

RESEARCH

# An elastic-net logistic regression approach to generate classifiers and gene signatures for types of immune cells and T helper cell subsets

Arezo Torang<sup>1</sup>, Paraag Gupta<sup>1</sup> and David J. Klinke II<sup>1,2\*</sup>

\*Correspondence:

david.klinke@mail.wvu.edu

<sup>1</sup>Department of Chemical and Biomedical Engineering, West Virginia University, 1306 Evansdale Dr, 26506 Morgantown, WV, USA  
Full list of author information is available at the end of the article

## Abstract

**Background:** Host immune response is coordinated by a variety of different specialized cell types that vary in time and location. While host immune response can be studied using conventional low-dimensional approaches, advances in transcriptomics analysis may provide a less biased view. Yet, leveraging transcriptomics data to identify immune cell subtypes presents challenges for extracting informative gene signatures hidden within a high dimensional transcriptomics space characterized by low sample numbers with noisy and missing values. To address these challenges, we explore using machine learning methods to select gene subsets and estimate gene coefficients simultaneously.

**Results:** Elastic-net logistic regression, a type of machine learning, was used to construct separate classifiers for ten different types of immune cell and for five T helper cell subsets. The resulting classifiers were then used to develop gene signatures that best discriminate among immune cell types and T helper cell subsets using RNA-seq datasets. We validated the approach using single-cell RNA-seq (scRNA-seq) datasets, which gave consistent results. In addition, we classified cell types that were previously unannotated. Finally, we benchmarked the proposed gene signatures against other existing gene signatures.

**Conclusions:** Developed classifiers can be used as priors in predicting the extent and functional orientation of the host immune response in diseases, such as cancer, where transcriptomic profiling of bulk tissue samples and single cells are routinely employed. Information that can provide insight into the mechanistic basis of disease and therapeutic response. The source code and documentation are available through GitHub: <https://github.com/KlinkeLab/ImmClass2019>.

**Keywords:** Immune Cells; Gene Signature; Machine Learning; Elastic-Net

1

2

### 3 **Background**

4 Host immune response is a coordinated complex system, consisting of different spe-  
5 cialized cell types of innate and adaptive immune cells that vary dynamically and  
6 in different anatomical locations. As shown in Fig. 1, innate immune cells comprise  
7 myeloid cells including eosinophils, neutrophils, basophils, monocytes, and mast  
8 cells. Adaptive immune cells are mainly B lymphocytes and T lymphocytes that  
9 specifically recognize different antigens [1]. Linking innate with adaptive immunity  
10 are antigen presenting cells, like macrophages and dendritic cells, and Natural Killer  
11 cells. Traditionally, unique cell markers have been used to characterize and separate  
12 different immune cell subsets from heterogeneous cell mixtures using flow cytom-  
13 etry [2, 3, 4]. However, flow cytometry measures on the order of 10 parameters  
14 simultaneously and relies on prior knowledge for selecting relevant molecular mark-  
15 ers, which could provide a biased view of the immune state within a sample [5].  
16 Recent advances in technology, like mass cytometry or multispectral imaging, have  
17 expanded the number of molecular markers, but the number of markers used for  
18 discriminating among cell types within a sample remains on the order of  $10^{1.5}$ .

19 In the recent years, quantifying tumor immune contexture using bulk transcrip-  
20 tomics or single-cell RNA sequencing data (scRNA-seq) has piqued the interest of  
21 the scientific community [6, 7, 8, 9, 10]. Advances in transcriptomics technology,  
22 like RNA sequencing, provide a much higher dimensional view of which genes are  
23 expressed in different immune cells (i.e., on the order of  $10^3$ ) instead of focusing  
24 on a small number of genes [11]. Conceptually, inferring cell types from data us-  
25 ing an expanded number of biologically relevant genes becomes more tolerant to  
26 non-specific noise and non-biological differences among samples and platforms. In  
27 practice, cell types can be identified using gene signatures, which are defined as  
28 sets of genes linked to common downstream functions or inductive networks that  
29 are co-regulated [12, 13], using approaches such as Gene Set Enrichment Analy-  
30 sis (GSEA) [12]. However, as microarray data can inflate detecting low abundance  
31 and noisy transcripts and scRNA-seq data can have a lower depth of sequencing,  
32 opportunities for refining methods to quantify the immune contexture using gene  
33 signatures still remain.

34 Leveraging transcriptomics data to identify types of immune cells presents an-  
35 alytic challenges for extracting informative gene signatures hidden within a high

36 dimensional transcriptomics space that is characterized by low sample numbers  
37 with noisy and missing values. Typically, the number of cell samples is in the range  
38 of hundreds or less, while the number of profiled genes is in the tens thousands [14].  
39 Yet, only a few number of genes are relevant for discriminating among immune cell  
40 subsets. Datasets with a large number of noisy and irrelevant genes decrease the  
41 accuracy and computing efficiency of machine learning algorithms, especially when  
42 the number of samples are very limited. Hence, it is essential to use feature selec-  
43 tion algorithms to reduce redundant genes [15]. The application of feature selection  
44 methods enables developing gene signatures in different biomedical fields of study  
45 [16]. There are many proposed feature selection methods to select gene sets with the  
46 properties that enable high accuracy classification. In recent years, regularization  
47 methods have become more popular, which efficiently select features [17] and also  
48 control for overfitting [18]. As a machine learning tool, logistic regression is consid-  
49 ered to be a powerful discriminative method [18]. However, logistic regression alone  
50 is not applicable for high-dimensional cell classification problems [19]. Regularized  
51 logistic regression, in the other hand, has been shown to be successfully applicable  
52 for high-dimensional problems [20]. Regularized logistic regression selects a small  
53 set of genes with strongest effects on the cost function [17]. A regularized logistic  
54 regression can be applied with different regularization terms. The most popular  
55 regularized terms are LASSO, Ridge [21], and elastic-net [22] which impose the  
56  $l_1$  norm,  $l_2$  norm, and linear combination of  $l_1$  norm and  $l_2$  norm regularization,  
57 respectively, to the cost function. It has been shown that, specially in very high  
58 dimensional problems, elastic-net outperforms LASSO and Ridge [17, 22].

59 In this study, we focused on two-step regularized logistic regression techniques to  
60 develop immune cell signatures and immune cell and T helper cell classifiers using  
61 RNA-seq data for the cells highlighted in bold in Fig. 1. The first step of the process  
62 included a pre-filtering phase to select the optimal number of genes and implemented  
63 an elastic-net model as a regularization method for gene selection in generating the  
64 classifiers. The pre-filtering step reduced computational cost and increased final  
65 accuracy by selecting the most discriminative and relevant set of genes. Finally, we  
66 illustrate the value of the approach in annotating gene expression profiles obtained  
67 from single-cell RNA sequencing. The second step generated gene signatures for

68 individual cell types using selected genes from first step and implemented a binary  
69 regularized logistic regression for each cell type against all other samples.

## 70 **Results**

71 We developed classifiers for subsets of immune cells and T helper cells separately  
72 with two main goals. First, we aimed to annotate RNA-seq data obtained from an  
73 enriched cell population with information as to the immune cell identity. Second, we  
74 developed gene signatures for different immune cells that could be used to quantify  
75 the prevalence from RNA-seq data obtained from a heterogeneous cell population.  
76 Prior to developing the classifiers, the data was pre-processed to remove genes that  
77 have low level of expression for most of samples (details can be found in Meth-  
78 ods section) and normalized to increase the homogeneity in samples from different  
79 studies and to decrease dependency of expression estimates to transcript length  
80 and GC-content. Genes retained that had missing values for some of the samples  
81 were assigned a values of -1. Next, regularized logistic regression (elastic-net) was  
82 performed and the optimal number of genes and their coefficients were determined.

### 83 **Generating and validating an immune cell classifier**

84 In development of the immune cell classifier, we determined the optimal number of  
85 genes in the classifier by varying the lambda value used in the regularized logistic  
86 regression of the training samples and assessing performance. To quantify the perfor-  
87 mance using different lambdas, a dataset was generated by combining true-negative  
88 samples, which were created by randomly scrambling associated genes and their  
89 corresponding value from the testing datasets, with the original testing data, which  
90 were untouched during training and provided true-positive samples. The accuracy  
91 of predicting the true-Positive samples were used to generate Receiver Operating  
92 Characteristic (ROC) curves (Fig. 2a). Performance using each lambda was quan-  
93 tified as the Area Under the ROC Curve (AUC).

94 The optimal lambda for immune cell classifier was the smallest value (i.e., highest  
95 number of genes) that maximized the AUC. Functionally, this lambda value repre-  
96 sents the trade-off between retaining the most possible number of informative genes  
97 (i.e., classifier signal) in the first step for developing the gene signature later, while  
98 not adding non-informative genes (i.e., classifier noise). Consequently, we selected

99 a lambda value of  $1e-4$  (452 genes) for the immune cell classifier, where the selected  
100 genes and their coefficients are shown in Table S1.

101 To explore correlations between the weights of selected genes with their expression  
102 level, we generated heatmaps shown in Fig. 2, panels b and c. A high level of gene  
103 expression is reflected as a larger positive coefficient in a classifier model, while  
104 low or absent expression results in a negative coefficient. This is interpreted as, for  
105 example, if gene A is not in cell type 1, the presence of this gene in a sample decreases  
106 the probability for that sample to be cell type 1. For instance, E-cadherin (CDH1)  
107 was not detected in almost all monocyte samples and thus has a negative coefficient.  
108 Conversely, other genes are only expressed in certain cell types, which results in a  
109 high positive coefficient. For instance, CYP27B1, INHBA, IDO1, NUPR1, and UBD  
110 are only expressed by M1 macrophages and thus have high positive coefficients.

111 The differential expression among cell types suggests that the set of genes in-  
112 cluded in the classifier model may also be a good starting point for developing  
113 gene signatures, which is highlighted in Fig. 2d. Here, we focused on the expres-  
114 sion of the 452 genes included in the classifier model and the correlations between  
115 samples clustered based on cell types. The off-diagonal entries in the correlation  
116 matrix are colored by euclidean distance values with the color indicating similarity  
117 between sample pairs (similar: pink versus dissimilar: blue) and color bars along the  
118 axes highlight the cell types for the corresponding RNA-seq samples. As expected,  
119 RNA-seq samples from the same cell type were highly similar. More interestingly,  
120 correlation between different cell types can also be seen, like high similarity between  
121 CD4+ and CD8+ T cell samples, CD8+ T cell and NK cell samples, and monocyte  
122 and dendritic cell samples. Collectively, these heatmaps illustrate that the selected  
123 genes are a highly condensed but still representative set of genes that include main  
124 characteristics of the immune cell types. It is also notable to compare the clustering  
125 result of cell types based on their coefficients in the classifier shown in Fig. 2b with  
126 similarity matrix in Fig. 2d. Since in the classifier coefficients are forcing the model  
127 to separate biologically close cell types (like CD4+ T cell and CD8+ T cell), the  
128 resulted clustering did not find them in close relationship (Fig. 2b). However, in the  
129 case of their expression values, their similarity is remains (Fig. 2d).

