior / pr(data)
  - normalizing constant  $\text{pr}(Y)$  often drops out of calculation



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

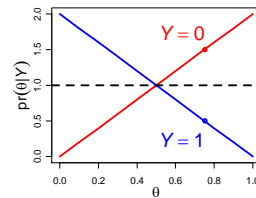
Bayes

NCSU QTL II: Yandell © 2005

3

## where is the second ball given the first?

- first ball  $\theta$
- second ball  $Y$
- prior  $\text{pr}(\theta) = 1$
- likelihood  $\text{pr}(Y|\theta) = \theta^{-Y}(1-\theta)^Y$
- posterior  $\text{pr}(\theta|Y) = ?$



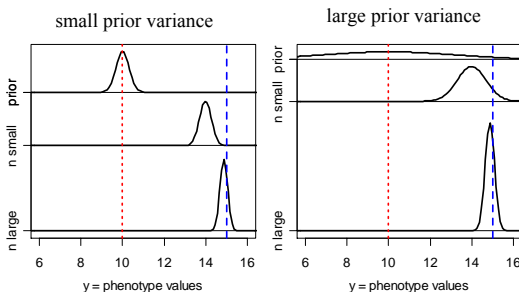
prior:  $\text{pr}(\theta) = 1$   
 likelihood:  $\text{pr}(Y|\theta) = \begin{cases} \theta & \text{if } Y=0 \\ 1-\theta & \text{if } Y=1 \end{cases}$   
 marginal:  $\text{pr}(Y) = \frac{1}{2}$   
 posterior:  $\text{pr}(\theta|Y) = \frac{\text{pr}(Y|\theta)\text{pr}(\theta)}{\text{pr}(Y)}$   
 $= \begin{cases} 2\theta & \text{if } Y=0 \\ 2(1-\theta) & \text{if } Y=1 \end{cases}$

Bayes

NCSU QTL II: Yandell © 2005

4

## Bayes posterior for normal data



Bayes

NCSU QTL II: Yandell © 2005

5

## Bayes posterior for normal data

- model  $Y_i = \mu + E_i$   
 environment  $E \sim N(0, \sigma^2)$ ,  $\sigma^2$  known  
 likelihood  $Y \sim N(\mu, \sigma^2)$   
 prior  $\mu \sim N(\mu_0, \kappa\sigma^2)$ ,  $\kappa$  known  
 posterior: mean tends to sample mean  
 single individual  $\mu \sim N(B_n \bar{Y}_n + (1-B_n)\mu_0, B_n \frac{\sigma^2}{n})$   
 sample of  $n$  individuals  $\mu \sim N\left(B_n \bar{Y}_n + (1-B_n)\mu_0, B_n \frac{\sigma^2}{n}\right)$   
 with  $\bar{Y}_n = \text{sum} \frac{Y_i}{n}$   
 fudge factor (shrinks to 1)  $B_n = \frac{\kappa n}{\kappa n + 1} \rightarrow 1$

Bayes

NCSU QTL II: Yandell © 2005

6

## 2. Bayesian inference for QTL

- develop priors on unknowns
  - unknowns:
    - missing genotypes  $Q$
    - effects  $\theta = (G_Q, \sigma^2)$
    - loci  $\lambda$  (see next section)
  - use empirical Bayes to set useful priors
- study posterior for unknowns given data
  - data:
    - phenotypes  $Y$
    - markers & linkage map  $X$
  - marginal posteriors for effects  $\theta$ , loci  $\lambda$

Bayes

NCSU QTL II: Yandell © 2005

7

## Bayesian priors for QTL

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

Bayes

NCSU QTL II: Yandell © 2005

8

## Bayesian model posterior

- augment data  $(Y, X)$  with unknowns  $Q$
- study unknowns  $(\theta, \lambda, Q)$  given data  $(Y, X)$ 
  - properties of posterior  $\text{pr}(\theta, \lambda, Q | Y, X)$
- sample from posterior in some clever way
  - multiple imputation or MCMC

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

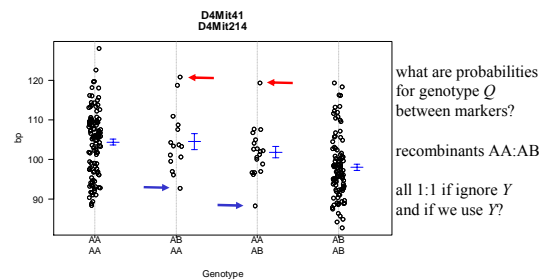
$$\text{pr}(\theta, \lambda | Y, X) = \sum_Q \text{pr}(\theta, \lambda, Q | Y, X)$$

Bayes

NCSU QTL II: Yandell © 2005

9

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



Bayes

NCSU QTL II: Yandell © 2005

10

## posterior on QTL genotypes

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

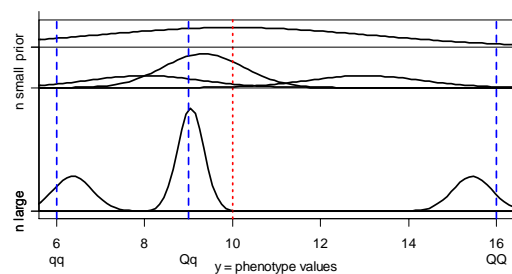
$$\text{pr}(Q | Y_i, X_i, \theta, \lambda) = \frac{\text{pr}(Q | X_i, \lambda) \text{pr}(Y_i | Q, \theta)}{\text{pr}(Y_i | X_i, \theta, \lambda)}$$

