bioRxiv preprint doi: https://doi.org/10.1101/2022.02.21.481356; this version posted February 22, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

# <sup>1</sup> Supervised Rank aggregation (SRA): A

<sup>2</sup> novel rank aggregation approach for

# ensemble-based feature selection

- 4 Rahi Jain<sup>1</sup>, Wei Xu<sup>2\*</sup>
- <sup>5</sup><sup>1</sup>Biostatistics Department, Princess Margaret Cancer Research Centre, Toronto, Ontario, Canada
- 6 <sup>2</sup>Dalla Lana School of Public Health, University of Toronto, Toronto, Ontario, Canada
- 7 <sup>\*</sup> Corresponding author
- 8 Telephone Number: (+1) 416-946-4497
- 9 Email: <u>wei.xu@</u>uhnres.utoronto.ca (WX)
- 10 Rahi Jain is a postdoctoral fellow at the Biostatistics Department, Princess Margaret Cancer
- 11 Research Centre, University Health Network. His research interests are primarily in feature selection
- 12 in high dimensional data.
- 13 Wei Xu is an associate professor at the Dalla Lana School of Public Health, University of Toronto. His
- 14 research interests focus on biostatistics and bioinformatics methodology and statistical genetics.

15

# 17 Abstract

Background: Feature selection (FS) is critical for high dimensional data analysis. Ensemble based feature selection (EFS) is a commonly used approach to develop FS techniques. Rank aggregation (RA) is an essential step of EFS where results from multiple models are pooled to estimate feature importance. However, the literature primarily relies on rule-based methods to perform this step which may not always provide an optimal feature set.

23 Method and Results: This study proposes a novel Supervised Rank Aggregation (SRA) approach to 24 allow RA step to dynamically learn and adapt the model aggregation rules to obtain feature 25 importance. The approach creates a performance matrix containing feature and model performance 26 value from all models and prepares a supervised learning model to get the feature importance. 27 Then, unsupervised learning is performed to select the features using their importance. We evaluate 28 the performance of the algorithm using simulation studies and implement it into real research 29 studies, and compare its performance with various existing RA methods. The proposed SRA method 30 provides better or at par performance in terms of feature selection and predictive performance of 31 the model compared to existing methods.

32 Conclusion: SRA method provides an alternative to the existing approaches of RA for EFS. While the 33 current study is limited to the continuous cross-sectional outcome, other endpoints such as 34 longitudinal, categorical, and time-to-event medical data could also be used.

### 35 Keywords

36 high dimensional data, supervised rank aggregation, artificial intelligence, machine learning,

37 ensemble feature selection, random forest

## 38 Introduction

39 A high dimensional data has challenges associated with model fitting, generalizability [1], and 40 computation complexity [2,3], which prevents modeling by many classic statistical techniques. 41 Feature selection is an important component in high dimensional data analysis domains like 42 genomics [4] and radiomics [5], as it helps reduce the dimensions of the dataset. Literature provides 43 many techniques to perform feature selection. However, these techniques could be categorized 44 based on their feature selection (FS) approach (Figure 1). One broad category of FS techniques uses 45 only expert or domain knowledge to perform feature selection [6,7]. These techniques work in 46 scenarios with few features without interaction among features and are well known in the research 47 domain [8]. Another broad category of FS techniques combines expert or domain knowledge with 48 data [9,10]. FS techniques designed in the Bayesian framework incorporate prior knowledge in the 49 feature selection process [9].

The third and major category of FS techniques relies on the dataset to perform feature selection and is referred to as data-based FS techniques in this paper. These techniques are sub-categorized into Filter, Wrapper, and Embedded FS techniques [11,12]. Filter methods select features based on internal data structures like association [13] and information gain [14,15]. Wrapper methods evaluate multiple subsets of features iteratively by building models to get the feature subset, which achieves the best performance [16–18]. Embedded methods build the model that simultaneously performs features selection [19–22].

Literature suggests different approaches to use the FS techniques for FS. These approaches can be categorized into base, hybrid, and ensemble approaches. In the base approach, a single FS technique is used. In the hybrid approach, multiple FS techniques are used in a sequence to perform feature selection [10,23]. Commonly, a filter based FS technique is used as coarse FS followed by a wrapper or an embedded based FS technique for final FS [23]. Some approaches create a sequence by combining expert based FS with other FS techniques [10]. bioRxiv preprint doi: https://doi.org/10.1101/2022.02.21.481356; this version posted February 22, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

-	-
6	~

#### Figure 1: Different feature selection approaches

64 In an ensemble approach, instead of a single model, multiple models are created from the same 65 dataset. The performance of features from these models is pooled and ranked based on their 66 relevance. Finally, the relevant features are selected based on the cut-off of importance. Two 67 approaches can be used to generate multiple models, namely homogenous ensemble approach and 68 heterogeneous ensemble approach [24-26]. In a homogenous ensemble approach, multiple 69 datasets are created from the same data by sub-setting the samples, features, or both followed by 70 using a single technique to build the model on each of these datasets [25]. In a heterogeneous 71 ensemble approach, a single dataset is modeled using different techniques to generate multiple 72 models [26]. An ensemble approach could perform better than single model approaches [27].

