

RubricOE: a learning framework for genetic epidemiology

Subrata Saha^{1,†}, Aldo Guzmán-Sáenz^{2,†}, Aritra Bose^{2,†}, Filippo Utro², Daniel E. Platt²,
and Laxmi Parida^{2,*}

¹Department of Systems Biology, Columbia University Medical Center, New York, NY 10032

²Computational Genomics, IBM T.J Watson Research Center, Yorktown Heights, NY 10598

[†]Equal contribution

*Corresponding author: parida@us.ibm.com

Running Title: RubricOE

Keywords: Association Studies; Machine Learning; Polygenic Risk Score

1 Abstract

2 Genetic epidemiology is a growing area of interest in the past years due to the availability of genetic
3 data with the decreasing cost of sequencing. Machine learning (ML) algorithms can be a very
4 useful tool to study the genetic factors on disease incidence or on different traits characterizing a
5 population. There are many challenges that plagues the field of genetic epidemiology including
6 the unbalanced case-control data sets, fallibility of standard genome wide association studies with
7 single marker analysis, heavily underdetermined systems with millions of markers in contrast of
8 a few thousands of samples, to name a few. Ensemble ML methods can be a very useful tool to
9 tackle many of these challenges and thus we propose RubricOE, a pipeline of ML algorithms with
10 error bar computations to obtain interpretable genetic and non-genetic features from genomic
11 or transcriptomic data combined with clinical factors in the form of electronic health records.
12 RubricOE is shown to be robust in simulation studies, detecting true associations with traits of
13 interest in arbitrarily structured multi-ethnic populations.

14 Introduction

15 In recent years, with the availability of large cohorts for complex diseases, identifying the asso-
16 ciated genetic variants as well as gene-environment interactions has become more challenging.
17 Genome-wide association studies (GWAS) detect associations between genetic variants (single
18 nucleotide polymorphism or SNP) and traits is a powerful tool but is sensitive to sample sizes,
19 population structure, rare variants, environmental factors, etc. However, GWAS has several chal-
20 lenges in disentangling environmental effect from true genetic associations of a disorder leading
21 to spurious associations. It also suffers from the curse of dimensionality with the availability of
22 millions of SNPs for several thousands individuals, leading to the “small n , large p problems”
23 (due to the under-determined nature of the genotype matrix). With the availability of more pop-
24 ulation controls than cases for binary traits due to low prevalence of many conditions/diseases
25 the genotype data sets are often unbalanced. The scale of the data as well as the unbalanced
26 nature of binary traits pose a substantial challenge for GWAS in large cohorts [1]. The SNPs are
27 often genotyped with linkage disequilibrium (LD) between them resulting in correlated variables,
28 which needs to be taken into account to find true genetic associations. Therefore, standard multi-

29 variable statistical approaches like linear or logistic regression are not well suited for genome-wide
30 data [2]. Due to these shortcomings there has been a need for sophisticated machine learning
31 approaches.

32 Machine learning (ML) algorithms can learn from the data without making any model as-
33 sumptions. It provides several alternatives for performing multi-SNP analyses away from the
34 single marker analysis as done in GWAS. It can employ techniques such as regularization and
35 cross-validation in regression to tackle overfitting issues for under-determined data sets. Ma-
36 chine learning methods have been applied on a broad range of problems in genomics including
37 how to recognize locations of transcriptions start sites, promoters, enhancers, etc in a genome
38 sequence [3]. It has also been employed for genomic selection in plant breeding [4] as well as in
39 GWAS [2] and genetic epidemiology [5]. ML can also elucidate non-linear SNP-SNP epistatic
40 interactions, such as testing for N^2 SNP-SNP combinations for N SNPs. Generally, with multi-
41 ple hypothesis Bonferroni corrected thresholds, this can be a difficult task with linear regression
42 based techniques. For binary traits of complex diseases, the challenge is to identify clusters of
43 alleles that suggest interactive pathways which are potential therapeutic targets.

44 In order to do that, it is important to understand how error analysis and statistical power
45 interact in ML as probes to characterize identified candidates for pathogenic SNPs or other
46 features. Therefore, we present a ML pipeline, called RubricOE, namely, a rubric for multi-omics
47 association studies and genetic epidemiology. It includes main stages for feature selection, and
48 for cross-validation applied to the selected feature set from the first stage. We consider the effect
49 of error analysis to feature selection and cross validation. Lastly, it also considers the possibility
50 of identifying SNPs that are purely epistatic. We find that RubricOE can correctly identify true
51 associations in simulated data. It observes a trade-off with genetic effect in simulated datasets
52 leading to minimal spurious associations for increased genetic effect and vice versa.

53 RubricOE is a robust ML pipeline that ranks disease associated SNPs and predicts the out-
54 come of different disorder or traits. It provides further interpretation of the selected associations
55 with their ranks and scores for studying their relative importance. The ranks and scores of SNPs
56 closely imitates Polygenic Risk Scores (PRS) and thus can be a very useful tool for studying her-
57 itability of different traits. RubricOE is less sensitive to noise than standard GWAS and provides
58 different optimizations in selecting hyper-parameters further elucidating the importance of each

59 selected SNP. It also provides a framework to integrate clinical Electronic Health Record (EHR)
60 information along with genomic data, providing further information on interactions between ge-
61 netic and non-genetic factors associated with a disease or trait.

