

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

**Use of partial least squares regression to predict single nucleotide polymorphism marker genotypes when some animals are genotyped with a low-density panel**

**This is the author's manuscript**

*Original Citation:*

*Availability:*

This version is available <http://hdl.handle.net/2318/1686986> since 2019-02-04T09:43:35Z

*Published version:*

DOI:10.1017/S1751731110002600

*Terms of use:*

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

**This is the author's final version of the contribution published as:**

Dimauro, C., Steri, R., Pintus, M., Gaspa, G., & Macciotta, N. (2011). Use of partial least squares regression to predict single nucleotide polymorphism marker genotypes when some animals are genotyped with a low-density panel. *Animal*, 5(6), 833-837.  
doi:10.1017/S1751731110002600

**The publisher's version is available at:**

**When citing, please refer to the published version.**

**Link to this full text:**

**<http://hdl.handle.net/hdl:2318/1686986>**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

**Use of partial least-squares regression to predict single-nucleotide polymorphism marker genotypes when some animals are genotyped with a low density panel**

C. Dimauro<sup>a</sup>, R. Steri, M. A. Pintus, G. Gaspa and N. P.P. Macciotta

*Dipartimento di Scienze Zootecniche, Università di Sassari, via De Nicola 9, 07100 Sassari, Italy*

<sup>a</sup>Corresponding author: Corrado Dimauro, Dipartimento di Scienze Zootecniche, Università di Sassari, via De Nicola 9, 07100 Sassari, Italy. Tel. +39 079229298; fax +39 079 229302. E-mail: [dimauro@uniss.it](mailto:dimauro@uniss.it)

Running head: Single nucleotide polymorphism prediction

23 **Abstract**

24

25 *High density SNP platforms are currently used in Genomic Selection (GS) programs to*  
26 *enhance the selection response. However, the genotyping of a large number of animals*  
27 *with high throughput platforms is rather expensive and may represent a constraint for a*  
28 *large-scale implementation of GS. The use of low density marker platforms could*  
29 *overcome this problem, but different SNP chips may be required for each trait and/or*  
30 *breed. In this paper a strategy of imputation independent from trait and breed, is proposed.*  
31 *A simulated population of 5,865 individuals with a genome of 6,000 SNP equally*  
32 *distributed on six chromosomes was considered. First, reference and prediction*  
33 *populations were generated by mimicking high and low density SNP platforms,*  
34 *respectively. Then, the partial least squares regression (PLSR) technique was applied to*  
35 *reconstruct the missing SNP in the low density chip. The proportion of SNP correctly*  
36 *reconstructed by the PLSR method ranged from 0.78 to 0.97 when 90% and 50% of*  
37 *genotypes were predicted, respectively. Moreover, data sets consisting of a mixture of*  
38 *actual and PLSR-predicted SNP or only actual SNP were used to predict genomic*  
39 *breeding values (GEBV). Correlations between GEBV and true breeding values varied*  
40 *from 0.74 to 0.76 respectively. Results of the study indicate that the PLSR technique can*  
41 *be considered a reliable computational strategy for predicting SNP genotypes in a low*  
42 *density marker platform with reasonable accuracies.*

43

44

45 **Keywords:** genomic selection; SNP prediction

46

47

48

49 **Implications**

50

51 In genomic selection programs, animals are genotyped with high-density SNP marker  
52 platforms with around 50-60K markers. However, being the number of phenotypes  
53 available markedly lower than the number of markers, several statistical shortcomings  
54 arise when data are analyzed. In this paper we propose the use of both high and low-  
55 density SNP marker platforms in combination with partial least squares regression (PLSR)  
56 technique to reconstruct the missing SNP in the low density chips. Savings obtained by  
57 using low density platforms could be used to enlarge the number of animals involved in the  
58 selection program.