73 In an ensemble approach, one of the essential steps is to pool together the performance of features 74 obtained from different models and is referred to as rank aggregation (RA) in this study. The 75 performance metric used for RA varies across the studies, like model estimates [8,21] and goodness 76 of fit [8]. Literature provides various techniques to aggregate the feature performance obtained 77 from different models, but these techniques mainly rely upon a pre-defined rule to aggregate the 78 performance of features, i.e., rule-based rank aggregation approaches. Commonly used methods to aggregate the performance of features is to find the mean, median or Robust rank aggregation (RRA) 79 80 performance of the feature across all models [28] [29]. However, they cannot learn from the data 81 about the RA rule dynamically and may even be sensitive towards extreme values like mean values.

In high-dimensional data analysis, the performances of rule-based analysis have been challenged by machine learning (ML) based approaches like supervised learning. ML-based approaches are considered effective in the dynamic and complex environment as compared to rule-based approaches because ML creates dynamic rules by learning and adapting to the existing environment [30]. In the case of ensemble FS, the data structure is dynamic and varies across datasets, so it may not always be possible for a predefined rule to give optimal results for all the scenarios [30,31]. Thus, it is desirable to explore the application of ML in all steps of ensemble FS owing to its dynamic learning characteristics. ML approaches like supervised learning are well established in the model
building step [32], but no supervised learning approach is designed for the RA step.

This study proposes a novel perspective to perform RA using the supervised learning approach of the ML called supervised rank aggregation (SRA). First, SRA creates a performance matrix that contains the performance of all features in all the models as the input and the performance of each model in achieving the final data analysis goal as the label. Then, supervised learning is used to find the relative rank or performance of features based on their potential to help achieve the best performance in the final data analysis.

97 SRA based ensemble feature selection (EFS) is highly innovative in many ways. Firstly, perspective is 98 unique as it pools and ranks features dynamically rather than using fixed rules for EFS. Secondly, it 99 provides a unique application of supervised learning models as they replace the static rule-based RA 100 approach with a dynamic rule-based RA approach. Thirdly, it is versatile, which allows its integration 101 with existing ensemble methods.

This paper provides the "Methodology" section to explain the SRA based EFS. Then, its performance is compared against existing rank aggregation methods used in EFS for simulations and real studies in the "Simulation Studies" and "Real Studies" sections. Finally, we summarize and provide future directions for research in the "Conclusion and Discussion" section.

# 106 Methodology

SRA methodology is developed to integrate the supervised learning in the rank aggregation step of the ensemble learning (Figure 2). A dataset of sample size, n, with given input feature space, p, and an outcome is fed into the EFS process, where multiple models are created either by creating multiple bootstrapped datasets from the original dataset (homogenous approach) or by using multiple modeling techniques (heterogeneous approach). Then a performance matrix is created from these multiple models by extracting feature performance and model performance. A supervised learning algorithm is trained on this performance matrix, and feature importance

obtained from the algorithm is used as the final feature ranking or importance. Finally, the features

115 are selected based on an importance cut-off obtained from a predefined threshold or an

- unsupervised ML algorithm. The proposed methodology is discussed below in more detail.
- 117

#### Figure 2: Graphical representation of SRA methodology based Ensemble feature selection

#### 118 Generate multiple models

119 From the original dataset D of feature space p, outcome y, and sample size n, k randomly sampled

datasets are generated by randomly sampling features without repeats  $q_i | i \in \{1, ..., k\}, 1 < q \le p$ .

121 All k sample datasets have a sample size of n by sampling with replacement from dataset D. A

122 model m is created for every k sample dataset using any modeling technique.

$$m_i: y_i = f(q_i) | i \in \{1, \dots, k\} \# (1)$$

123 where, modeling technique used to prepare the  $i^{th}$  dataset model  $m_i$  will determine the function f.

124 In this study, RIDGE regression is used as the modeling technique for building the models. Optimal

125 hyperparameter values for each model are obtained from 10-fold cross-validation.

#### 126 Create Performance Matrix

127 A performance matrix C is prepared from m models containing k rows and p + 1 columns. The 128 matrix contains feature performance, FP as the input features and model performance for study 129 objective, MP as the outcome or label for all m models.

$$C = |FP_{ij} \quad MP_i| \mid i \in \{1, \dots, k\}, j \in \{1, \dots, p\} \# (2)$$

In the current study, model estimates are used as *FP* metric and predictive performance of a model
on the left out bootstrap samples as *MP*. Accordingly, *MP* metric used in the study is inverse of root
mean square error (RMSE).

#### 133 Supervised Rank Aggregation

134 A supervised learning model (SLM) is created from the performance matrix with *FP* of *p* features as

135 predictors and *MP* as the outcome.

$$SLM:MP = g(FP)\#(3)$$

136 where, machine learning technique used for SLM will determine the function g. Currently, only ML

137 techniques like penalized regression and decision trees which could provide feature importance,

138 *fimp* in achieving the model performance could be used.