62 Results

63 Simulated Data

64 We applied RubricOE on six different simulation data sets with 1,000 individuals across 10,000
65 SNPs, related to two different phenotype distribution cases for each of the three population
66 structure simulation scenarios related to genotypes. We observe that in the most challenging
67 scenario with very low genetic effect (10%) consisting of only ten causal SNPs, RubricOE detects
68 30% of causal SNPs in the real-world TGP scenario. The Youden’s J statistic on the validation
69 data in this case is understandably very low due to the overwhelming amount of noise in the
70 dataset. Alternatively, when the genetic effect (70%) is substantially more than environmental
71 effect and noise, RubricOE is able to discover accurately around 90% of the true causal SNPs
72 in PSD and 70% in TGP, respectively. The Youden’s J statistic are also relatively higher with
73 around 0.45 in PSD model, even though there are only ten true causal SNPs contributing to the
binary phenotype. The score curves of the above experiments show that the peak scores of each

	BN		PSD		TGP	
	Causal SNPs	J statistic	Causal SNPs	J statistic	Causal SNPs	J statistic
10-20-70	0	0.038	10%	0.056	30%	0.061
70-20-10	70%	0.505	90%	0.446	70%	0.261

Table 1: RubricOE performance on validation data in simulation studies evaluated by the percentage of causal SNPs detected and the mean Youden’s J statistic observed.

74
75 feature in RubricOE is reached by using a very small set of features. As RubricOE computes the
76 rank of the features by the descending order of their score and then computes Youden’s J statistic
77 iteratively on the dataset, it gives an opportunity to understand how much a subset of feature
78 contributes to the total classification accuracy of the learning algorithm. Thus, we observe, in
79 case of PSD which accurately identifies 90% of causal associations, it needs only the first few
80 hundred features to reach its peak cumulative score (Figures 4-6 in Appendix), deeming the rest
81 of the SNPs non-informative. We see similar patterns in other simulation scenarios such as TGP

82 and BN as well. We can leverage this information in making RubricOE faster by asking it to
83 compute curves only with a few thousand features instead of the large number of SNPs, genes and
84 clinical features, sometimes amounting to a few millions. This provides a significant scale-up in
85 computational time of the pipeline. As SVM has quadratic complexity, this feature of RubricOE
86 is particularly useful.

87 **Discussion**

88 RubricOE provides a ML framework with which a stable set of features discriminating between the
89 healthy controls and diseased patients, can be obtained. This differs from a single marker analysis
90 as performed in GWAS. ML algorithms allow multi-SNP analysis along with interpretation of
91 the subset of SNPs classifying different diseases. It ranks the SNPs along with their relative
92 scores as calculated by Youden’s J statistic. One of the advantages of ML algorithms is in
93 applying criteria for accepting features that focus on aggregated predictive power. For regression,
94 the fraction of variation predicted by a feature is closely related, and for standard regression
95 equal, to the square of the z-score for the coefficient. The downside is that large models may be
96 underdetermined, and feature selection may identify false features due to overfitting. Regulation
97 (Lasso and Ridge) mitigates that, but complicates the relationship between feature significance
98 and how much variation is accounted for by the feature. PRS is comprised by most of the
99 features of ML algorithms explored here, including susceptibility to overfitting. So the diffusion
100 of functional associations may be spurious.

101 Associations can be spurious due to environmental effects as well as genotyping errors, etc.
102 We explored robustness of RubricOE to such effects in the simulation studies and observed that
103 the framework was able to correctly detect about 70% of true causal associations in the simulation
104 scenarios closely emulating real-world population structure when the genetic effect is significantly
105 more than non-genetic factors. RubricOE is also able to detect about 20-30% of causal associations
106 when the genetic effect is only 10%, among which only 0.1% is causal. This demonstrates the
107 prowess of RubricOE in detecting rare associations in the presence of noise.

108 The estimate of an individual’s genetic liability to a trait or disease of interest is captured by
109 the PRS computation. It returns a score which is parallel to the Youden’s J statistic described
110 here. Another factor that RubricOE can provide a solution to is the infamous “missing heritabil-

ity” issue. Single marker analysis as done in GWAS does not account for the heritability of a trait or disease across generations. PRS aims to solve that problem, but, RubricOE provides a straight-forward one-pass solution for doing multi-SNP tests. One of the advantages of ML algorithms is in applying criteria for accepting features that focus on aggregated predictive power. For regression, the fraction of variation predicted by a feature is closely related, and for standard regression equal, to the square of the z-score for the coefficient. The downside is that large models may be underdetermined, and feature selection may identify false features due to overfitting. Regularization (Lasso and Ridge) mitigates that, but complicates the relationship between feature significance and how much variation is accounted for by the feature. We showed that, for random sampling variation, cross-validation does not really boost power. However, cross-validation is an effective way to recognize biased variations (e.g. batch effects). PRS is comprised by most of the features of machine learning algorithms explored here, including susceptibility to overfitting. So the diffusion of functional associations may be spurious. A detailed study on the parallels between PRS and Youden’s J statistic as well as demonstrate how it approaches the “missing heritability” issue across multi-ethnic cohorts would highlight the importance of studying other ML algorithms such as neural network or deep learning in disorder trait prediction and interpretable feature selection. RubricOE, or a rubric for omics and epidemiology can be useful in stable feature selection from integrated multi-omic and clinical data, elucidating their interactions.