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

## 75 **Introduction**

76

77 Traditional genetic evaluations for livestock combine phenotypic data with pedigree  
78 relationships to estimate the probability that genes are transferred to the next generations.  
79 Genomic selection (GS), on the contrary, exploits dense marker information represented  
80 by single nucleotide polymorphism (SNP) to evaluate genomic breeding values (GEBV) by  
81 estimating the effect of chromosome segments on phenotypes (Hayes and Goddard,  
82 2008). Advances in high throughput technologies have led to the construction of dense  
83 SNP platforms that could trace the inheritance of individual genes. High density marker  
84 (HDM) platforms with 50 – 60 K SNP are currently used in GS programmes. However, the  
85 number of genotyped animals is considerably smaller than the number of markers. In dairy  
86 cattle, the ratio number of animals vs. number of markers is, on average, between 0.08-  
87 0.15, apart from USA and Canada where it is around 0.45 (VanRaden *et al.*, 2009). Such a  
88 data asymmetry results in several statistical shortcomings, as collinearity among predictors  
89 and issues in multiple testing procedures. Furthermore, the well known curse of multi-  
90 dimensionality should become now more relevant, due to the recent commercial  
91 availability of the 777 K SNP Illumina Bead-chip.

92 The use of low density marker platforms (LDM) may represent a interesting technical  
93 option to reduce the genotyping costs and enlarge the number of animals involved in GS  
94 programmes. However, the reduction of SNP density is expected to decrease GEBV  
95 accuracy. Weigel *et al.* (2009) reported a loss of about one-third in the gain of reliability of  
96 GEBV for lifetime profit in cattle when a low-density assays with 750-1,000 SNP was used.  
97 In this study, SNP were chosen either on the basis of their chromosomal location (evenly  
98 spaced) or for their relevance on the considered trait. Habier *et al.* (2007) combined the  
99 use of evenly spaced SNP and co-segregation information from LDM to track HDM

100 inheritance within families. On simulated data, they found a reduction in GEBV accuracy  
101 ranging from 1 to about 25%, depending on the considered scenario.

102 The use of the above mentioned methodologies can be useful to reduce the number of  
103 SNP but, separate chips for each trait and/or breed may be required. In this paper an  
104 alternative strategy, independent from trait or breed, is proposed. The method starts by  
105 creating a reference (REF) and a prediction (PRED) population of animals genotyped with  
106 HDM (containing  $N$  SNP) and LDM ( $n$  SNP) platforms, respectively ( $N > n$ ). Missing  $k$ -  
107 markers ( $k = N-n$ ) in PRED population are reconstructed by using a suitable mathematical  
108 tool and, as a final result, a PRED population with  $N$  SNP as in HDM is obtained. These  
109 markers are a mixture of actual and predicted SNP.

110 The most straightforward computational method for predicting unknown SNP markers in  
111 the LDM platform is the multivariate multiple regression. However, considering that  
112 adjacent SNP are highly correlated, the predictive capability of the model could be  
113 compromised by the multicollinearity among predictors (Draper and Smith, 1981). Partial  
114 least squares regression (PLSR), originally developed in the computational chemistry  
115 context (Hoeskuldsson, 1988), has become an established tool for modeling linear  
116 relations between multivariate measurements. It is characterized by an higher prediction  
117 efficiency compared to ordinary multivariate regression or principal component regression  
118 (Macciotta *et al.*, 2006). PLSR has been already used in GS studies by Solberg *et al.*  
119 (2009) for reducing the dimensionality of predictors in the calculation of GEBV. In the  
120 present study, the PLSR technique is applied to predict missing SNP when animals are  
121 genotyped with a LDM platform. Actually, this statistical technique is particularly useful  
122 when a set of correlated dependent variables ( $\mathbf{Y}$ ) have to be predicted from a set of  
123 correlated independent variables ( $\mathbf{X}$ ). PLSR maximizes the correlation structures between  
124  $\mathbf{Y}$  and  $\mathbf{X}$  and overcomes the multicollinearity problems by combining features of principal  
125 components analysis and multiple regression (Abdi, 2003).

126 The aim of this work is to test the ability of PLSR for predicting missing SNP genotypes  
127 when a PRED population is created by using a LDM platform of SNP markers.