#### 139 Feature Selection

140 The importance for each feature is used to select target features  $q_{best}$ . It is assumed that the 141 features with more importance should be target features as they are more relevant in achieving 142 higher model performance. In literature, the cut-off value for features is obtained by using a pre-143 defined threshold [8,21], rule-based threshold estimation [33], or unsupervised learning based 144 threshold estimation [21]. A predefined threshold may require the tuning step to arrive at an 145 appropriate cut-off value, which will give optimal results for a given scenario [8,21]. Rule-based 146 methods may not always provide optimal results [30,31]. Thus, in this study, the K-means based 147 unsupervised learning technique is used for obtaining the threshold cut-off as it will eliminate the 148 need for tuning and dynamically adapt to the given scenario. Since clustering will be happening on a 149 single dimension, hence high dimension limitation of K-mean clustering is avoided. K-means is used 150 to cluster the features into two groups, and the features in the cluster with a higher mean fimp 151 value are selected as final features  $q_{best}$ . Pseudo Algorithm summarizes the complete SRA based 152 ensemble feature selection algorithm.

 Pseudo Algorithm: SRA based ensemble feature selection

 Input:
 Feature data X (p × n)

 Target feature Y (1 × n)

 Number of sample datasets k

 Performance matrix C = {empty}

 Output:
 Final Feature set q<sub>best</sub>

Begin: # Step I: Generate multiple models for i=1 to k Generate  $q_i$  random features from pGenerate samples  $(X^i, Y^i \in R^{n \times (q_i+1)})$ Build embedded model  $m_i$  (like RIDGE) from  $(X^i, Y^i)$ end for # Step II: Prepare Performance Matrix for i=1 to k Compute feature performance estimate  $FP_i$  of the model Compute model performance estimate  $MP_i$  of the model Add  $(FP_i, MP_i)$  to C end for # Step III: Supervised Rank Aggregation Build a supervised learning model (like LASSO, RIDGE, Random Forest), SLM from C Compute feature importance estimate (like feature coefficient, feature importance score) fimp from SLM model for p# Step IV: Feature Selection Generate two clusters  $(c_1, c_2)$  from fimp of p features using unsupervised learning like K-means Compute mean *fimp* of  $c_1 m c_1 = \mu(c_1)$ Compute mean *fimp* of  $c_2 m c_2 = \mu(c_2)$ if  $mc_1 = MAXIMUM(mc_1, mc_2)$ Add  $c_1$  features to  $q_{best}$  to get final feature selection else if  $mc_2$  = MAXIMUM( $mc_1$ ,  $mc_2$ ) Add  $c_2$  features to  $q_{best}$  to get final feature selection else Add p features to  $q_{best}$  to get final feature selection End

# 153 Simulation Studies

154 We perform simulation studies to evaluate the proposed RA method and compare its performance 155 with multiple other RA methods for EFS. The study generates high-dimensional feature space for 156 marginal models using multivariate normal distributions. The study uses regression model  $y = \beta_0 + \sum_{i=1}^p \beta_i x_i + \epsilon$  to provide a continuous outcome variable of simulated data with sample 157 158 size, n for marginal models.  $\beta$  represents the effect of different features and intercept term on the outcome,  $\varepsilon \sim N(0, \sigma^2)$  is the normally distributed error term and  $x_i \sim N(0, 1)$  are normally 159 160 distributed input features, p. Multi-collinearity is added between features using the covariance 161 matrix as given below:

$\int x_1 x_1$	•	$x_1 x_{15}$	•	•	$x_1 x_p$		۲1		5		ן0	
	•		•	•	•				5			
$x_{15}x_{1}$	•	$x_{15}x_{15}$	•	÷	$x_{15}x_{p}$	=	5	5	1		0	
$x_p x_1$	•	$x_{p}x_{15}$	•	•	$x_p x_p$		L <sub>0</sub>		0		1 1	

162 Multiple scenarios are simulated by changing  $p, n, \beta$ , the number of target features, and k (Table 1). 163 Only true features are assigned a non-zero  $\beta$  value. We prepare homogenous ensemble models for 164 feature selection. The dataset for each model is generated by randomly sampling two to p features 165 from p feature space and sub-setting n samples from the original dataset with replacement. RIDGE is 166 used to build models for each dataset. A penalized effect size of each feature obtained from the 167 RIDGE models is scaled using the absolute maximum value, which is used as a feature performance 168 metric.

169 Implementation of SRA is shown using three different supervised learning algorithms, namely LASSO 170 (SRA-Lasso), RIDGE (SRA-Ridge), random forest (SRA-RF). While LASSO and RIDGE perform 171 supervised learning using a linear combination of features, random forest performs supervised 172 learning using a non-linear combination of features. A supervised learning model in each SRA is 173 prepared using optimized hyperparameter values.

SRA performance is compared with existing rule-based RA methods, namely, mean based RA (MeRA), maximum based RA (MaRA), minimum based RA (MiRA), median based RA (MedRA), coefficient of variation based RA (CVRA), standard deviation based RA (SDRA), robust rank aggregation (RRA), t-test based RA (tRA), and Wilcoxon signed-rank test based RA (WRA). R 4.0.3 is used for the analysis. The study has used some inbuilt packages in statistical language R for the analysis like *glmnet* package [34] for LASSO and RIDGE, *randomForest* package [35] for random forest, and *RobustRankAggreg* package [29] for RRA.

181 The different RA methods are evaluated for their ability to select target features, discriminate 182 between target and noise features, and predictive performance of the models built using selected 183 features. We use the F1 score for feature discrimination ability evaluation and inverse RMSE for the 184 test data for the predictive performance evaluation. RIDGE is used to build the final model from the 185 selected features for predictive performance evaluation. Ten trials are performed for each scenario.