Bayes

NCSU QTL II: Yandell © 2005

11

## posterior genotypic means $G_Q$



Bayes

NCSU QTL II: Yandell © 2005

12



## genetic effect posterior given $Q$

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

prior:  $G_Q \sim N(\bar{Y}_*, \kappa\sigma^2)$

posterior:  $G_Q \sim N\left(B_Q \bar{Y}_Q + (1-B_Q)\bar{Y}_*, B_Q \frac{\sigma^2}{n_Q}\right)$

$$n_Q = \text{count}\{Q_i = Q\}, \bar{Y}_Q = \frac{\sum_{i:Q_i=Q} Y_i}{n_Q}$$

fudge factor:  $B_Q = \frac{\kappa n_Q}{\kappa n_Q + 1} \rightarrow 1$

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

- sample variance is proportional to chi-square
  - $ns^2/\sigma^2 \sim \chi^2(n)$
  - likelihood of sample variance  $s^2$  given  $n, \sigma^2$
- conjugate prior is inverse chi-square
  - $v\tau^2/\sigma^2 \sim \chi^2(v)$
  - prior of population variance  $\sigma^2$  given  $v, \tau^2$
- posterior is weighted average of likelihood and prior
  - $(v\tau^2 + ns^2)/\sigma^2 \sim \chi^2(v+n)$
  - posterior of population variance  $\sigma^2$  given  $n, s^2, v, \tau^2$
- empirical choice of hyper-parameters
  - $\tau^2 = s^2/3, v=6$
  - $E(\sigma^2/v, \tau^2) = s^2/2, \text{Var}(\sigma^2/v, \tau^2) = s^4/4$

## 3. Markov chain sampling of architectures

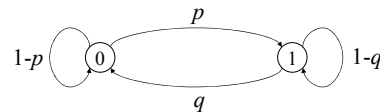
- construct Markov chain around posterior
  - want posterior as stable distribution of Markov chain
  - in practice, the chain tends toward stable distribution
    - initial values may have low posterior probability
    - burn-in period to get chain mixing well
- hard to sample  $(\lambda, Q, \theta, m)$  from joint posterior
  - update  $(\lambda, Q, \theta)$  from full conditionals for  $m$ -QTL model
  - update  $m$  using reversible jump technology

$$(\lambda, Q, \theta, m) \sim \text{pr}(\lambda, Q, \theta, m | Y, X)$$

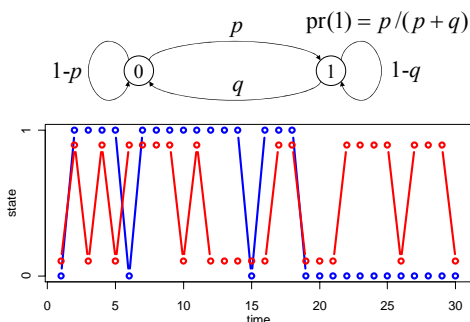
$$(\lambda, Q, \theta, m)_1 \rightarrow (\lambda, Q, \theta, m)_2 \rightarrow \dots \rightarrow (\lambda, Q, \theta, m)_N$$

## 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
- toy problem
  - two states (0,1)
  - move chances depend on current state  $\text{pr}(1) = p/(p+q)$
  - what is the chance of being in state 1?



## Markov chain idea



## Gibbs sampler idea

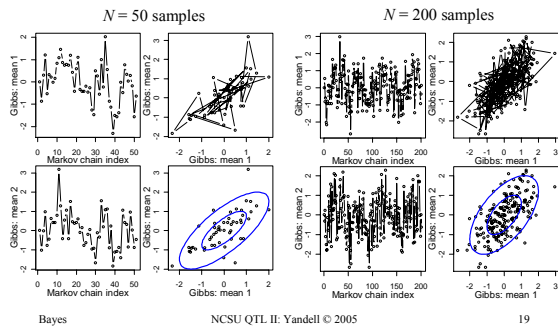
- toy problem
  - want to study two correlated effects
  - could sample directly from their bivariate distribution
- instead use Gibbs sampler:
  - sample each effect from its full conditional given the other
  - pick order of sampling at random
  - repeat many times

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

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

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

## Gibbs sampler samples: $\rho = 0.6$

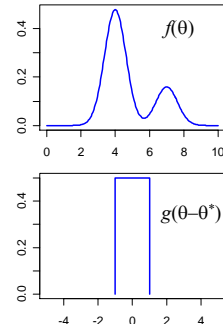


## MCMC sampling of $(\lambda, Q, \theta)$

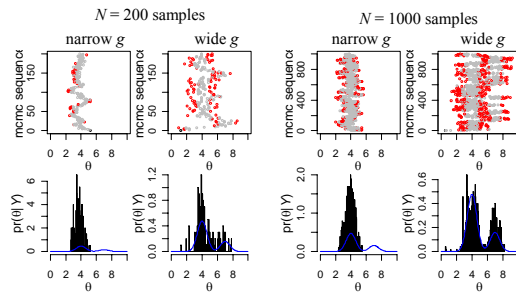
- Gibbs sampler
  - effects  $\theta = (G, \sigma^2)$   $\lambda \sim \frac{\text{pr}(Q | X, \lambda) \text{pr}(\lambda | X)}{\text{pr}(Q | X)}$
  - genotypes  $Q$   $Q \sim \frac{\text{pr}(Q | Y, X, \theta, \lambda)}{\text{pr}(Y | Q)}$
  - not loci  $\lambda$   $\theta \sim \frac{\text{pr}(Y | Q, \theta) \text{pr}(\theta)}{\text{pr}(Y | Q)}$
- extension of Gibbs sampler
  - Metropolis-Hastings sampler
  - does not require normalization
  - loci  $\lambda$ :  $\text{pr}(Q | X)$  difficult to compute

## Metropolis-Hastings idea

- want to study distribution  $f(\theta)$ 
  - take Monte Carlo samples
    - unless too complicated
    - take samples using ratios of  $f$
- 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:  $A = 1$