128

## 129 **Materials and methods**

130

### 131 *The data*

132 Data were extracted from an archive generated for the XII QTLs – MAS workshop, freely  
133 available at: <http://www.computationalgenetics.se/QTLMAS08/QTLMAS/DATA.html>. The  
134 base population consisted of 100 individuals (50 males and 50 females). A genome of six  
135 chromosomes (total length 6 M) with 6,000 biallelic SNP, equally spaced in the genome at  
136 a distance of 0.1 cM, was generated. A total of 48 biallelic QTLs were included, with  
137 positions sampled from the genetic map of the mouse genome and effects derived from a  
138 gamma distribution (Hayes and Goddard, 2001). Initial allelic frequencies of both SNP and  
139 QTL were set to 0.5. Then 50 generations of random mating followed. Generations from  
140 51 to 57 were used to create the definitive archive of 5,865 individuals. For each  
141 generation 15 males and 150 females were randomly selected to be parents of the next  
142 generation. Each male had 100 sons and was mated to 10 females (10 sons for female).  
143 Animals belonging to the generations from 51 to 54 had pedigree, phenotype, and marker  
144 information available. For the last 3 generations only pedigree and marker information  
145 were available. These animals constituted the PRED population and were obtained by  
146 randomly selecting 400 animals for each generation (a total of 1200 individuals). True  
147 breeding values (TBV) were created as the sum of all QTL effects across the entire  
148 genome. Phenotypes were generated by adding to the TBV an environmental noise drawn  
149 from a normal distribution with mean zero and variance equal to the residual variance  
150 defined to obtain a heritability of 0.30. For further details on the data generation see Lund  
151 *et al.* (2009).



152

153 *The PLSR technique*

154 PLSR is a multivariate extension of the multiple regression analysis. It is particularly useful  
155 when (i) the number of predictor variables is similar to or higher than the number of  
156 observations and/or (ii) predictors are highly correlated (i.e. there is strong collinearity).

157 The basic model is:

158  $\mathbf{Y}=\mathbf{XB}+\mathbf{E}$

159 where  $\mathbf{Y}$  is a  $n \times m$  response matrix,  $\mathbf{X}$  is a  $n \times p$  design matrix,  $\mathbf{B}$  is a  $n \times m$  regression  
160 coefficient matrix, and  $\mathbf{E}$  is a  $n \times m$  error term. In PLSR, matrices  $\mathbf{X}$  and  $\mathbf{Y}$  are  
161 simultaneously decomposed into a set of new variables (called latent factors). Factors are  
162 extracted in order to explain as much as possible of the covariance between  $\mathbf{X}$  and  $\mathbf{Y}$  and  
163 to minimize the covariance between variables inside each matrix. Extracted latent factors  
164 account for successively lower proportions of original variance and are defined as linear  
165 combinations of predictor and response variables (Hubert and Branden, 2003). Key  
166 elements in the different calculation steps of the PLSR are: the scores, i.e. values of the  
167 extracted latent factors both for the dependent ( $\mathbf{U}$ ) and independent variables ( $\mathbf{T}$ ), and  
168 factor loadings ( $\mathbf{Q}$ ) expressing correlations between extracted factors and original  
169 dependent variables. Considering a REF and a PRED population, latent factor scores ( $\mathbf{T}_{\text{ref}}$ )  
170 extracted from  $\mathbf{X}_{\text{ref}}$ , are used to predict scores of latent factors extracted from  $\mathbf{Y}_{\text{ref}}$  ( $\mathbf{U}_{\text{ref}}$ )

171  $\mathbf{U}_{\text{ref}} = \mathbf{BT}_{\text{ref}} \quad (1)$

172 Then, the estimated regression coefficients  $\mathbf{B}$  are used to predict values of  $\mathbf{Y}_{\text{pred}}$  in the  
173 PRED population as:

174  $\hat{\mathbf{Y}}_{\text{pred}} = \mathbf{BT}_{\text{pred}}\mathbf{Q}'_{\text{ref}} \quad (2)$

175 where  $\mathbf{Q}'_{\text{ref}}$  is the transposed matrix of factor loadings extracted from  $\mathbf{Y}_{\text{ref}}$ .

