

Tomic et al, SIMON: open-source knowledge discovery platform

32 **Main text:**

33 Over the past several years, due to the technological breakthroughs in genome sequencing¹, high-
34 dimensional flow cytometry²⁻⁴, mass cytometry^{5,6}, and multi-parameter microscopy^{7,8} the amount and
35 complexity of biological data has become increasingly intractable and it is no longer feasible to extract
36 knowledge without using sophisticated computer algorithms. Therefore, researchers are in need of novel
37 computational approaches that can cope with complexity and heterogeneity of data in an objective and
38 unbiased way. Machine learning (ML), a subset of artificial intelligence, is a computational approach
39 developed to identify patterns from the data in order to make predictions on new data⁹. ML has a profound
40 impact on biological research¹⁰⁻¹², including genomics¹³, proteomics¹⁴⁻¹⁶, cell image analysis¹⁷, drug
41 discovery and development¹⁸, and cell phenotyping^{6,19,20} which revolutionized our understanding of
42 biological complexity. Recently, using systems-level analysis of genetic, transcriptional, and proteomic
43 signatures to predict patients' response to vaccines^{21,22}, therapies and disease progression²³⁻²⁷, ML has
44 become primary computational approach used in the 'precision medicine'²⁸.

45 The biggest challenge is the proper application of ML methods and the translation of the results into
46 meaningful insights. The analysis of massive datasets and extraction of knowledge using ML requires
47 knowledge of many different computational libraries for data pre-processing and cleaning, data partitioning,
48 model building and tuning, evaluation of the performance for the model and minimizing overfitting¹¹. Tools
49 to achieve these tasks have been mainly developed either in R (<https://www.r-project.org/>)^{29,30} or Python
50 (www.python.org/)³¹, which have today become leading statistical programming languages in data science.
51 Because R and Python are free and open-source, they have been quickly adopted by a large community of
52 programmers who are building new libraries and improving existing ones. As of May 2020, there are 15,658
53 R packages available in the CRAN package repository (<https://cran.r-project.org/>). Many of the packages
54 offer different modeling functions and have different syntax for model training, predictions and
55 determination of variable importance. Due to the lack of a unified method for proper application of ML
56 process, even experienced bioinformaticians struggle with these time-consuming ML tasks. To provide a
57 uniform interface and standardize the process of building predictive models, ML libraries were developed,
58 for example mlr3³² (<https://mlr3.mlr-org.com>), the classification and regression training (caret)^{30,33}
59 (<https://rdrr.io/cran/caret>), scikit-learn³⁴ (<https://scikit-learn.org>), mlPy³⁵ (<https://mlpy.fbk.eu>), SciPy
60 (<https://www.scipy.org/>) including also ones for deep learning, such as TensorFlow
61 (<https://www.tensorflow.org/>), PyTorch (<https://pytorch.org/>) and Keras (<https://keras.io/>). Since those
62 libraries do not have a graphical user interface, usage requires extensive programming experience and
63 general knowledge of R or Python making it inaccessible for many life science researchers. Therefore, there

Tomic et al, SIMON: open-source knowledge discovery platform

64 is an increased effort to harmonize those libraries and develop a software that will facilitate application of
65 ML in life sciences.

66 The software should provide a standardized ML method for data pre-processing, data partitioning,
67 building predictive models, evaluation of model performance and selection of features. Moreover, such
68 software should be adapted for biological datasets that have high percentage of missing values³⁶, have
69 imbalanced participant distributions (i.e. having a high number of infected subjects, but only a relatively
70 small number of healthy controls)³⁷ or suffer from a “*curse of dimensionality*”, i.e. poor predictive power,
71 as can be observed when the number of features is much greater than the number of samples³⁸. Additionally,
72 beyond ML process, the software should support exploratory analysis and visualization of the results using
73 user-friendly graphical interface. The fast-paced technological development dramatically increased size of
74 biological datasets and computational power needed for analysis. Therefore, open-source web-based
75 software supporting cloud processing architecture is essential. Additionally, software should support an
76 automated ML³⁹ (autoML) process that rapidly builds high-performance predictive models by identifying
77 optimal ML method, including selection of an appropriate algorithm, optimization of model
78 hyperparameters and evaluation of the best-performing models. AutoML improves the efficiency of ML
79 process and resulting models often outperform hand-designed models^{39,40}.

