





# From sequence data to genomic prediction







#### **Course overview**

- Day 1
  - Introduction
  - Generation, quality control, alignment of sequence data
  - Detection of variants, quality control and filtering
- Day 2
  - Imputation from SNP array genotypes to sequence data
- Day 3
  - Genome wide association studies with SNP array and sequence variant genotypes

#### Day 4 & 5

- Genomic prediction with SNP array and sequence variant genotypes (BLUP and Bayesian methods)
- Use of genomic selection in breeding programs

- Introduction to genomic selection
- Genomic prediction with BLUP
- Genomic prediction with Bayesian methods
- Examples in real data

 Problem marker assisted selection is only a proportion of genetic variance is tracked with markers

- Eg. 10 QTL << 5% of the genetic variance

- Alternative is to trace all segments of the genome with markers
  - Divide genome into chromosome segments based on marker intervals?
  - Capture all QTL = all genetic variance

#### **Genomic selection**

M M M M M M M M M M M M M





Effect of "2" allele

+0.3 L milk

 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M
 M

 Predict genomic breeding values as sum of effects over all SNP



 Predict genomic breeding values as sum of effects over all SNP



Number of SNP

- Genomic prediction exploits linkage disequilibrium
  - Assumption is that markers picking up QTL and will have same effect across the whole population
- Possible within dense marker maps now available

 Genomic prediction avoids bias in estimation of effects due to multiple testing, as all effects fitted simultaneously

#### Genomic selection



- First step is to predict the chromosome segment effects in a reference population
- Number of effects >>> than number of records
- Eg. 50,000 SNPs
- From ~ 2000 records?
- Need methods that can deal with this

BLUP = best linear unbiased prediction
Model:

$$\mathbf{y} = \boldsymbol{\mu} \mathbf{1}_{\mathbf{n}} + \sum_{i=1}^{p} \mathbf{X}_{i} \mathbf{g}_{i} + \mathbf{e}$$

• In BLUP we assume SNP effects come from normal distribution with same variance  $E(\mathbf{g}) \sim N(0, \sigma_g^2)$ 

BLUP assumes normal distribution of SNP effects



- **BLUP** = best linear unbiased prediction
- Then we can estimate segment effects as:

$$\begin{bmatrix} \land \\ \mu \\ \land \\ g \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{n}'\mathbf{1}_{n} & \mathbf{1}_{n}'X \\ \mathbf{X'1}_{n} & \mathbf{X'X} + \mathbf{I}\lambda \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1}_{n}'y \\ \mathbf{X'y} \end{bmatrix}$$

• 
$$\lambda = \sigma_e^2 / \sigma_g^2$$

- Example
- A "simulated" data set
- Single chromosome, with 10 markers
- Phenotypes "simulated"
  - overall mean of 1
  - an effect for SNP 1 of 2 allele of 1
  - normally distributed error term with mean 0 and variance
     1.

#### • Example

		X										
Animal	Y		1	2	3	4	5	6	7	8	9	10
1	0.1	19	0	0	0	0	0	0	1	2	0	2
2	2 1.2	23	1	0	0	1	1	1	2	1	0	1
3	0.8	36	1	0	0	1	0	0	1	1	1	1
4	- 1.2	23	1	1	1	1	0	1	2	1	1	1
5	5 0.4	45	0	1	1	1	1	1	2	1	0	1

- 10 SNPs
- Only 5 phenotypic records.

• Example

		Х									
Animal	Y	1	2	3	4	5	6	7	8	9	10
1	0.19	0	0	0	0	0	0	1	2	0	2
2	1.23	1	0	0	1	1	1	2	1	0	1
3	0.86	1	0	0	1	0	0	1	1	1	1
4	1.23	1	1	1	1	0	1	2	1	1	1
5	0.45	0	1	1	1	1	1	2	1	0	1

- Assume value of 1 for  $\lambda$
- $1_n = [1 \ 1 \ 1 \ 1 \ 1]$



#### • Example

Mean	0.47
SNP1	0.29
SNP2	-0.05
SNP3	-0.05
SNP4	0.08
SNP5	-0.02
SNP6	0.13
SNP7	0.13
SNP8	-0.08
SNP9	0.11
SNP10	-0.08

 Now we want to predict GEBV for a group of young animals without phenotypes.

$$\mathbf{GEBV} = \mathbf{X} \mathbf{g}^{\wedge}$$

• We have the g\_hat, and we can get **X** from their haplotypes (after genotyping).....

Progeny	X									
1	1	1	1	1	1	1	2	1	0	1
2	1	0	0	1	1	1	2	1	0	1
3	1	0	0	1	1	1	2	1	0	1
4	1	0	0	1	1	1	2	1	0	1
5	0	0	0	0	0	0	1	2	0	2

• GEBV

$\hat{\mathbf{GEBV}} = \mathbf{X}\hat{\mathbf{g}}$								
X		∧ <b>g</b>	GEBV					
1111	112101	0.29	0.47					
1001	112101	-0.05	0.58					
1001	112101	-0.05	0.58					
1001	112101	0.08	0.58					
0000	0 1 2 0 2	-0.02	-0.20					
		0.13						
		0.13						
		-0.08						
		0.11						
		-0.08						

