

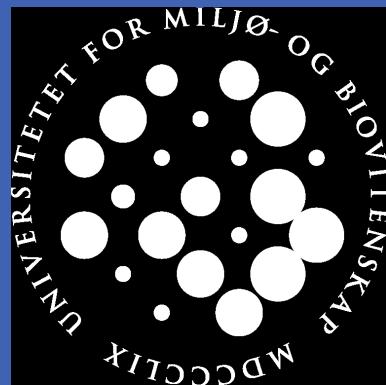
# **Using partial least square regression (PLSR) and principal component regression (PCR) in prediction of genome wide breeding values**

T. R. Solberg<sup>1</sup>, A.K. Sonesson<sup>2</sup>, J. A. Woolliams<sup>1,3</sup>, T. H. E. Meuwissen<sup>1</sup>

<sup>1</sup>*University of Life Sciences, Dept. of Animal and Aquacultural Sciences, P.O. Box 5003, N-1430 Aas, Norway.*

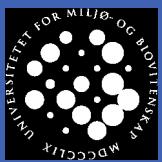
<sup>2</sup>*AKVAFORSK (Institute of Aquaculture Research Ltd.), P. O. Box 5010, N-1432 Aas, Norway.*

<sup>3</sup>*Roslin Institute (Edinburgh), Roslin, Midlothian EH25 9PS, UK.*



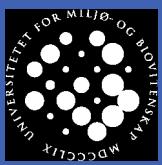
## Introduction

- Dense molecular data available
  - Genomic selection (e.g. Meuwissen *et al.*, 2001; Gianola *et al.*, 2006; Schaeffer, 2006; Solberg *et al.*, 2006)
- SNP data is a typical case where the number of predictors can be larger than the number of records
- Partial Least Square Regression (PLSR) and Principal Component Regression (PCR) are methods designed for such situations
- Test the hypothesis that genome wide evaluation can be obtained using models of reduced dimensionality.
- Compared to the 'bayesB' method (Solberg *et al.*, 2006).



## Introduction (cont.)

- PLSR and PCR mostly used within econometrics, chemometrics and social sciences (e.g. Wold, 1981; Wold, 1985; Martens & Næs, 1991)
- Basic idea to reduce the number of predictors (principal components, PC)



## Materials and methods

### Population structure/genome

- Stochastic simulation study,  $N_e=100$ . Random mating for 1000 generations



Increased to 1000 animals in  $t = 1001$  and  $t = 1002$

Phenotypic values simulated by  $P_i = TBV_i + \varepsilon_i$  ( $\varepsilon \sim N(0, \sigma^2)$ )

Animals genotyped in generation  $t = 1001$  and  $t = 1002$

BVs predicted on animals in generation  $t = 1002$ , using marker information in generation  $t = 1001$ .

## M & M (cont.)

- Genome simulated with 10 chromosomes, each with a length of 100 cM
- 4 different marker densities evaluated, 1 cM, 0.5 cM, 0.25 cM and 0.125 cM spacing between the markers

(1010 mrk  8080 mrk)

## M & M (cont.)

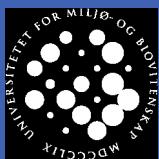
- Marker genotype data organised as  $m \times p$  matrix ( $\mathbf{X}$ )  
 $(1000 \times 1010 \longrightarrow 1000 \times 8080)$
- $\mathbf{y}$  vector and  $\mathbf{X}$  matrix centered by subtracting their mean and scaled by deviding by their SD
- PCR create the PC by linear combinations among predictors
- PLSR create the PC by maximising the covariance between the predictors and the response

### PLSR

- SIMPLS algorithm performed (de Jong, 1993)
  - Compute dominant eigenvectors of the matrix  $\mathbf{X} (\mathbf{X}'\mathbf{y})$ , since the data is univariate)
  - Deflating  $\mathbf{X}$  matrix and  $\mathbf{y}$  vector to compute the next PC