80 To address these challenges, we developed SIMON (Sequential Iterative Modeling “Over Night”), a free
81 and open-source software for application of ML in life sciences that facilitates production of high-
82 performing ML models and allows researchers to focus on knowledge discovery process. SIMON provides
83 a user-friendly, uniform interface for building and evaluating predictive models using a variety of ML
84 algorithms. Currently, there are 182 different ML algorithms available (**Supplementary table 1**). The entire
85 ML process which is based on the caret³³ library, from model building and evaluation to feature selection
86 in SIMON is fully automated. This allows advanced ML users to focus on other important aspects necessary
87 to build highly accurate models, such as data preprocessing, feature engineering and model deployment. It
88 also makes the entire ML process more accessible to domain-knowledge experts that formulated the
89 research hypothesis and collected the data, but lack programming ML skills. Additionally, to prevent
90 optimistic accuracy estimates and to optimize the model for generalization to unseen data, SIMON
91 introduces unified process for model training, hyperparameter tuning and model evaluation by generation
92 of training, validation and test sets. Training set is used for building models, validation set is used for
93 hyperparameter tuning and finally, models are evaluated in an unbiased way using the test, also known as
94 holdout set that has never been used in training. Beside the standardized ML process, the initial install
95 version offers a set of core components specifically suited for analysis of biomedical data, such as multi-
96 set intersection function for integration of data with many missing values⁴¹ (<https://cran.r->

Tomic et al, SIMON: open-source knowledge discovery platform

97 project.org/web/packages/mulset/index.html), method for identifying differentially expressed genes using
98 significance analysis in microarrays (SAM)⁴², graphical representation of the clustering analysis important
99 for detection of batch effects, graphical display of the correlation analysis and graphical visualizations of
100 the ML results that can be downloaded as publication-ready figures in scalable vector graphics (SVG)
101 format. Finally, SIMON is available in two versions as a single mode and a server version. The single mode
102 is developed as a SIMON Docker container (<https://www.docker.com/>), ensuring code reproducibility and
103 solving installation compatibility issues across major operating systems (Windows, MacOS and Linux). In
104 both versions parallel computing is supported which is essential for more efficient ML analysis by
105 distributing the workload across several processors. To promote collaboration, data sharing and support
106 distributed cloud processing, SIMON is also available as a server version. The server version can be
107 installed on a private or a public Linux cloud service. Distributed cloud processing (multiNode) is
108 implemented utilizing OpenStack, a free and open source cloud computing platform
109 (<https://www.openstack.org/>). The advantage of the server versions is that it has multiNode capability
110 which allows users to distribute workload on multiple computers simultaneously to optimize SIMON
111 performance. The multiNode process can be used to horizontally scale analysis to large infrastructure, such
112 as high performance computing clusters to meet the computational needs and accommodate parallel
113 processing of large amounts of data. Additionally, in the server version, users can configure data storage
114 either on a local server or in a cloud using service which is interoperable with Amazon Web Services S3
115 application programming interface (AWS S3 API)⁴³. SIMON is also translated into multiple languages by
116 collaborative open-source effort. SIMON source code is regularly updated, and both source code and
117 compiled software are available from the project's website at <http://www.genular.org/>. Overall, SIMON is
118 designed to provide a uniform knowledge discovery interface adaptable to the increasing size of biomedical
119 datasets allowing data scientists, bioinformaticians and domain-knowledge experts to solve biological
120 research questions.

121 We demonstrate the accuracy, ease of use and power of SIMON on five different biomedical datasets and
122 build predictive models for arboviral infection severity (SISA)⁴⁴, the identification of the cellular immune
123 signature associated with a high-level of physical activity (Cyclists)⁴⁵, the determination of the humoral
124 responses that mediate protection against *Salmonella* Typhi infection (VAST)⁴⁶, early-stage detection of
125 colorectal cancer from microbiome data (Zeller)^{47,48}, and for the detection of liver hepatocellular carcinoma
126 cells (LIHC)⁴⁹ (**Fig. 1 b, c, d, e, Supplementary protocol**). To build models using the SISA dataset
127 (described in the **Supplementary methods** and available as **Supplementary table 2**), 11 ML algorithms
128 were used, five from the original publication⁴⁴ (treebag, k nearest neighbors, random forest, stochastic
129 generalized boosting model and neural network) and additionally, 'sda', shrinkage discriminant analysis;