- Where do we get  $\sigma_q^2$  from?
- Can estimate total additive genetic variance and divide by number of segments, eg  $\sigma_g^2 = \sigma_a^2 / p$
- If using single markers take account of heterozygosity

$$\sigma_g^2 = \sigma_a^2 / 2 \sum_{i=1}^p q_i (1 - q_i)$$

Ridge regression (Bayesian approach)Cross validation

- An equivalent model
- If there are many QTLs whose effects are normally distributed with constant variance,
- Then genomic selection equivalent to replacing the expected relationship matrix with the realised or genomic relationship matrix (G) estimated from DNA markers in normal BLUP equations.
  - $G_{ij}$  = proportion of genome that is IBD between animals i and j

- An equivalent model
- Rescale X to account for allele frequencies
   -w<sub>ii</sub> = x<sub>ii</sub> 2p<sub>i</sub>
- Then breeding values are  $-\mathbf{v} = \mathbf{W}\mathbf{g}$  (GEBV =  $\mathbf{X}\hat{\mathbf{g}}$ )
- And

**G** = **WW'** 
$$/2\sum_{j=1}^{p} p_{j}(1-p_{j})$$

• Then



• An equivalent model

$$\mathbf{y} = \mathbf{1}_{\mathbf{n}}\,\boldsymbol{\mu} + \mathbf{Z}\mathbf{v} + \mathbf{e}$$

$$\begin{bmatrix} \land \\ \mu \\ \land \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{1_n'1_n} & \mathbf{1_n'Z} \\ \mathbf{Z'1_n} & \mathbf{Z'Z} + \mathbf{G}^{-1} \frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1_n'y} \\ \mathbf{Z'y} \end{bmatrix}$$

Genomic prediction with BLUP
An equivalent model

Model 1.

$$\mathbf{y} = \mathbf{1}_{\mathbf{n}} \boldsymbol{\mu} + \sum_{i=1}^{p} \mathbf{X}_{i} \mathbf{g}_{i} + \mathbf{e} \begin{bmatrix} \hat{\mu} \\ \mu \\ \hat{g} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{\mathbf{n}}' \mathbf{1}_{\mathbf{n}} & \mathbf{1}_{\mathbf{n}}' \mathbf{X} \\ \mathbf{X}' \mathbf{1}_{\mathbf{n}} & \mathbf{X}' \mathbf{X} + \mathbf{I} \frac{\sigma_{e}^{2}}{\sigma_{g}^{2}} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1}_{\mathbf{n}}' \mathbf{y} \\ \mathbf{X}' \mathbf{y} \end{bmatrix} \quad \mathbf{GEBV} = \mathbf{X} \mathbf{g}$$

- Model 2.

Genomic prediction with BLUP
An equivalent model

Model 1.

$$\mathbf{y} = \mathbf{1}_{\mathbf{n}} \boldsymbol{\mu} + \sum_{i=1}^{p} \mathbf{X}_{i} \mathbf{g}_{i} + \mathbf{e} \begin{bmatrix} \hat{\boldsymbol{\mu}} \\ \hat{\boldsymbol{\mu}} \\ \hat{\mathbf{g}} \end{bmatrix} =$$

$$\begin{bmatrix} \mathbf{1_n'1_n} & \mathbf{1_n'X} \\ \mathbf{X'1_n} & \mathbf{X'X} + \mathbf{I}\frac{\sigma_e^2}{\sigma_g^2} \end{bmatrix}$$

$$\mathbf{GEBV} = \mathbf{X} \mathbf{g}^{\wedge}$$

 $\begin{bmatrix} \mathbf{1}_{n} & \mathbf{y} \\ \mathbf{X} & \mathbf{y} \end{bmatrix}$ 

#### - Model 2.

$$\mathbf{y} = \mathbf{1}_{\mathbf{n}}\,\boldsymbol{\mu} + \mathbf{Z}\mathbf{v} + \mathbf{e}$$

$$\begin{bmatrix} \land \\ \mu \\ \land \\ \mathbf{V} \end{bmatrix} = \begin{bmatrix} \mathbf{1_n'1_n} & \mathbf{1_v'Z} \\ \mathbf{Z'1_n} & \mathbf{Z'Z} + \mathbf{G}^{-1} \frac{\sigma_e^2}{\sigma_v^2} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1_n'y} \\ \mathbf{Z'y} \end{bmatrix}$$

Holstein reference	n = 781
Jersey reference	n=287
Holstein validation	n = 400
Jersey validation	n = 77



- An equivalent model
- Why use model 2 (GBLUP).
  - If number of markers >>> large than number of animals, more computationally efficient
  - Can be integrated into national evaluations more readily?
  - Calculate accuracy of GEBV from inverse coefficient matrix



- Introduction to genomic selection
- Genomic prediction with BLUP
- Genomic prediction with Bayesian methods
- Examples in real data

- Alternative assumptions regarding the distribution of SNP effects
- Introduction to Bayesian methods
- Genomic prediction with Bayesian methods
- Comparison of accuracy of methods

#### Genomic selection



#### Alternative prior assumptions for SNP effects

- BLUP assumes normally distributed QTL effects
- Does not match prior knowledge of distributions of QTL effects for some traits


#### Alternative prior assumptions for SNP effects

- Students t distribution?
  - BayesA
- Many zero effects and proportion Students t distribution?
  - BayesB
- Many zero effect and rest normal distribution
  - BayesCpi
- Double exponential effects
  - BayesianLASSO
- Multiple normal distributions
  - BayesMulti, BayesR

• Bayes theorem

 $P(x \mid y) \propto P(y \mid x)P(x)$ 

• Bayes theorem

 $P(x \mid y) \propto P(y \mid x)P(x)$ 

Probability of parameters x given the data y (posterior)

