





Comparison of Bayesian Models for genomic selection using real dairy data

K.L. Verbyla*, P.J. Bowman, B.J. Hayes, H. Raadsma,

M. Khatkar and M.E. Goddard



Prepared for Session 28 EAAP 2009 *Email: klara.verbyla@dpi.vic.gov.au







Project Aim: To find the best method for GENOMIC SELECTION for implementation and use by the Australian dairy industry





State Covernment of Victoria Primary Industries



Generally, the methods to predict of GEBV face 2 statistical issues:

•Number of SNP (p) is greater than the number of individuals or records (n) i.e **p>n** problem

- oversaturated or overparameterised model

•Large number of SNP effects that are zero or close to zero (need a sparse model).

WHAT APPROACH SHOULD BE USED??

-Dimension Reduction (PCA, PLS)

-Machine learning (SVM)

–Shrinkage models – penalised methods exploring sparsity (LASSO)

-Variable Dimension Model approaches (SSVS)

-Variable Selection (reduced set of SNPs)

Project: Tested 20 Methods!

Bayesian Inference





Data

Reference Population

- 1098 Holstein Friesian bulls progeny tested ≤ 2003
 Validation Population
- 400 Holstein Friesian bulls progeny tested > 2003

Phenotypes

• deregressed breeding values for protein, fat, milk volume, protein%, fat%, fertility, ASI (Australian Selection Index), APR (Australian Profit Ranking) and Overall type.

Genotypes

•39,048 markers

Evaluate methods on

- r(GEBV,ABV)
- Correlation of Predicted GEBV with Australian breeding value (ABV) •MSE
- •Regression Coefficient











Bayesian Models Statistical Model:

$$y = \mu 1_n + \sum_{j=1}^{p} X_j \beta_j + Zu + e$$

- *y* is the vector of phenotypes of the trait for n individuals
- μ is the mean
- 1_n is a vector of ones of length n
- X'_{j} is a vector of indicator variables representing the genotypes of the jth marker for all individuals ($x_{ij}=0,1,2$)
- β_i is the size of the SNP effect associated with marker *j*
- \vec{u} is the vector of random polygenic effects of length n (Z is the associated design matrix)

$$u \sim N(0, \sigma_u^2 A)$$

• e is the residual error

$$e \sim N(0, \sigma_e^2 \mathbf{I}_n)$$

$$GEBV = \hat{u} + X\hat{\beta}$$



State Covernment of Victoria Department of Primary Industries



Bayes A

Prior Distributions SNP effects

 $\boldsymbol{\beta}_i \mid \boldsymbol{\sigma}_i^2 \sim N(0, \boldsymbol{\sigma}_i^2)$

$$\sigma_i^2 \sim \chi^{-2}(r,S) \sim \gamma^{-1}\left(\frac{r}{2},\frac{rS}{2}\right)$$

r degrees of freedom and scale parameter S

SNP EFFECTS

- Normal-inverse scaled chi square (t distribution)
- unequal variance
- assumes that all SNPs have an effect Gibbs Sampler







Bayes BLUP

Prior Distributions SNP effects

 $\beta_{i} \mid \sigma_{\beta}^{2} \sim N(0, \sigma_{\beta}^{2})$ $\sigma_{\beta}^{2} \sim \chi^{-2}(r, S) \sim \gamma^{-1}\left(\frac{r}{2}, \frac{rS}{2}\right)$

r degrees of freedom and scale parameter S

SNP EFFECTS

- Normal
- equal variance
- Infinitesimal assumptions
- assumes that all SNPs have an effect
- •Gibbs Sampler





Bayes C Stochastic Search Variable Selection (SSVS)

(George and McCulloch, 1993)

Use latent variable γ_i (0,1)

Prior Distributions SNP effects $\beta_i / \gamma_i, \sigma_i^2 \sim (1 - \gamma_i) N(0, \sigma_i^2 / 100) + \gamma_i N(0, \sigma_i^2)$ $\sigma_i^2 \sim \chi^{-2}(r, S)$ $\gamma_i \sim bernoulli(p_i)$ $1 - p(\gamma_i = 0) = p(\gamma_i = 1) = p_i$ SNI

SNPs with γ_i =0, posterior values limited to values close to 0 (but not removed from the model- NO changing dimensionality) – GIBBS SAMPLER

SAME ASSUMPTIONS AS BAYES B

SNP EFFECTS

Mixture of two Normalinverse scaled chi square distributions (t distributions)
unequal variance
assumes that a few SNP have an significant effect







Methods

- Bayes A
 - All SNPs
 - Selected SNPs with weights
 - Selected SNPs without weights
- Bayes BLUP
 - All SNPs
 - Selected SNPs with weights
 - Selected SNPs without weights
- Bayes C
 - All SNPs