Tomic et al, SIMON: open-source knowledge discovery platform

130 *'hdda'*, high dimensional discriminant analysis; *'svmLinear2'*, support vector machine with linear kernel;
131 *'pcaNNet'*, neural networks with feature extraction; *'LogitBoost'*, boosted logistic regression and naïve
132 Bayes. Due to the unified ML process for training, tuning and evaluating predictive models, users can test
133 a variety of ML algorithms in SIMON. Since the same training and test sets are used by different algorithms,
134 resulting models can be compared and the best performing models can be selected. After manually setting
135 initial parameters for data partitioning, predictor and outcome variables, exploratory classes, pre-processing
136 and selecting ML algorithms (**Fig. 1a**), SIMON automatically performs all necessary ML analysis steps to
137 build, tune and evaluate predictive models. The process of building all 11 models on the SISA dataset in
138 SIMON finished in 59 sec on a standard laptop (Intel® Core™ i7 Processor 7700HQ and 16 GB of RAM).
139 In SIMON, users can evaluate model performance using standard performance measurements such as
140 accuracy, sensitivity, specificity, precision, recall, area under the receiver operating characteristic curve
141 (AUROC), precision-recall area under curve (prAUC), and logarithmic loss (LogLoss) on training and
142 holdout, test sets (**Fig. 1b**). The shrinkage discriminant analysis model (*'sda'*) had the highest AUROC of
143 0.97 on the training set and also performed well on the holdout, test set (test AUROC 0.96) (**Fig. 1c**,
144 **Supplementary table 3**, the model is available as the **Supplementary data 1**).

145 To demonstrate SIMON's capabilities for analyzing biomedical datasets with missing data, we applied
146 SIMON to (i) the Cyclists dataset studying the impact of physical activity on the immune system in
147 adulthood⁴⁵ (**Supplementary table 4**) and (ii) the VAST dataset collected from a clinical trial which was
148 undertaken to evaluate typhoid vaccine efficacy⁵⁰ (**Supplementary table 5**). Description of both datasets
149 is available in **Supplementary methods**. The percentage of missing values was 8% in the Cyclists dataset
150 and 21% in the VAST dataset either due to the exclusion of samples not passing quality control criteria or
151 the lack of sample volume to repeat experiments and obtain reportable data. To build models using the
152 datasets with missing values, we used the multi-set intersection (*mulset*) function⁴¹ to identify shared
153 features between donors and generate resamples (**Supplementary protocol**). Because *mulset* function
154 generates multiple resamples from the initial dataset based on shared features, it is useful for removal of
155 missing values and can be used for integration of data collected from different assays and across clinical
156 studies⁴¹. For the Cyclists dataset, the *mulset* function generated 146 resamples. The models were built for
157 each of the 146 resamples using five ML algorithms (naïve Bayes, svmLinear2, pcaNNet, logistic
158 regression and hdda) to identify immune cell subsets enriched in the cohort of master cyclists. The analysis
159 finished in 41 min and 24 sec. The model with the highest performance measures was built with naïve
160 Bayes on the resample with 96 donors that shared 31 features (train AUROC 0.99 and test AUROC 1) (**Fig.**
161 **1d, Supplementary table 6 and Supplementary data 2**). The *mulset* function generated 206 resamples
162 from the initial VAST dataset with varying number of donors and features. Resamples with less than 10

Tomic et al, SIMON: open-source knowledge discovery platform

163 donors in the test set were removed prior ML process to prevent too optimistic predictive estimates using
164 the holdout set. Therefore, the ML analysis was performed on 58 resamples using same five ML algorithms
165 as for the Cyclists dataset. The entire analysis finished in 31 min and 1 sec. The top performing model was
166 built on the resample with 47 donors that shared 13 features with the naïve Bayes algorithm (train AUROC
167 0.73 and test AUROC 0.71) (**Fig. 1d, Supplementary table 7 and Supplementary data 3**).

168 We also applied SIMON to (i) a dataset with a large number of features measured using whole-
169 metagenome shotgun sequencing of fecal samples (Zeller dataset, **Supplementary table 8**), and (ii) the
170 liver hepatocellular carcinoma dataset from TCGA with an imbalanced sample distribution of tumor and
171 adjacent normal tissue samples (LIHC dataset, **Supplementary table 9**). Both datasets are described in
172 **Supplementary methods**. For the Zeller dataset, models were built using ML algorithms known to perform
173 well in the situations where more features were measured than individuals, such as shrinkage discriminant
174 analysis⁵¹, high dimensional discriminant analysis⁵² and neural network with feature extraction⁵³. Two
175 additional algorithms were included, svmLinear2 and LogitBoost. The complete analysis was performed in
176 less than 1 min (0:38 min). The sda algorithm built the model with the highest performance (train AUROC
177 0.86 and test AUROC 0.81) having a higher performance measure than the published LASSO linear
178 regression model⁴⁷ (train AUROC 0.84 and test AUROC 0.85) (**Fig. 1e, Supplementary table 10 and**
179 **Supplementary data 4**). For the LIHC dataset we used same five ML algorithms as for Zeller dataset and
180 analysis finished in 11 min and 30 sec. For such highly imbalanced dataset the precision-recall AUC
181 (prAUC)⁵⁴ is a much better performance measurement than AUROC that reported near-perfect performance
182 (**Fig. 1e**). The prAUC provides information how well the model correctly detects cancer cells, while it is
183 less stringent on the evaluation of healthy cells. To avoid obtaining overly optimistic prediction results
184 (often observed on imbalanced datasets), we ranked models based on the prAUC of the training set
185 (**Supplementary table 11**). The model that had the best performance was built using the svmLinear2
186 algorithm (train prAUC 0.83) and it also performed well on the holdout, test set (prAUC 0.73)
187 (**Supplementary data 5**).