Table 2 results suggest that all methods can select some target features under all scenarios, but SRA-Ridge consistently outperformed rule-based RA methods. SRA-Ridge selected almost all the target features in all scenarios. The performance of the other two SRA methods is at par with existing RA methods. Further, the results suggest that SRA-Ridge has a better or at par feature discriminative ability than other methods. Thus, SRA-Ridge not only selects target features but is also good in rejecting noise features as compared to other methods. The results suggest that SRA could be a good candidate to select target features.

Further, SRA-Ridge based selected features can build good predictive models and consistently outperformed rule-based RA methods (Table 2). These findings suggest that the SRA may provide better or at par prediction performance than existing methods. Further, SRA could enhance the performance of ensemble-based approaches in high-dimensional settings.

# 197 Real Studies

Three real studies are analyzed to compare the performance of SRA and existing RA methods. Study I is Community Health Status Indicators (CHSI) study that collected US county data (n=3141) containing 578 features to understand non-communicable diseases [36]. Study II is National Social Life, Health and Aging Project (NSHAP) study that collected aged Americans data (n=4377) containing 1470 features to understand their health and well-being [37]. Study III is the DNA methylation data (n=27578) containing 108 samples to understand its relationship with human age [38,39].

Table 3 shows the final cleaned dataset for these three studies used for analysis. Features and samples are filtered to remove highly correlated features, non-continuous features, missing values, and very low standard deviation. The final cleaned dataset is randomly split into training and test dataset. The test dataset is used to evaluate the predictive performance of the features selected by
different RA methods. The study uses inverse RMSE as the predictive performance metric. The mean
performance of ten trials is used for comparison between RA methods. In the cases of Study I and
Study II, 100 ensemble models are created, while in Study III, 1000 ensemble models are created.

212 The results from Table 4 suggest that SRA methods provided better or at par predictive performance 213 than existing RA methods. The better performance of the SRA method suggests that it may be more 214 reliable than existing RA methods in identifying the target features. Further, unlike the simulated 215 data results, different SRA methods have shown different performances. In the case of Study I and 216 Study II, SRA-Ridge has the best predictive performance, but in Study III, SRA-Lasso has the best 217 predictive performance, which suggests that SRA methods performance may change with dataset 218 and ensemble models. In general, the variation in performance of feature selection techniques with 219 dataset has been reported in the literature and could be attributed to data characteristics [40].

220 In the current study, Study III data is also used to compare the performance of SRA based selected 221 methylated features with state-of-art literature based selected features [41,42]. SRA-Lasso is used to 222 obtain the target features. The complete dataset is used for the FS step rather than the training 223 data. SRA-Lasso identified 484 methylation sites compared to 353 methylation sites identified by the 224 literature, but only ten methylation sites are shared between the two approaches (Supplementary 225 File 1). The selected methylation sites from the two approaches are compared for their predictive 226 performance on the test data. Accordingly, the Study III dataset is split into training (80%) and test 227 (20%) data. RIDGE model is prepared using training data followed by predictive performance 228 measurement on test data. It is found that SRA-Lasso based selected features provided a marginally 229 better predictive performance (RMSE<sup>-1</sup>(95% Cl): 0.06 (0.05-0.07)) compared to literature 230 recommended selected features (RMSE<sup>-1</sup>(95% CI): 0.05 (0.04-0.05)).

Further, we identified the differentially expressed genes associated with selected methylated sites
using *BioMethyl* package [43]. SRA-Lasso based selected methylated sites are linked with 288 genes,

233 but literature based selected methylated sites are linked with only 136 genes (Supplementary File 2). 234 Only ten genes, namely SFRP1, STRA6, BNC1, CSPG5, DCHS1, DIRAS3, TCF15, ERG, PIPOX, and 235 MCAM, are shared between the two approaches. Literature also provides a database, GenAge, of 236 307 genes commonly associated with age [44]. Among the 136 genes linked with literature-based 237 methylation sites, only 1 out of 308 genes is found (Supplementary File 3). However, among the 288 238 genes linked with SRA-Lasso based methylation sites, 9 out of 308 genes are found (Supplementary 239 File 3). Thus, SRA-Lasso may be relevant in identifying target features that have both biological 240 importance and good predictive performance.

# 241 Conclusion and Discussion

This paper proposes SRA, an innovative approach, to perform rank aggregation in ensemble models for feature selection. The approach allows dynamic learning of feature performance pooling strategy, which current rule-based rank aggregation methods do not perform. The approach is flexible and could be incorporated into any ensemble technique. The SRA could identify target features while retaining very few noise features compared to other methods. The simulated data studies showed that SRA outperforms existing methods in feature selection and prediction performance. Similar performance in real datasets also demonstrates the practical relevance of SRA.

The proposed method has certain limitations. The scope of the current study is limited to concept testing. Consequently, the robustness of the approach on different data types and modeling techniques could be the focus of future research. The ensemble model used in the study assumes a linear combination of features. Thus, future research could study SRA for algorithms designed to explore the non-linear combinations of features.

## 254 Key Points

Supervised Rank Aggregation (SRA) methods are better than rule-based rank aggregation
 methods for ensemble-based feature selection (EFS).