• Bayes theorem

 $P(x \mid y) \propto P(y \mid x)P(x)$ 

Probability ofIs proportional toparameters x given+the data y (posterior)

• Bayes theorem

 $P(x \mid y) \propto P(y \mid x)P(x)$ 

Probability ofIsparameters x giventhe data y (posterior)

Is proportional to Probability of data y given the x (likelihood of data)

• Bayes theorem

 $P(x \mid y) \propto P(y \mid x)P(x)$ 

Probability of<br/>parameters x givenIs proportional to<br/>data y given the<br/>x (likelihood of<br/>data)Prior

- Consider an experiment where we measure height of 10 people to estimate average height
- We want to use prior knowledge from many previous studies that average height is 174cm with standard error 5cm

y=average height + e

• Bayes theorem

$$P(x \mid y) \propto P(y \mid x)P(x)$$

#### Prior probability of x (average height)



• Bayes theorem

$$P(x \mid y) \propto P(y \mid x)P(x)$$

#### From the data.....

$$\overline{x} = 178$$
  
s.e = 5

#### 0.09 0.08 0.07 0.06 Density 0.05 0.04 0.03 0.02 0.01 0 165 170 175 180 185 160 190 Height

Prior probability of x (average height)

• Bayes theorem

$$P(x \mid y) \propto P(y \mid x)P(x)$$

# Likelihood of data (y) given height x, most likely x = 178cm Prior probability of x (average height)



• Bayes theorem



- Bayes theorem
- Less certainty about prior information? Use *less* informative (flat) prior



Height

185

190

- Bayes theorem
- Less certainty about prior information? Use *less* informative (flat) prior



- Bayes theorem
- More certainty about prior information? Use *more* informative prior



190

- Bayes theorem
- More certainty about prior information? Use *more* informative prior



## Genomic prediction

- Alternative assumptions regarding the distribution of SNP effects
- Introduction to Bayesian methods
- Genomic prediction with Bayesian methods
- Comparison of accuracy of methods

### Genomic selection

- For some traits prior knowledge suggests tdistribution of effects
- How to incorporate this into our predictions?



### Genomic selection

- The **t distribution** can be presented as a two level hierarchical model
- Allow different variances between chromosome segments
- Assume a distribution of these variances
- Computationally easier to deal with than original form



 Now lets allow different variances of chromosome segment effects



#### 

#### Distribution of g<sub>j</sub>



- Now lets allow different variances of chromosome segment effects
- Need two levels of models

– Data

$$P(\mathbf{g}, \mu \mid y) \propto P(y \mid \mathbf{g}, \mu) P(\mathbf{g}, \mu)$$

- Variances of chromosome segment effects

$$P(\sigma_{gi}^2 \mid g_i) \propto P(g_i \mid \sigma_{gi}^2) P(\sigma_{gi}^2)$$

- Now lets allow different variances of chromosome segment effects
- Data



 $\begin{bmatrix} \land \\ \mu \\ \land \\ \mathbf{g}_1 \\ \vdots \\ \mathbf{g}_p \end{bmatrix} = \begin{bmatrix} \mathbf{1}_n'\mathbf{1}_n & \mathbf{1}_n'\mathbf{X}_1 & \cdots & \mathbf{1}_n'\mathbf{X}_p \\ \mathbf{X}_1'\mathbf{1}_n & \mathbf{X}_1'\mathbf{X}_1 + \mathbf{I}\frac{\sigma_e^2}{\sigma_{g1}^2} & \cdots & \mathbf{X}_1'\mathbf{X}_p \\ \vdots & \ddots & \ddots & \ddots \\ \mathbf{X}_p'\mathbf{1}_n & \mathbf{X}_p'\mathbf{X}_1 & \cdots & \mathbf{X}_p'\mathbf{X}_p + \mathbf{I}\frac{\sigma_e^2}{\sigma_{gp}^2} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1}_n'y \\ X_1'y \\ \vdots \\ X_p'y \end{bmatrix}$ 

Variances of chromosome segments

$$P(\sigma_{gi}^2 \mid g_i) \propto P(g_i \mid \sigma_{gi}^2) P(\sigma_{gi}^2)$$

- Note that these variance components are not the parameters of interest
- However they are useful intermediates to arrive at better inferences for the g<sub>i</sub>
- Amount of shrinkage of effects varies between segments

Variances of chromosome segments

$$P(\sigma_{gi}^2 \mid g_i) \propto P(g_i \mid \sigma_{gi}^2) P(\sigma_{gi}^2)$$

• Prior?

Inverted chi square convenient for variances

#### • Prior?

- Inverted chi square convenient for variances
- An inverted chi square with v degrees of freedom and scaled by  $S^2$ , eg.

$$S^2/\chi_v^2$$

Describes a distribution with

• mean

$$\frac{vS^2}{(v-2)}$$

variance

$$\frac{2v^2S^4}{(v-2)^2(v-4)}$$

 Larger v, more informative prior = more belief about variance

v=2



v=2

v=20



Variances of chromosome segments

$$P(\sigma_{gi}^2 \mid g_i) \propto P(g_i \mid \sigma_{gi}^2) P(\sigma_{gi}^2)$$

• Prior?

$$S^2/\chi_v^2$$

 We can choose v and S<sup>2</sup> so that the prior reflects our knowledge that there are many QTL of small effect and few of large effect