176 The standard algorithms for computing latent factors are nonlinear and iterative (NIPALS  
177 and SIMPLS algorithms, for example) and require the use of dedicated software (for more  
178 details see Wold *et al.*, 2001; de Jong, 1993). In this work, the PLS procedure of SAS-  
179 STAT software (SAS Institute INC, Cary, NC) was used.

180

### 181 *The PLSR method for SNP genotypes prediction*

182 To simulate a PRED population genotyped with a LDM platform, the first  $k$ -SNP were  
183 assumed to be not known. SNP from  $k+1$  to 1,000 represented the predictors (i.e.  $\mathbf{X}_{\text{ref}}$  and  
184  $\mathbf{X}_{\text{pred}}$ ) and were known both for REF and PRED population. SNP from 1 to  $k$  were known in  
185 REF ( $\mathbf{Y}_{\text{ref}}$ ) and were used to calculate the matrix of regression coefficients  $\mathbf{B}$  (equation 1).

186 Then, using the equation (2), the  $\hat{\mathbf{Y}}_{\text{pred}}$  matrix was predicted. Being that the genotype at  
187 each SNP is coded as the number of allele 1 copies, i.e. 0, 1 or 2, results (columns  
188 in  $\hat{\mathbf{Y}}_{\text{pred}}$  each containing the predicted SNP genotype) were rounded to the nearest integer.

189 The goodness of SNP prediction was evaluated by calculating correlations between real  
190 ( $\mathbf{Y}_{\text{pred}}$ ) and PLSR predicted ( $\hat{\mathbf{Y}}_{\text{pred}}$ ) SNP genotypes. Considering that for  $k$  predicted SNP  
191  $k$  correlations were calculated, the average value of these correlations, for each prediction  
192 scenario, was considered. Moreover, percentage of correct predictions across SNP and  
193 mean percentage of corrected SNP predictions for each animal were calculated.

194 A crucial point in PLSR modeling is how many latent factors should be retained to  
195 correctly define the complexity of one experiment. When several and correlated predictors  
196 are used, the risk of obtaining a model able to fit data well but with a very poor predictive  
197 power is rather high. This problem is known as model “over-fitting”. It is usually handled by

198 testing the predictive significance of the successive extracted factors. Cross-validation in  
199 combination with PRESS statistics is commonly used to this purpose (Wold *et al.*, 2001).  
200 However, in the present study several scenarios involving a great number of predictors are  
201 compared and, therefore, the use of the above cited tests become problematic in terms of  
202 computation time and resources. For these reasons, the best number of extracted latent  
203 factors in each scenario was fixed empirically by comparing the obtained results with real  
204 data (the procedure will be explained in the next section).

205

#### 206 *Setup of the PLSR method*

207 Location of missing SNP along the chromosome, number of latent factors to be extracted  
208 for each scenario, number of SNP to be predicted and the minimum number of genotyped  
209 animals to use as REF population are relevant aspects for the method be efficiently  
210 performed in practice. They were tested in successive steps during the development of the  
211 PLSR method. All the computations were done separately per chromosome .

212 Step 1: four scenarios of chromosome location of SNP to be predicted ( $k=100$ ) in PRED  
213 population were tested: at the beginning (SNP1 – SNP100), in the middle (SNP451-  
214 SNP550), at the end (SNP901 – SNP1,000), or evenly spaced in the chromosome.

215 Step 2: once the best SNP location was assessed, the optimum number of latent factors to  
216 be extracted was evaluated. In PLSR procedure, the number of factors can not exceed the  
217 number of the independent variables. Therefore, for each chromosome, several  
218 simulations were performed where 100 SNP were predicted with a number of factors  
219 ranging from 10 to 900.

220 Step 3: prediction accuracy for different number of SNP to be predicted was investigated  
221 using the following proportions for missing SNP in PRED population: 10%, 25%, 50%,

222 75% and 90%. At the end of the PLSR procedure, a series of new data sets for PRED  
223 population, each containing 10%, 25%, 50%, 75% and 90% of PLSR predicted SNP, were  
224 produced.