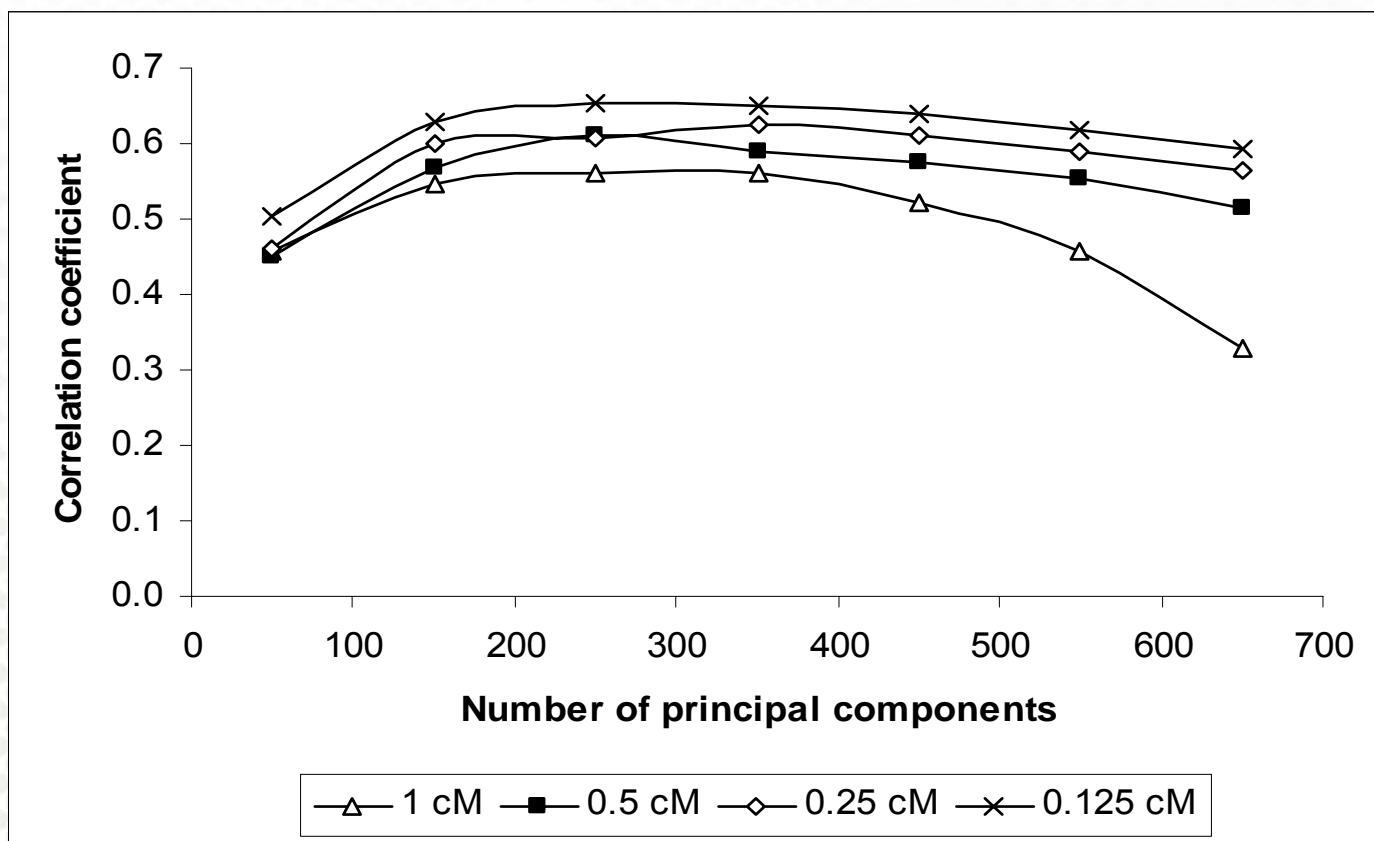
### PCR

- Singular value decomposition of the  $\mathbf{X}$  matrix to find the PC (Press *et al.*, 2003)



