

1 Simultaneous SNP selection and adjustment for
2 population structure in high dimensional prediction
3 models

4 Sahir R Bhatnagar^{1,2}, Yi Yang⁴, Tianyuan Lu², Erwin Schurr⁶,
5 JC Lored-Osti⁷, Marie Forest², Karim Oualkacha³, and
6 Celia MT Greenwood^{1,2,5}

7 ¹Department of Epidemiology, Biostatistics and Occupational Health,
8 McGill University

9 ²Lady Davis Institute, Jewish General Hospital, Montréal, QC

10 ³Département de Mathématiques, Université de Québec À Montréal

11 ⁴Department of Mathematics and Statistics, McGill University

12 ⁵Departments of Oncology and Human Genetics, McGill University

13 ⁶Department of Medicine, McGill University

14 ⁷Department of Mathematics and Statistics, Memorial University

15 August 19, 2019

16 **Abstract**

17 Complex traits are known to be influenced by a combination of environmental fac-

tors and rare and common genetic variants. However, detection of such multivariate associations can be compromised by low statistical power and confounding by population structure. Linear mixed effects models (LMM) can account for correlations due to relatedness but have not been applicable in high-dimensional (HD) settings where the number of fixed effect predictors greatly exceeds the number of samples. False positives or false negatives can result from two-stage approaches, where the residuals estimated from a null model adjusted for the subjects' relationship structure are subsequently used as the response in a standard penalized regression model. To overcome these challenges, we develop a general penalized LMM framework called `ggmix` for simultaneous SNP selection and adjustment for population structure in high dimensional prediction models. Our method can accommodate several sparsity-inducing penalties such as the lasso, elastic net and group lasso, and also readily handles prior annotation information in the form of weights. We develop a blockwise coordinate descent algorithm which is highly scalable, computationally efficient and has theoretical guarantees of convergence. Through simulations and two real data examples, we show that `ggmix` leads to better sensitivity and specificity compared to the two-stage approach or principal component adjustment with better prediction accuracy. `ggmix` can be used to construct polygenic risk scores and select instrumental variables in Mendelian randomization studies. Our algorithms are available in an R package (<https://github.com/greenwoodlab/ggmix>).

1 Author Summary