130 *Evaluating the Immune Cell classifier using scRNA-seq datasets*

131 To evaluate the proposed classifier in immune cell classification, two publicly ac-  
132 cessible datasets generated by scRNA-seq technology were used [23, 24]. The first  
133 dataset reported by [23] included malignant, immune, stromal and endothelial cells  
134 from 15 melanoma tissue samples. We focused on the immune cell samples, which  
135 included 2761 annotated samples of T cells, B cells, *Mphi* and NK cells, and 294  
136 unresolved samples. The immune cells in this study were recovered by flow cytom-  
137 etry by gating on CD45 positive cells. Annotations were on the basis of expressed  
138 marker genes while unresolved samples were from the CD45-gate and classified as  
139 non-malignant based on inferred copy number variation (CNV) patterns (i.e., CNV  
140 score  $< 0.04$ ).

141 Following a pre-processing step to filter and normalize the samples similar to  
142 the training step, the trained elastic-net logistic regression model was used to clas-  
143 sify cells into one of the different immune subsets based on the reported scRNA-seq  
144 data with the results summarized in Fig. 3a. The inner pie chart shows the prior cell  
145 annotations reported by [23] and the outer chart shows the corresponding cell anno-  
146 tation predictions by our proposed classifier. Considering T cells as either CD4+ T  
147 cell or CD8+ T cell, the overall similarity between annotations provided by [23] and  
148 our classifier prediction is 96.2%. The distribution in cells types contained within  
149 the unresolved samples seemed to be slightly different than the annotated samples  
150 as we predicted the unresolved samples to be mainly CD8+ T cells and B cells.

151 The only cell type with low similarity between our classifier predictions and prior  
152 annotations was NK cells, where we classified almost half of samples annotated  
153 previously as NK cells as CD8+ T cell. Discriminating between these two cell types  
154 is challenging as they share many of the genes related to cytotoxic effector function  
155 and can also be subclassified into subsets, like CD56bright and CD56dim NK subsets  
156 [25]. To explore this discrepancy, we compared all annotated samples based on their  
157 CD8 score and NK score provided by the classifier, as shown in Fig. 3b. Although  
158 the number of NK cell samples are relatively low, it seems that the NK samples  
159 consist of two groups of samples: one with a higher likelihood of being a NK cell  
160 and a second with almost equal likelihood for being either CD8+ T cell or NK cell.  
161 We applied principal component analysis (PCA) to identify genes associated with  
162 this difference and used Enrichr for gene set enrichment [26, 27]. Using gene sets

163 associated with the Human Gene Atlas, the queried gene set was enriched for genes  
164 associated with CD56 NK cells, CD4+ T cell and CD8+ T cell. Collectively, the  
165 results suggests that the group of cells with similar score for NK and CD8 in the  
166 classifier model are Natural Killer T cells.

167 We also analyzed a second dataset that included 317 epithelial breast cancer  
168 cells, 175 immune cells and 23 non-carcinoma stromal cells, from 11 patients di-  
169 agnosed with breast cancer [24]. We only considered samples annotated previously  
170 as immune cells, which were annotated as T cells, B cells, and myeloid samples  
171 by clustering the gene expression signatures using non-negative factorization. The  
172 scRNA-seq samples were similarly pre-processed and analyzed using the proposed  
173 classifier, with the results shown in Fig. 4. The inner pie chart shows the prior cell  
174 annotations reported by [24] and the outer chart shows the corresponding predicted  
175 cell annotation by our proposed classifier. Considering T cells as either CD4+ T  
176 cell or CD8+ T cell, 94.4% of reported T cells are predicted as the same cell type  
177 and other 5.6% is predicted to be DC or NK cells. However, for reported B cells  
178 and myeloid cells, we predicted relatively high portion of samples to be T cells (  
179 15.7% of B cells and 40% of myeloid cells). The rest of the myeloid samples were  
180 predicted to be macrophages or dendritic cells. Collectively, our proposed classifier  
181 agreed with many of the prior cell annotations and annotated many of the samples  
182 that were previously unresolved.

### 183 Developing a classifier for T Helper cell subsets

184 Similar to the immune cell classifier, we next wanted to generate a classifier to dis-  
185 tinguish among T helper cells and applied regularized logistic regression to corre-  
186 sponding training samples. We explored different values of the regression parameter  
187 lambda to find the optimal number of genes. To visualize the performance of differ-  
188 ent lambdas, we generated True-Negative samples by randomly scrambling testing  
189 datasets. Original testing data that were completely untouched during training were  
190 used as True-Positive samples. The True-Negative and True-Positive samples were  
191 used to generate ROC curves (Fig. 5a) and the AUC was used to score each lambda  
192 value. Generally, the lambda values for T helper cell classifier represents the trade-  
193 off between retaining genes and keeping the AUC high. However, there appeared to  
194 be an inflection point at a lambda value of 0.05 whereby adding additional genes,

195 by increasing lambda, reduced the AUC. Consequently, we selected a lambda value  
196 equal to 0.05 (72 genes) for T helper classifier. The selected genes and their coeffi-  
197 cients are listed in Table S1. The gene list was refined subsequently by developing  
198 a gene signature.

199 Similar to the immune cell classifier, the coefficients of the selected genes for the T  
200 helper cell classifier correlated with their expression levels, as seen by comparing the  
201 heatmaps shown in Fig. 5, panels b and c. For instance, FUT7 has been expressed in  
202 almost all T helper cell samples except for iTreg that result in a negative coefficient  
203 for this cell type. In addition, there are sets of genes for each cell type that have large  
204 coefficients only for certain T helper cell subsets, like ALPK1, TBX21, IL12RB2,  
205 IFNG, RNF157 for Th1 that have low expression in other cells. As illustrated in  
206 Fig. 5d, the genes included in the classifier don't all uniquely associate with a  
207 single subset but collectively enable discriminating among T helper cell subsets.  
208 Interestingly, the T helper subsets stratified into two subgroups where naive T  
209 helper cells (Th0) and inducible T regulatory (iTreg) cells were more similar than  
210 effector type 1 (Th1), type 2 (Th2), and type 17 (Th17) T helper cells. Similar  
211 to the immune cell classifier, we also noted that the clustering of the classifier  
212 coefficients is different from what similarity matrix shows in Fig. 5d because the  
213 classifier coefficients aim to create a "classifying distance" among closely related  
214 cell types.

215 Finally by comparing the results of immune cell classifier with that of the T helper  
216 classifier, the intensity of differences among cell types can be seen in Fig. 2c and  
217 Fig. 5c. In the first figure you can find completely distinct set of genes in each cell  
218 type while in the second figure the gene sets are not as distinct which could be due  
219 to either the few number of samples or high biological similarity between T helper  
220 cell types.

## 221 Application of the Classifiers

222 Clinical success of immune checkpoint inhibitors (ICI) for treating cancer coupled  
223 with technological advances in assaying the transcriptional signatures in individual  
224 cells, like scRNA-seq, has invigorated interest in characterizing the immune contex-  
225 ture within complex tissue microenvironments, like cancer. However as illustrated  
226 by the cell annotations reported by [24], identifying immune cell types from noisy



227 scRNA-seq signatures using less biased methods remains an unsolved problem. To  
228 address this problem, we applied our newly developed classifiers to characterize the  
229 immune contexture in melanoma and explored differences in immune contexture  
230 that associated with immune checkpoint response. Of note, some melanoma pa-  
231 tients respond to ICIs durably but many others show resistance [28]. Specifically,  
232 we annotated immune cells in the melanoma scRNA-seq datasets [23, 29] using  
233 our classifiers separately for each patient sample and ordered samples based on the  
234 treatment response, with the results shown in Fig. 6a, b. We used the percentage  
235 of cell type for each tumor samples as it is more informative and meaningful than  
236 using absolute cell numbers. It is notable that untreated and NoInfo samples likely  
237 include both ICI-resistant and ICI-sensitive tumors.

238 In comparing samples from resistant tumors to untreated tumors, we found in-  
239 terestingly that there are samples with high prevalence of NK in untreated tumors  
240 (Mel53, Mel81, and Mel82) while no samples in resistant tumors have a high preva-  
241 lence of NK cells. The mentioned untreated tumors also have no or very low number  
242 of Th2 cells in their populations. In addition, untreated tumors have a more uni-  
243 form distribution of immune cell types in contrast to ICI-resistant ones, which could  
244 reflect a therapeutic bias in immune cell prevalence in the tumor microenvironment  
245 due to ICI treatment.

246 Next, we combined the annotation data from both classifiers and applied PCA  
247 and clustering analysis, as shown in Fig. 6, panels c and d. Using scrambled data  
248 to determine principal components and their associated eigenvalues that are not  
249 generated by random chance (i.e., a negative control), we kept the first and second  
250 principal components that capture 68% and 21% of the total variance, respectively,  
251 and neglected other components that fell below the negative control of 8.4%. As it  
252 shown in 6c, resistant samples mainly located in lowest value of second principal  
253 component (PC2). Upon closer inspection of the cell loadings within the eigen-  
254 vectors, the low values of PC2 corresponds to a low prevalence of  $M\phi$  or high  
255 percentage of B cells. In addition, based on the first principal component (PC1),  
256 resistant samples have either lowest values of PC1 (Mel74, Mel75, Mel58, Mel 78)  
257 which correspond to higher than average prevalence of CD8+ T cells or highest  
258 values of PC1 (Mel60, Mel72, Mel94) that show higher than average prevalence of  
259 B cells.

260 In hierarchical clustering, the optimal number of clusters was selected based on cal-  
261 culation of different cluster indices using the NbClust R package [30] which mainly  
262 identified two or three clusters as the optimal number. In considering three group-  
263 ings of the hierarchical clustering results shown in 6d, seven out of eight ICI-resistant  
264 samples clustered in first two clusters while the third cluster mainly contained un-  
265 treated samples. The comparison of results from PCA and clustering analyses shows  
266 that the first cluster contained samples with extreme low value of PC1 which itself  
267 divided into two groups; one with extreme low value of PC2 and the other with  
268 higher amount of PC2. The second cluster located in highest amount of PC1 and  
269 lowest amount of PC2. All remained samples were clustered as third group, which  
270 were predominantly untreated samples. The difference in clustering suggests dissim-  
271 ilarities between ICI-resistant and untreated samples and the possibility of having  
272 ICI-sensitive tumors in untreated samples.