188 The *drowsiness* contributed the most to the top-performing SISA model, confirming the findings from
189 the original study⁴⁴ (**Supplementary table 12**). To standardize the process for evaluation of the features
190 and their contribution to the models, we implemented the variable importance score evaluation functions
191 from the caret library³³. This allows users to compare features selected across models. In the case of SISA
192 dataset, *drowsiness* contributed the most in all of the models built (**Supplementary table 13**), indicating
193 the importance of this symptom and its correlation with hospitalization. The features that contributed the
194 most to the Cyclists model were total memory, unswitched memory and naïve B cells, recent thymic
195 emigrants, CD8+ T cells with TEMRA phenotype, and regulatory T cells (CD25+ Foxp3+ CD4+ T cells)

Tomic et al, SIMON: open-source knowledge discovery platform

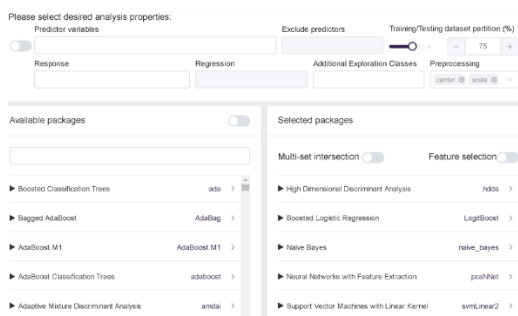
196 **(Supplementary table 14)**. In comparison to age-matched physically inactive individuals (non-cyclists),
197 the master cyclists had increased frequencies of recent thymic emigrants, naïve B cells and CD3 cells, and
198 decreased frequencies of memory B cells and CD8 T cells with TEMRA phenotype, confirming that ageing
199 of immune systems i.e. immunosenescence can be reduced by high levels of physical activity⁴⁵ (**Fig. 1f,**
200 **Supplementary figure 1**). To further explore the relationship between selected features, users can perform
201 correlation analysis to reveal highly correlated features (**Fig. 1g, Supplementary protocol**). Naïve and
202 memory B cells were identified as being highly correlated (**Fig. 1g**), as expected since these subsets were
203 determined from the same flow cytometry plots and their relationship is inversely correlated. Removal of
204 those highly correlated features can help to build more accurate models. Removal of naïve B cells resulted
205 in building predictive model with the same performance measurements as the model built on the entire
206 dataset (train AUROC 0.99 and test AUROC 1) (**Supplementary table 15**), while removal of total memory
207 B cells lowered the accuracy estimates (train AUROC 0.98 and test AUROC 1) (**Supplementary table 16**),
208 indicating the importance of memory B cells to discriminate between master cyclists and non-cyclists. In
209 the VAST dataset, individuals with higher IgA, IgA1, IgA2 and IgG2 titers against native Vi polysaccharide
210 (nViPS) antigen and higher IgA and IgG3 titers against biotinylated Vi polysaccharide (ViBiot) on the day
211 of the challenge were protected against the typhoid challenge supporting the data from univariate analysis⁴⁶
212 (**Supplementary table 17 and Supplementary figure 2**). Moreover, using the clustering function of
213 SIMON's exploratory analysis module, we can quickly identify that the IgA2 signature dominates the
214 responses after vaccination with a purified Vi polysaccharide (Vi-PS), while the IgG2 signature was
215 dominant for the Vi tetanus toxoid conjugate (Vi-TT) vaccine⁴⁶ (**Fig. 1h, Supplementary protocol**). For
216 the Zeller dataset, the same features as originally reported⁴⁰ contributed the most to the model, including
217 *Fusobacterium nucleatum* and *Peptostreptococcus stomatis* (**Supplementary table 18**). The features that
218 contributed the most to the LIHC model were well-known genes identified to be upregulated in LIHC such
219 as *GABRD* and *PLVAP*⁵⁵ and genes enriched in adjacent normal tissue samples *ANGPTL6*⁵⁶, *VIPRI*⁵⁷ and
220 *OIT3*⁵⁸ as a typical signature for healthy liver tissue (**Supplementary table 19, Supplementary figure 3**).