225 Step 4: the effect of the SNP reduction in the estimation of genomic breeding values was  
226 tested by evaluating GEBV's either in original and in five data sets, generated in step 3,  
227 which contain the mixture of actual and PLSR predicted SNP. Effects of SNP markers on  
228 phenotypes in the REF population were estimated with a mixed linear model that included  
229 the fixed effects of mean, sex (1,2) and generation (1,2,3,4), and the random effects of  
230 SNP genotypes (Meuwissen *et al.* 2001). Overall mean and effects of SNP genotypes  
231 were then used to predict GEBV in PRED population (Macciotta *et al.*, 2010). Accuracies  
232 were evaluated by calculating Pearson correlations between GEBV and true breeding  
233 values.

234 Step 5: finally, considering a possible application of the method on real data, accuracy of  
235 the PLSR predictions were tested for different sizes of the REF population, from 5,000 to  
236 600 individuals. In all the simulations, the size of PRED population was kept constant  
237 (600).

238

## 239 **Results and discussion**

240

241 Step 1: the effects of SNP location on prediction accuracy can be observed in Table 1  
242 where average correlations between actual and PLSR-predicted SNP genotypes for  
243 different scenarios are reported. Lowest correlations were obtained when markers to be  
244 predicted are located at the beginning or at the end of the chromosome. A slight increase  
245 of accuracy can be observed when SNP are located in the middle of the chromosome. The  
246 highest value was found for evenly spaced missing SNP. These results were expected,

247 considering the decaying pattern of correlation between loci for increasing distances, and  
248 are in agreement with figures reported by Habier *et al.* (2009) who had already used  
249 evenly spaced SNP to simulate low density marker panels. In any case, the value of the  
250 mean correlation for the best scenario is notably high and may represent a useful  
251 indication for constructing a LDM platform without trait or breed constraints.

252 Step 2 : Figure 1 displays pattern of mean correlations between 100 actual and PLSR  
253 predicted SNP for increasing number of extracted latent factors for the first chromosome.  
254 There is a rapid increase of prediction accuracy from 10 up to 100 factors (from 47% to  
255 93%). A plateau of 98% is then reached when about 150 - 200 factors are extracted.  
256 These results indicate that the number of latent factors to be extracted should be higher or,  
257 at least, equal to the number of predicted SNP.

258 Step 3: the variation of prediction accuracy for different number of SNP to be predicted is  
259 reported in Table 2. Moving from 10% to 75% missing SNP, there is small decrease (about  
260 6%) in the average correlation between actual and predicted genotypes. In any case,  
261 prediction accuracy is higher than 90% even when two-third of the SNP are predicted. It  
262 slightly falls below 0.80 when 90% of SNP have to be predicted. However, even in this  
263 case, the accuracy can be considered satisfactory. If confirmed on real data, results of the  
264 present study may indicate that a chip with 5.4 K SNP evenly spaced across the genome  
265 could represent a suitable base for reconstructing, with a reasonable accuracy, the profile  
266 of an high density platform of 54 K SNP (i.e. the one currently used for cattle). In a recent  
267 study carried out with the bovine 54 K SNP, Weigel *et al.* (2010) using the algorithm  
268 implemented in fastPHASE 1.2 software (University of Washington TechTransfer Digital  
269 Ventures Program, Seattle, WA), reported a proportion of correctly reconstructed missing  
270 SNP of about 0.88 when 90% SNP were predicted. Druet and Georges (2010) combined  
271 fastPHASE and Beagle (Browning and Browning, 2007) algorithms to take into account  
272 both population (linkage disequilibrium) and familial (Mendelian segregation and linkage)

273 information to predict missing genotypes. They found, with 50% missing genotypes, an  
274 imputation error of 3% and 1% for sparse and dense marker map, respectively. In the  
275 present work, the proportion of correctly reconstructed SNP for 90% and 50% missing  
276 genotypes was 0.86 and 0.98, respectively (Table 2).

277 The SNP genotype profile of each animal was also well reconstructed by the PLSR  
278 method. When 90% SNP were predicted, more than 84% of animals presented a  
279 percentage of corrected SNP reconstruction ranging from 80 to 100%. Moreover, when  
280 predicted SNP were lower than 75%, all animals had a proportion of corrected  
281 reconstructed SNP ranging from 95 to 100%.