### 273 Developing Gene Signatures

274 While classifiers are helpful for annotating scRNA-seq data as the transcriptomic  
275 signature corresponds to a single cell, gene signatures are commonly used to deter-  
276 mine the prevalence of immune cell subsets within transcriptomic profiles of bulk  
277 tissue samples using deconvolution methods. Leveraging the classifier results, we  
278 generated corresponding gene signatures using binary elastic-net logistic regression.  
279 Specifically, classifier genes with non-zero coefficients were used as initial features of  
280 the models, which were regressed to the same training and testing datasets as used  
281 for developing the classifiers. Lambda values were selected for each immune and T  
282 helper cell subset based on similar method of lambda selection for classifiers and  
283 their values and corresponding AUC are shown in Table S2. Finally, all generated  
284 signatures are summarized in Table S3.

285 We visualized the expression levels of remained set of genes, which at least occur  
286 in one gene signature, in Fig. 7. The expression of genes retained in immune cell sig-  
287 natures (Fig. 7a) and T helper cell signatures (Fig. 7b) were clustered by similarity  
288 in expression (rows) and by similarity in sample (columns). For both immune and  
289 T helper cell subsets, samples of same cell type were mainly clustered together. The  
290 only exception is for macrophages ( $M\phi$  and M2) which can be attributed to high  
291 biological similarity and a low number of technical replicates for these cell types.

292 In general, the gene set generated from the logistic regression model performed  
293 well with far fewer requisite genes in the testing set, a desirable result for a gene  
294 set intended to be used for immunophenotyping. In Fig. 8, the results of the bench-  
295 marking are shown separated by comparative gene set. Both the CIBERSORT and  
296 Single-Cell derived gene sets contain an average of 64 and 135 genes, respectively,  
297 while the logistic regression gene set contains an average of just 19. The new lo-  
298 gistic regression gene set performed comparably to the existing contemporary gene  
299 sets and far exceeded the performance of the manually curated gene set used previ-  
300 ously [6]. The benchmarking results indicate that logistic regression gene set is an  
301 improvement in efficacy over compact gene sets, such as those that are manually  
302 annotated or hand-picked. Meanwhile, the logistic regression gene set also demon-  
303 strates an optimization of broader gene sets that contain too many genes for deep  
304 specificity when used in further analysis. The inclusion of too many genes in the set  
305 can dilute the real data across a constant level of noise, while including too few lacks  
306 the power to draw conclusions with high confidence. The logistic regression gene  
307 set demonstrates a balance of these two issues through its highly refined selection  
308 of genes that can be fine-tuned using its lambda parameter.

## 309 Discussion

310 Recent developments in RNA sequencing enable a high fidelity view of the tran-  
311 scriptomic landscape associated with host immune response. Despite considerable  
312 progress in parsing this landscape using gene signatures, gaps remain in developing  
313 unbiased signatures for individual immune cell types from healthy donors using high  
314 dimensional RNA-seq data. Here, we developed two classifiers - one for immune cell  
315 subsets and one for T helper cell subsets - using elastic-net logistic regression with  
316 cross validation. The features of these classifiers have been used as starting point for  
317 generation of gene signatures captured with fifteen binary elastic-net logistic regres-  
318 sion models as the most relevant gene sets to distinguish among different immune  
319 cell types without making too much noise.

320 Gene signatures in previous studies have been developed and used mainly as  
321 a base for deconvolution of tumor microenvironment and to find the fractions of  
322 existing immune cells. Therefore, as the first step, determining cell-specific gene  
323 signatures critically influences the results of deconvolution methods [31]. Newman

324 et al. defined gene signatures for immune cells using two-sided unequal variances  
325 t-test as base matrix for CIBERSORT [8]. In another study, Li et al. in developing  
326 TIMER, generated gene signatures for six immune cell types with selecting genes  
327 with expression levels that have a negative correlation with tumor purity [9]. More  
328 recently, Racle et al. developed a deconvolution tool based on RNA-seq data (EPIC)  
329 by pre-selecting genes based on ranking by fold change and then selected genes  
330 by manually curating and comparing the expression levels in blood and tumor  
331 microenvironment [10]. Finally, quanTIseq (the most recently developed tool for  
332 deconvolution) has been developed for RNA-seq data based on the gene signatures  
333 generated by quantizing the expression levels into different bins and selecting high  
334 quantized genes for each cell type that have low or medium expression in other cell  
335 types [7]. Although all methods obtained high accuracy based on their developed  
336 signatures, a more rigorous and unbiased gene signature developed by RNA-seq  
337 data and precise feature selection methods can be used to improve the accuracy  
338 even further and validate the process for downstream analyses.

339 In addition, to identify cell types based on their transcriptome, clustering tech-  
340 niques have been used in many studies [32, 33]. However, there are high variability  
341 levels of gene expression even in samples from the same cell type. Moreover, tran-  
342 scriptomics data has high dimensions (tens of thousands) and this is too complicated  
343 for clustering techniques specially because only few number of genes are discrimi-  
344 native. To overcome these problems some studies used supervised machine learning  
345 methods like Support Vector Machine (SVM) [34, 35]. However, to the best of our  
346 knowledge, this paper is the first to apply two-step regularized logistic regression  
347 on RNA-seq transcriptomic of immune cells. This method increases the chance to  
348 capture the most discriminative set of genes for each cell type based on the power  
349 of an elastic-net [22]. In addition, using a two-step elastic net logistic regression  
350 enabled eliminating the most irrelevant genes while keeping the most possible sig-  
351 nificant genes in the first step and more deeply selecting among them in the second  
352 step to generate robust gene signatures for immune cells.

353 Moreover, contemporary methods have only considered a limited number of im-  
354 mune cell types, and specifically T helper subsets as individual cell types have been  
355 neglected [23, 29, 24] in comprehensive studies. Therefore, the other novel aspect  
356 of this study is the separation of models for immune cells and T helper cells and

357 development of gene signatures for vast number of immune cell types (fifteen differ-  
358 ent immune cell types) including different T helper cell subsets. This can be used  
359 to study immune system in different diseases in more depth. As we used publicly  
360 available RNA-seq datasets for immune cells and T helper cells, we acknowledge  
361 that our developed classifiers and gene signatures may be still constrained by the  
362 limited number of samples specifically for T helper cells. As more data describing  
363 the transcriptome of for immune cells will become accessible, one can update the  
364 classifiers and gene signatures. Despite the limited number of samples used in the  
365 approach, the developed classifiers can even be applied to completely untouched  
366 and large datasets [23, 24] that have been generated using scRNA-Seq technology  
367 which creates noisier data.

## 368 **Conclusions**

369 Here, we developed an immune cell classifier and classifier for T helper cell subsets  
370 along with gene signatures to distinguish among fifteen different immune cell types.  
371 Elastic-net logistic regression was used to generate classifiers with 10-fold cross-  
372 validation after normalizing and filtering two separate RNA-seq datasets that were  
373 generated using defined homogeneous cell populations. Subsequently, we generated  
374 gene signatures using a second step of binary regularized logistic regression applied  
375 to the RNA-seq data using previously selected classifier genes. As an external val-  
376 idation, the resulting classifiers accurately identified the type of immune cells in  
377 scRNA-seq datasets. Our classifiers and gene signatures can be considered for a  
378 different downstream applications. First, the classifiers may be used to detect the  
379 type of immune cells in under explored bulks and to verify uncertainly annotated  
380 immune cells. Second, the gene signatures could be used to study tumor micro-  
381 environments and the connections of immune systems with cancer cells, which is  
382 emerging to be an important clinical question.

## 383 **Methods**

### 384 **Data Acquisition**

385 RNA-seq datasets for 15 different immune cell types including T helper cells, were  
386 obtained from ten different studies [36, 37, 38, 39, 40, 41, 42, 43, 44, 45] which were  
387 publicly accessible as part of *Gene Expression Omnibus* [46]. The list of samples  
388 is provided as Supplementary Table S1. Cell types divided into two groups: the

389 immune cells includes B cells, CD4+ and CD8+ T cells, monocytes (Mono), neu-  
390 trophils (Neu), natural killer (NK) cells, dendritic cells (DC), macrophage ( $M\phi$ ),  
391 classically (M1) and alternatively (M2) activated macrophages, and the T helper  
392 cells includes Th1, Th2, Th17, Th0, and Regulatory T cells (Treg). The goal was  
393 to train the gene selection model on immune cell types, and CD4+ T cell subsets  
394 (T helper cells), separately. As if these two groups of cells are analyzed together,  
395 many of the genes that potentially could be used to discriminate among T helper  
396 cell subsets might be eliminated as they overlap with genes associated with CD4+  
397 T cells.

398 In short, a total of 233 samples were downloaded and divided into two sets of  
399 185 and 48 samples, for immune cells and T helper cells, respectively. Moreover,  
400 immune cell samples have been further divided into 108 training and 77 testing  
401 samples. Numbers for T helper samples are 31 and 17, respectively. Training and  
402 testing data include samples from all studies. For a verification dataset, scRNA-  
403 seq data derived from CD45+ cell samples obtained from breast cancer [24] and  
404 melanoma [23] were used with GEO accession numbers of GSE75688 and GSE72056,  
405 respectively.

#### 406 Data Normalization

407 The expression estimates provided by the individual studies were used, regardless  
408 of the underlying experimental and data processing methods (Table S1). For devel-  
409 oping individual gene signatures and cell classification models, we did not use raw  
410 data due to sample heterogeneity such as different experimental methods and data  
411 processing techniques used by different studies as well as differences across biolog-  
412 ical sources. Rather, we applied a multistep normalization process before training  
413 models. To eliminate obvious insignificant genes from our data, for immune cell  
414 samples, genes with expression values higher than or equal to five, in at least five  
415 samples have been kept, otherwise, they were eliminated from the study. However,  
416 for T helper samples, due to fewer number of samples, four samples with values  
417 higher than or equal to five were enough to be considered in the study. After first  
418 step of filtering, the main normalization step was used to decrease dependency of  
419 expression estimates to transcript length and GC-content[47, 48]. For all four sets  
420 of samples, including training and testing samples for immune cells and for T helper

421 cells, expression estimates were normalized separately by applying *withinLaneNor-*  
422 *malization* and *betweenLaneNormalization* functions from EDASeq package [49] in  
423 R programming language (R 3.5.3), to remove GC-content biases and between-lane  
424 differences in count distributions [49]. After normalization, the second step of filtra-  
425 tion, just similar to the first step, was applied to eliminate genes with insignificant  
426 expression.