221 Overall, SIMON is a powerful software platform for data mining that facilitates pattern recognition and
222 knowledge extraction from high-quality, heterogenous biological and clinical data, especially where there
223 is missing data, an imbalanced distribution and/or high dimensionality. It can be used for identification of
224 genetic, microbial and immunological correlates of protection and help guiding further analysis of the
225 biomedical data.

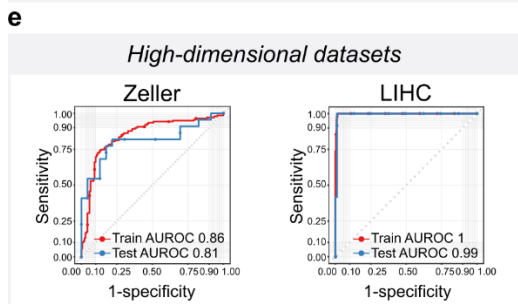
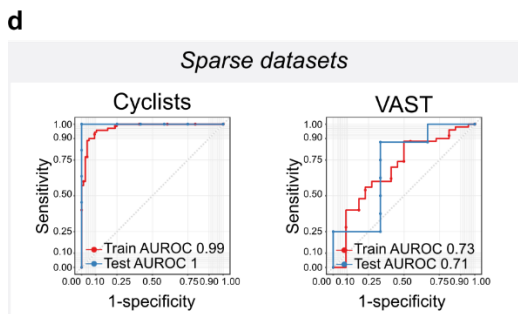
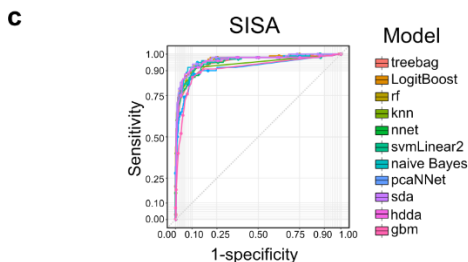
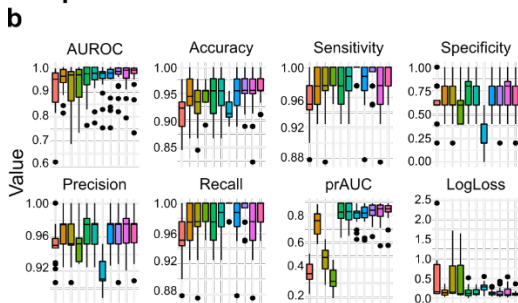
226 **Figure**

Tomic et al, SIMON: open-source knowledge discovery platform

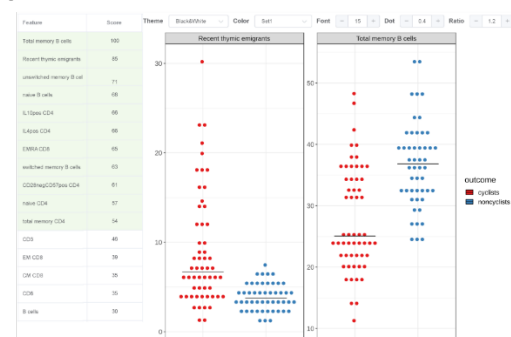
a Step 1. Building predictive models



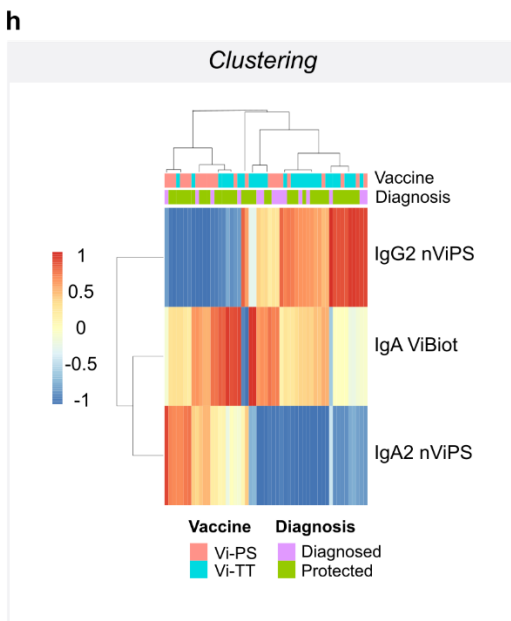
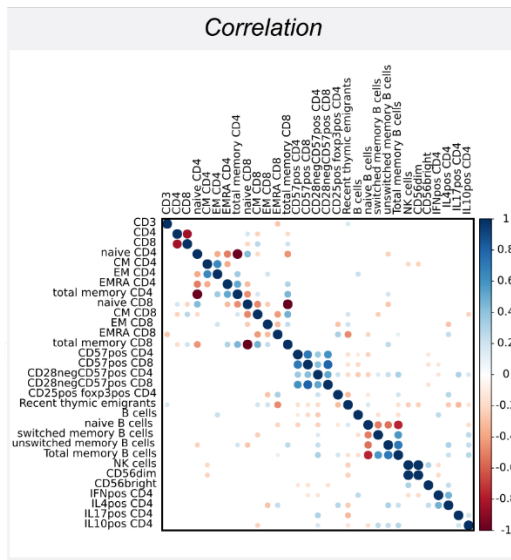
b Step 2. Model evaluation and selection