#### 

#### Distribution of g<sub>j</sub>





Variances of chromosome segments

$$P(\sigma_{gi}^2 | \mathbf{g_i}) \propto P(\mathbf{g_i} | \sigma_{gi}^2) P(\sigma_{gi}^2)$$

#### Posterior?

- An advantage of choosing the inverse chi-square distribution for the prior is that the posterior will also be an inverse chi-square distribution
  - Degrees of freedom = prior + data
  - Scaling factor = sums of squares prior (S<sup>2</sup>) + sums of squares from data

Variances of chromosome segments

$$P(\sigma_{gi}^2 | \mathbf{g_i}) \propto P(\mathbf{g_i} | \sigma_{gi}^2) P(\sigma_{gi}^2)$$

Posterior?
- n<sub>i</sub> = number of haplotype effects

$$\chi^{-2}_{(v+n_i,S^2+\mathbf{g_i'g_i})}$$

Variances of chromosome segments

$$P(\sigma_{gi}^2 | \mathbf{g_i}) \propto P(\mathbf{g_i} | \sigma_{gi}^2) P(\sigma_{gi}^2)$$

• Posterior?

$$\chi^{-2}_{(4.012+n_i,0.002+\mathbf{g_i'g_i})}$$

 But posterior cannot be estimated directly, dependent on g<sub>i</sub>!!

- Solution is to use Gibbs sampling
  - Draw samples from the posterior distributions of parameters conditional on all other effects
  - The average of these samples can be used as the estimates of the parameters

# **Bayesian methods** Gibbs sampling scheme - Parameters to estimate and their posteriors Series1 $\chi^{-2}_{(4.012+n_i,0.002+\mathbf{g_i'g_i})}$ $-P(\sigma_{qi}^2|g_i)$ $\chi^{-2}_{(n-2,\mathbf{e'e})}$ $-P(\sigma_e^2|\mathbf{e})$ $-\mathsf{P}(\mu|\mathbf{y},\mathbf{e},\mathbf{g},\sigma_{e}^{N}) \left(\frac{1}{n}(\mathbf{1}'_{n}\mathbf{y}-\mathbf{1}'_{n}\mathbf{X}\mathbf{g}),\sigma_{e}^{2}/n\right)$ $-\mathsf{P}(\mathsf{g}_{ij}|\mathbf{y},\mu,\mathbf{g}\neq ij,\sigma_{gi}^{2},\sigma_{e}^{2})^{N}\left(\frac{\mathbf{X}_{ij}^{'}\mathbf{y}-\mathbf{X}_{ij}^{'}\mathbf{X}_{g}}{\mathbf{X}_{ij}^{'}+\sigma_{e}^{2}/\sigma_{gi}^{2}},\sigma_{e}^{2}/(\mathbf{X}_{ij}^{'}+\mathbf{X}_{ij}^{2}+\sigma_{e}^{2}/\sigma_{gi}^{2}})\right)$

# **Bayesian methods** Gibbs sampling scheme Parameters to estimate and their posteriors Series1 $-P(\sigma_{qi}^2|g_i)$ $\chi^{-2}_{(4.012+n_i,0.002+\mathbf{g_i'g_i})}$ $\chi^{-2}_{(n-2,\mathbf{e'e})}$ Series $-P(\sigma_{e}^{2}|e)$ $-\mathsf{P}(\mu|\mathbf{y},\mathbf{e},\mathbf{g},\sigma_{e}^{N})\left(\frac{1}{n}(\mathbf{1}'_{n}\mathbf{y}-\mathbf{1}'_{n}\mathbf{X}\mathbf{g}),\sigma_{e}^{2}/n\right)$

 $- \mathsf{P}(\mathsf{g}_{ij}|\mathsf{y},\mu,\mathsf{g}\neq ij,\sigma_{gi}^{2},\sigma_{e}^{2})^{N} \left(\frac{\mathbf{X}_{ij}^{'}\mathbf{y}-\mathbf{X}_{ij}^{'}\mathbf{X}_{g_{ij}-0}^{-}-\mathbf{X}_{ij}^{'}\mathbf{1}_{n}^{\mu}}{\mathbf{X}_{ij}^{'}\mathbf{X}_{ij}^{2}+\sigma_{e}^{2}/\sigma_{gi}^{2}},\sigma_{e}^{2}\right)^{N} \left(\frac{\mathbf{X}_{ij}^{'}\mathbf{y}-\mathbf{X}_{ij}^{'}\mathbf{X}_{g_{ij}-0}^{-}-\mathbf{X}_{ij}^{'}\mathbf{1}_{n}^{\mu}}{\mathbf{X}_{ij}^{'}\mathbf{X}_{ij}^{-}+\sigma_{e}^{2}/\sigma_{gi}^{2}},\sigma_{e}^{2}\right)^{N}$ 


# The Gibbs chain Step 1. Initialise value of g, eg. g=0.01 and μ, eg μ=0.01

- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$ 



$$\chi^{-2}_{(4.012+n_i,0.002+\mathbf{g_i'g_i})}$$

# The Gibbs chain Step 1. Initialise value of g, eg. g=0.01 and μ, eg μ=0.01

- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$ 



$$\chi^{-2}_{(4.012+n_i,0.002+\mathbf{g_i'g_i})}$$

• 
$$\sigma_{g1}^2 = 0.95$$

The Gibbs chain

- Step 1. Initialise value of **g**, eg. **g**=0.01 and  $\mu$ , eg  $\mu$ =0.01
- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$
- Step 3. Draw a sample from  $P(\sigma_e^2|e)$ First calculate the **e** as

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{g} - \mathbf{1}_{n}^{'} \boldsymbol{\mu}$$

The Gibbs chain

- Step 1. Initialise value of **g**, eg. **g**=0.01 and  $\mu$ , eg  $\mu$ =0.01
- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$
- Step 3. Draw a sample from  $P(\sigma_e^2|e)$ First calculate the **e** as

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{g} - \mathbf{1}_{n}' \boldsymbol{\mu}$$

- Then sample...