## Metropolis-Hastings samples



## full conditional for locus

- cannot easily sample from locus full conditional
 
$$\text{pr}(\lambda | Y, X, \theta, Q) = \text{pr}(\lambda | X, Q)$$

$$= \text{pr}(\lambda) \text{pr}(Q | X, \lambda) / \text{constant}$$
- to explicitly determine constant, must average
  - over all possible genotypes
  - over entire map
- Gibbs sampler will not work in general
  - but can use method based on ratios of probabilities
  - Metropolis-Hastings is extension of Gibbs sampler

## Metropolis-Hastings Step

- pick new locus based upon current locus
  - propose new locus from some distribution  $g(\cdot)$ 
    - pick value near current one? (usually)
    - pick uniformly across genome? (sometimes)
  - accept new locus with probability  $A$ 
    - otherwise stick with current value

$$A(\lambda_{old}, \lambda_{new}) = \min \left( 1, \frac{\text{pr}(\lambda_{new}) \text{pr}(Q | X, \lambda_{new}) g(\lambda_{new}, \lambda_{old})}{\text{pr}(\lambda_{old}) \text{pr}(Q | X, \lambda_{old}) g(\lambda_{old}, \lambda_{new})} \right)$$

## Brassica napus data

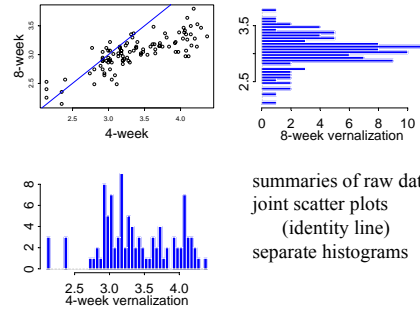
- 4-week & 8-week vernalization effect
  - $\log(\text{days to flower})$
- genetic cross of
  - Stellar (annual canola)
  - Major (biennial rapeseed)
- 105 F1-derived double haploid (DH) lines
  - homozygous at every locus (*QQ* or *qq*)
- 10 molecular markers (RFLPs) on LG9
  - two QTLs inferred on LG9 (now chromosome N2)
  - corroborated by Butruille (1998)
  - exploiting synteny with *Arabidopsis thaliana*

Bayes

NCSU QTL II: Yandell © 2005

25

## Brassica 4- & 8-week data

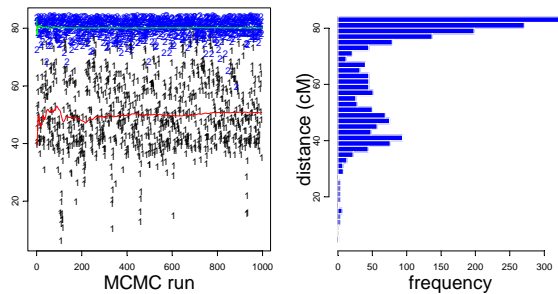


Bayes

NCSU QTL II: Yandell © 2005

26

## Brassica 8-week data locus MCMC with $m=2$



Bayes

NCSU QTL II: Yandell © 2005

27

## 4-week vs 8-week vernalization

- |   |  |
|---|--|
| <p>4-week vernalization</p> <ul style="list-style-type: none"> <li>• longer time to flower</li> <li>• larger LOD at 40cM</li> <li>• modest LOD at 80cM</li> <li>• loci well determined</li> </ul> | <p>8-week vernalization</p> <ul style="list-style-type: none"> <li>• shorter time to flower</li> <li>• larger LOD at 80cM</li> <li>• modest LOD at 40cM</li> <li>• loci poorly determined</li> </ul> |
|---|--|

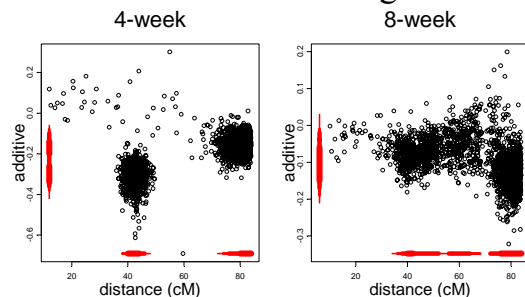
cM	add	cM	add
40	.30	40	.06
80	.16	80	.13

Bayes

NCSU QTL II: Yandell © 2005

28

## Brassica credible regions



Bayes

NCSU QTL II: Yandell © 2005

29

## 4. sampling across architectures

- search across genetic architectures  $M$  of various sizes
  - allow change in  $m = \text{number of QTL}$
  - allow change in types of epistatic interactions
- compare architectures
  - Bayes factors: previous talk
- methods for search
  - reversible jump MCMC
  - Gibbs sampler with loci indicators
- complexity of epistasis
  - Fisher-Cockerham effects model
  - general multi-QTL interaction & limits of inference

Bayes

NCSU QTL II: Yandell © 2005

30

## reversible jump issues

- use reversible jump MCMC to change  $m$ 
  - adjust to change of variables between models
    - bookkeeping helps in comparing models
  - Green (1995); Richardson Green (1997)
- think model selection in multiple regression
  - but regressors (QTL genotypes) are unknown
  - linked loci = collinear regressors = correlated effects
  - consider only additive genetic effects here
    - genotype coding  $Q = -1, 0, 1$  centered on average genotype

$$G(Q) = \mu + \beta(Q) \text{ with } \beta(Q) = \alpha \times (Q - \bar{Q})$$

Bayes

NCSU QTL II: Yandell © 2005

31

## model selection in regression

- consider known genotypes  $Q$  at 2 known loci  $\lambda$ 
  - models with 1 or 2 QTL
- jump between 1-QTL and 2-QTL models
  - adjust parameters when model changes
  - $\alpha$  and  $\alpha_1$  differ due to collinearity of QTL genotypes

$$m = 1 : Y = \mu + \alpha(Q_1 - \bar{Q}_1) + e$$

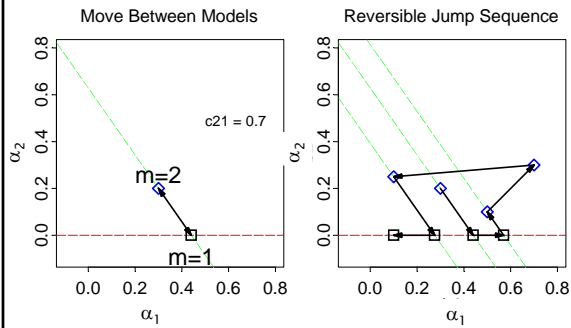
$$m = 2 : Y = \mu + \alpha_1(Q_1 - \bar{Q}_1) + \alpha_2(Q_2 - \bar{Q}_2) + e$$