- SRA Ridge could give much better discrimination between true and noise features as well as
- 258 predictive performance than rule-based rank aggregation methods
- SRA could be useful in detecting the genomic features like methylation sites which could
- 260 have biological relevance

# 261 Declarations

- 262 Ethics approval and consent to participate
- 263 Not Applicable
- 264 Consent for publication
- 265 Not Applicable
- 266 Availability of data and materials
- All the datasets and code are in the github link: https://github.com/rahijaingithub/SRA.
- 268 Competing interests
- 269 The authors declare that they have no competing interests
- 270 Funding
- 271 This work was supported by the Natural Sciences and Engineering Research Council of Canada
- 272 [NSERC Grant RGPIN-2017-06672 to W.X.]; and the Prostate Cancer Canada [Translation Acceleration
- 273 Grant 2018 to R.J. and W.X.].
- 274 Author Contributions
- 275 ALL AUTHORS HAVE READ AND APPROVED THE MANUSCRIPT.
- 276 Conceptualisation: RJ, WX
- 277 Formal Analysis: RJ

bioRxiv preprint doi: https://doi.org/10.1101/2022.02.21.481356; this version posted February 22, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

- 278 Investigation: RJ
- 279 **Methodology:** RJ, WX
- 280 Software: RJ
- 281 Supervision: RJ, WX
- 282 Validation: RJ, WX
- 283 Writing-original draft: RJ
- 284 Writing-review & editing: RJ, WX
- 285 Acknowledgements
- 286 Not Applicable

## 287 Reference

- 288 1. Bellman R. Dynamic Programming. Math. Sci. Eng. 1967; 40:101–137
- 289 2. Fan J, Li R. Statistical challenges with high dimensionality<sup>2</sup>: feature selection in knowledge
- 290 discovery. Proc. Int. Congr. Math. Madrid, August 22–30, 2006 2007; 595–622
- 3. Ayesha S, Hanif MK, Talib R. Overview and comparative study of dimensionality reduction
   techniques for high dimensional data. Inf. Fusion 2020; 59:44–58
- 4. Piles M, Bergsma R, Gianola D, et al. Feature Selection Stability and Accuracy of Prediction Models for Genomic Prediction of Residual Feed Intake in Pigs Using Machine Learning. Front. Genet. 2021;
- 295 12:
- 5. Healy G, Salinas-Miranda E, Jain R, et al. Pre-operative radiomics model for prognostication in
   resectable pancreatic adenocarcinoma with external validation. Eur. Radiol. 2021; Online:
- 298 6. Walter S, Tiemeier H. Variable selection: Current practice in epidemiological studies. Eur. J.
  299 Epidemiol. 2009; 24:733–736
- 7. Heinze G, Wallisch C, Dunkler D. Variable selection A review and recommendations for the
   practicing statistician. Biometrical J. 2018; 60:431–449
- 302 8. Jain R, Xu W. HDSI: High dimensional selection with interactions algorithm on feature selection
   303 and testing. PLoS One 2021; 16:1–17
- 304 9. Mitchell TJ, Beauchamp JJ. Bayesian variable selection in linear regression. J. Am. Stat. Assoc.
  305 1988; 83:1023-1032
- 306 10. Zycinski G, Barla A, Squillario M, et al. Knowledge Driven Variable Selection (KDVS) a new
- 307 approach to enrichment analysis of gene signatures obtained from high-throughput data. Source

308 Code Biol. Med. 2013; 8:1–14

- 309 11. Yang P, Huang H, Liu C. Feature selection revisited in the single-cell era. Genome Biol. 2021;
  310 22:1–17
- 12. Dhal P, Azad C. A comprehensive survey on feature selection in the various fields of machine
   learning. Appl. Intell. 2021; 51:1–39
- 313 13. Chormunge S, Jena S. Correlation based feature selection with clustering for high dimensional
  314 data. J. Electr. Syst. Inf. Technol. 2018; 5:542–549
- 14. Dash M, Liu H, Yao J. Dimensionality reduction of unsupervised data. Proc. Ninth IEEE Int. Conf.
  Tools with Artif. Intell. 1997; 532–539
- 317 15. Zhang R, Nie F, Li X, et al. Feature selection with multi-view data: A survey. Inf. Fusion 2019;
  318 50:158–167
- 319 16. Kohavi R, John GH. Wrappers for feature subset seelction. Artif. Intell. 1997; 97:273–324
- 17. Tarkhaneh O, Nguyen TT, Mazaheri S. A novel wrapper-based feature subset selection method
   using modified binary differential evolution algorithm. Inf. Sci. (Ny). 2021; 565:278–305
- 18. Alweshah M, Alkhalaileh S, Al-betar MA. Coronavirus herd immunity optimizer with greedy
   crossover for feature selection in medical diagnosis. Knowledge-Based Syst. 2020; 235:107629
- 19. Tibshirani R. Regression shrinkage and selection via the lasso: A retrospective. J. R. Stat. Soc. Ser.
  B Stat. Methodol. 2011; 73:273–282
- 20. Chun H, Keleş S. Sparse partial least squares regression for simultaneous dimension reduction
  and variable selection. J. R. Stat. Soc. Ser. B Stat. Methodol. 2010; 72:3–25
- 328 21. Jain R, Xu W. RHDSI: A novel dimensionality reduction based algorithm on high dimensional
   329 feature selection with interactions. Inf. Sci. (Ny). 2021; 574:590–605
- 330 22. Lal TN, Chapelle O, Weston J. Embedded Methods. Featur. Extr. Found. Appl. 2006; 165:137–165
- 331 23. Hancer E, Xue B, Zhang M. A survey on feature selection approaches for clustering. Artif. Intell.
   332 Rev. 2020; 53:4519-4545
- 333 24. Seijo-Pardo B, Porto-Díaz I, Bolón-Canedo V, et al. Ensemble feature selection: Homogeneous
   334 and heterogeneous approaches. Knowledge-Based Syst. 2017; 118:124–139
- 25. Hosni M, Idri A, Abran A. On the value of filter feature selection techniques in homogeneous
  ensembles effort estimation. J. Softw. Evol. Process 2021; 33:e2343
- 26. Mera-Gaona M, López DM, Vargas-Canas R, et al. Framework for the ensemble of feature
  selection methods. Appl. Sci. 2021; 11:1–16
- 27. Tsai CF, Sung YT. Ensemble feature selection in high dimension, low sample size datasets: Parallel
   and serial combination approaches. Knowledge-Based Syst. 2020; 203:106097
- 28. Noureldien N, Mohmoud S. The Efficiency of Aggregation Methods in Ensemble Filter Feature
  Selection Models. Trans. Mach. Learn. Artif. Intell. 2021; 9:39–51
- 343 29. Kolde R, Laur S, Adler P, et al. Robust rank aggregation for gene list integration and meta344 analysis. Bioinformatics 2012; 28:573–580
- 345 30. van Ginneken B. Fifty years of computer analysis in chest imaging: rule-based, machine learning,
- deep learning. Radiol. Phys. Technol. 2017; 10:23–32