 $\chi^{-2}_{(n-2,\mathbf{e'e})}$ 

The Gibbs chain

- Step 1. Initialise value of **g**, eg. **g**=0.01 and  $\mu$ , eg  $\mu$ =0.01
- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$
- Step 3. Draw a sample from  $P(\sigma_e^2|e)$ First calculate the **e** as

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{g} - \mathbf{1}_n \,\boldsymbol{\mu}$$

- Then sample...



$$\chi^{-2}_{(n-2,\mathbf{e'e})}$$

The Gibbs chain

- Step 1. Initialise value of **g**, eg. **g**=0.01 and  $\mu$ , eg  $\mu$ =0.01
- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$
- Step 3. Draw a sample from  $P(\sigma_e^2|e)$
- Step 4. Draw a sample from  $P(\mu|y,g,\sigma_e^2)$

The Gibbs chain

- Step 1. Initialise value of **g**, eg. **g**=0.01 and  $\mu$ , eg  $\mu$ =0.01
- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$
- Step 3. Draw a sample from  $P(\sigma_e^2|e)$
- Step 4. Draw a sample from  $P(\mu|y,g,\sigma_e^2)$



 $-\mu = -0.1$ 

The Gibbs chain

- Step 1. Initialise value of **g**, eg. **g**=0.01 and  $\mu$ , eg  $\mu$ =0.01
- Step 2. For each *i*, draw from  $P(\sigma_{gi}^2|g_i)$
- Step 3. Draw a sample from  $P(\sigma_e^2|e)$
- Step 4. Draw a sample from  $P(\mu|y,g,\sigma_e^2)$
- Step 5. For each  $g_{ij}$ , draw from  $P(g_{ii}|y,\mu,g,\sigma_{qi}^2,\sigma_e^2)$

 $-g_{11} = 0.5$ 



• The Gibbs chain

 Repeat steps 2-5 many times to build up samples from posterior distributions of the parameters

The Gibbs chain

- Repeat steps 2-5 many times to build up samples from posterior distributions of the parameters
- Finally, take estimates of parameters as average over many cycles
- Discard first ~ 100 cycles as dependent on starting values

#### • Example

- Consider a data set with three markers. The data set was simulated as:
- the effect of a 2 allele at the first marker is 3, the effect of a 2 allele at the second marker is 0, and the effect of a 2 allele at the third marker was -2.
- the  $\mu$  was 3
- $\sigma_e^2$  was 0.23. The data set was:

#### • Example

Animal	Phenotype	Marker1 allele 1	Marker1 allele 2	Marker2 allele 1	Marker 2 allele 2	Marker3 allele 1	Marker 3 allele 2
1	9.68	2	2	2	1	1	1
2	2 5.69	2	2	2	2	2	2
3	3 2.29	1	2	2	2	2	2
4	3.42	1	1	2	1	1	1
5	5.92	2	1	1	1	1	1
6	5 2.82	2	1	2	1	2	2
7	<b>5.07</b>	2	2	2	1	2	2
8	8 8.92	2	2	2	2	1	1
9	) 2.4	1	1	2	2	1	2
10	9.01	2	2	2	2	1	1
11	4.24	1	2	1	2	2	1
12	6.35	2	2	1	1	1	2
13	8 8.92	2	2	1	2	1	1
14	-0.64	1	1	2	2	2	2
15	5 5.95	2	1	1	1	1	1
16	6.13	1	2	2	1	1	1
17	6.72	2	1	2	1	1	1
18	4.86	1	2	2	1	1	2
19	6.36	2	2	2	2	2	2
20	0.81	1	1	2	1	1	2
21	9.67	2	2	1	2	1	1
22	2. 7.74	2	2	2	1	1	2
23	3 1.45	1	1	2	2	2	1
24	1.22	1	1	2	1	2	1
25	-0.52	1	1	2	2	2	2

#### • Example

- The Bayesian approach was applied, fitting single marker effects
- X matrix
  - Number of copies of two allele for each animal, eg. 2 1 0 for animal 1.

The Gibbs chain

 Step 1. Initialise value of **g**, μ
 g1=0.01, g2=0.01, g3=0.01, μ=0.1

The Gibbs chain

 Step 1. Initialise value of **g**, μ
 g1=0.01, g2=0.01, g3=0.01, μ=0.1
 Step 2. For *i=1,2,3*, draw from P(σ<sub>ai</sub><sup>2</sup>|g<sub>i</sub>)

 $\chi^{-2}_{(4.012+n_i,0.002+\mathbf{g_i'g_i})}$ 

The Gibbs chain

 Step 1. Initialise value of **g**, μ
 g1=0.01, g2=0.01, g3=0.01, μ=0.1
 Step 2. For *i=1,2,3*, draw from P(σ<sub>ai</sub><sup>2</sup>|g<sub>i</sub>)



• σ<sub>g1</sub><sup>2</sup>=0.002, σ<sub>g2</sub><sup>2</sup>=0.06, σ<sub>g3</sub><sup>2</sup>=0.009



$$\chi^{-2}_{(n-2,\mathbf{e'e})}$$

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{g} - \mathbf{1}_n \,\boldsymbol{\mu}$$