f Step 3. Feature selection



g Step 4. Exploratory analysis



Tomic et al, SIMON: open-source knowledge discovery platform

228 **Figure 1. SIMON machine learning workflow. Step 1. Building predictive models.** (a) Screenshot of the SIMON
229 graphical user interface demonstrating input selection for machine learning analysis, such as predictors and response
230 (outcome) variables, additional exploration classes, training/test split, preprocessing functions and desired machine
231 learning algorithms. **Step 2. Model evaluation and selection.** Comparison of (b) box plots of performance
232 measurements calculated for 11 predictive models and (c) ROC curves built on the SISA dataset. Comparison of ROC
233 curves calculated from the training and test sets on (d) datasets with missing values (Cyclists and VAST) and (e) high-
234 dimensional datasets (Zeller and LIHC). **Step 3. Feature selection.** (f) Screenshot of the SIMON interface showing
235 variable importance scores calculated for each feature and graphical visualization of the selected features from the
236 Cyclists dataset. **Step 4. Exploratory analysis.** (g) Correlation analysis on the Cyclists dataset. (h) Clustering analysis
237 on the VAST dataset.

Tomic et al, SIMON: open-source knowledge discovery platform

238 **References**

- 239 1 Stuart, T. & Satija, R. Integrative single-cell analysis. *Nature reviews. Genetics* **20**, 257-272 (2019).
- 240 2 Nolan, J. P. & Condello, D. Spectral flow cytometry. *Current protocols in cytometry* **Chapter 1**,
241 Unit1 27 (2013).
- 242 3 Gregori, G. *et al.* Hyperspectral cytometry at the single-cell level using a 32-channel photodetector.
243 *Cytometry. Part A : the journal of the International Society for Analytical Cytology* **81**, 35-44
244 (2012).
- 245 4 Futamura, K. *et al.* Novel full-spectral flow cytometry with multiple spectrally-adjacent fluorescent
246 proteins and fluorochromes and visualization of in vivo cellular movement. *Cytometry. Part A : the*
247 *journal of the International Society for Analytical Cytology* **87**, 830-842 (2015).
- 248 5 Bandura, D. R. *et al.* Mass cytometry: technique for real time single cell multitarget immunoassay
249 based on inductively coupled plasma time-of-flight mass spectrometry. *Analytical chemistry* **81**,
250 6813-6822 (2009).
- 251 6 Bendall, S. C. *et al.* Single-cell mass cytometry of differential immune and drug responses across
252 a human hematopoietic continuum. *Science* **332**, 687-696 (2011).
- 253 7 Angelo, M. *et al.* Multiplexed ion beam imaging of human breast tumors. *Nature medicine* **20**, 436-
254 442 (2014).
- 255 8 Giesen, C. *et al.* Highly multiplexed imaging of tumor tissues with subcellular resolution by mass
256 cytometry. *Nature methods* **11**, 417-422 (2014).
- 257 9 Bishop, C. M. *Pattern Recognition and Machine Learning*. (Springer-Verlag New York, 2006).
- 258 10 Yip, K. Y., Cheng, C. & Gerstein, M. Machine learning and genome annotation: a match meant to
259 be? *Genome biology* **14**, 205 (2013).
- 260 11 Chicco, D. Ten quick tips for machine learning in computational biology. *BioData mining* **10**, 35
261 (2017).
- 262 12 Deo, R. C. Machine Learning in Medicine. *Circulation* **132**, 1920-1930 (2015).
- 263 13 Libbrecht, M. W. & Noble, W. S. Machine learning applications in genetics and genomics. *Nature*
264 *reviews. Genetics* **16**, 321-332 (2015).
- 265 14 Bonetta, R. & Valentino, G. Machine learning techniques for protein function prediction. *Proteins*
266 **88**, 397-413 (2020).
- 267 15 Jurtz, V. *et al.* NetMHCpan-4.0: Improved Peptide-MHC Class I Interaction Predictions Integrating
268 Eluted Ligand and Peptide Binding Affinity Data. *Journal of immunology* **199**, 3360-3368 (2017).
- 269 16 Lin, H. H., Ray, S., Tongchusak, S., Reinherz, E. L. & Brusica, V. Evaluation of MHC class I peptide
270 binding prediction servers: applications for vaccine research. *BMC immunology* **9**, 8 (2008).