129 **Methods**

130 **Datasets**

131 **Simulation Study**

132 We generated an extensive set of simulations with real-world challenging scenarios to demonstrate
133 the robustness to different scenario and power to detect true disorder or trait associated SNPs. We
134 simulated and analyzed six different data sets with 1,000 individuals and 10,000 SNPs, pertaining
135 to three population structure scenarios and for each, two phenotype simulation scenarios. The
136 data was simulated based on a binary trait model using the Odds Ratio (OR) as the classifier for

137 disease status from a quantitative trait model described in previous work [6].

$$y_j = \alpha + \sum_{i=1}^m \beta_i \mathbf{X}_{i,j} + \lambda_j + \epsilon_j \quad (1)$$

138 where β_i is the genetic effect of SNP i on the trait, λ_j is the random non-genetic effect and ϵ_j
 139 is the random noise variation for individual j . $\mathbf{X}_{i,j}$ is the i^{th} marker for the j^{th} individual and
 140 $y \in \mathbb{R}^m$ is the trait response variable (binary or continuous).

141 For the genotype data, we simulated allele frequencies using (i) Balding-Nichols (BN) model [7]
 142 based on allele-frequency and F_{ST} estimates calculated on the HapMap data set, (ii) three dif-
 143 ferent levels of admixture by varying the parameter α from $\{0.01, 0.1, 0.5\}$ in Pritchard-Stephens-
 144 Donnelly model (PSD) [8] and (iii) structure estimated from 1000 Genomes Project (TGP) [9].
 The three scenarios allow us to simulate genotypes with varying degree of population struc-

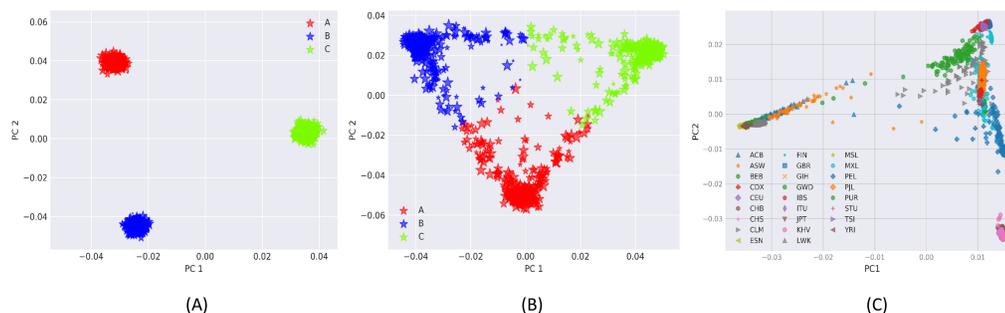


Figure 1: Projection of the samples from three populations simulated from (i) BN (ii) PSD ($\alpha = \{0.1, 0.1, 0.1\}$) and (iii) TGP model on the top two axes of variation.

145
 146 ture [6, 10]. BN provides three arbitrarily structured group of populations who are not mixing
 147 with each other. This provides an ideal case of structure to evaluate the model. PSD model
 148 accounts for admixture between these arbitrarily structured clusters of individuals (such as in-
 149 dividuals between each group share alleles between them) and the degree of admixture can be
 150 altered by the parameter α . Lastly, TGP model provides the most real-world like scenario where
 151 the individuals are sampled from allele frequencies corresponding to different population across
 152 the world as represented in the 1000 Genomes data set [9].

153 We used different variance ratio for the phenotype simulation scenarios. For the first scenario,
 154 we allowed 70% genetic effect, 20% environmental effect and 10% noise effect. This allowed us

155 to study whether the model can accurately detect the true associations when the genetic effect is
156 substantial. In comparison, we reversed the variances of the genetic effect (10%) and noise (70%)
157 to simulate a more challenging scenario. The causal genetic effect is simulated by independently
158 simulating the first 0.001 of the total number of SNPs, n (10,000 SNPs). $\beta_i \sim \text{Normal}(0, 0.5)$ for
159 $i = 1, 2, \dots, 0.001 * n$. For all other $i > 0.001 * n$ we set $\beta_i = 0$.

160 **RubricOE**

161 The ML pipeline, “RubricOE” outlined in Figure 2, has 3 major components: Quality Control
162 (QC) of input omics data, Iterative Feature Selection (IFS), and Stable Set Construction (SSC).
163 Optional selection of machine learning components with hyper-parameters may be employed in
164 each to optimize performance. The pipeline employs a nested test-training set configuration.
165 The outer test-training set split reserves the test set (“validation”) for final SNP evaluation
166 and denotes the training set as “working” data. Within the latter, further train-test splits are
167 performed to rank the features.