#### 427 Missing Values

428 In contrast to previous studies that only considered intersection genes [50], in order  
429 to avoid of deletion of discriminative genes, we tried to keep genes with high ex-  
430 pression, as much as possible. However, for most of genes, values for some samples  
431 were not estimated. Hence, to deal with these missing values, we used an imputa-  
432 tion method [51] and instead of mean imputation we set a dummy constant since  
433 mean imputation in this case is not meaningful and can increase error. Specifically,  
434 we generated a training set for each group of cell types, by duplicating the original  
435 training set 100 times and randomly eliminating ten percent of expression values.  
436 We next set -1 for all these missing values (both original missing values and those  
437 we eliminated) as a dummy constant because all values are positive and it is easy to  
438 be learned by the system as noise. This approach makes the system learn to neglect  
439 specific value (-1) and treat it like noise, instead of learning it as a feature of the  
440 samples.

#### 441 Classifier Training and Testing

442 Considering the few number of training samples in comparison with the high di-  
443 mensions (15453 genes in immune cell samples and 9146 genes in the T helper  
444 samples) and to avoid both over fitting the model and adding noise to the pre-  
445 diction model, we used regularization with logistic regression to decrease the total  
446 number of genes and select the most discriminative set of genes. To perform gene  
447 selection, we trained a lasso-ridge logistic regression (elastic-net) model, which au-  
448 tomatically sets the coefficients of a large number of genes to zero and pruned  
449 the number of genes as features of the classifier. We cross-validated the model by  
450 implementing `cv.glmnet` function with `nfold=10` from `glmnet` package [21] in R pro-  
451 gramming language, using training sets for both groups of cell types. We normalized  
452 the gene expression values using a  $\log_2$  transform over training sets to decrease the

453 range of values that can affect the performance of the model ( $\log_2(\text{counts}+1)$ ).  
454 In order to find the optimal number of genes, we tried 7 different lambdas and  
455 tested the results over the testing samples (*cv.glmnet(family="multinomial", al-*  
456 *pha=0.93, thresh=1e-07, lambda=c(0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001),*  
457 *type.multinomial="grouped", nfolds=10)*). To select the optimal value for lambda,  
458 True-Negative samples were generated by randomly scrambling testing datasets,  
459 then we generated ROC curves and considered original testing datasets as True-  
460 Positive samples.

#### 461 Developing Gene Signatures

462 Genes selected by the classifier models were used as initial point to build gene  
463 signatures. In this case, we trained a new binary elastic-net model for each cell type  
464 by considering a certain cell type as one class and all other cell types as another class.  
465 The training and testing samples used to build gene signatures were the training  
466 and testing samples used in developing the classifiers with the difference being  
467 that they only contained the selected genes. Similar steps including dealing with  
468 missing values, applying  $\log_2$  and visualization by ROC to select optimal number  
469 of genes were applied for each cell type. This two-step gene selection approach has  
470 the advantage that it eliminates a large number of undiscriminating genes at the  
471 first and finally select few number of genes for each cell type.

#### 472 Benchmarking

473 Fisher exact testing was used for each gene set to characterize true and system-  
474 atically scrambled data as a measure of performance of the gene set as a means  
475 of distinguishing between cell subtypes. Data was scrambled by randomly redis-  
476 tributing expression values by gene as well as patient in order to establish negative  
477 control values for determining specificity. The threshold for expression binarization  
478 for Fisher exact testing was selected based on gene expression histograms of the  
479 data to separate the measured expression from background noise levels, with 2.48  
480 being used as the threshold (after  $\log_2$  normalization). One-thousand iterations  
481 were processed and compiled in order to produce ROC curves with 95% confidence  
482 intervals shaded about the averaged ROC curve for each gene set's performance.  
483 The tested gene sets were the logistic regression gene set, the CIBERSORT gene set



484 [8], the single cell gene set [29], and the manually curated gene set that had been  
485 used previously.

## 486 **List of abbreviations**

487 ROC: receiver-operator curves  
488 scRNA-seq: single-cell RNA-seq  
489 AUC: area under the ROC curve  
490 CNV: copy number variation  
491 PCA: principal component analysis  
492 ICI: immune checkpoint inhibitor  
493 SVM: support vector machine

494

## 495 **Declarations**

### 496 **Ethics approval and consent to participate**

497 The results described in this manuscript consist of secondary analyses of existing data and was determined by the  
498 West Virginia University IRB to qualify for an exemption from human subject research under U.S. HHS regulations  
499 45 CFR 46.101(b)(4).

### 500 **Consent for publication**

501 All of the authors have read the final manuscript and consent for publication.

### 502 **Availability of data and material**

503 The datasets supporting the conclusions of this article are available in Gene Expression Omnibus repository  
504 [<https://www.ncbi.nlm.nih.gov>] with the following GEO accession numbers: GSE60424, GSE64655, GSE36952,  
505 GSE84697, GSE74246, GSE70106, GSE55536, GSE71645, GSE66261, GSE96538, GSE75688, GSE72056. R scripts  
506 used in the analyses can be found on GitHub [<https://github.com/KlinkeLab/ImmClass2019>].

### 507 **Competing interests**

508 The authors declare that they have no competing interests.

### 509 **Funding**

510 This work was supported by the National Science Foundation (NSF) (CBET-1644932 to DJK) and the National  
511 Cancer Institute (NCI) (R01CA193473 to DJK). The content is solely the responsibility of the authors and does not  
512 necessarily represent the official views of the NCI, the National Institutes of Health, or the National Science  
513 Foundation.

### 514 **Authors' contributions**

515 Designed study: AT and DK; performed analyses and interpreted results: AT, PG, and DK; and drafted initial  
516 manuscript: AT, PG, and DK. All authors edited and approved the final version of the manuscript.

### 517 **Acknowledgements**

#### 518 **Author details**

519 <sup>1</sup>Department of Chemical and Biomedical Engineering, West Virginia University, 1306 Evansdale Dr, 26506  
520 Morgantown, WV, USA. <sup>2</sup>Department of Microbiology, Immunology, and Cell Biology, West Virginia University, 1  
521 Medical Center Drive, 26506 Morgantown, WV, USA.