- 347 31. Cronin RM, Fabbri D, Denny JC, et al. A comparison of rule-based and machine learning
- 348 approaches for classifying patient portal messages. Int. J. Med. Inform. 2017; 105:110–120
- 349 32. Lopez-Rincon A, Mendoza-Maldonado L, Martinez-Archundia M, et al. Machine learning-based
- ensemble recursive feature selection of circulating mirnas for cancer tumor classification. Cancers
   (Basel). 2020; 12:1–27
- 33. Seijo-Pardo B, Bolón-Canedo V, Alonso-Betanzos A. Testing Different Ensemble Configurations
   for Feature Selection. Neural Process. Lett. 2017; 46:857–880
- 34. Friedman JH, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models via
   Coordinate Descent. J. Stat. Softw. 2010; 33:1–22
- 356 35. Liaw A, Wiener M. Classification and Regression by randomForest. R News 2002; 2:18–22
- 36. [Dataset] Centers for Disease Control and Prevention. Community Health Status Indicators (CHSI)
   to Combat Obesity, Heart Disease and Cancer. Healthdata.gov 2012;
- 359 37. [DATASET] Waite LJ, Laumann EO, Levinson WS, et al. National Social Life, Health, and Aging
- Project (NSHAP): Wave 1, [United States], 2005-2006 (ICPSR 20541). Inter-university Consort. Polit.
  Soc. Res. 2019;
- 362 38. Numata S, Ye T, Hyde TM, et al. DNA methylation signatures in development and aging of the
   363 human prefrontal cortex. Am. J. Hum. Genet. 2012; 90:260–272
- 364 39. [Dataset] Akalin A. compGenomRData. Github 2019;
- 40. Parmezan ARS, Lee HD, Spolaôr N, et al. Automatic recommendation of feature selection
   algorithms based on dataset characteristics. Expert Syst. Appl. 2021; 185:115589
- 367 41. Pelegi-Siso D, De Prado P, Ronkainen J, et al. Methylclock: A Bioconductor package to estimate
- 368 DNA methylation age methylclock: A Bioconductor package to estimate DNA methylation age.
   369 Bioinformatics 2021; 37:1759–1760
- 42. Horvath S. DNA methylation age of human tissues and cell types. Genome Biol. 2015; 16:1–19
- 43. Wang Y, Franks JM, Whitfield ML, et al. BioMethyl: An R package for biological interpretation of
  DNA methylation data. Bioinformatics 2019; 35:3635–3641
- 44. Tacutu R, Thornton D, Johnson E, et al. Human Ageing Genomic Resources: New and updated
- databases. Nucleic Acids Res. 2018; 46:D1083–D1090
- 375





			Sample Size			
Scenario	$oldsymbol{eta}$ (Non-Zero coefficients)	р			σ	k
			Train ( <i>n</i> )	Test		
A	$\{ \beta_i   i = \{1,, 10\} \} = \{0.9,, 0.9\}$	75	100	500	0.25	300
_		4.00	100		<del>-</del>	4.0.0
В	$\{ \beta_i   i = \{1,, 10\} \} = \{0.5,, 0.5\}$	100	100	500	0.25	100
C	$(R \mid i - (1  15)) = (0.4 - 0.8  0.4 - 0.8  0.4)$	175	275	500	0.25	100
C	$\{p_i \mid i = \{1,, 15\}\} = \{0.4, -0.0, 0.4, -0.0,, 0.4\}$	1/5	275	500	0.25	100
D	$\{\beta_i \mid i = \{1, \dots, 15\}\} = \{0.4, -0.8, 0.4, -0.8, \dots, 0.4\}$	75	275	500	0.25	100
E	$\{\beta_i   i = \{1,, 15\}\} = \{0.4, -0.8, 0.4, -0.8,, 0.4\}$	75	225	500	0.25	200
F	$\{\beta_i \mid i = \{1, \dots, 20\}\} = \{0.4, -0.8, 0.4, -0.8, \dots, -0.8\}$	125	225	500	0.25	200