• 
$$\sigma_e^2 = 53.38$$



- Step 4. Draw a sample from  $P(\mu|y,g,\sigma_e^2)$ 

$$N\left(\frac{1}{n}\left(\mathbf{1'_n y} - \mathbf{1'_n Xg}\right), \sigma_e^2 / n\right)$$

• µ=3.25

- The Gibbs chain - Step 1. Initialise value of  $\mathbf{g}$ ,  $\mu$ • g1=0.01, g2=0.01, g3=0.01, μ=0.1 -Step 2. For i=1,2,3, draw from P( $\sigma_{qi}^2|g_i$ ) •  $\sigma_{q1}^2 = 0.002$ ,  $\sigma_{q2}^2 = 0.06$ ,  $\sigma_{q3}^2 = 0.09$ - Step 3. Draw a sample from  $P(\sigma_e^2|e)$ •  $\sigma_{e}^{2}$  = 53.38 - Step 4. Draw a sample from  $P(\mu|y,g,\sigma_e^2)$ • μ=3.25
  - Step 5. Draw a sample from  $P(g_{ij}|y,\mu,\mathbf{g}\neq ij,\sigma_{gi}^{2},\sigma_{e}^{2}) = N\left[\frac{X'_{ij}y-X'_{ij}Xg_{(ij=0)}-X'_{ij}\mathbf{1}_{n}\mu}{X'_{ii}X_{ii}+\sigma_{e}^{2}/\sigma_{gi}^{2}},\sigma_{e}^{2}/(X_{ij}'X_{ij}+\sigma_{e}^{2}/\sigma_{gi}^{2})\right]$

- The Gibbs chain - Step 1. Initialise value of  $\mathbf{g}$ ,  $\mu$ • g1=0.01, g2=0.01, g3=0.01, µ=0.1 -Step 2. For i=1,2,3, draw from P( $\sigma_{qi}^2|g_i$ ) •  $\sigma_{a1}^2 = 0.002$ ,  $\sigma_{a2}^2 = 0.06$ ,  $\sigma_{a3}^2 = 0.009$ - Step 3. Draw a sample from  $P(\sigma_e^2|e)$ •  $\sigma_{a}^{2}$  = 53.38 - Step 4. Draw a sample from  $P(\mu|y,g, \sigma_e^2,e)$ • µ=3.25
  - Step 5. Draw a sample from P(g<sub>ij</sub>|y,μ,g≠ij,σ<sub>gi</sub><sup>2</sup>,σ<sub>e</sub><sup>2</sup>)
    g1=-0.02, g2=-0.81,g3=-0.005

- P(g<sub>1</sub>|y,μ,**g**≠1,
$$\sigma_{g1}^2,\sigma_e^2$$
)



- P(g<sub>1</sub>|y,μ,**g**≠1,
$$\sigma_{g1}^2, \sigma_e^2$$
)

"Burn in"



$$\hat{g}_1 = 2.97$$

density.default(x = gStore[200:10000, 1])



$$\hat{g}_1 = 2.97$$
  $\hat{g}_2 = 0.002$   $\hat{g}_1 = -1.81$ 



#### Vector of SNP effects for calculating GEBV





 Alternative priors for variance of segment haplotype/snp effects

- Meuwissen BayesA

 $\chi^{-2}_{(4.012,0.002)}$ 

$$\chi^{-2}_{(4.012+n_i,0.002+\mathbf{g_i'g_i})}$$

- Xu (2003)• Uninformative

$$\chi^{-2}_{(0,0)}$$

$$\chi^{-2}_{(1,\mathbf{g'g})}$$

-Ter Braak (2006)  $p(\sigma_{gi}^2) \propto (\sigma_{gi}^2)^{-1+\alpha}$ 

$$g_i'g_i/\chi_{1-2a}^{-2}$$

- Meuwissen BayesB

 $\sigma_{gi}^{2} = 0 \text{ with probability } \pi,$  $\sigma_{gi}^{2} \sim \chi^{-2} \left(\nu, S\right) \text{ with probability } \left(1 - \pi\right),$ 

#### Meuwissen BayesB

- BayesA prior information is many QTL with small effects and few with moderate effects
- But we have more prior knowledge than this – some chromosome segments will have no effect at all (*contain no QTL*)
  - σ<sub>gi</sub><sup>2</sup>=0,g<sub>i</sub>=0
- How to sample from the posterior?

$$\sigma_{gi}^2 = 0$$
 with probability  $\pi$ ,  
 $\sigma_{gi}^2 \sim \chi^{-2} (\nu, S)$  with probability  $(1 - \pi)$ ,



#### • Meuwissen BayesB – If we sample $\sigma_{gi}^2$ from $\chi^{-2}_{(4.012+n_i,0.002+g_i'g_i)}$ – We will never sample 0, as the distribution has no mass at zero.