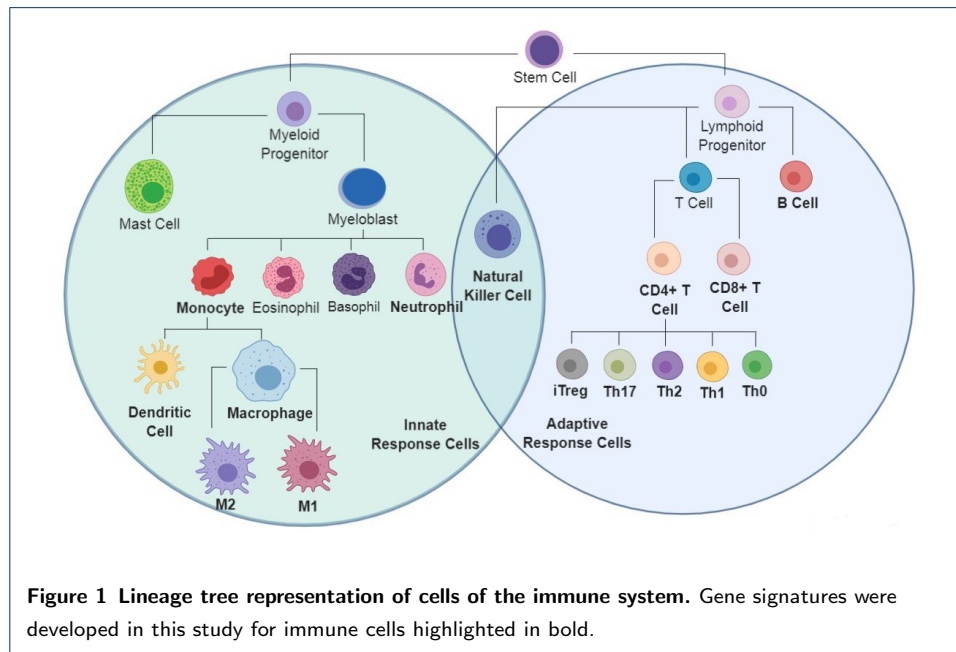
522 **References**

- 523 1. Carmona, S.J., Teichmann, S.A., Ferreira, L., Macaulay, I.C., Stubbington, M.J., Cvejic, A., Gfeller, D.:  
524 Single-cell transcriptome analysis of fish immune cells provides insight into the evolution of vertebrate immune  
525 cell types. *Genome research*, 207704 (2017)
- 526 2. Bendall, S.C., Simonds, E.F., Qiu, P., El-ad, D.A., Krutzik, P.O., Finck, R., Bruggner, R.V., Melamed, R.,  
527 Trejo, A., Ornatsky, O.I., *et al.*: Single-cell mass cytometry of differential immune and drug responses across a  
528 human hematopoietic continuum. *Science* **332**(6030), 687–696 (2011)
- 529 3. Shay, T., Kang, J.: Immunological genome project and systems immunology. *Trends in immunology* **34**(12),  
530 602–609 (2013)
- 531 4. Kinter, A.L., Hennessey, M., Bell, A., Kern, S., Lin, Y., Daucher, M., Planta, M., McGlaughlin, M., Jackson,  
532 R., Ziegler, S.F., *et al.*: Cd25+ cd4+ regulatory t cells from the peripheral blood of asymptomatic hiv-infected  
533 individuals regulate cd4+ and cd8+ hiv-specific t cell immune responses in vitro and are associated with  
534 favorable clinical markers of disease status. *Journal of Experimental Medicine* **200**(3), 331–343 (2004)
- 535 5. Vegh, P., Haniffa, M.: The impact of single-cell rna sequencing on understanding the functional organization of  
536 the immune system. *Briefings in functional genomics* (2018)
- 537 6. Kaiser, J.L., Bland, C.L., Klinke, D.J.: Identifying causal networks linking cancer processes and anti-tumor  
538 immunity using bayesian network inference and metagene constructs. *Biotechnology progress* **32**(2), 470–479  
539 (2016)
- 540 7. Finotello, F., Mayer, C., Plattner, C., Laschober, G., Rieder, D., Hackl, H., Krogsdam, A., Posch, W.,  
541 Wilflingseder, D., Sopper, S., *et al.*: quantiseq: quantifying immune contexture of human tumors. *bioRxiv*,  
542 223180 (2017)
- 543 8. Newman, A.M., Liu, C.L., Green, M.R., Gentles, A.J., Feng, W., Xu, Y., Hoang, C.D., Diehn, M., Alizadeh,  
544 A.A.: Robust enumeration of cell subsets from tissue expression profiles. *Nature methods* **12**(5), 453 (2015)
- 545 9. Li, B., Severson, E., Pignon, J.-C., Zhao, H., Li, T., Novak, J., Jiang, P., Shen, H., Aster, J.C., Rodig, S., *et*  
546 *al.*: Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome biology* **17**(1),  
547 174 (2016)
- 548 10. Racle, J., de Jonge, K., Baumgaertner, P., Speiser, D.E., Gfeller, D.: Simultaneous enumeration of cancer and  
549 immune cell types from bulk tumor gene expression data. *Elife* **6**, 26476 (2017)
- 550 11. Kidd, B.A., Peters, L.A., Schadt, E.E., Dudley, J.T.: Unifying immunology with informatics and multiscale  
551 biology. *Nature immunology* **15**(2), 118 (2014)
- 552 12. Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A.,  
553 Pomeroy, S.L., Golub, T.R., Lander, E.S., *et al.*: Gene set enrichment analysis: a knowledge-based approach for  
554 interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences* **102**(43),  
555 15545–15550 (2005)
- 556 13. Liberzon, A., Subramanian, A., Pinchback, R., Thorvaldsdóttir, H., Tamayo, P., Mesirov, J.P.: Molecular  
557 signatures database (msigdb) 3.0. *Bioinformatics* **27**(12), 1739–1740 (2011)
- 558 14. Zheng, C.-H., Chong, Y.-W., Wang, H.-Q.: Gene selection using independent variable group analysis for tumor  
559 classification. *Neural Computing and Applications* **20**(2), 161–170 (2011)
- 560 15. Wu, M.-Y., Dai, D.-Q., Shi, Y., Yan, H., Zhang, X.-F.: Biomarker identification and cancer classification based  
561 on microarray data using laplace naive bayes model with mean shrinkage. *IEEE/ACM transactions on*  
562 *computational biology and bioinformatics* **9**(6), 1649–1662 (2012)
- 563 16. Cui, Y., Zheng, C.-H., Yang, J., Sha, W.: Sparse maximum margin discriminant analysis for feature extraction  
564 and gene selection on gene expression data. *Computers in biology and medicine* **43**(7), 933–941 (2013)
- 565 17. Algamil, Z.Y., Lee, M.H.: Regularized logistic regression with adjusted adaptive elastic net for gene selection in  
566 high dimensional cancer classification. *Computers in biology and medicine* **67**, 136–145 (2015)
- 567 18. Liang, Y., Liu, C., Luan, X.-Z., Leung, K.-S., Chan, T.-M., Xu, Z.-B., Zhang, H.: Sparse logistic regression  
568 with a  $l_{1/2}$  penalty for gene selection in cancer classification. *BMC bioinformatics* **14**(1), 198 (2013)
- 569 19. Bielza, C., Robles, V., Larrañaga, P.: Regularized logistic regression without a penalty term: An application to  
570 cancer classification with microarray data. *Expert Systems with Applications* **38**(5), 5110–5118 (2011)
- 571 20. Cawley, G.C., Talbot, N.L.: Gene selection in cancer classification using sparse logistic regression with bayesian  
572 regularization. *Bioinformatics* **22**(19), 2348–2355 (2006)

- 573 21. Friedman, J., Hastie, T., Tibshirani, R.: Regularization paths for generalized linear models via coordinate  
574 descent. *Journal of statistical software* **33**(1), 1 (2010)
- 575 22. Zou, H., Hastie, T.: Regularization and variable selection via the elastic net. *Journal of the Royal Statistical*  
576 *Society: Series B (Statistical Methodology)* **67**(2), 301–320 (2005)
- 577 23. Tirosh, I., Izar, B., Prakadan, S.M., Wadsworth, M.H., Treacy, D., Trombetta, J.J., Rotem, A., Rodman, C.,  
578 Lian, C., Murphy, G., *et al.*: Dissecting the multicellular ecosystem of metastatic melanoma by single-cell  
579 rna-seq. *Science* **352**(6282), 189–196 (2016)
- 580 24. Chung, W., Eum, H.H., Lee, H.-O., Lee, K.-M., Lee, H.-B., Kim, K.-T., Ryu, H.S., Kim, S., Lee, J.E., Park,  
581 Y.H., *et al.*: Single-cell rna-seq enables comprehensive tumour and immune cell profiling in primary breast  
582 cancer. *Nature communications* **8**, 15081 (2017)
- 583 25. Caligiuri, M.A.: Human natural killer cells. *Blood* **112**(3), 461–469 (2008)
- 584 26. Chen, E.Y., Tan, C.M., Kou, Y., Duan, Q., Wang, Z., Meirelles, G.V., Clark, N.R., Ma'ayan, A.: Enrichr:  
585 interactive and collaborative html5 gene list enrichment analysis tool. *BMC bioinformatics* **14**(1), 128 (2013)
- 586 27. Kuleshov, M.V., Jones, M.R., Rouillard, A.D., Fernandez, N.F., Duan, Q., Wang, Z., Koplev, S., Jenkins, S.L.,  
587 Jagodnik, K.M., Lachmann, A., *et al.*: Enrichr: a comprehensive gene set enrichment analysis web server 2016  
588 update. *Nucleic acids research* **44**(W1), 90–97 (2016)
- 589 28. Sharma, P., Hu-Lieskovan, S., Wargo, J.A., Ribas, A.: Primary, adaptive, and acquired resistance to cancer  
590 immunotherapy. *Cell* **168**(4), 707–723 (2017)
- 591 29. Jerby-Arnon, L., Shah, P., Cuoco, M.S., Rodman, C., Su, M.-J., Melms, J.C., Leeson, R., Kanodia, A., Mei, S.,  
592 Lin, J.-R., *et al.*: A cancer cell program promotes t cell exclusion and resistance to checkpoint blockade. *Cell*  
593 **175**(4), 984–997 (2018)
- 594 30. Charrad, M., Ghazzali, N., Boiteau, V., Niknafs, A., Charrad, M.M.: Package 'nbclust'. *Journal of statistical*  
595 *software* **61**, 1–36 (2014)
- 596 31. Finotello, F., Trajanoski, Z.: Quantifying tumor-infiltrating immune cells from transcriptomics data. *Cancer*  
597 *Immunology, Immunotherapy* **67**(7), 1031–1040 (2018)
- 598 32. Xu, C., Su, Z.: Identification of cell types from single-cell transcriptomes using a novel clustering method.  
599 *Bioinformatics* **31**(12), 1974–1980 (2015)
- 600 33. Grün, D., Lyubimova, A., Kester, L., Wiebrands, K., Basak, O., Sasaki, N., Clevers, H., van Oudenaarden, A.:  
601 Single-cell messenger rna sequencing reveals rare intestinal cell types. *Nature* **525**(7568), 251 (2015)
- 602 34. Hu, Y., Hase, T., Li, H.P., Prabhakar, S., Kitano, H., Ng, S.K., Ghosh, S., Wee, L.J.K.: A machine learning  
603 approach for the identification of key markers involved in brain development from single-cell transcriptomic  
604 data. *BMC genomics* **17**(13), 1025 (2016)
- 605 35. Yao, F., Zhang, C., Du, W., Liu, C., Xu, Y.: Identification of gene-expression signatures and protein markers for  
606 breast cancer grading and staging. *PloS one* **10**(9), 0138213 (2015)
- 607 36. Linsley, P.S., Speake, C., Whalen, E., Chaussabel, D.: Copy number loss of the interferon gene cluster in  
608 melanomas is linked to reduced t cell infiltrate and poor patient prognosis. *PloS one* **9**(10), 109760 (2014)
- 609 37. Hoek, K.L., Samir, P., Howard, L.M., Niu, X., Prasad, N., Galassie, A., Liu, Q., Allos, T.M., Floyd, K.A., Guo,  
610 Y., *et al.*: A cell-based systems biology assessment of human blood to monitor immune responses after  
611 influenza vaccination. *PloS one* **10**(2), 0118528 (2015)
- 612 38. Beyer, M., Mallmann, M.R., Xue, J., Staratschek-Jox, A., Vorholt, D., Krebs, W., Sommer, D., Sander, J.,  
613 Mertens, C., Nino-Castro, A., *et al.*: High-resolution transcriptome of human macrophages. *PloS one* **7**(9),  
614 45466 (2012)
- 615 39. Şenbabaoğlu, Y., Gejman, R.S., Winer, A.G., Liu, M., Van Allen, E.M., de Velasco, G., Miao, D., Ostrovskaya,  
616 I., Drill, E., Luna, A., *et al.*: Tumor immune microenvironment characterization in clear cell renal cell carcinoma  
617 identifies prognostic and immunotherapeutically relevant messenger rna signatures. *Genome biology* **17**(1), 231  
618 (2016)
- 619 40. Corces, M.R., Buenrostro, J.D., Wu, B., Greenside, P.G., Chan, S.M., Koenig, J.L., Snyder, M.P., Pritchard,  
620 J.K., Kundaje, A., Greenleaf, W.J., *et al.*: Lineage-specific and single-cell chromatin accessibility charts human  
621 hematopoiesis and leukemia evolution. *Nature genetics* **48**(10), 1193 (2016)
- 622 41. Kumar, N.A., Cheong, K., Powell, D.R., da Fonseca Pereira, C., Anderson, J., Evans, V.A., Lewin, S.R.,  
623 Cameron, P.U.: The role of antigen presenting cells in the induction of hiv-1 latency in resting cd4+ t-cells.