## Table 1: Description of the scenarios of simulation studies

1

#### 1 Table 1: Comparison of model performance between SRA methods and Existing methods under six

### 2 scenarios in terms of target feature selection, feature discrimination ability (F1 Score) and outcome

3

#### prediction (1/RMSE)

				Scen	arios		
RA	technique	А	В	С	D	Ε	F
			-	Target Features	s (%) [µ(95%Cl)	]	
		100	100	46	46	47	51
	CVKA	(100-100)	(100-100)	(45-47)	(45-47)	(47-47)	(48-53)
	ΜΔΒΔ	100	100	87	97	85	66
	100 (10) (	(100-100)	(100-100)	(81-93)	(95-99)	(78-93)	(56-76)
	MeRA	100	100	47	47	47	53
		(100-100)	(100-100)	(47-47)	(47-47)	(47-47)	(51-55)
	MedRA	100	100	4/	4/	4/	53
bu		(100-100)	(100-100)	(4/-4/)	(46-49)	(4/-4/)	(30-30)
isti	MIRA	6Z (48.76)	94 (89.99)	95 (01 08)	// (72 22)	62 (92 99)	88 (96 90)
EX		(48-76)	(09-99)	(91-98) 47	(72-02)	(02-05) 17	(00-05)
	RRA	(97-100)	(97-101)	(47-47)	(46-50)	(46-49)	(50-54)
		78	71	34	40	35	39
	SDRA	(67-89)	(67-75)	(29-39)	(35-45)	(27-42)	(32-46)
		<b>100</b>	<b>100</b>	46	45	47	51
	tRA	(100-100)	(100-100)	(45-47)	(44-47)	(47-47)	(48-53)
		100	100	49	53	53	37
	WKA	(100-100)	(100-100)	(45-53)	(53-53)	(53-53)	(34-40)
	Lasso RF	92	37	41	58	63	46
		(87-97)	(18-56)	(37-45)	(51-65)	(57-69)	(38-54)
RA		98	53	63	67	65	61
S		(95-100)	(43-63)	(54-73)	(57-76)	(56-75)	(53-69)
	Ridge	100	99	95 (80,100)	95	100	92
P۸	technique	(100-100)	(97-100)	(89-100) E1 Score[	(92-99) u(95%CI)]	(100-100)	(88-90)
na.	leunnque	1.00	0.93	0.63	μ(95%Cl)]	0.64	0.67
	CVRA	(1 00-1 00)	(0.95 (0.89-0.97)	(0.65-0.64)	(0.65	(0.64-0.64)	(0.65-0.69)
		0.58	0.70	0.60	0.83	0.65	0.44
	MARA	(0.55-0.61)	(0.68-0.72)	(0.55-0.64)	(0.8-0.86)	(0.61-0.69)	(0.39-0.49)
		0.83	0.81	0.64	0.64	0.64	0.69
	MeRA	(0.8-0.86)	(0.80-0.82)	(0.64-0.64)	(0.64-0.64)	(0.64-0.64)	(0.67-0.71)
	ModRA	0.85	0.81	0.64	0.64	0.64	0.69
6	IVIEUNA	(0.82-0.87)	(0.80-0.82)	(0.64-0.64)	(0.63-0.65)	(0.64-0.64)	(0.67-0.71)
tin	MIRA	0.41	0.60	0.79	0.75	0.70	0.76
Exis		(0.33-0.49)	(0.53-0.67)	(0.73-0.85)	(0.72-0.78)	(0.67-0.73)	(0.74-0.79)
	RRA	0.90	0.82	0.64	0.65	0.64	0.68
		(0.87-0.93)	(0.80-0.83)	(0.64-0.64)	(0.63-0.66)	(0.63-0.65)	(0.67-0.7)
	SDRA	0.30	0.22	0.07	U.16	0.16	0.12
		(0.27-0.32)	(0.20-0.24)	(0.06-0.07)	(0.12-0.18)	(0.14-0.17)	(0.10-0.14)
	tRA	1.00 (1.00)			U.62		
		0.40	(0.00-0.20)	(0.02-0.04) 0.14	(U.01-U.04) A 29	(0.04-0.04) 0 20	(80.0-00) 0 1 Q
	WRA	(0.38-0.42)	0.33 (0.32-0.34)	(0.14 (0.13-0.16)	(0.25 (0.27-0.3)	(0.28-0.32)	(0 17-0 2)
<u> </u>	lasso	0.96	0.32	0.52	0.73	0.77	0.62
	Lasso	0.90	0.55	0.56	0.75	0.77	0.02

bioRxiv preprint doi: https://doi.org/10.1101/2022.02.21.481356; this version posted February 22, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