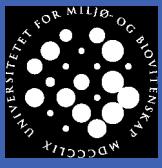
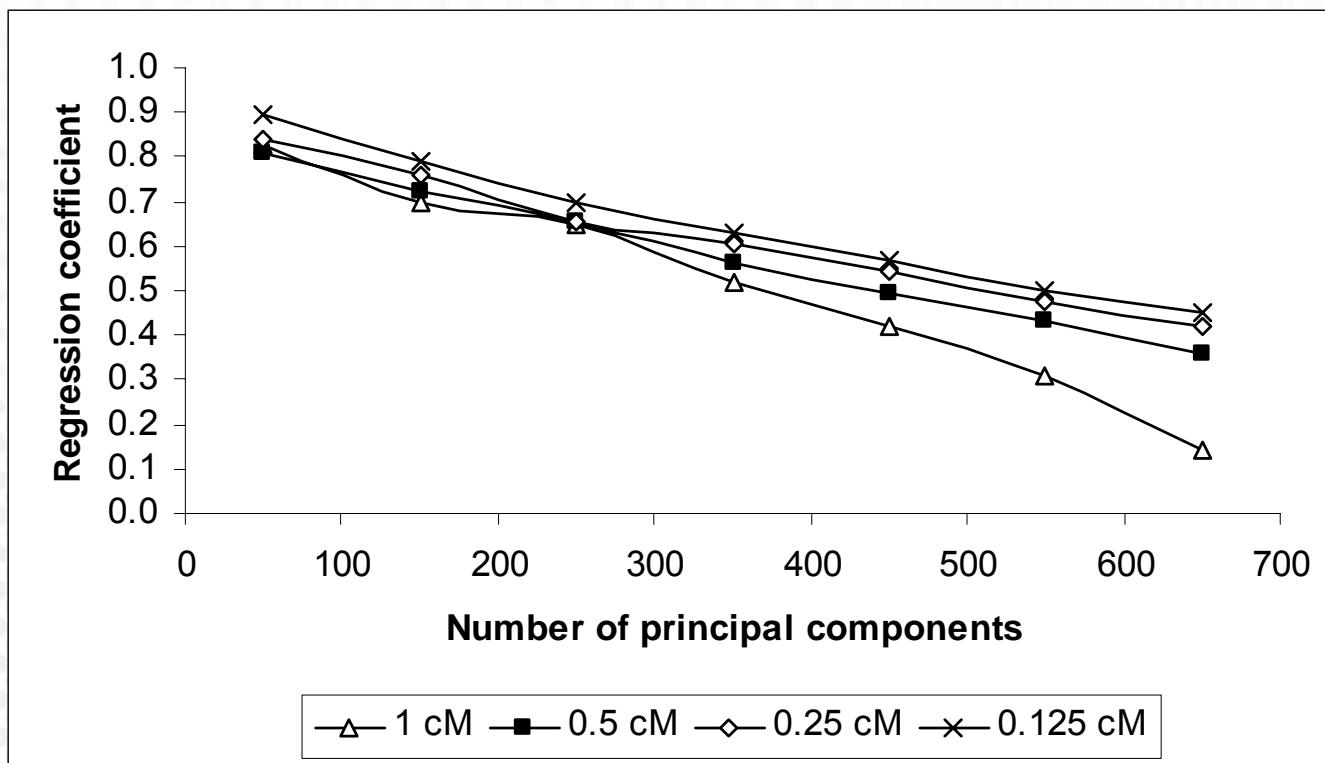
## Results

- Optimum number of principal components using PCR (correlation)