- 624 Retrovirology **12**(1), 76 (2015)
- 625 42. Zhang, H., Xue, C., Shah, R., Bermingham, K., Hinkle, C.C., Li, W., Rodrigues, A., Tabita-Martinez, J., Millar,  
626 J.S., Cuchel, M., et al.: Functional analysis and transcriptomic profiling of ipsc-derived macrophages and their  
627 application in modeling mendelian disease. *Circulation research*, 114 (2015)
- 628 43. Kanduri, K., Tripathi, S., Larjo, A., Mannerström, H., Ullah, U., Lund, R., Hawkins, R.D., Ren, B.,  
629 Lähdesmäki, H., Lahesmaa, R.: Identification of global regulators of t-helper cell lineage specification. *Genome  
630 medicine* **7**(1), 122 (2015)
- 631 44. Spurlock III, C.F., Tossberg, J.T., Guo, Y., Collier, S.P., Crooke III, P.S., Aune, T.M.: Expression and functions  
632 of long noncoding rnas during human t helper cell differentiation. *Nature communications* **6**, 6932 (2015)
- 633 45. Schmidt, A., Marabita, F., Kiani, N.A., Gross, C.C., Johansson, H.J., Éliás, S., Rautio, S., Eriksson, M.,  
634 Fernandes, S.J., Silberberg, G., et al.: Time-resolved transcriptome and proteome landscape of human  
635 regulatory t cell (treg) differentiation reveals novel regulators of foxp3. *BMC biology* **16**(1), 47 (2018)
- 636 46. Edgar, R., Domrachev, M., Lash, A.E.: Gene expression omnibus: Ncbi gene expression and hybridization array  
637 data repository. *Nucleic acids research* **30**(1), 207–210 (2002)
- 638 47. Oshlack, A., Wakefield, M.J.: Transcript length bias in rna-seq data confounds systems biology. *Biology direct*  
639 **4**(1), 14 (2009)
- 640 48. Robinson, M.D., Oshlack, A.: A scaling normalization method for differential expression analysis of rna-seq  
641 data. *Genome biology* **11**(3), 25 (2010)
- 642 49. Risso, D., Schwartz, K., Sherlock, G., Dudoit, S.: Gc-content normalization for rna-seq data. *BMC*  
643 *bioinformatics* **12**(1), 480 (2011)
- 644 50. Schwalie, P.C., Ordóñez-Morán, P., Huelsken, J., Deplancke, B.: Cross-tissue identification of somatic stem and  
645 progenitor cells using a single-cell rna-sequencing derived gene signature. *Stem Cells* **35**(12), 2390–2402 (2017)
- 646 51. García-Laencina, P.J., Sancho-Gómez, J.-L., Figueiras-Vidal, A.R.: Pattern classification with missing data: a  
647 review. *Neural Computing and Applications* **19**(2), 263–282 (2010)

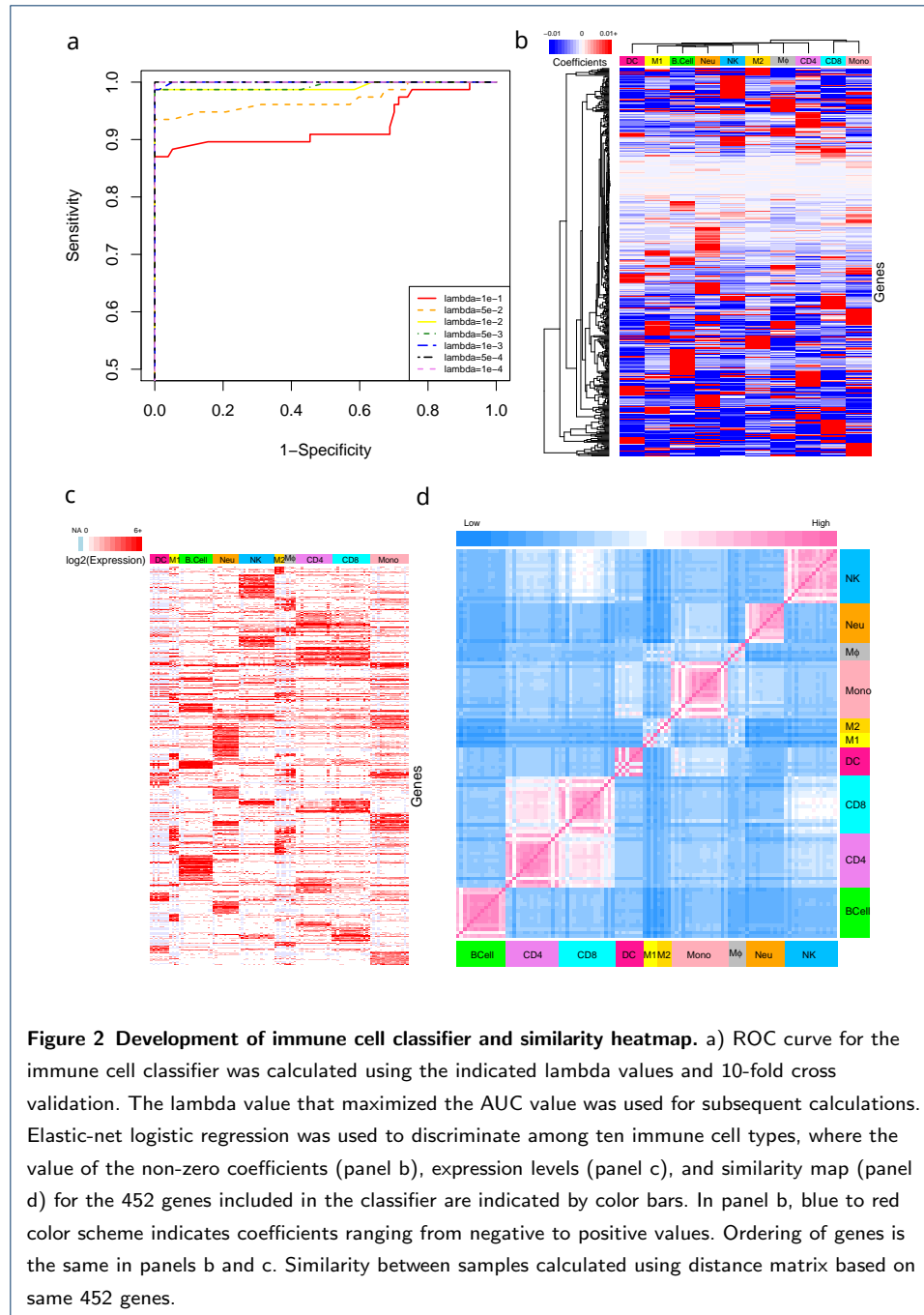
648 **Figures**



649 **Additional Files**

650 Table S1. — Coefficients of immune cell classifier and T helper cell classifier

651 Coefficients of immune cell classifier were located in the first sheet and coefficients of T helper cells were located in  
652 the second sheet.



653 Table S2. — Lambda Selection by AUC Values

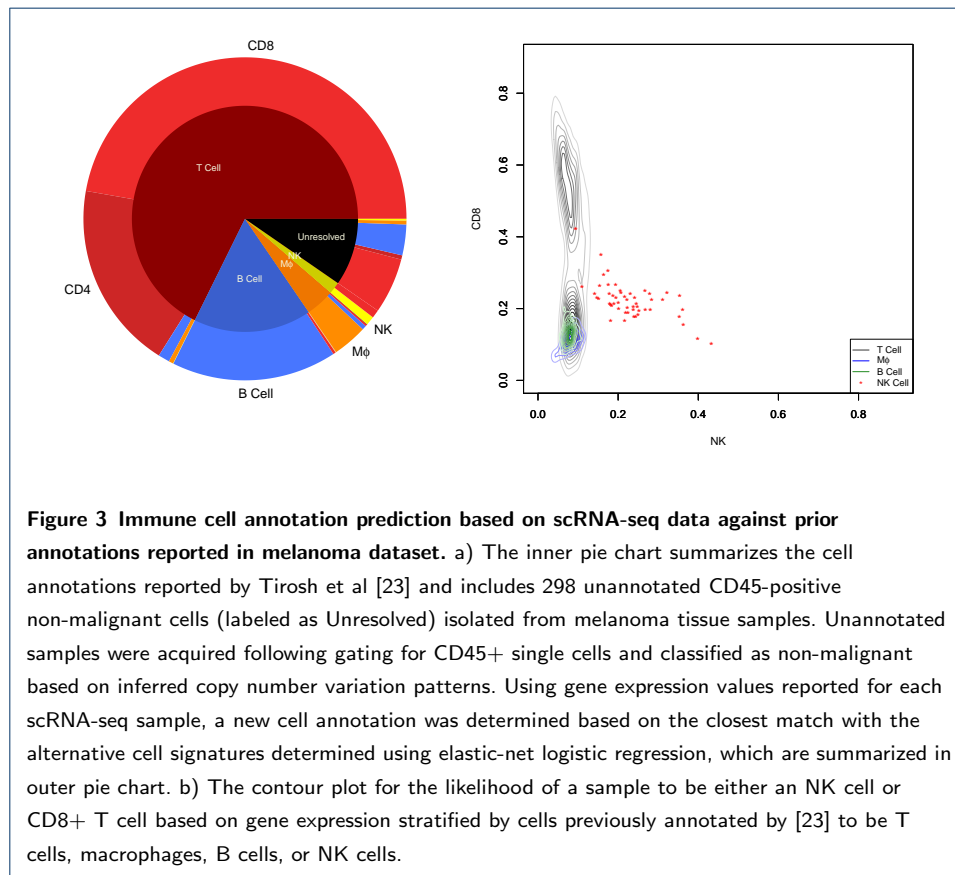
654 Lambdas with corresponding calculated AUC. The final column shows the selected lambdas

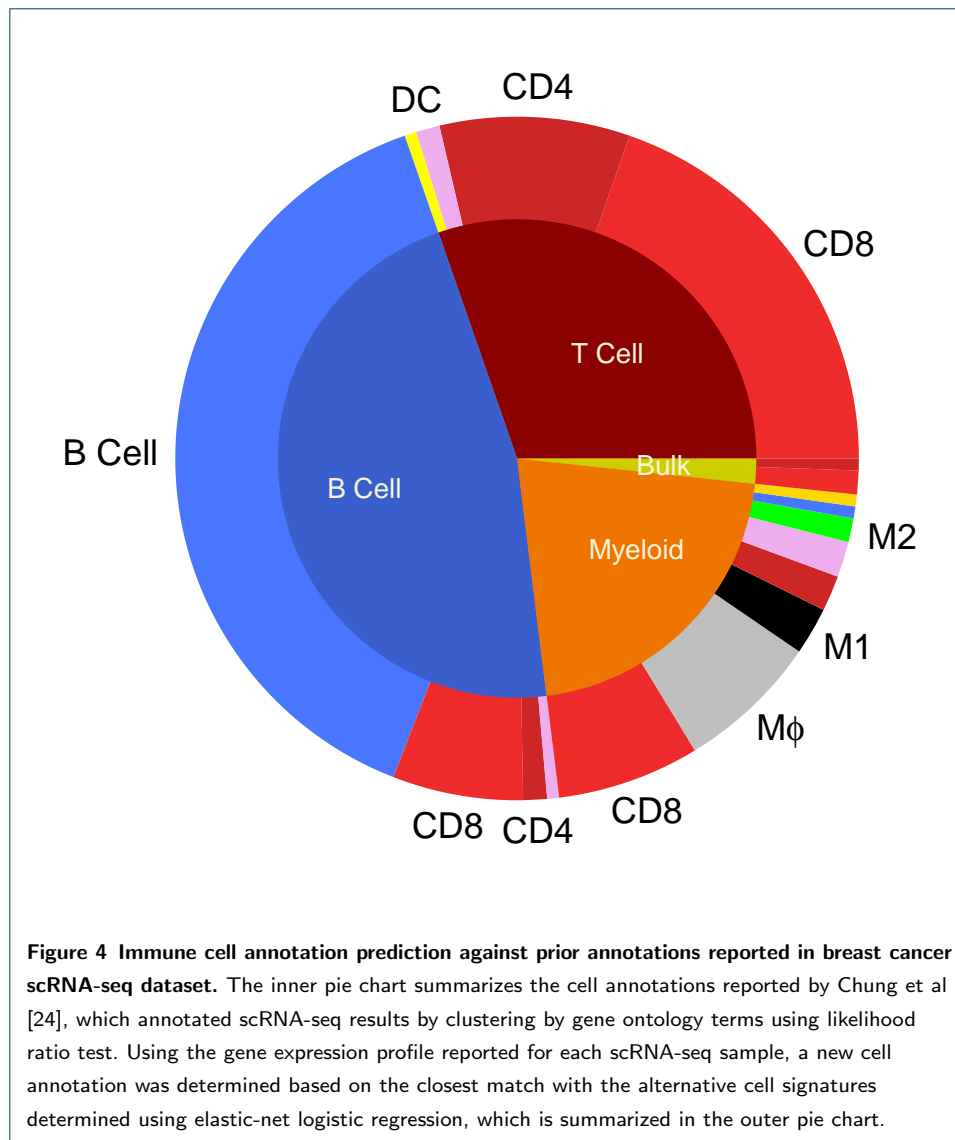
655 Table S3. — Genes in developed gene signature for immune and T helper cells