		(0.93-0.98)	(0.16-0.50)	(0.53-0.62)	(0.67-0.79)	(0.72-0.81)	(0.55-0.7)
	DE	0.75	0.29	0.64	0.73	0.77	0.70
	ĸr	(0.69-0.81)	(0.24-0.33)	(0.57-0.71)	(0.66-0.8)	(0.71-0.83)	(0.64-0.76)
	Didee	1.00	0.99	0.97	0.98	1.00	0.95
	Ridge	(1.00-1.00)	(0.98-1.00)	(0.94-1.00)	(0.96-0.99)	(1.00-1.00)	(0.93-0.98)
RA	technique		Predicti	ve Performanc	æ (1/RMSE) [μ(	95%CI)]	
		3.50	3.82	0.81	0.84	0.83	0.75
	CVRA	(3.29-3.71)	(2.90-4.75)	(0.79-0.84)	(0.81-0.86)	(0.80-0.85)	(0.72-0.77)
		2.67	3.65	1.73	2.42	1.43	0.58
	IVIARA	(2.43-2.90)	(2.77-4.54)	(1.43-2.03)	(1.97-2.86)	(1.01-1.85)	(0.51-0.65)
		2.94	3.67	0.82	0.84	0.83	0.76
	IVIERA	(2.56-3.31)	(2.83-4.51)	(0.80-0.84)	(0.82-0.87)	(0.80-0.85)	(0.73-0.79)
	ModPA	2.96	3.67	0.82	0.85	0.83	0.76
D	WIEUKA	(2.55-3.36)	(2.83-4.51)	(0.80-0.84)	(0.82-0.89)	(0.80-0.85)	(0.73-0.79)
ting	MIRA	0.80	2.58	2.45	1.45	1.74	1.31
six		(0.27-1.34)	(2.00-3.17)	(1.97-2.93)	(1.29-1.61)	(1.57-1.91)	(1.25-1.37)
-	DDA	2.92	3.54	0.82	0.87	0.84	0.75
		(2.42-3.42)	(2.61-4.46)	(0.80-0.84)	(0.83-0.91)	(0.81-0.87)	(0.73-0.77)
	SDPA	1.10	1.03	0.68	0.77	0.71	0.53
	JUNA	(0.53-1.68)	(0.96-1.10)	(0.66-0.7)	(0.74-0.8)	(0.66-0.76)	(0.47-0.58)
	+Ρ.Δ	3.50	3.79	0.81	0.84	0.83	0.75
		(3.29-3.71)	(2.91-4.67)	(0.79-0.83)	(0.81-0.86)	(0.80-0.85)	(0.72-0.77)
	\ <b>M/R</b> Δ	2.18	2.62	0.45	0.47	0.46	0.37
	WIG	(1.96-2.39)	(2.39-2.85)	(0.44-0.45)	(0.45-0.48)	(0.45-0.46)	(0.36-0.38)
	lasso	2.17	0.77	0.79	1.00	1.09	0.72
	Labbo	(1.32-3.02)	(0.51-1.03)	(0.76-0.82)	(0.85-1.16)	(0.96-1.21)	(0.66-0.79)
RA	RF	2.62	0.88	0.70	0.87	0.73	0.55
S		(2.07-3.17)	(0.77-0.99)	(0.62-0.78)	(0.46-1.27)	(0.55-0.9)	(0.50-0.59)
	Ridge	3.50	3.83	2.58	2.65	2.98	1.87
RIU	1	(3.29-3.71)	(2.91-4.75)	(2.07-3.08)	(2.02-3.28)	(2.51-3.44)	(1.46-2.28)

1

## Table 1: Summary of the real datasets

Real Marginal		Outcome feature	Sar	nple size	k	
Studies	Features (p)	Outcome leature	Total	Train	Test	
Study I	45	Height	1035	207	828	100
Study II	74	Number of unhealthy days	177	141	36	100
Study III	2289	Age	108	86	22	1000

#### 1 Table 1: Comparison of SRA methods with Existing methods for three real studies in terms of outcome

## 2

#### prediction (1/RMSE)

		Study						
<b>RA technique</b>		I	11	<i>III</i>				
	-	Predictiv	/e Performance (1/RMSE) [μ(95%	6CI)]				
	CVRA	1.08(1.07-1.1)	1.28(1.22-1.35)	2.36(2.15-2.56)				
	MARA	1.14(1.12-1.16)	1.25(1.2-1.31)	2.14(1.95-2.34)				
	MeRA	3.14(2.96-3.31)	1.28(1.23-1.33)	2.35(2.17-2.54)				
bı	MedRA	3.16(2.98-3.35)	1.28(1.23-1.33)	2.38(2.19-2.57)				
aistii	MIRA	2.96(2.76-3.17)	1.22(1.16-1.27)	1.76(1.67-1.86)				
Ê	RRA	3.13(2.96-3.3)	1.28(1.19-1.36)	2.39(2.18-2.61)				
	SDRA	1.06(1.03-1.09)	1.16(1.09-1.23)	1.05(1.02-1.09)				
	tRA	1.08(1.07-1.1)	1.28(1.2-1.36)	2.35(2.15-2.55)				
	WRA	1.13(1.11-1.14)	1.25(1.19-1.3)	2.05(1.87-2.23)				
	Lasso	1(0.99-1.01)	1.23(1.14-1.31)	2.72(2.25-3.19)				
SRA	RF	2.51(1.86-3.16)	1.23(1.14-1.32)	1.88(1.75-2)				
	Ridge	3.21(3.03-3.39)	1.28(1.22-1.34)	2.30(2.12-2.47)				