Bayes

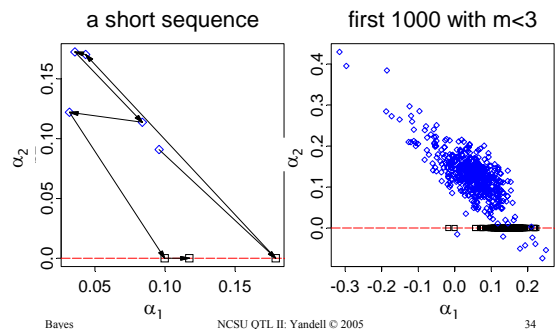
NCSU QTL II: Yandell © 2005

32

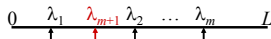
## geometry of reversible jump



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



## reversible jump MCMC



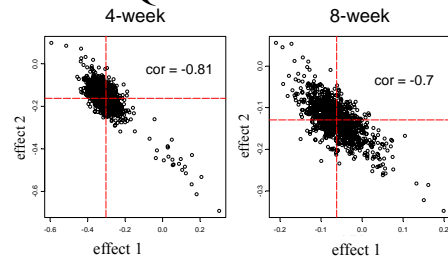
- Metropolis-Hastings updates: draw one of three choices
- update  $m$ -QTL model with probability  $1-b(m+1)-d(m)$ 
    - update current model using full conditionals
    - sample  $m$  QTL loci, effects, and genotypes
  - add a locus with probability  $b(m+1)$ 
    - propose a new locus along genome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

Bayes

NCSU QTL II: Yandell © 2005

35

## collinear QTL = correlated effects



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

Bayes

NCSU QTL II: Yandell © 2005

36

## R/bim: our RJ-MCMC software

- R: [www.r-project.org](http://www.r-project.org)
  - freely available statistical computing application R
  - library(bim) builds on Broman's library(qtl)
- QTLCart: [statgen.ncsu.edu/qtlcart](http://statgen.ncsu.edu/qtlcart)
- [www.stat.wisc.edu/~yandell/qtl/software/Bmapqtl](http://www.stat.wisc.edu/~yandell/qtl/software/Bmapqtl)
- genesis
  - initially designed by JM Satagopan (1996)
  - major revision and extension by PJ Gaffney (2001)
    - whole genome
    - multivariate update of effects; long range position updates
    - substantial improvements in speed, efficiency
    - pre-burnin: initial prior number of QTL very large
  - incorporated into QTLCart (S Wang 2003)
  - built as official R library (H Wu, Yandell, Gaffney, CF Jin 2003)

Bayes

NCSU QTL II: Yandell © 2005

37

## Gibbs sampler with loci indicators

- partition genome into intervals
  - at most one QTL per interval
  - interval = marker interval or large chromosome region
- use loci indicators in each interval
  - $\delta = 1$  if QTL in interval
  - $\delta = 0$  if no QTL
- Gibbs sampler on loci indicators
  - still need to adjust genetic effects for collinearity of  $Q$
  - see work of Nengjun Yi (and earlier work of Ina Hoeschele)

$$Y = \mu + \delta_1 \alpha_1 (Q_1 - \bar{Q}_1) + \delta_2 \alpha_2 (Q_2 - \bar{Q}_2) + e$$

Bayes

NCSU QTL II: Yandell © 2005

38

## epistatic interactions

- model space issues
  - 2-QTL interactions only?
  - Fisher-Cockerham partition vs. tree-structured?
  - general interactions among multiple 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 Xu (2000) *Genetics*; Yi, Xu, Allison (2003) *Genetics*; Yi (2004)

Bayes

NCSU QTL II: Yandell © 2005

39

## limits of epistatic inference

- power to detect effects
  - epistatic model size grows exponentially
    - $|M| = 3^m$  for general interactions
  - power depends on ratio of  $n$  to model size
    - want  $n/|M|$  to be fairly large (say  $> 5$ )
    - $n = 100, m = 3, n/|M| \approx 4$
- empty cells mess up adjusted (Type 3) tests
  - missing  $q_1 Q_2 / q_1 Q_2$  or  $q_1 Q_2 q_3 / q_1 Q_2 q_3$  genotype
  - null hypotheses not what you would expect
  - can confound main effects and interactions
  - can bias AA, AD, DA, DD partition

Bayes

NCSU QTL II: Yandell © 2005

40

## Multiple Traits & Microarrays

1. why study multiple traits together? 2-13
  - diabetes case study
  - central dogma via microarrays
2. design issues for expensive phenotypes 14-21
  - selective phenotyping
3. why are traits correlated? 22-26
  - close linkage or pleiotropy?
4. how to handle high throughput? 27-40
  - dimension reduction: multivariate stats
  - principal components on phenotypes