- Meuwissen BayesB – If we sample  $\sigma_{gi}^2$  from  $\chi^{-2}_{(4.012+n_i,0.002+g_i'g_i)}$ 
  - We will never sample 0 if  $g_i'g_i > 0$ , as the distribution has no mass at zero.
  - But if  $\sigma_{gi}^2 > 0$ , then sampling  $g_i = 0$  has infinitesimal (basically zero) probability

Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} | y^{*}) = p(\sigma_{gi}^{2} | y^{*}) \times p(g_{i} | \sigma_{gi}^{2}, y^{*})$$

We want to sample from this

Can do it by sampling from these two distributions

Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} \mid y^{*}) = p(\sigma_{gi}^{2} \mid y^{*}) \times p(g_{i} \mid \sigma_{gi}^{2}, y^{*})$$

We want to sample from this

 $P(g_i|y,\mu,g,{\sigma_{gi}}^2,{\sigma_e}^2)$ 



Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} | y^{*}) = p(\sigma_{gi}^{2} | y^{*}) \times p(g_{i} | \sigma_{gi}^{2}, y^{*})$$

??

Sample  $\sigma_{gi}^2$  without conditioning on  $g_i$ 

Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} \mid y^{*}) = p(\sigma_{gi}^{2} \mid y^{*}) \times p(g_{i} \mid \sigma_{gi}^{2}, y^{*})$$

Cannot be expressed as a known distribution = cannot use Gibbs for this bit
Use a Metropolis Hastings algorithm

Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} \mid y^{*}) = p(\sigma_{gi}^{2} \mid y^{*}) \times p(g_{i} \mid \sigma_{gi}^{2}, y^{*})$$

- Step1 Sample  $\sigma_{g_{new}}^2$ , from prior( $\sigma_{g_{new}}^2$ )

$$\sigma_{gi}^2 = 0$$
 with probability  $\pi$ ,  
 $\sigma_{gi}^2 \sim \chi^{-2} (\nu, S)$  with probability (1 –



Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} \mid y^{*}) = p(\sigma_{gi}^{2} \mid y^{*}) \times p(g_{i} \mid \sigma_{gi}^{2}, y^{*})$$

- Step1 Sample  $\sigma_{g_{new}}^2$ , from prior( $\sigma_{g_{new}}^2$ )

 $-\sigma_{g_{new}}^2=0$ 


Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} \mid y^{*}) = p(\sigma_{gi}^{2} \mid y^{*}) \times p(g_{i} \mid \sigma_{gi}^{2}, y^{*})$$

- Step1 Sample  $\sigma_{g_{new}}^2$ , from prior( $\sigma_{g_{new}}^2$ )

 $-\sigma_{g_{new}}^2=0.5$ 



Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} \mid y^{*}) = p(\sigma_{gi}^{2} \mid y^{*}) \times p(g_{i} \mid \sigma_{gi}^{2}, y^{*})$$

- Step 1 Sample  $\sigma_{g_new}^2$ , from prior( $\sigma_{g_new}^2$ ) - Step 2 Evaluate p(y\*|  $\sigma_{g_new}^2$ ) (Likelihood)

$$L(\mathbf{y}^* | \sigma_{ginew}^2 = \frac{1}{2\pi^{1/2n} |\mathbf{V}|^{1/2}} e(-0.5^* (\mathbf{y}^* \mathbf{V}^{-1} \mathbf{y}^*)) \quad \mathbf{V} = \mathbf{X}(\mathbf{I}\sigma_{ignew}^2) \mathbf{X}' + \mathbf{I}\sigma_{\mathbf{e}}^2)$$

Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} \mid y^{*}) = p(\sigma_{gi}^{2} \mid y^{*}) \times p(g_{i} \mid \sigma_{gi}^{2}, y^{*})$$

 $\begin{array}{l} - \mbox{Step 1 Sample } \sigma_{g_new}^2, \mbox{from } prior(\sigma_{g_new}^2) \\ - \mbox{Step 2 Evaluate } p(y^* | \ \sigma_{g_new}^2) \ (\mbox{Likelihood}) \\ - \mbox{Step 3 Replace } \sigma_{gi}^2 \ \mbox{with } \sigma_{g_new}^2 \ \mbox{probability } \\ min[p(y^* | \ \sigma_{g_new}^2) / \ p(y^* | \ \sigma_{gi}^2):1] \end{array}$ 

Meuwissen BayesB

– Solution: sample  $\sigma_{gi}^2$ ,  $g_i$  simultaneously from the distribution:

$$p(\sigma_{gi}^{2}, g_{i} | y^{*}) = p(\sigma_{gi}^{2} | y^{*}) \times p(g_{i} | \sigma_{gi}^{2}, y^{*})$$

- Step 1 Sample  $\sigma_{g_new}^2$ , from prior( $\sigma_{g_new}^2$ ) - Step 2 Evaluate p(y\*|  $\sigma_{g_new}^2$ ) (Likelihood) - Step 3 Replace  $\sigma_{gi}^2$  with  $\sigma_{g_new}^2$  probability min[p(y\*|  $\sigma_{q_new}^2$ )/ p(y\*|  $\sigma_{qi}^2$ ):1]
- Step 4 Repeat ~ 100 cycles

- Comparison of accuracy of methods (Meuwissen et al. 2001)
  - Genome of 1000 cM simulated, marker spacing of 1 cM.
  - Markers surrounding each 1-cM region combined into haplotypes.
  - Due to finite population size (Ne = 100), marker haplotypes were in linkage disequilibrium with QTL between markers.
  - Effects of haplotypes predicted in one generation of 2000 animals
  - Breeding values for progeny of these animals predicted based on marker genotypes

• Comparison of accuracy of methods (Meuwissen et al. 2001)

 $r_{\text{TBV};\text{EBV}} + \text{SE} \ b_{\text{TBV},\text{EBV}} + \text{SE}$ 

LS  $0.318 \pm 0.018 \ 0.285 \pm 0.024$ BLUP  $0.732 \pm 0.030 \ 0.896 \pm 0.045$ BayesA  $0.798 \ 0.827$ BayesB  $0.848 \pm 0.012 \ 0.946 \pm 0.018$ 

- Comparison of accuracy of methods (Meuwissen et al. 2001)
  - The least squares method does very poorly, primarily because the haplotype effects are over-estimated.