168 The working data is used within the IFS stage. In this stage, replications of data test-training
169 splits are each subject to iterative feature ranking, using machine learning algorithms (Ridge
170 Regression in this case), and scored against the replicated test sets using a scoring metric (e.g.
171 Youden or h^2) using another machine learning algorithm (Support Vector Machines or SVM with
172 linear kernel in this case) to identify a subset of SNPs. If the ML is non-linear in character,
173 feature selection is based on recursive feature elimination (RFE) [11] is an option we have used.
174 These rankings are combined to form a final candidate set. Models constructed from the final
175 set are applied to the “unseen” (test) set for final scoring. Iterative replications are performed
176 subsequently, and the set of SNPs that persistently survive in the intersection are called the
177 “Stable Set”.

178 **Iterative Feature Selection**

179 **Robust SNP ranking** Some of the ranking algorithms that were applied here, besides ap-
180 plication of simple ridge regression, are described here. We attempt to make our SNP ranking
181 model robust by employing multiple linear SVMs (LSVMs) in a concerted way. Each LSVM runs
182 on slightly different data sets randomly sampled from original dataset D . In the experiment, we

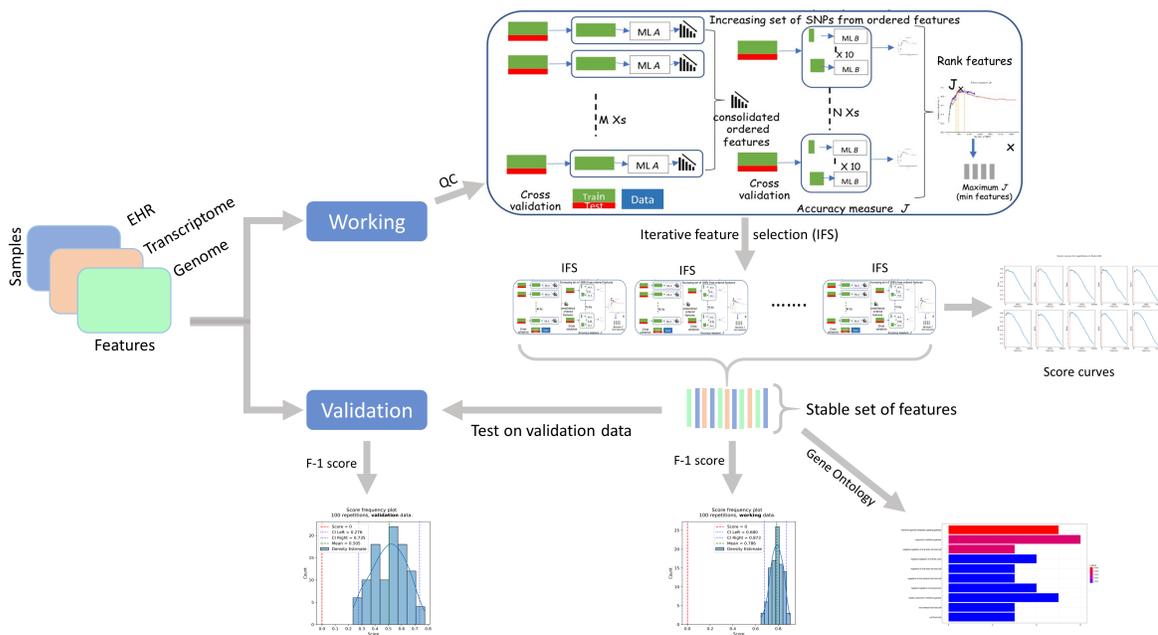


Figure 2: Components of RubricOE, the ensemble ML pipeline producing a stable set of features from multi-omics data.

183 randomly picked 90% unique samples from D for each LSVM to build a robust learning model.
 184 Due to this small change, the weight associated with each SNP is not identical for different runs.
 185 By averaging the weights across a set of weight vectors from LSVMs of the same configuration,
 186 we can make the weight vector robust. Since weights are directly associated with the importance
 187 of the SNPs, ranking should also be robust.

188 Let the number of LSVMs to rank a set S of SNPs of interest be L . Each SVM $l \in L$ is
 189 trained on a slightly smaller dataset D' to build an inductive model and we get a corresponding
 190 weight vector w_l ($1 \leq l \leq |L|$). According to [11], the significance of the i^{th} SNP $I_i = (w_l^i)^2$ where
 191 w_l^i represents the i^{th} weight component of w_l^{th} weight vector. At the end of this step, we get a set
 192 of $|L|$ weight vectors. Since, in each run, we introduce slightly different dataset D' by randomly
 193 sampling the original dataset D , the weight vectors will be different from each other. To make it
 194 stable we normalize each weight vector and average each component. Each weight vector w_l will
 195 have n components and is normalized as follows: $w_l' = \frac{w_l}{\sum_{i=1}^n |w_l^i|} \dots (2)$. The i^{th} component of the
 196 final weight vector W is formed as follows: $W^i = \sum_{l=1}^{|L|} |(w_l^i)| \dots (3)$. Now, the importance of the
 197 i^{th} SNP is defined as $I_i = \frac{W^i}{n}$. We sort the SNPs based on their significance in non-decreasing
 198 order. So, the highest significant SNP will have lowest rank, the 2^{nd} highest significant SNP will

199 have 2^{nd} lowest rank, and so forth.