Traits

NCSU QTL II: Yandell © 2004

1

## 1. why study multiple traits together?

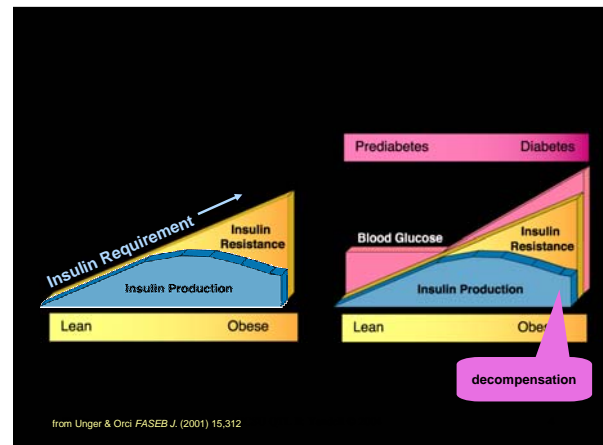
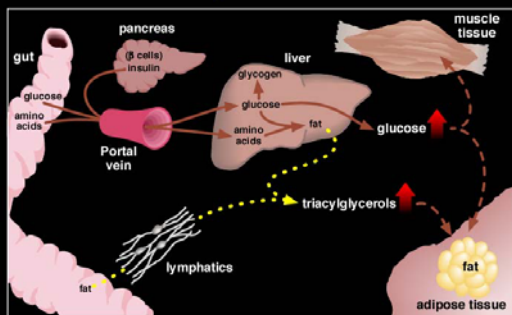
- avoid reductionist approach to biology
  - address physiological/biochemical mechanisms
  - Schmalhausen (1942); Falconer (1952)
- separate close linkage from pleiotropy
  - 1 locus or 2 linked loci?
- identify epistatic interaction or canalization
  - influence of genetic background
- establish QTL x environment interactions
- decompose genetic correlation among traits
- increase power to detect QTL

Traits

NCSU QTL II: Yandell © 2004

2

## Type 2 Diabetes Mellitus



from Unger & Orci *FASEB J.* (2001) 15,312

## Insulin Resistant Mice



Bill Dove



BTBR strain

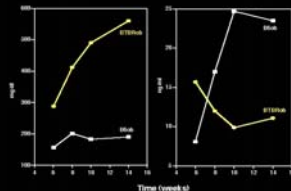


glucose

insulin



(courtesy AD Attie)



## studying diabetes in an F2

- segregating cross of inbred lines
  - B6.ob x BTBR.ob → F1 → F2
  - selected mice with ob/ob alleles at leptin gene (chr 6)
  - measured and mapped body weight, insulin, glucose at various ages (Stoehr et al. 2000 *Diabetes*)
  - sacrificed at 14 weeks, tissues preserved
- gene expression data
  - Affymetrix microarrays on parental strains, F1
    - (Nadler et al. 2000 *PNAS*; Ntambi et al. 2002 *PNAS*)
  - RT-PCR for a few mRNA on 108 F2 mice liver tissues
    - (Lan et al. 2003 *Diabetes*; Lan et al. 2003 *Genetics*)
  - Affymetrix microarrays on 60 F2 mice liver tissues
    - design (Jin et al. 2004 *Genetics* tent. accept)
    - analysis (work in prep.)

Traits

NCSU QTL II: Yandell © 2004

6

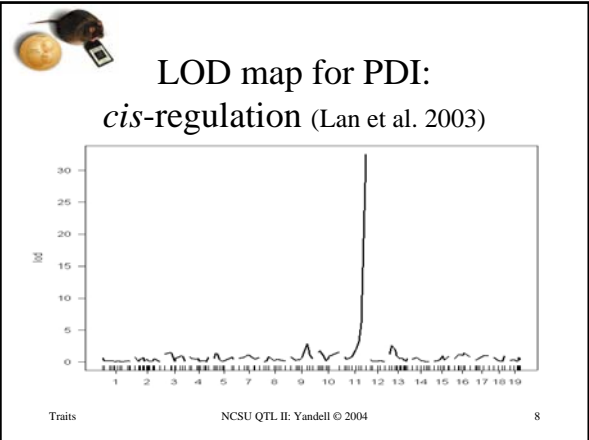
## why map gene expression as a quantitative trait?

- *cis*- or *trans*-action?
  - does gene control its own expression?
  - or is it influenced by one or more other genomic regions?
  - evidence for both modes (Brem et al. 2002 *Science*)
- simultaneously measure all mRNA in a tissue
  - ~5,000 mRNA active per cell on average
  - ~30,000 genes in genome
  - use genetic recombination as natural experiment
- mechanics of gene expression mapping
  - measure gene expression in intercross (F2) population
  - map expression as quantitative trait (QTL)
  - adjust for multiple testing

Traits

NCSU QTL II: Yandell © 2004

7



Traits

NCSU QTL II: Yandell © 2004

8

## mapping microarray data

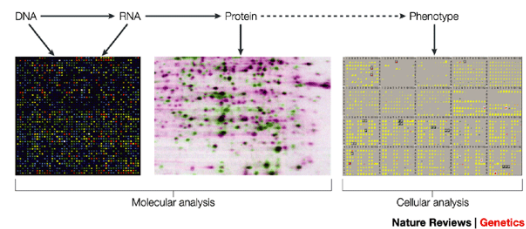
- single gene expression as trait (single QTL)
  - Dumas et al. (2000 *J Hypertens*)
- overview, wish lists
  - Jansen, Nap (2001 *Trends Gen*); Cheung, Spielman (2002); Doerge (2002 *Nat Rev Gen*); Bochner (2003 *Nat Rev Gen*)
- microarray scan via 1 QTL interval mapping
  - Brem et al. (2002 *Science*); Schadt et al. (2003 *Nature*); Yvert et al. (2003 *Nat Gen*)
  - found putative *cis*- and *trans*- acting genes
- multivariate and multiple QTL approach
  - Lan et al. (2003 *Genetics*)