Tomic et al, SIMON: open-source knowledge discovery platform

- 271 17 Kan, A. Machine learning applications in cell image analysis. *Immunology and cell biology* **95**,
272 525-530 (2017).
- 273 18 Vamathevan, J. *et al.* Applications of machine learning in drug discovery and development. *Nature*
274 *reviews. Drug discovery* **18**, 463-477 (2019).
- 275 19 Newell, E. W., Sigal, N., Bendall, S. C., Nolan, G. P. & Davis, M. M. Cytometry by time-of-flight
276 shows combinatorial cytokine expression and virus-specific cell niches within a continuum of
277 CD8+ T cell phenotypes. *Immunity* **36**, 142-152 (2012).
- 278 20 Horowitz, A. *et al.* Genetic and environmental determinants of human NK cell diversity revealed
279 by mass cytometry. *Science translational medicine* **5**, 208ra145 (2013).
- 280 21 Chaudhury, S. *et al.* Identification of Immune Signatures of Novel Adjuvant Formulations Using
281 Machine Learning. *Scientific reports* **8**, 17508 (2018).
- 282 22 Chaudhury, S. *et al.* Combining immunoprofiling with machine learning to assess the effects of
283 adjuvant formulation on human vaccine-induced immunity. *Human vaccines &*
284 *immunotherapeutics* **16**, 400-411 (2020).
- 285 23 Warsinske, H. C. *et al.* Assessment of Validity of a Blood-Based 3-Gene Signature Score for
286 Progression and Diagnosis of Tuberculosis, Disease Severity, and Treatment Response. *JAMA*
287 *network open* **1**, e183779 (2018).
- 288 24 Robinson, M. *et al.* A 20-Gene Set Predictive of Progression to Severe Dengue. *Cell reports* **26**,
289 1104-1111 e1104 (2019).
- 290 25 Sweeney, T. E. *et al.* A community approach to mortality prediction in sepsis via gene expression
291 analysis. *Nature communications* **9**, 694 (2018).
- 292 26 Mayhew, M. B. *et al.* A generalizable 29-mRNA neural-network classifier for acute bacterial and
293 viral infections. *Nature communications* **11**, 1177 (2020).
- 294 27 Kourou, K., Exarchos, T. P., Exarchos, K. P., Karamouzis, M. V. & Fotiadis, D. I. Machine learning
295 applications in cancer prognosis and prediction. *Computational and structural biotechnology*
296 *journal* **13**, 8-17 (2015).
- 297 28 Beckmann, J. S. & Lew, D. Reconciling evidence-based medicine and precision medicine in the
298 era of big data: challenges and opportunities. *Genome medicine* **8**, 134 (2016).
- 299 29 R: A language and environment for statistical computing (R Foundation for Statistical Computing,
300 Vienna, Austria, 2013).
- 301 30 Kuhn, M. Building Predictive Models in R Using the caret Package. *Journal of Statistical Software*
302 **28**, 1-26 (2008).
- 303 31 Guttag, J. V. *Introduction to Computation and Programming Using Python: With Application to*
304 *Understanding Data*. Second edn, 472 / 466 (The MIT Press, 2016).
- 305 32 Lang, M. *et al.* mlr3: A modern object-oriented machine learning framework in R. *Journal of Open*
306 *Source Software* **4**, 1903 (2019).