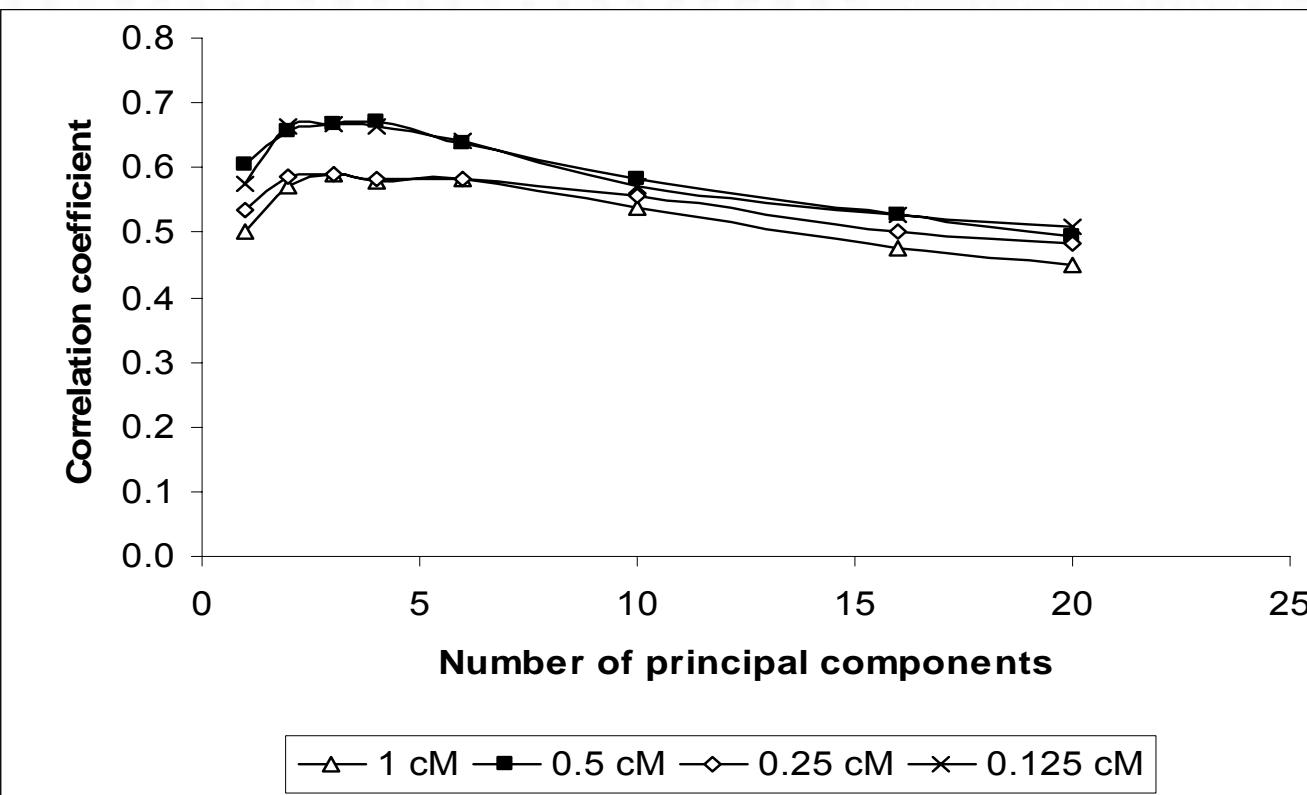
## Results (cont.)

- Optimum number of principal components using PCR (regression)