282 Step 4: accuracies displayed in Table 3 indicate that the use of PLSR-predicted SNP does  
283 not affect the estimation of genomic breeding values. Correlations between true breeding  
284 values and GEBV remain basically the same moving from the scenario where all used  
285 SNP are actual to the one where 90% of marker genotypes are PLSR-predicted (Table 3).  
286 These results are similar to those obtained by Habier *et al.* (2009) who reported a  
287 reduction in GEBV accuracy of about 4% moving from a SNP panel density of 0.05 cM to  
288 10cM.

289 Step 5: finally, Figure 2 displays accuracies of SNP prediction obtained with different sizes  
290 of REF population. As the number of fully genotyped animals becomes smaller,  
291 correlations between actual and predicted SNP slowly decrease reaching a value of 93%  
292 when the number of REF animals is twice (2,000) the total number of SNP per  
293 chromosome. Correlations dramatically drop (<70%) for a number of fully genotyped  
294 animals equal to 600. Considering that on real data each bovine chromosome has on  
295 average 1000-1200 SNP after data editing, a minimum number of 2,000-2,500 fully  
296 genotyped animals could be enough to obtain reliable predictions from the PLSR method.

297

298 **Conclusions**

299

300 The use of LDM platforms in combination with a suitable computational algorithm able to  
301 predict the missing genotypes with respect to HDM chips is an option for reducing  
302 genotyping costs in GS programs. Savings could be used to enlarge the genotyped  
303 population thus enhancing the efficiency of the breeding scheme. In this paper, the ability  
304 of PLSR technique for predicting missing SNP genotypes in LDM platforms was tested.  
305 The method correctly assigned from 86 to 98% of missing genotypes, when 90 and 50%  
306 SNP were predicted, respectively. Moreover, only a slight difference (2%) in GEBV  
307 accuracies was observed using actual SNP or a mixture of actual and predicted SNP.  
308 Finally, a size of around 2,000-2,500 fully genotyped animals with a 54 K SNP chip was  
309 found to be a reliable REF population to reconstruct the SNP profile of a PRED population  
310 of animals genotyped with a LDM chip containing 5,4 K evenly spaced SNP.

311

312

313 **Acknowledgments**

314 Research was funded by the Italian Ministry of Agriculture (Rome, Italy), grant SELMOL.

315



316 **References**

- 317 Abdi H 2003. Partial least squares (PLS) regression. In *Encyclopaedia of social sciences research*  
318 *methods* (eds M Lewis–Beck, A Bryman and T Futing) pp. 1-7. Sage Publication, Thousand Oaks,  
319 CA.
- 320 Browning SR and Browning BL 2007. Rapid and accurate haplotype phasing and missing-data  
321 inference for whole-genome association studies by use of localized haplotype clustering. *American*  
322 *Journal of Human Genetics* 81, 1084-1097.
- 323 De Jong S 1993. SIMPLS: an alternative approach to partial least squares regression.  
324 *Chemometrics and Intelligent Laboratory Systems* 18, 251–263.
- 325 Draper NR and Smith H 1981. *Applied regression analysis*. John Wiley and Sons, New York.
- 326 Druet T and Georges M 2010. Hidden Markov model combining linkage and linkage disequilibrium  
327 information for haplotype reconstruction and quantitative trait locus fine mapping. *Genetics* 184,  
328 789-798.
- 329 Habier D, Fernando RL and Dekkers JCM 2007. The impact of genetic relationship information on  
330 genome-assisted breeding values. *Genetics* 177, 2389-2397.
- 331 Habier D, Fernando RL and Dekkers JCM 2009. Genomic selection using low-density marker  
332 panels. *Genetics* 182, 343-353.
- 333 Hayes BJ and Goddard M E 2001. The distribution of the effects of genes affecting quantitative  
334 traits in livestock. *Genetics Selection Evolution* 33, 209-229.
- 335 Hayes BJ and Goddard ME 2008. Technical note: prediction of breeding values using marker-  
336 derived relationship matrices. *Journal of Animal Science* 86, 2089-2092.
- 337 Hoeskuldsson A 1988. Partial least squares PLS methods. *Journal of Chemometrics* 88, 211-228.
- 338 Hubert M and Branden KV 2003. Robust methods for partial least squares regression. *Journal of*  
339 *Chemometrics* 17, 537-549.