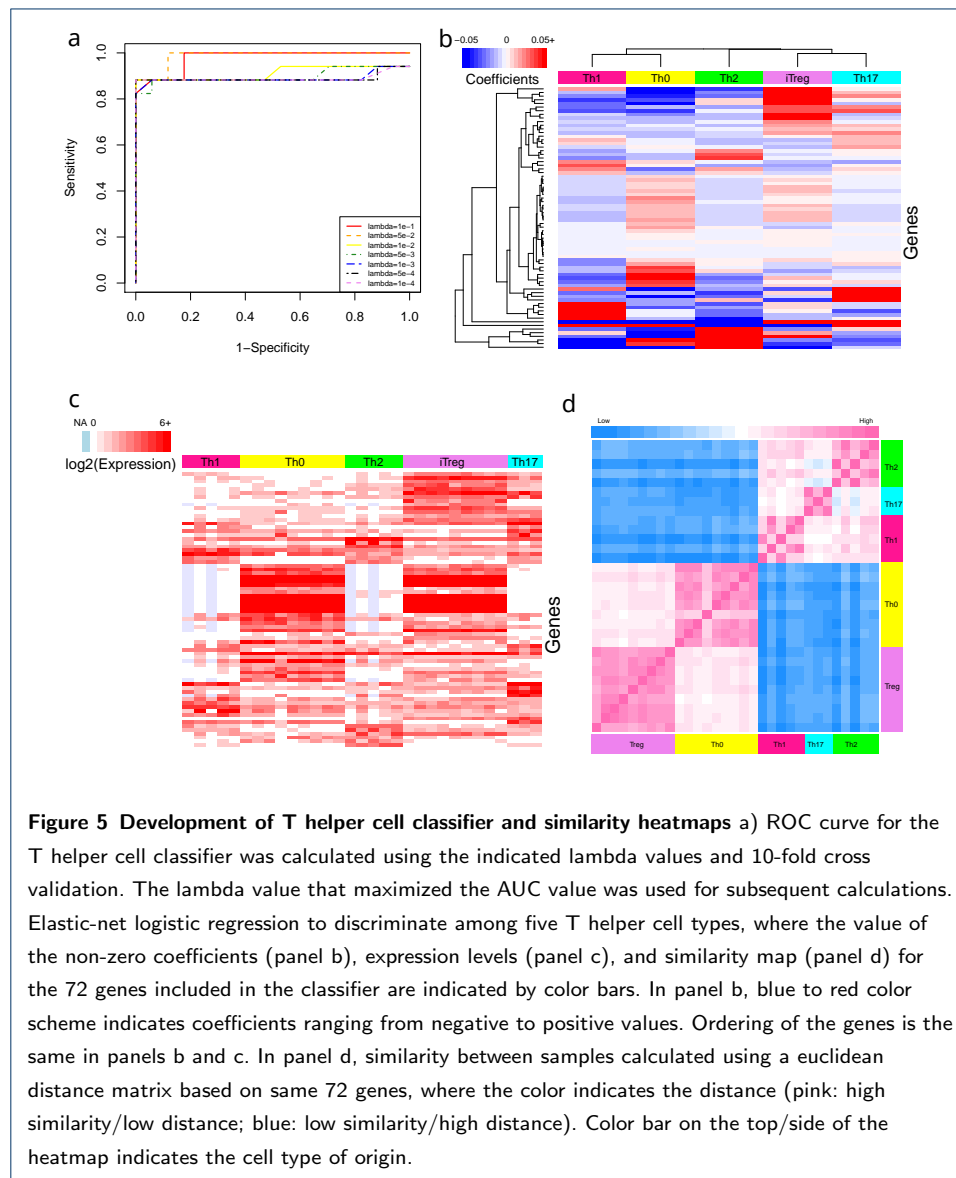
656 Yellow boxes show genes with negative impact in possibility of being related cell type.

657 Table S4. — Data information used in training models.

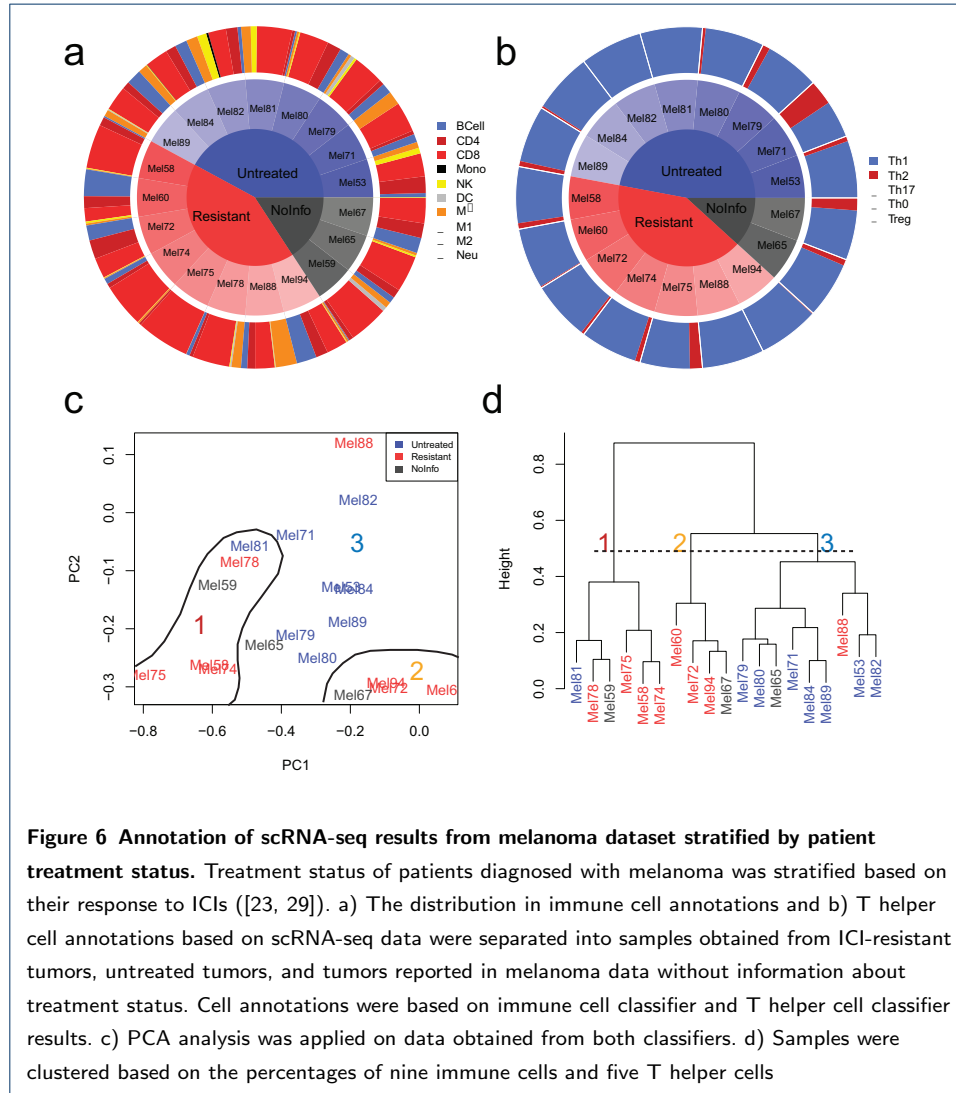
658 The second sheet shows names that used in creating datasets.

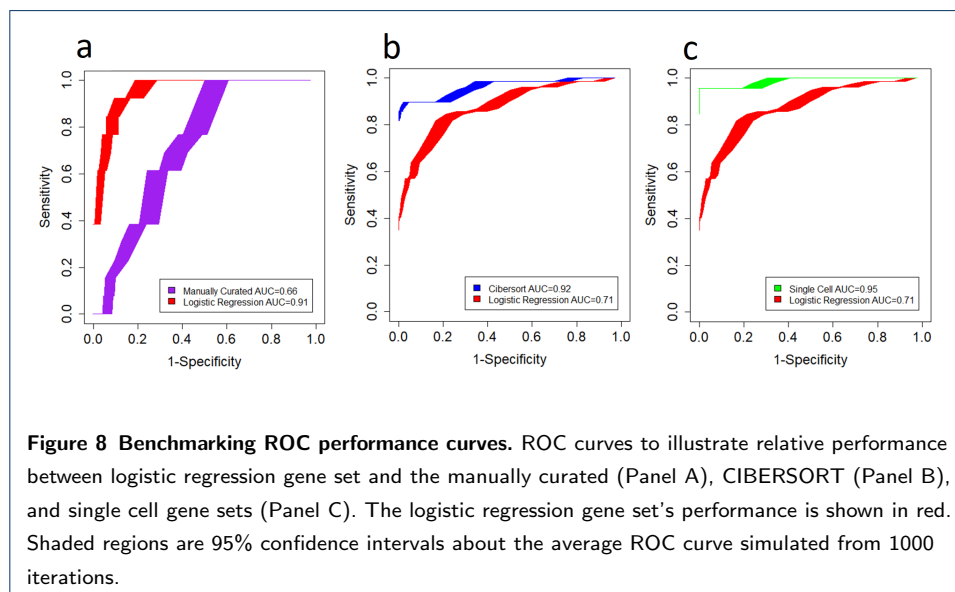
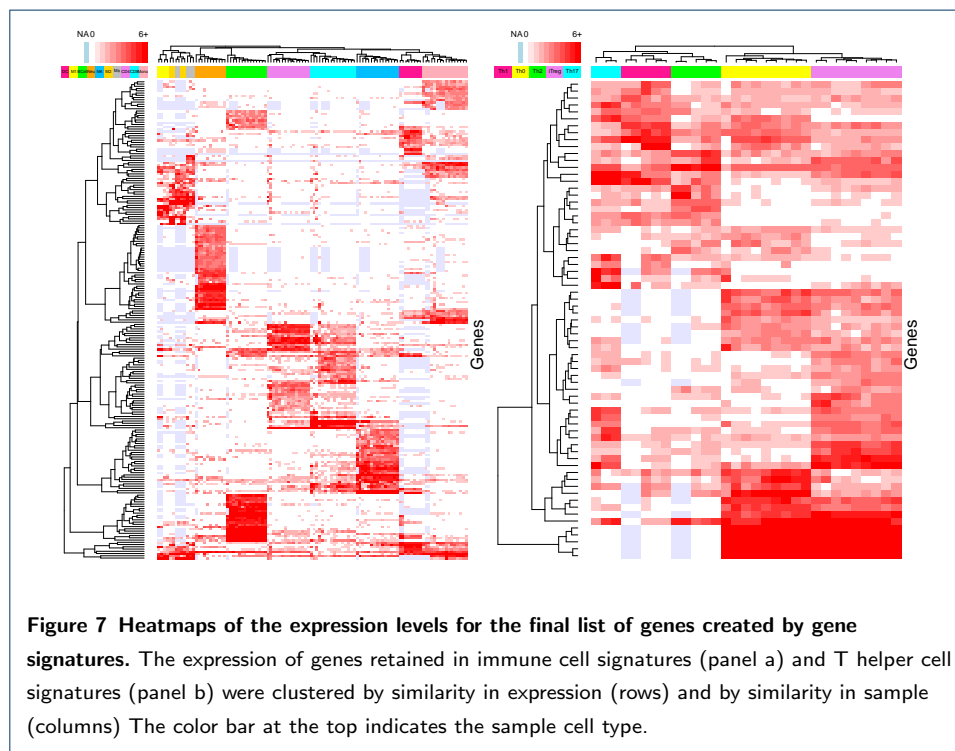