Tomic et al, SIMON: open-source knowledge discovery platform

- 307 33 caret: Classification and Regression Training v. 6.0-80 (R package, 2018).
- 308 34 Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **12**, 2825-2830 (2011).
309
- 310 35 Albanese, D. *et al.* mlp: Machine Learning Python. *arXiv* **1202.6548** (2012).
311 <<https://arxiv.org/abs/1202.6548>>.
- 312 36 Bell, M. L., Fiero, M., Horton, N. J. & Hsu, C. H. Handling missing data in RCTs; a review of the
313 top medical journals. *BMC medical research methodology* **14**, 118 (2014).
- 314 37 Pes, B. Handling Class Imbalance in High-Dimensional Biomedical Datasets. *8th International
315 Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE)*.
316 150-155 (2019 IEEE).
- 317 38 Bellman, R. E. *Dynamic programming*. (Princeton University Press, 1957).
- 318 39 Automated Machine Learning: Methods, Systems, Challenges in *The Springer Series on
319 Challenges in Machine Learning* (eds F. Hutter, L. Kotthoff, & J. Vanschoren) (Springer, 2018).
320 <<http://automl.org/book>>.
- 321 40 Thornton, C., Hutter, F., Hoos, H. H. & Leyton-Brown, K. Auto-WEKA: combined selection and
322 hyperparameter optimization of classification algorithms. *KDD: Knowledge Discovery and Data
323 Mining 2013*. 847–855.
- 324 41 Tomic, A. *et al.* SIMON, an Automated Machine Learning System, Reveals Immune Signatures of
325 Influenza Vaccine Responses. *J Immunol* **203**, 749-759 (2019).
- 326 42 Tusher, V. G., Tibshirani, R. & Chu, G. Significance analysis of microarrays applied to the ionizing
327 radiation response. *Proceedings of the National Academy of Sciences of the United States of
328 America* **98**, 5116-5121 (2001).
- 329 43 Murty, J. *Programming Amazon Web Services: S3, EC2, SQS, FPS, and SimpleDB*. 604 (O'Reilly
330 Media, 2009).
- 331 44 Sippy, R. *et al.* Severity Index for Suspected Arbovirus (SISA): Machine learning for accurate
332 prediction of hospitalization in subjects suspected of arboviral infection. *PLoS neglected tropical
333 diseases* **14**, e0007969 (2020).
- 334 45 Duggal, N. A., Pollock, R. D., Lazarus, N. R., Harridge, S. & Lord, J. M. Major features of
335 immunesenescence, including reduced thymic output, are ameliorated by high levels of physical
336 activity in adulthood. *Aging cell* **17** (2018).
- 337 46 Dahora, L. C. *et al.* IgA and IgG1 Specific to Vi Polysaccharide of Salmonella Typhi Correlate
338 With Protection Status in a Typhoid Fever Controlled Human Infection Model. *Frontiers in
339 immunology* **10**, 2582 (2019).
- 340 47 Zeller, G. *et al.* Potential of fecal microbiota for early-stage detection of colorectal cancer.
341 *Molecular systems biology* **10**, 766 (2014).

Tomic et al, SIMON: open-source knowledge discovery platform

- 342 48 Pasolli, E. *et al.* Accessible, curated metagenomic data through ExperimentHub. *Nature methods*
343 **14**, 1023-1024 (2017).
- 344 49 Geistlinger, L. *et al.* Toward a gold standard for benchmarking gene set enrichment analysis.
345 *Briefings in bioinformatics* (2020).
- 346 50 Jin, C. *et al.* Efficacy and immunogenicity of a Vi-tetanus toxoid conjugate vaccine in the
347 prevention of typhoid fever using a controlled human infection model of Salmonella Typhi: a
348 randomised controlled, phase 2b trial. *The Lancet* **390**, 2472–2480 (2017).
- 349 51 Mkhadri, A. Shrinkage parameter for the modified linear discriminant analysis. *Pattern*
350 *Recognition Letters* **16**, 267-275 (1995).
- 351 52 Bouveyron, C., Girard, S. & Schmid, C. High-Dimensional Discriminant Analysis.
352 *Communications in Statistics - Theory and Methods* **36**, 2607–2623 (2007).
- 353 53 Ripley, B. D. *Pattern Recognition and Neural Networks*. (Cambridge University Press, 1996).
- 354 54 Davis, J. & Goadrich, M. The Relationship Between Precision-Recall and ROC Curves.
355 *Proceedings of the 23rd International Conference on Machine Learning*. (2006).
- 356 55 Sarathi, A. & Palaniappan, A. Novel significant stage-specific differentially expressed genes in
357 hepatocellular carcinoma. *BMC cancer* **19**, 663 (2019).
- 358 56 Oike, Y. *et al.* Angiopoietin-related growth factor antagonizes obesity and insulin resistance. *Nat*
359 *Med* **11**, 400-408 (2005).
- 360 57 Lu, S., Lu, H., Jin, R. & Mo, Z. Promoter methylation and H3K27 deacetylation regulate the
361 transcription of VIPR1 in hepatocellular carcinoma. *Biochemical and biophysical research*
362 *communications* **509**, 301-305 (2019).
- 363 58 Xu, Z. G. *et al.* A novel liver-specific zona pellucida domain containing protein that is expressed
364 rarely in hepatocellular carcinoma. *Hepatology* **38**, 735-744 (2003).
- 365