Traits

NCSU QTL II: Yandell © 2004

9

## central dogma via microarrays (Bochner 2003)

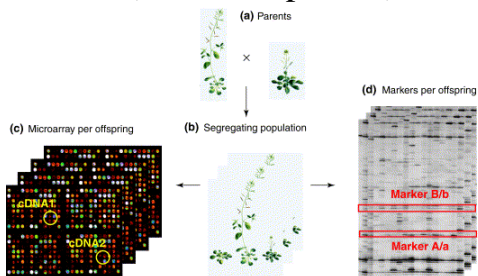


Traits

NCSU QTL II: Yandell © 2004

10

## idea of mapping microarrays (Jansen Nap 2001)

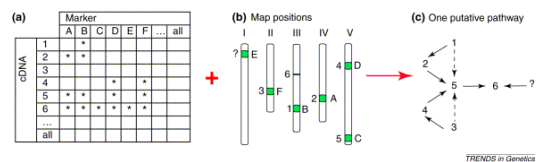


Traits

NCSU QTL II: Yandell © 2004

11

## goal: unravel biochemical pathways (Jansen Nap 2001)



Traits

NCSU QTL II: Yandell © 2004

12



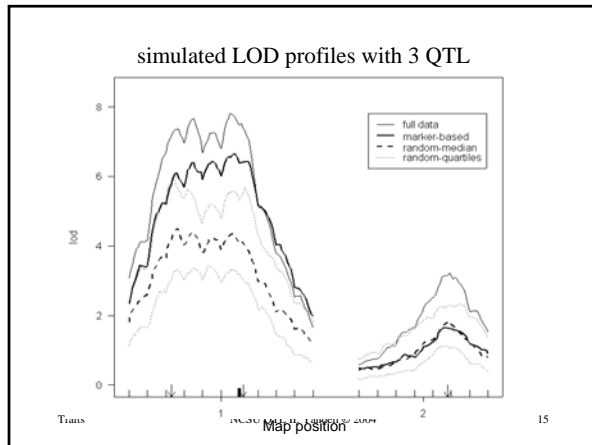
## 2. design issues for expensive phenotypes (thanks to CF “Amy” Jin)

- microarray analysis ~ \$1000 per mouse
  - can only afford to assay 60 of 108 in panel
  - wish to not lose much power to detect QTL
- selective phenotyping
  - genotype all individuals in panel
  - select subset for phenotyping
  - previous studies can provide guide

Traits

NCSU QTL II: Yandell © 2004

14



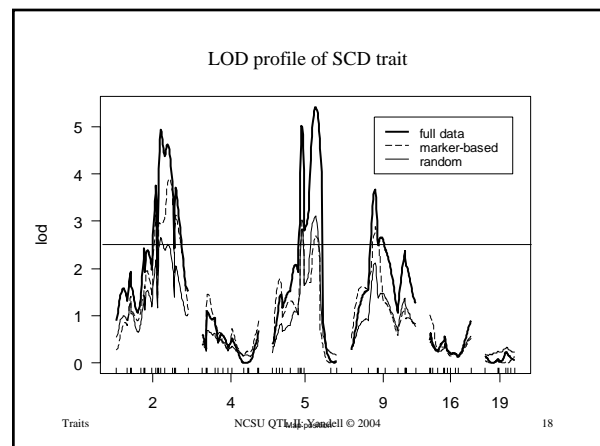
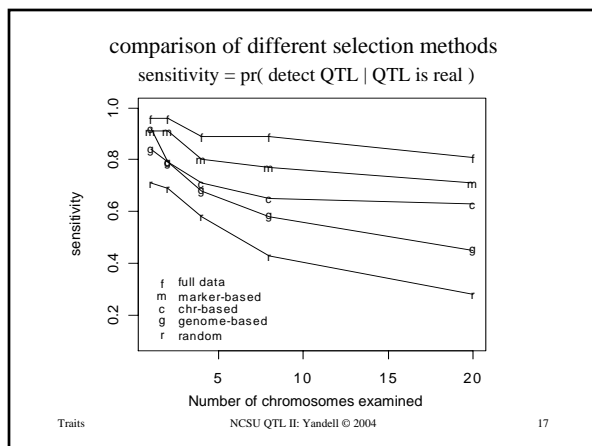
## selective phenotyping

- genotype all individuals in panel
  - whole genome or selected genomic regions?
  - maintain high power in selected regions
  - sensitivity similar to random sample in other regions
- select subset for phenotyping
  - select individuals with large genetic distance
  - use experimental design concepts (Jin et al. 2004)
- previous studies: key regions of chr 2,4,5,9,16,19
  - QTL for important physiological traits

Traits

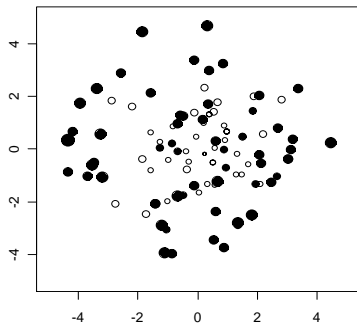
NCSU QTL II: Yandell © 2004

16





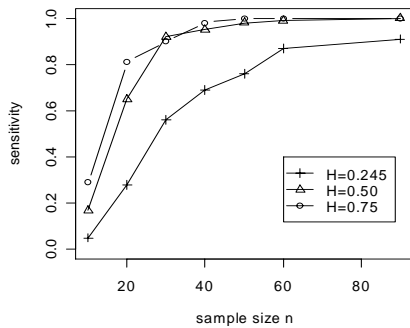
multidimensional scaling of mice selection  
(close points have similar genotypes)



is this relevant to large QTL studies?

- why not phenotype entire mapping panel?
  - selectively phenotype subset of 50-67%
  - may capture most effects
  - with little loss of power
- two-stage selective phenotyping?
  - genotype & phenotype subset of 100-300
    - could selectively phenotype using whole genome
  - QTL map to identify key genomic regions
  - selectively phenotype subset using key regions

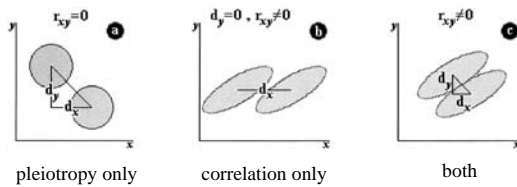
sensitivity = pr( detect QTL | QTL is real )  
depends on heritability and proportion sampled (of  $N=100$ )



3. why are traits correlated?