200 **Robust SNP subset selection** Finding optimal set of features from n features will take $\mathcal{O}(2^n)$
201 time which is exponential in the number of features. To reduce the search space, we follow two
202 step procedure. At first, we employ robust SNP ranking algorithm (as described above) to rank
203 SNPs based on their significance. In the second step, we linearly search through the ranked SNP
204 space to get the subset of SNPs maximizing the classification accuracy (such as Youden index).
205 Next, we describe the procedure in detail.

206 Let $r_0, r_1, r_2, \dots, r_{n-1}$ be the ranks of the given SNPs where r_0 is the smallest rank, r_1 is the
207 second smallest rank, etc. Next, we take top x SNPs (i.e. r_0, r_1, \dots, r_x where $x \ll n$) and compute
208 Youden index using 10-*fold* cross validation based on those x SNPs. Subsequently, we compute
209 Youden indices based on top $2x$ SNPs, top $3x$ SNPs, etc. As soon as there is no improvement
210 over the maximum Youden index seen so far in succeeding iterations, we stop the linear search
211 and return back those SNPs having maximum Youden index. These SNPs constitute our robust
212 subset of SNPs.

213 **Stable Set Construction**

214 Finally, the whole procedure (i.e. SNP ranking and subset selection) is repeated multiple times to
215 obtain stable set of features. For an illustrative example, let's assume we run the entire procedure
216 p times. Consequently, we get p rankings of SNPs and p subsets of robust SNPs (through ranking
217 and subset selection procedure). Finally, we extract common SNPs among the p subsets. Suppose,
218 t_i represent a robust subset of SNPs from i^{th} repetition where $1 \leq i \leq p$. Then the stable subset
219 of SNPs will be $T = t_1 \cap t_2 \cap t_3 \cap \dots \cap t_p$.

220 **Cross Validation**

221 Cross validation splits data into a training, or discovery, set D used to determine a set of SNPs P_D
222 that passes a threshold predicting a phenotype, and a test set V independent of the training set,
223 used to test, or validate, the predictive power of the test P_V on the training set's best variants.
224 This is used in the "Feature Set Selection" and "Stable Set Construction" stages of our nested
225 training-test set configuration. Therefore, a major component of power analysis of these methods
226 revolve around evaluating factors impacting cross validation.

227 The discovery set parameters are $\mathbb{P}(P_D|\bar{A}) = \alpha_D$ and $\mathbb{P}_{\mathbb{D}}(P_D|A) = 1 - \beta_D$, and the validation
 228 set parameters are $\mathbb{P}(P_V|\bar{A}) = \alpha_V$ and $\mathbb{P}_{\mathbb{V}}(P_V|A) = 1 - \beta_V$. $\mathbb{P}(A) = f$ in both groups. P_D and
 229 P_V are assumed to be independent. That is $\mathbb{P}(P_V|P_D) = \mathbb{P}(P_V)$, etc, and $\mathbb{P}(P_V \cap A|P_D \cap A) =$
 230 $\mathbb{P}(P_V|\cap A)$, Following arguments above,

$$\mathbb{P}(P_D \cap P_V) = (1 - \beta_D)(1 - \beta_V)f + \alpha_D\alpha_B(1 - f)$$

$$\mathbb{P}(P_D \cap P_{\bar{V}}) = (1 - \beta_D)\beta_V f + \alpha_D(1 - \alpha_B)(1 - f)$$

$$\mathbb{P}(P_{\bar{D}} \cap P_V) = \beta_D(1 - \beta_V)f + (1 - \alpha_D)\alpha_B(1 - f)$$

$$\mathbb{P}(P_{\bar{D}} \cap P_{\bar{V}}) = \beta_D\beta_V f + (1 - \alpha_D)(1 - \alpha_B)(1 - f)$$

231 These relations satisfy $\mathbb{P}(P_D \cap P_V) + \mathbb{P}(P_D \cap P_{\bar{V}}) + \mathbb{P}(P_{\bar{D}} \cap P_V) + \mathbb{P}(P_{\bar{D}} \cap P_{\bar{V}}) = 1$.

232 We apply the above to estimate numbers of SNPs expected to be shared between training and
 233 test sets in some of the steps of our ML algorithm.

An odds ratio may be defined

$$r = \frac{\mathbb{P}(P_D \cap P_V)\mathbb{P}(\bar{P}_D \cap \bar{P}_V)}{\mathbb{P}(P_D \cap \bar{P}_V)\mathbb{P}(\bar{P}_D \cap P_V)}$$

234 If there were no biologically active SNPs, then $f = 0$, and $r = 1$. So this can test whether the
 235 number of SNPs in both test and training sets is more than expected by chance.

236 Model selection

237 We use two different ML algorithms for the two steps outlined in the RubricOE pipeline (Figure 2),
 238 IFS and SSC. We use Ridge Regression for ranking the features in IFS and LSVM to score the
 239 Youden curves for optimization.

240 One standard option for selecting the initial feature ranking has been to use regression coeffi-
 241 cients. In our approach, one method for ranking in the ‘IFS stage is ridge regression ranking by
 242 magnitude of the coefficients. Another alternative was to consider the variation in the phenotype
 243 covered by the coefficient.

