



## Genome Information Assisted Prediction

Theo Meuwissen

*Norw. Univ. Life Sci.*

*Ås, Norway*



This presentation represents the views of the Authors, not the EC.  
The EC is not liable for any use that may be made of the information



## Introduction

- GS adopted in many countries
- Increasingly dense SNPchips
  - Moving to sequence data
  - Expectations are high
- However, in cattle : 50k SNPs => 700k
  - Little improvement
  - $G_{50k} \approx G_{700k} \approx G_{\text{sequence}}$



## Sequence Data

- A wealth of data, but:
  - what about the signal to noise ratio?
  - 1000 QTN are hidden amongst ~20 million SNPs
    - As needles in a hay stack
- Yet, we need to find them in some way
  - GBLUP and BayesB/C/R seem to use too little info
- Use genomic info beyond SNP genotypes?



## genomic info on SNP beyond genotype?

- In/near exon
- Conserved region (GERP score)
- In/near differentially expressed gene
- In/near eQTL (eQTL data base)
- In/near methylationQTL (mQTL)
- In/near genes of network known to affect trait
- gain/loss of stop codon
- Non-synonymous SNPs...



## AIM:

- Develop a GEBV prediction model that utilises extra genomic info about the SNP
- Genome Information Assisted Prediction
- Not the aim:
  - Determine what this extra genomic info actually is



## Hierarchical model: 1st layer

- Like BayesC:

$$y_i = \mu + \sum_j I_j X_{ij} b_j + e_i$$

$X_{ij}$  = standardised genotype

$b_j$  = SNP effect, with prior  $N(0, \sigma_b^2)$

$I_j$  = indicator variable whether SNP has effect (0/1)



## 2<sup>nd</sup> layer:

- Threshold model for  $I_j$ :

$$E(I_j) = \pi_j$$

$$\text{Probit}(\pi_j) = S\beta$$

$S$  = a matrix containing extra genomic information  
(eg: distance to gene; belong to pathway X (0/1))

$\beta$  = regression coefficients  
(unknown; are estimated)



## Extension to random $\beta$

- Prior  $\beta \sim N(0, \mathbf{A}_R \sigma_\beta^2)$   
 $\mathbf{A}_R$  = Autoregressive correlation matrix  
Auto correlation  $\rho$  is assumed known
- Thus:  $\pi_j$  depends on  $\pi_{j+1}, \pi_{j-1}, \dots$





## Implementation

- By Gibbs-sampling
  - EM-type of implementation seems possible
- Still fixed parameters:
  - $\rho$
  - $\sigma_{\beta}^2$ 
    - Determines how much  $\pi_j$  can deviate from its mean
- No results yet...



## Conclusions

- Model was build to include extra genomic info
  - Distance to gene; pathway info; conserved region...
- Hierarchical model:
  - 1st layer: models  $y$  as a function of  $b$
  - 2nd : models  $\pi_j$  as a function of extra genomic info
- 2nd model may be autoregressive
  - $\pi_j$  depends on neighboring SNPs