State Covernment of Victoria Department of Primary Industries



SNP Pre-selection

The Cooperative Research Centre for

BEEF GENETIC TECHNOLOGIES

Single SNP analysis (ASReml)

$$y = 1_n \mu + X\beta + Z_1 u_1 + Z_2 u_2 + e$$

- X is a vector of indicator variables representing the genotypes of the current SNP marker for all individuals (X_n=0,1,2) and β is the associated effect of the SNP
- u_1 is the random sire effect (Z₁ associated design matrix)

 $u_1 \sim N(0, \sigma_{u_1}^2 I)$

u₂ is the random maternal grand sire effect (Z₂ associated design matrix)

 $u_2 \sim N(0, \sigma_{u_2}^2 A)$

Fitted with and without weights

-Weights = Number of Effective Records

-SNPs with p-value <0.1 included in predictive set.







Results r(GEBV,ABV)

Method	Fat	Fat%	Milk	Protein	Protein%
Bayes BLUP - All SNPs	0.528	0.630	0.648	0.613	0.660
Bayes BLUP - Selected SNPs (Unweighted)	0.527	0.689	0.646	0.596	0.678
Bayes BLUP - Selected SNPs (Weighted)	0.543	0.643	0.659	0.610	0.661
Bayes A - All SNPs	0.538	0.700	0.631	<mark>0.572</mark>	0.645
Bayes A - Selected SNPs - (Unweighted)	0.543	0.712	0.639	0.579	0.667
Bayes A - Selected SNPs - (Weighted)	0.538	0.704	0.635	0.583	0.648
Bayes C	0.557	0.728	0.644	0.588	0.670



(

State Covernment Victoria Department of Primary Industries



Fat % - DGAT1

The Cooperative Research Centre for

BEEF GENETIC TECHNOLOGIES

 Table 1. Effect of the DGAT1 K232A Mutation on Sires' Daughter Yield Deviations (DYDs) for Milk Yield and Composition

Trait	$\alpha/2 \pm 2$ std.err.	r _{QTL}	Р value _{qтL}	r ² _{polygenic}	r ² error
Milk yield (kgs)	-158 ± 24.5	0.18	5.00E - 35	0.49	0.32
Fat yield (kgs)	5.23 ± 0.9	0.15	1.57E – 29	0.55	0.30
Protein vield (kgs)	-2.82 ± 0.7	0.08	1.70E – 15	0.65	0.26
Fat (%)	0.17 ± 0.012	0.51	4.33E - 122	0.29	0.19
Protein (%)	0.04 ± 0.006	0.14	5.05E – 28	0.66	0.20

(i) $\alpha/2$: QTL allele substitution effect on DYD (\approx halve breeding value), corresponding in the mixed model to the regression coefficient on the number of *K* alleles in the *DGAT1 K232A* genotype, and to $\alpha/2$, where α is defined according to Falconer and Mackay 1996. (ii) r_{QTL}^2 : proportion of the trait variance explained by the *DGAT1 K232A* polymorphism. (iii) *P* value_{QTL}: statistical significance of the *DGAT1 K232A* effect. (iv) $r_{polygenic}^2$: proportion of the trait variance explained by the random, polygenic effect in the mixed model. (v) r_{error}^2 : proportion of the trait variance unexplained by the model.

- QTL explains > 50% of genetic variance in fat%
- QTL allele is common and acts additively
- \rightarrow Major violation of BLUP assumptions





Results r(GEBV,ABV)







State Covernment Victoria Department of Primary Industries



RESULTS - IN CONTEXT:

A SUBSET OF ALL METHODS USING ALL SNPs

	Average Correlation
Method	Across all traits
Bayes BLUP	0.589
Bayes A	0.578
Bayes C	0.597
LASSO	0.595
SVR	0.587
GBLUP	0.588
PLS	0.592



State Covernment Victoria Department of Primary Industries

Conclusions



- Only small differences in accuracy and bias of GEBV from different methods
- Method by trait interaction. Better results when priors matched "real" distribution of QTL effects

 \rightarrow best method is trait dependent!

- Pre-selection of SNP neither reduces or increases the accuracy of predicted GEBV
- Bayesian BLUP performs as well as or better than the other methods EXCEPT for traits with QTL that explain large amount of genetic variance eg. Fat % with DGAT1
- Still a need to find a method that produces equally accurate GEBV across traits with different genetic architecture.



State Covernment Victoria Department of Primary Industries



Recognitions

Personally:

- Ph.D Supervisors: Mike Goddard, Ben Hayes and Richard Huggins.
- Melbourne University, Beef CRC and DPI Victoria for their ongoing support.
- Marie Curie EST Sabretrain Fellowship Supervisors: Roel Veerkamp, Mario Calus and Han Mulder (Wageningen University and Research)

Project:

B. Hayes

H. Raadsma

P. Bowman

M. Khatkar

M.E. Goddard

All other people involved in the project