In standard regression, given $\text{cov}(y, y^T) = \Sigma$, a χ^2 distributed statistic measuring goodness of fit of a model $y = Xb$ is

$$\mathcal{E}^2 = (y - Xb)^T \Sigma^{-1} (y - Xb).$$

This may be rewritten

$$\mathcal{E}^2 = \mathcal{E}_{res}^2 + (b - \beta)^T (X^T \Sigma^{-1} X) (b - \beta)$$

where

$$\mathcal{E}_{res}^2 = y^T \Sigma^{-1} y - (X\beta)^T \Sigma^{-1} (X\beta)$$

with

$$\beta = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} y,$$

and

$$\text{cov}(\beta, \beta^T) = (X^T \Sigma^{-1} X)^{-1}.$$

Given this, it is notable that the residual may be written

$$\mathcal{E}_{res}^2 = y^T \Sigma^{-1} y - \beta^T (\text{cov}(\beta, \beta^T))^{-1} \beta,$$

244 noting that X and b were augmented by a feature identical to unity in X , and equal to an
 245 offset in b . This gives a measure of variability to the variability in the mean for y . Therefore,
 246 the multivariate equivalent of sums of squares of the z-scores, including effects of correlations,
 247 represents the contribution of the predicted values to the residual. We identify the unexplained
 248 variability to be $h^2 = \mathcal{E}_{res}^2$, with the proportion explained by individual features h_i^2 for β_i as
 249 $h_i^2 = \beta_i^T (\text{cov}(\beta, \beta^T))^{-1} \beta_i$ for each component by itself. We note that $y^T \Sigma^{-1} y$ is not centered.
 250 The contribution from offsets to the variability are included in the offset variable included in the
 251 augmented β 's offset estimation. Centering y shifts the offset, but does not change the coefficients.
 252 The residual error \mathcal{E}_{res}^2 represents the proportion of phenotype not predicted by the regression
 253 model, including genetics and adjustment.

Feature selection regressions tend to be underdetermined, so a regularization that preserves the form of \mathcal{E}_{res}^2 , and that provides error bars and covariances for the coefficients is desirable. This argues for an $L2$ regularization. In this case, it is assumed that the coefficients b in the regression are themselves distributed, and have terms with variance C^{-1} and mean 0. Then for $\text{cov}(y, y^T) = \Sigma$, and $\text{cov}(b, b^T) = C^{-1}I$ for a number C , a χ^2 distributed statistic measuring

goodness of fit of a model $y = Xb$ is

$$\mathcal{E}^2 = (y - Xb)^T \Sigma^{-1} (y - Xb) + Cb^T b.$$

This may be rewritten

$$\mathcal{E}^2 = \mathcal{E}_{res}^2 + (b - \beta)^T (X^T \Sigma^{-1} X + C \cdot I) (b - \beta),$$

where

$$\mathcal{E}_{res}^2 = y^T \Sigma^{-1} y - \beta^T (X^T \Sigma^{-1} X + C \cdot I) \beta,$$

with

$$\beta = (X^T \Sigma^{-1} X + C \cdot I)^{-1} X^T \Sigma^{-1} y,$$

and

$$\text{cov}(\beta, \beta^T) = (X^T \Sigma^{-1} X + C \cdot I)^{-1} X^T \Sigma^{-1} X (X^T \Sigma^{-1} X + C \cdot I)^{-1}.$$

In the above, the

$$(X^T \Sigma^{-1} X + C \cdot I)^{-1}$$

term induced by regularization tempers the response of β to variations in y . But the contribution to \mathcal{E}_{res}^2 that reflects the $\text{cov}(\beta, \beta^T)$ contribution is the $(X^T \Sigma^{-1} X)$, while the $C \cdot I$ does not reflect predictive information. Expanding, this becomes

$$\mathcal{E}_{res}^2 = y^T \Sigma^{-1} y - (X\beta)^T \Sigma^{-1} (X\beta) - C\beta^T \beta.$$

Given this, note that

$$(X\beta)^T \Sigma^{-1} (X\beta) \neq \beta^T (\text{cov}(\beta, \beta^T))^{-1} \beta$$

254 So the amount of variation due to the y variation accounted for in the residual is not equivalent
 255 to the z-score squared of the coefficients in the way that it is in non-regularized linear regression.
 256 Note that the part of the contribution to \mathcal{E}_{res}^2 that responds to y is $(X\beta)^T \Sigma^{-1} (X\beta)$. However,
 257 while the direct relationship between variability is lost in ridge regression, the expression capturing
 258 the dependency on y variation should be applied for ranking, using the $h_i^2 = (X\beta_i)^T \Sigma^{-1} (X\beta_i)$

259 as in standard regression described above, but with the ridge-regression estimated coefficients.
260 Note also that $X\beta$ is still the estimator for y , so $(X\beta)^T \Sigma^{-1} (X\beta)$ is still the component that
261 describes contributions to the proportion of variability in $y^T \Sigma^{-1} y$. We apply ridge regression to
262 feature selection as described above, computing the regression in the singular-value decomposition
263 (SVD) basis. The value of the regularization parameter C is chosen to be the median singular
264 value obtained from the SVD.