- environmental correlation
  - non-genetic, controllable by design
  - historical correlation (learned behavior)
  - physiological correlation (same body)
- genetic correlation
  - pleiotropy
    - one gene, many functions
    - common biochemical pathway, splicing variants
  - close linkage
    - two tightly linked genes
    - genotypes  $Q$  are collinear

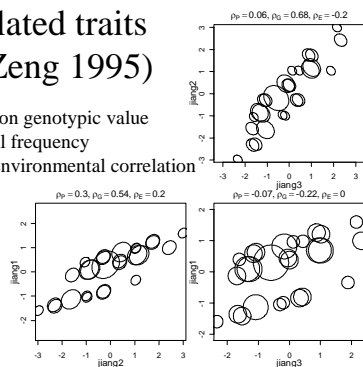
interplay of pleiotropy & correlation

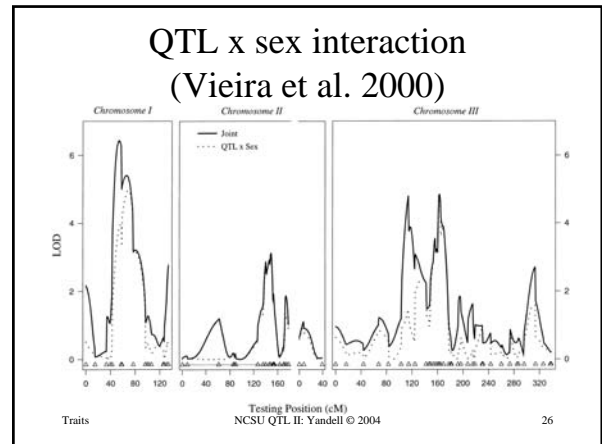
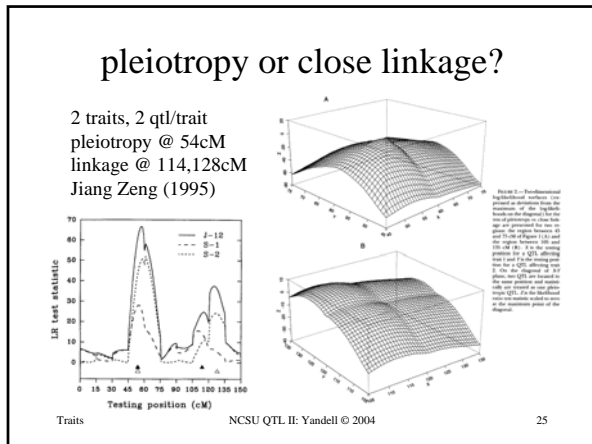


3 correlated traits  
(Jiang Zeng 1995)

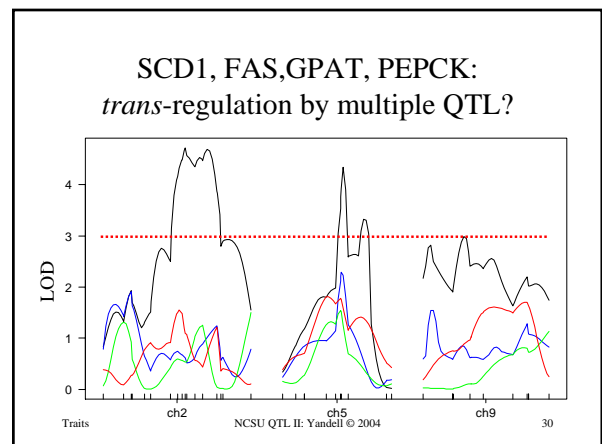
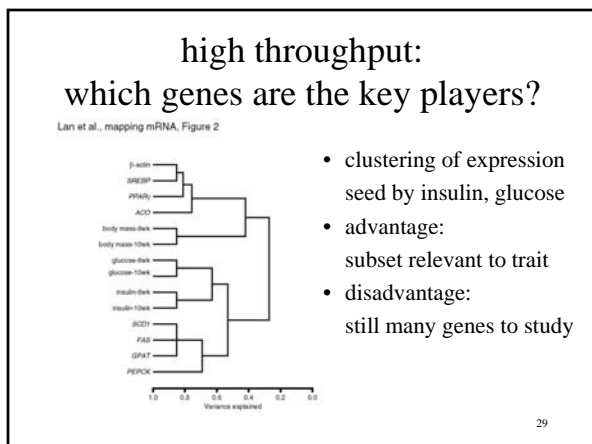
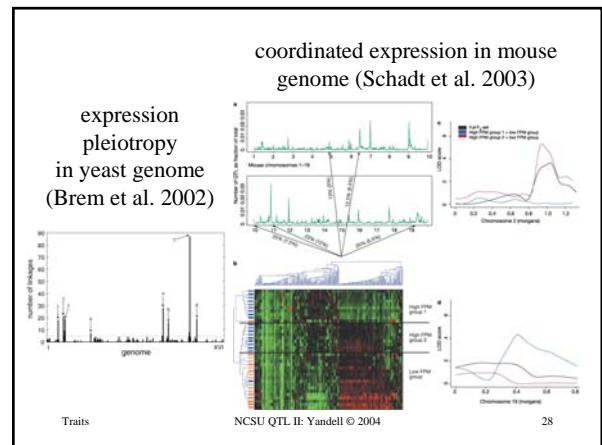
ellipses centered on genotypic value  
width for nominal frequency  
main axis angle environmental correlation  
3 QTL, F2  
27 genotypes

note signs of  
genetic and  
environmental  
correlation





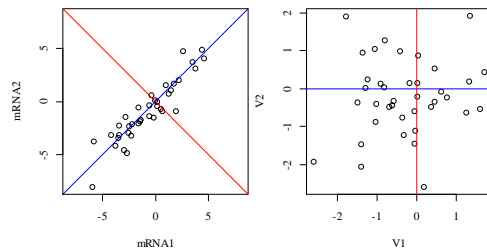
- ### 4. high throughput dilemma
- want to focus on gene expression network
    - ideally capture pathway in a few dimensions
    - allow for complicated genetic architecture
  - may have multiple controlling loci
    - could affect many genes in coordinated fashion
    - could show evidence of epistasis
    - quick assessment via interval mapping may be misleading
  - try mapping principle components as super-traits
    - capture key multivariate features of multiple traits
    - elicit biochemical pathways (Henderson et al. Hoeschele 2001; Ong Page 2002)
- Traits
- NCSU QTL II: Yandell © 2004
- 27