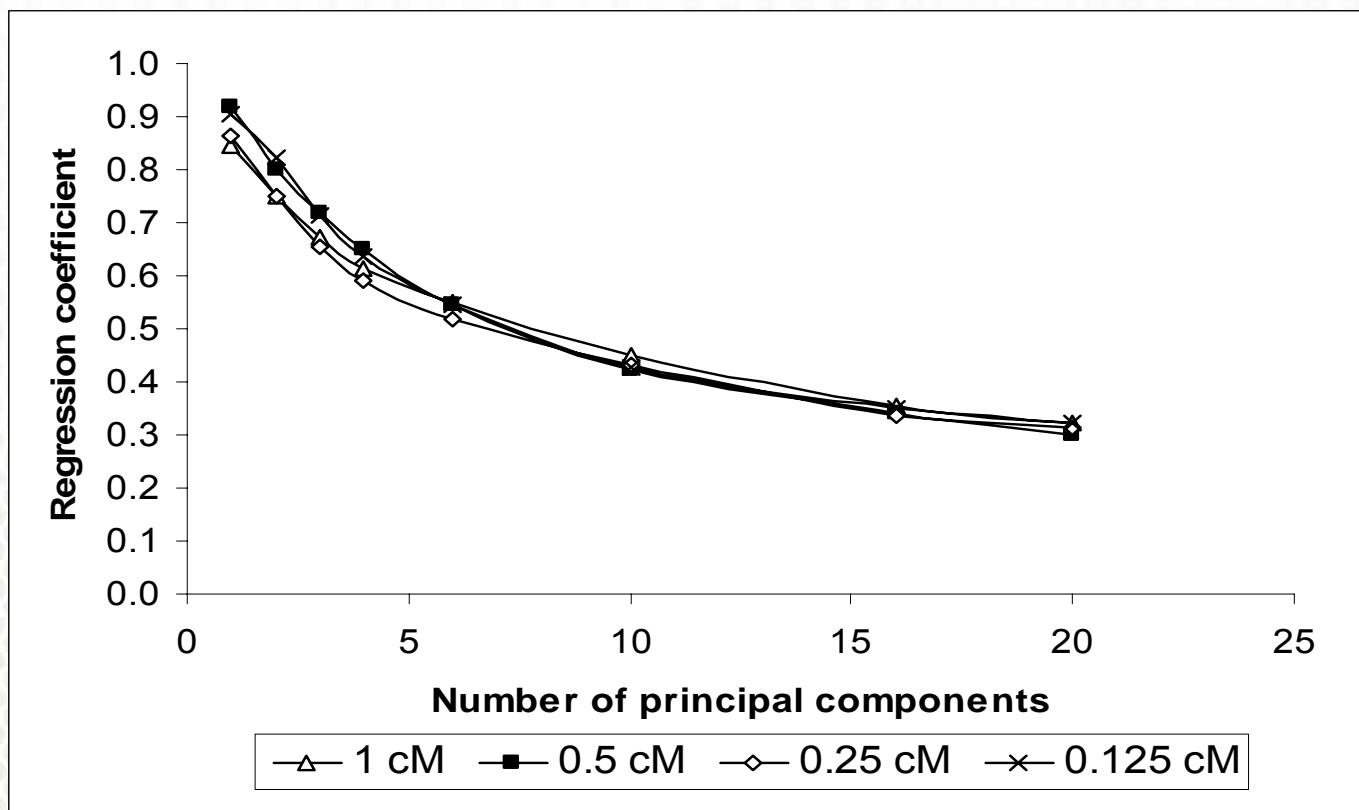
## Results (cont)

- Optimum number of principal components using PLSR (correlation)



## Results (cont.)

- Optimum number of principal components using PLSR (regression)