340 Lund M S, Sahana D, De Koning DJ, Su G and Carlborg Ö 2009. Comparison of analyses of  
341 QTLMAS XII common dataset. I: genomic selection. BMC proceedings 3 (suppl. 1), S1.

342 Macciotta NPP, Dimauro C, Bacciu N, Fresi P and Cappio-Borlino A 2006. Use of a partial least-  
343 squares regression model to predict test day of milk, fat and protein yields in dairy goats. Animal  
344 Science 82, 463-468.

345 Macciotta NPP, Gaspa G, Steri R, Nicolazzi E, Dimauro C, Pieramati C and Cappio-Borlino A  
346 2010. Use of principal component analysis to reduce the number of predictor variables in the  
347 estimation of Genomic Breeding Values. Journal of Dairy Science 93, 2765-2774.

348 Meuwissen THE, Hayes BJ and Goddard ME 2001. Prediction of total genetic values using  
349 genome-wide dense marker maps. Genetics 157, 1819-1829.

350 Solberg TR, Sonesson AK, Woolliams J and Meuwissen THE 2009. Reducing dimensionality for  
351 prediction of genome-wide breeding values. Genetics Selection Evolution 41, 29.

352 VanRaden PM, Van Tassell CP, Wiggans GR, Sonstengard TS, Schnabel RD *et al* 2009.  
353 Reliability of genomic predictions for north American Holstein bulls. Journal of Dairy Science 92,  
354 4414-4423.

355 Weigel KA, De Los Campos G, González-Recio O, Naya H, Wu L, Long N, Rosa GJ and Gianola,  
356 D 2009. Predictive ability of direct genomic values for lifetime net merit of Holstein sires using  
357 selected subsets of single nucleotide polymorphism markers. Journal of Dairy Science 92, 5248-  
358 5257.

359 Weigel KA, Van Tassell CP, O'Connell JR, VanRaden PM and Wiggans GR 2010. Prediction of  
360 unobserved single nucleotide polymorphism genotypes of Jersey cattle using reference panels and  
361 population-based imputation algorithms. Journal of Dairy Science 93, 2229-2238.

362 Wold S, Michael Sjöström M, Eriksson L 2001. PLS-regression: a basic tool of chemometrics.  
363 Chemometrics and Intelligent Laboratory Systems 58,109–130.



366 **Table 1** Mean correlations (and related standard deviations) between 100 actual and predicted

367 SNP in each chromosome

Missing SNP position	Correlations	
	Mean	St. Dev.
First 100	0.57	0.17
Middle 100	0.75	0.11
Last 100	0.68	0.14
One every 10	0.93	0.09

368

369

370 **Table 2** Mean correlations (and related standard deviations) between actual and predicted SNP for  
371 increasing percentage of predicted SNP. Proportions of correct SNP prediction are also reported

Percentage of predicted SNP	Correlations		Proportion of correct SNP prediction
	Mean	St. Dev.	
10%	0.98	0.07	0.99
25%	0.98	0.07	0.99
50%	0.97	0.08	0.98
75%	0.92	0.08	0.95
90%	0.78	0.13	0.86

372

373

374

375 **Table 3** *GEBV accuracies for different ratio of available/predicted SNP.*

Real SNP	Predicted SNP	GEBV accuracy
100%	0%	0.76
75%	25%	0.76
50%	50%	0.76
25%	75%	0.75
10%	90%	0.74

376

377

378

379

380

381 Figure captions:

382

383 **Figure 1** Pattern of the mean correlations between actual and predicted SNP for increasing  
384 number of extracted factors during the PLSR procedure

385

386 **Figure 2** Mean correlations between actual and predicted SNP for different numbers of fully  
387 genotyped animals

388

389