



**Application of SNP reduction approaches and random forest for the identification of population informative markers in cosmopolitan and local cattle breeds**

Journal:	<i>Italian Journal of Animal Science</i>
Manuscript ID	TJAS-2017-0288
Manuscript Type:	Abstract Submission
Date Submitted by the Author:	27-Feb-2017
Complete List of Authors:	Bertolini, Francesca; University of Bologna, Department of Agricultural and Food Sciences Galimberti, Giuliano; University of Bologna MASTRANGELO, Salvatore; Università degli Studi di Palermo, Scienze Agrarie e Forestali Di Gerlando, Rosalia; Università degli Studi di Palermo, Scienze Agrarie e Forestali Bagnato, Alessandro; Università degli Studi di Milano, Strillacci, Maria; University of Milano Portolano, Baldassare; University of Palermo Fontanesi, Luca; University of Bologna, Dept. of Agricultural and Food Sciences - Division of Animal Sciences
Abstract:	In livestock, single nucleotide polymorphism genotyping arrays have been used to differentiate breeds and populations for several downstream applications, including breed allocation of individuals, breeds of origin of crossbred animals, authentication of mono breed products, comparative analyses of selection signatures among several other uses. We already tested a combination of principal component analysis (PCA), used as pre-selection method, and random forest (RF) used as classification method to

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

	<p>assign cosmopolitan Italian breeds with no or very low error rate. In this work, we increased the number of breeds and approaches, to have a more comprehensive view of the strategies available and the applicability to local Italian breeds. The most common cosmopolitan dairy or dual purpose breeds (Holstein, Brown, Simmental) and 3 local breeds subjected to limited or no breeding programs (Reggiana, Modicana and Cinisara) were analyzed comparing several methods of SNPs pre-selection (Delta, Fst and PCA) in addition to RF classifications. From these classifications, two panels of 96 and 48 SNPs that contained the most discriminant SNPs were created for each pre-selection method. The results showed that the 96-SNP panels were generally more able to discriminate all breeds, while for the 48- SNP panels the error rate increased mainly for autochthonous breeds, particularly for Cinisara. This was probably a consequence of limited selection pressure, admixed origin, and ascertain bias on the construction of the SNP chip that was not designed considering these breeds. Several selected SNPs are located nearby genes affecting breed-specific traits (e.g. coat color and stature) or associated to production traits. The 96-SNP panel obtained after a preselection chromosome by chromosome, and used in the previous work with cosmopolitan breeds only, could identify informative SNPs that were particularly useful for the assignment of minor breeds. This panel reached the lowest value of out of bag (OOB) error in the RF test even in the Cinisara, whose value was quite high in all other panels. Moreover, this panel contained also the lowest number of SNPs in linkage disequilibrium. Our results showed the usefulness and power of the combination of PCA pre-selection and RF also for the discrimination of local cattle breeds.</p> <p>Acknowledgements This work was funded by Innovagen MiPAAF and by PON01_02249 - R&amp;C 2007-2013 - MIUR projects.</p>



# Application of SNP reduction approaches and random forest for the identification of population informative markers in cosmopolitan and local cattle breeds

Francesca Bertolini<sup>1,2</sup>, Giuliano Galimberti<sup>3</sup>, Salvatore Mastrangelo<sup>4</sup>, Rosalia Di Gerlando<sup>4</sup>, Maria Giuseppina Strillacci<sup>5</sup>, Alessandro Bagnato<sup>5</sup>, Baldassare Portolano<sup>4</sup>, Luca Fontanesi<sup>1</sup>

<sup>1</sup>*Dipartimento di Scienze e Tecnologie Agroalimentari, Università di Bologna, Bologna, Italy*

<sup>2</sup>*Department of Animal Science, Iowa State University, Ames, Iowa, USA*

<sup>3</sup>*Dipartimento di Scienze Statistiche "Paolo Fortunati", Università di Bologna, Italy*

<sup>4</sup>*Dipartimento Scienze Agrarie e Forestali, Università di Palermo, Palermo, Italy*

<sup>5</sup>*Dipartimento di Medicina Veterinaria, Università degli Studi di Milano, Milano, Italy*

Corresponding author: [luca.fontanesi@unibo.it](mailto:luca.fontanesi@unibo.it)

In livestock, single nucleotide polymorphism genotyping arrays have been used to differentiate breeds and populations for several downstream applications, including breed allocation of individuals, breeds of origin of crossbred animals, authentication of mono breed products, comparative analyses of selection signatures among several other uses. We already tested a combination of principal component analysis (PCA), used as pre-selection method, and random forest (RF) used as classification method to assign cosmopolitan Italian breeds with no or very low error rate. In this work, we increased the number of breeds and approaches, to have a more comprehensive view of the strategies available and the applicability to local Italian breeds. The most common cosmopolitan dairy or dual purpose breeds (Holstein, Brown, Simmental) and 3 local breeds subjected to limited or no breeding programs (Reggiana, Modicana and Cinisara) were analyzed comparing several methods of SNPs pre-selection (Delta,  $F_{st}$  and PCA) in addition to RF classifications. From these classifications, two panels of 96 and 48 SNPs that contained the most discriminant SNPs were created for each pre-selection method. The results showed that the 96-SNP panels were generally more able to discriminate all breeds, while for the 48-SNP panels the error rate increased mainly for autochthonous breeds, particularly for Cinisara. This was probably a consequence of limited selection pressure, admixed origin, and ascertain bias on the construction of the SNP chip that was not designed considering these breeds. Several selected SNPs are located nearby genes affecting breed-specific traits (e.g. coat color and stature) or associated to production traits. The 96-SNP panel obtained after a preselection chromosome by chromosome, and used in the previous work with cosmopolitan breeds only, could identify informative SNPs that were particularly useful for the assignment of minor breeds. This panel reached the lowest value of out of bag (OOB) error in the RF test even in the Cinisara, whose value was quite high in all other panels. Moreover, this panel contained also the lowest number of SNPs in linkage disequilibrium. Our results showed the usefulness and power of the combination of PCA pre-selection and RF also for the discrimination of local cattle breeds.

## Acknowledgements

This work was funded by Innovagen MiPAAF and by PON01\_02249 - R&C 2007-2013 - MIUR projects.