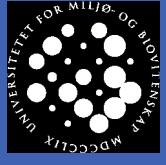


## Results (cont.) Correlation

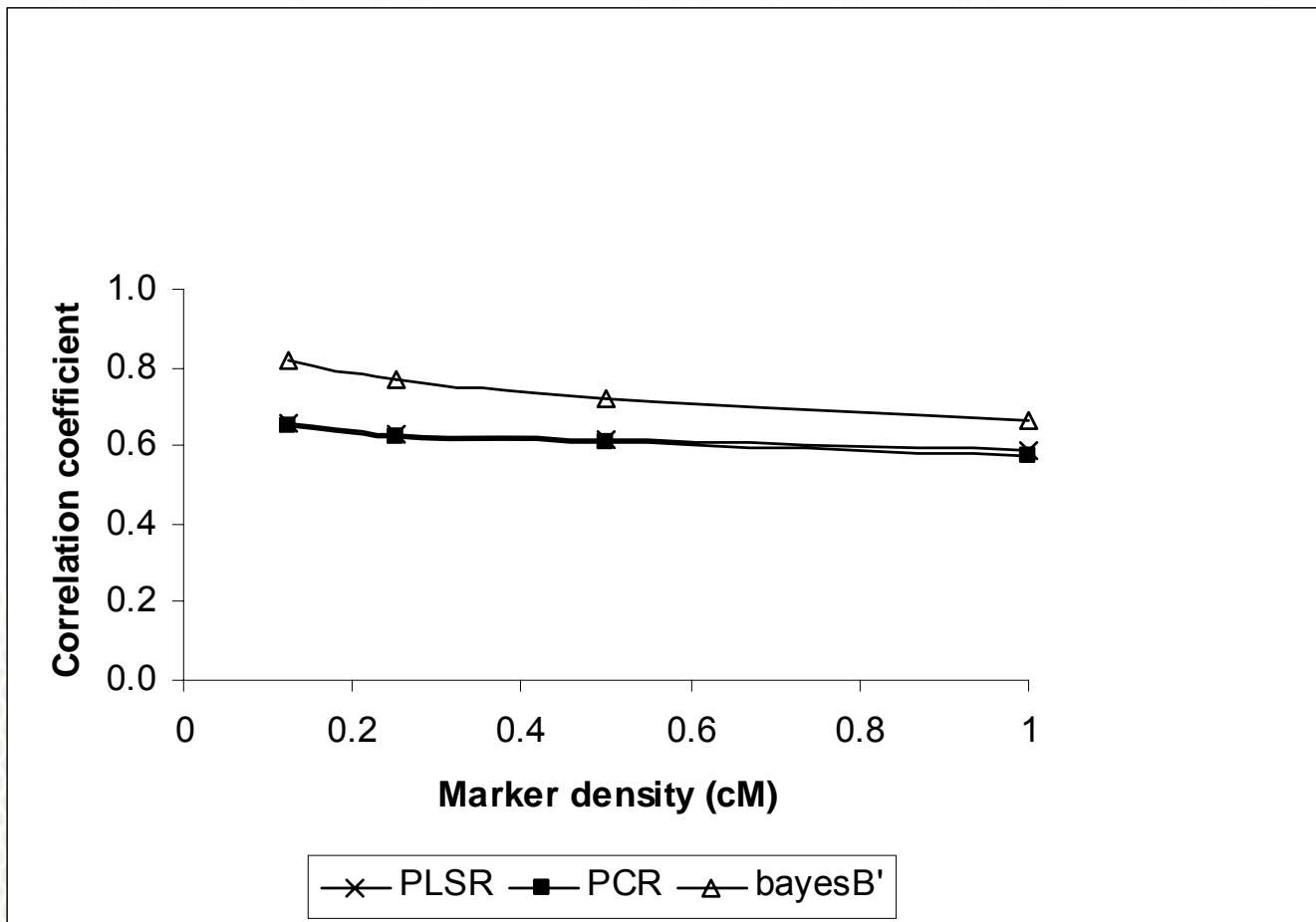
	PCR	PLSR	'bayesB'
Marker density	$r_{TBV;EBV} \pm SE$	$r_{TBV;EBV} \pm SE$	$r_{TBV;EBV} \pm SE$
1 cM	$0.570 \pm 0.014$	$0.588 \pm 0.014$	$0.663 \pm 0.019$
0.5 cM	$0.606 \pm 0.009$	$0.612 \pm 0.011$	$0.717 \pm 0.014$
0.25 cM	$0.620 \pm 0.010$	$0.630 \pm 0.009$	$0.771 \pm 0.010$
0.125 cM	$0.649 \pm 0.013$	$0.658 \pm 0.012$	$0.820 \pm 0.017$

## Results (cont.) Regression

	PCR	PLSR	'bayesB'
Marker density	$b_{TBV;EBV} \pm SE$	$b_{TBV;EBV} \pm SE$	$b_{TBV;EBV} \pm SE$
1 cM	$0.645 \pm 0.011$	$0.749 \pm 0.012$	$0.880 \pm 0.016$
0.5 cM	$0.642 \pm 0.011$	$0.754 \pm 0.012$	$0.894 \pm 0.018$
0.25 cM	$0.668 \pm 0.013$	$0.782 \pm 0.011$	$0.909 \pm 0.012$
0.125 cM	$0.692 \pm 0.010$	$0.801 \pm 0.010$	$0.910 \pm 0.012$



## Results (cont.)



## Conclusions

- PLSR and PCR computationally very fast (hours)
- Simpler and less assumptions compared to 'bayesB'
- Possible to predict breeding values for next generation (reasonable high accuracy)
- Accuracy increased as the marker density increased
- Why use it?
- Simple, computationally fast, non-parametric method with no assumptions