265 **Characterizing Stable Features**

266 Logistic regression is evaluated for the stable set applied to the training sets and to the final
267 test set to understand how cross-validated selection of final features relates to their estimated
268 uncertainties, and to understand more clearly the power analysis situation in this problem.

269 **Code Availability.** RubricOE executable is available at <https://github.com/ComputationalGenomics/>
270 RubricOE.

271 **Data Availability.** All simulated data was used in this manuscript.

272 **Acknowledgements.** This work was conducted by SS as part of his Postdoctoral Research at
273 IBM T.J Watson Research Center.

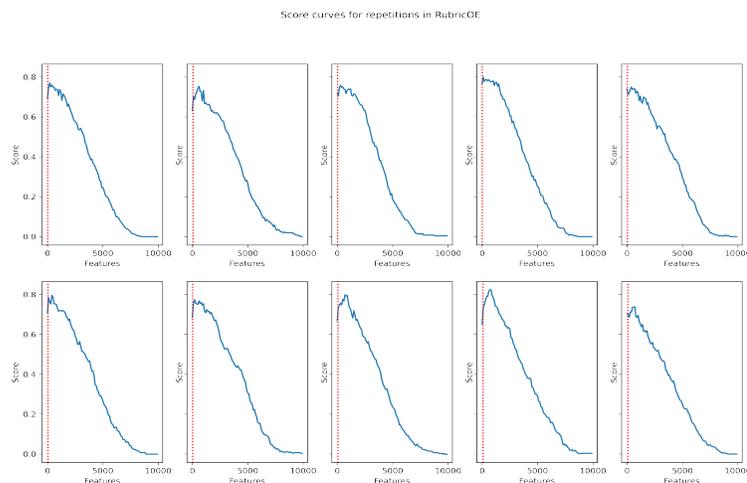
274 References

- 275 [1] Zhou, W., Nielsen, J. B., Fritsche, L. G., Dey, R., Gabrielsen, M. E., Wolford, B. N., LeFaive,
276 J., VandeHaar, P., Gagliano, S. A., Gifford, A., et al. (2018). Efficiently controlling for case-
277 control imbalance and sample relatedness in large-scale genetic association studies. *Nature*
278 *genetics* **50**(9), 1335–1341.
- 279 [2] Szymczak, S., Biernacka, J. M., Cordell, H. J., González-Recio, O., König, I. R., Zhang,
280 H., and Sun, Y. V. (2009). Machine learning in genome-wide association studies. *Genetic*
281 *epidemiology* **33**(S1), S51–S57.
- 282 [3] Libbrecht, M. W. and Noble, W. S. (2015). Machine learning applications in genetics and
283 genomics. *Nature Reviews Genetics* **16**(6), 321–332.
- 284 [4] González-Camacho, J. M., Ornella, L., Pérez-Rodríguez, P., Gianola, D., Dreisigacker, S.,
285 and Crossa, J. (2018). Applications of machine learning methods to genomic selection in
286 breeding wheat for rust resistance. *The plant genome* **11**(2), 170104.
- 287 [5] Dasgupta, A., Sun, Y. V., König, I. R., Bailey-Wilson, J. E., and Malley, J. D. (2011).
288 Brief review of regression-based and machine learning methods in genetic epidemiology: the
289 Genetic Analysis Workshop 17 experience. *Genetic epidemiology* **35**(S1), S5–S11.
- 290 [6] Song, M., Hao, W., and Storey, J. D. (2015). Testing for genetic associations in arbitrarily
291 structured populations. *Nature genetics* **47**(5), 550.
- 292 [7] Balding, D. J. and Nichols, R. A. (1995). A method for quantifying differentiation between
293 populations at multi-allelic loci and its implications for investigating identity and paternity.
294 *Genetica* **96**(1-2), 3–12.
- 295 [8] Pritchard, J. K., Stephens, M., and Donnelly, P. Jun (2000). Inference of population structure
296 using multilocus genotype data. *Genetics* **155**(2), 945–959. 10835412[pmid].
- 297 [9] Auton, A., Abecasis, G. R., Altshuler, D. M., Durbin, R. M., Abecasis, G. R., Bentley, D. R.,
298 Chakravarti, A., Clark, A. G., Donnelly, P., Eichler, E. E., et al. (2015). A global reference
299 for human genetic variation. *Nature* **526**(7571), 68–74.

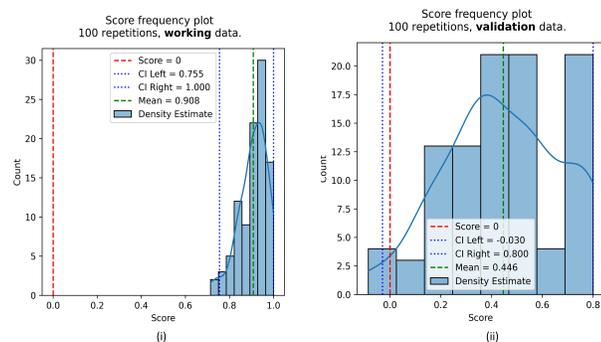
- 300 [10] Hao, W., Song, M., and Storey, J. D. (2015). Probabilistic models of genetic variation in
301 structured populations applied to global human studies. *Bioinformatics* **32**(5), 713–721.
- 302 [11] Guyon, I., Weston, J., Barnhill, S., and Vapnik, V. January (2002). Gene Selection for
303 Cancer Classification using Support Vector Machines. *Machine Learning* **46**(1), 389–422.