### from gene expression to super-genes

- PC or SVD decomposition of multiple traits
  - $Y = t \text{ traits} \times n \text{ individuals}$
  - decompose as  $Y = UDW^T$ 
    - $U, W$  = ortho-normal transforms (eigen-vectors)
    - $D$  = diagonal matrix with singular values
- transform problem to principal components
  - $W_1$  and  $W_2$  uncorrelated "super-traits"
- interval map each PC separately
  - $W_1 = \mu^*_1 + G^*_{1Q} + e^*_1$
- may only need to map a few PCs

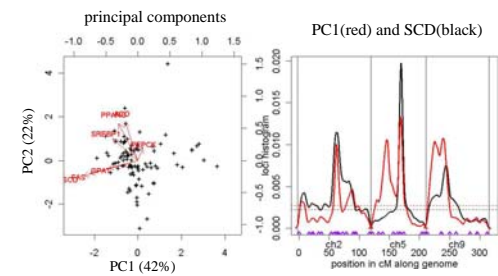
### PC simply rotates & rescales to find major axes of variation



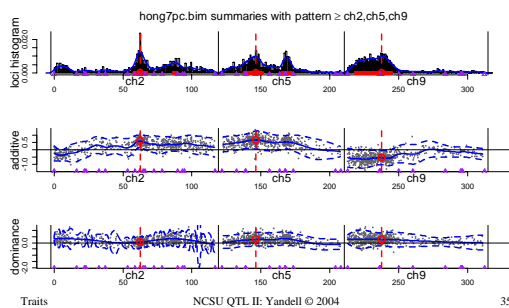
### QTL via Principal Components

- *Drosophila* gonad shape
  - Liu et al. (1996); Zeng et al. (2000)
- other refs of interest
  - Weller et al. (1996); Mangin et al. (1998); Olson et al. (1999); Mahler et al. (2002)
- problems
  - PC may have no relation to genetics!
  - residuals from QTL correlated across PCs
  - PC is descriptive summary, not interpretive

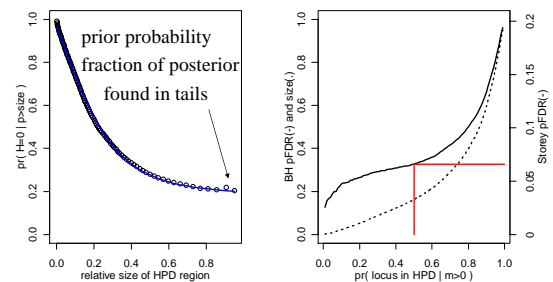
### multivariate screen for gene expressing mapping



### mapping first diabetes PC as a trait



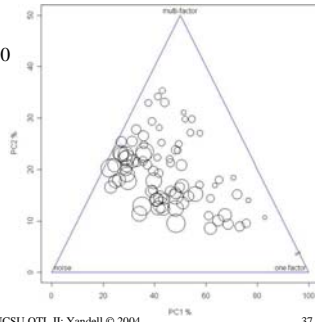
### pFDR for PC1 analysis



## PC across microarray functional groups

1500+ mRNA of 30,000  
85 functional groups  
60 mice  
2-35 mRNA / group  
which are interesting?

examine PC1, PC2  
size = # unique mRNA

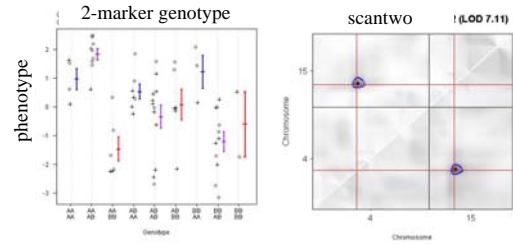


Traits

NCSU QTL II: Yandell © 2004

37

## PC-guided search of mRNA (red lines at main QTL for PC1)



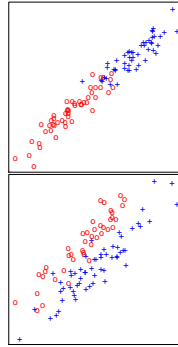
Traits

NCSU QTL II: Yandell © 2004

38

## improvements on PC?

- what is our goal?
  - reduce dimensionality
  - focus on QTL
- PC reduces dimensionality
  - but may not relate to genetics
- canonical discriminant analysis
  - rotate to improve discrimination
  - redo at each putative QTL
  - Gilbert and le Roy (2003,2004)



Traits

NCSU QTL II: Yandell © 2004

39

## how to map multiple traits?

- WinQTL/QTL Cartographer: IM & CIM
  - Jiang Zeng (1995); statgen.ncsu.edu/qtlcart
- MultiQTL: 1-2 QTL with PC on residuals
  - Korol et al. (2001); www.multiqtl.com
- 1-2 QTL with DA across traits
  - Gilbert and le Roy (2003, 2004)
- QTL Express: Haley-Knott regression
  - Knott Haley (2000); qtl.cap.ed.ac.uk
- SOLAR: outbred pedigrees
  - Almasy Blangero (1997); Williams et al. (1999)

Traits

NCSU QTL II: Yandell © 2004

40