- Comparison of accuracy of methods (Meuwissen et al. 2001)
  - The least squares method does very poorly, primarily because the haplotype effects are over-estimated.
  - Increased accuracy of the Bayesian approach because method sets many of the effects of the chromosome segments close to zero in BayesA, or zero in BayesB

- Comparison of accuracy of methods (Meuwissen et al. 2001)
  - The least squares method does very poorly, primarily because the haplotype effects are over-estimated.
  - Increased accuracy of the Bayesian approach because method sets many of the effects of the chromosome segments close to zero in BayesA, or zero in BayesB
  - Also "shrinks" estimates of effects of other chromosome segments based on a prior distribution of QTL effects.

- Comparison of accuracy of methods (Meuwissen et al. 2001)
  - The least squares method does very poorly, primarily because the haplotype effects are over-estimated.
  - Increased accuracy of the Bayesian approach because method sets many of the effects of the chromosome segments close to zero in BayesA, or zero in BayesB
  - Also "shrinks" estimates of effects of other chromosome segments based on a prior distribution of QTL effects.
  - Accuracies were very high, as high as following progeny testing for example

- Introduction to genomic selection
- Genomic prediction with BLUP
- Genomic prediction with Bayesian methods
- Examples in real data

- 1500 Australian dairy bulls
- genotyped for 56000 genome wide SNPs
- Phenotypes average of daughters milk production



- Split data into two sub-populations
   Reference: Bulls born < 2003</li>
  - Validation: Bulls born >= 2003

- Split data into two sub-populations
  - Reference: Bulls born < 2003</p>
  - Validation: Bulls born >= 2003
- Accuracy
  - Correlation of genomic breeding values with EBVs (which include daughter information) in validation set

Table 3 MEBV- Correlation between predicted MEBV and ABV in the validation data set (Bulls proven in years 2005, 2006, 2007)

Method	Protein kg	Fat kg	Protein %	Fat %
Bayes B	0.55	0.51	0.68	0.73
Bayes A	0.53	0.48	0.66	0.70
BLUP	0.60	0.48	0.66	0.64





• Bayesian C∏ (Habier et al 2011)

- Two criticisms of BayesB
  - Posterior of locus-specific variance has only one additional degree of freedom, compared to its prior regardless of the number of genotypes, so
    - Degree of shrinkage of depends strongly on prior
    - Little information coming from data
  - ∏ is treated as known, not estimated
     from the data

• Bayesian C∏ (Habier et al 2011)

Use a common σ<sub>gi</sub><sup>2</sup> across all SNP

Many degrees of freedom from data
A "BLUP" for SNP in model

Estimate ∏ from data

Sample from
Beta(K - m(t) + 1, m(t) + 1).

 Where K is number of SNP, m(t) is the number of SNP in the model at iteration t (eg. Those not set to zero)

- Bayesian C∏ (Habier et al 2011)
  - Accuracy in German Holstein Friesian data set

Trait	GBLUP	BayesA	BayesB	BayesCpi
Milk Yield	0.48	0.48	0.40	0.43
Fat Yield	0.51	0.56	0.52	0.54
Protein Yield	0.21	0.22	0.17	0.21
Somatic cells	0.17	0.17	0.12	0.14

• Can draw inferences about trait architecture?



- BayesR -> variants belong to one of 4 normal distributions, with zero, very small, small, medium variance
- Posterior proportion of variants in each distribution



# Real Data, 800K

#### • Reference

- Holstein = 3049 bulls, 8478 cows
- Jersey = 770 bulls, 3917 cows
- Validation
  - Holstein = 262 bulls
  - Jersey = 105 bulls
  - Australian Reds = 114 bulls
- GEBV with GBLUP, BayesR
- (Kemper et al GSE, 2014)







# Real Data, 800K

#### • r(GEBV,DTD)

	Fat	Milk	Protein	Fat%	Protein%	Average
Holstein						
GBLUP	0.60	0.59	0.58	0.72	0.83	0.66
BAYESR	0.64	0.62	0.57	0.81	0.84	0.69
Jersey						
GBLUP	0.56	0.62	0.67	0.64	0.76	0.65
BAYESR	0.56	0.69	0.71	0.76	0.79	0.70
Australian Reds						
GBLUP	0.20	0.16	0.11	0.32	0.34	0.22
BAYES	0.26	0.21	0.13	0.44	0.36	0.28







- Methods for deriving prediction equation differ in assumptions about distribution of QTL effects
  - BLUP = normal distribution with known variance
  - Ridge regression = normal distribution with prior assumption about variance
  - BayesA = t-distribution, degree of shrinkage known apriori, or sampled
  - BayesB = mixture distribution, many effects zero
  - BayesianLASSO, double exponential distribution of effects
  - Bayesian C∏, estimate ∏ from data, common variance across SNP
  - BayesR = multiple normal distributions

- Bayesian methods can have an advantage when:
- QTL of moderate to large effect on the trait (eg Fat%, DGAT1)
- Very large numbers of SNP (eg 800K) (but need large reference sets) – set some SNP effects to zero
- Multi-breed, across breed genomic predictions