304 **Appendix**

305 **Simulated Data**

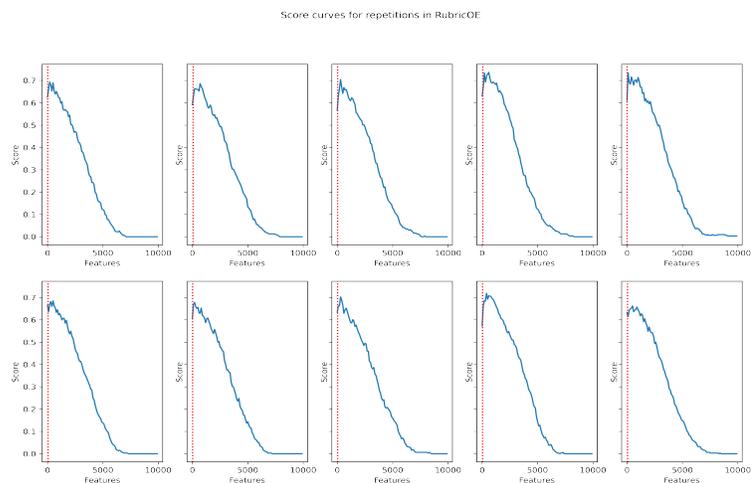


(a) Youden's J statistic score curves if IFS stage of RubricOE. The optimal number of SNPs to reach the highest peak cumulative score is marked by the red dotted line.

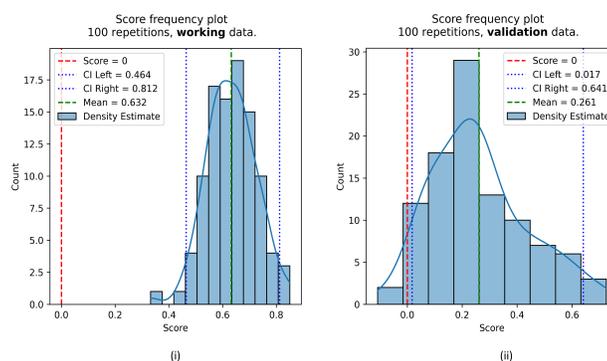


(b) Histograms of F-1 score on working and validation data.

Figure 3: Performance of RubricOE on PSD scenario with 70% genetic effect.

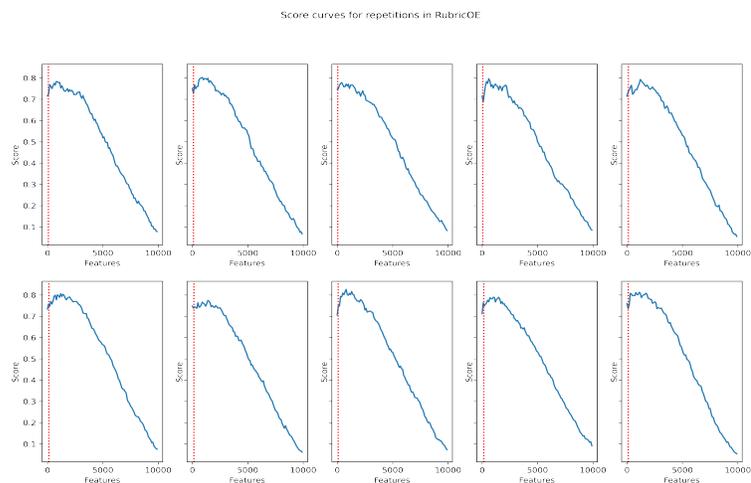


(a) Youden's J statistic score curves if IFS stage of RubricOE. The optimal number of SNPs to reach the highest peak cumulative score is marked by the red dotted line.

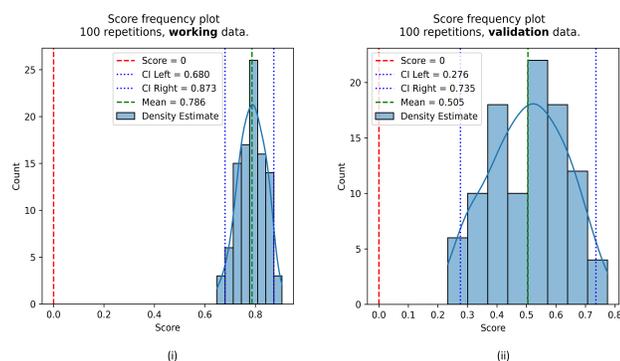


(b) Histograms of F-1 score on working and validation data.

Figure 4: Performance of RubricOE on TGP scenario with 70% genetic effect.



(a) Youden's J statistic score curves if IFS stage of RubricOE. The optimal number of SNPs to reach the highest peak cumulative score is marked by the red dotted line.



(b) Histograms of F-1 score on working and validation data.

Figure 5: Performance of RubricOE on BN scenario with 70% genetic effect.