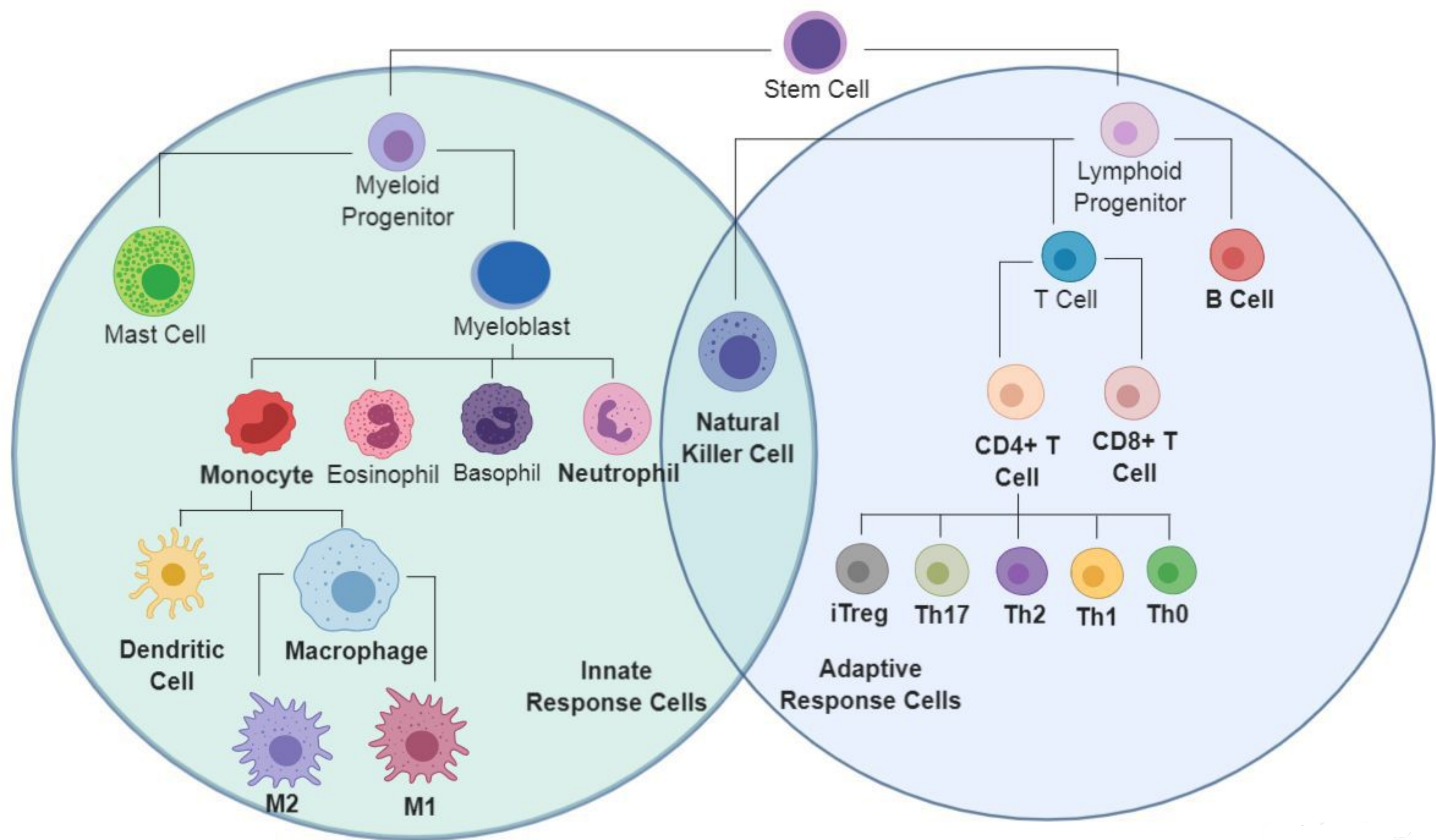


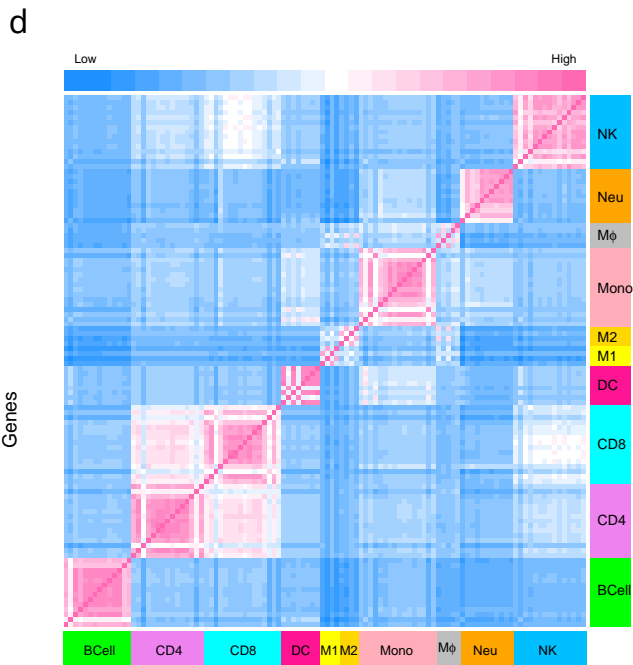
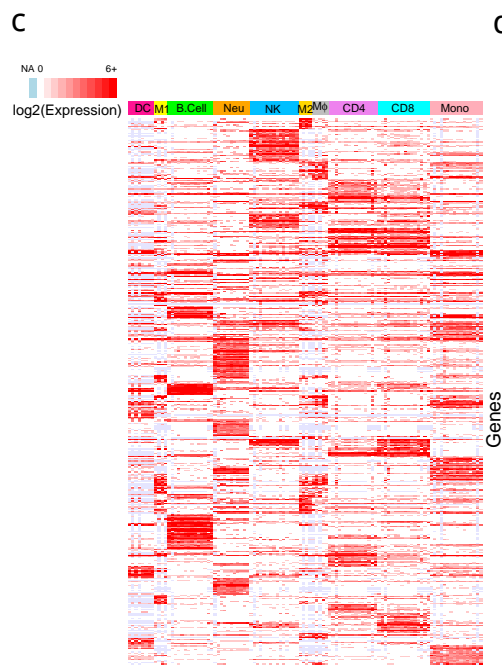
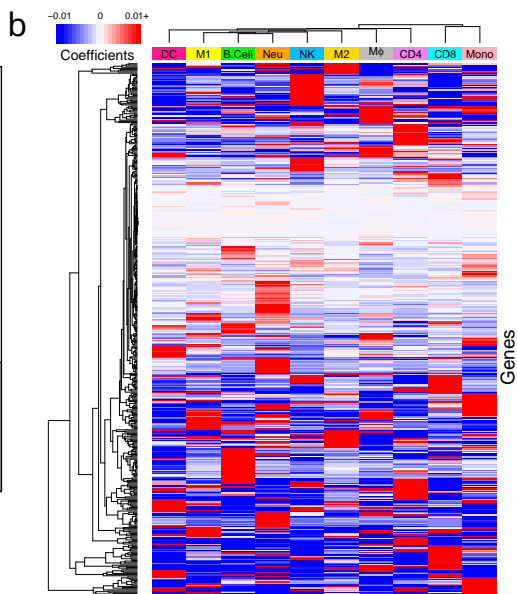
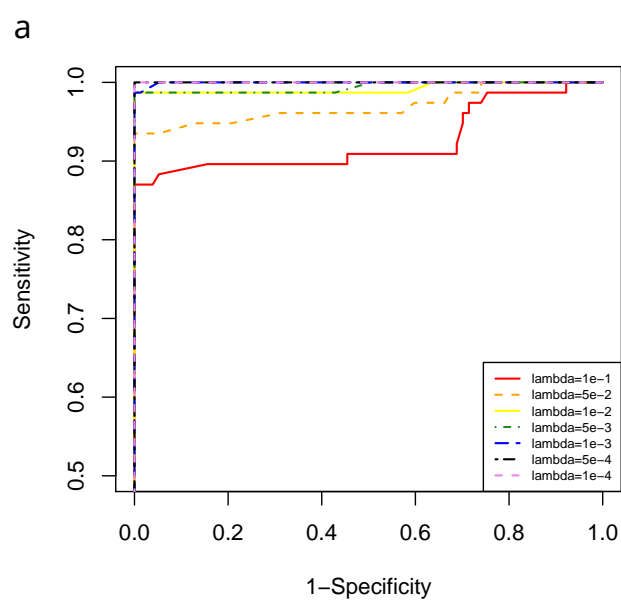


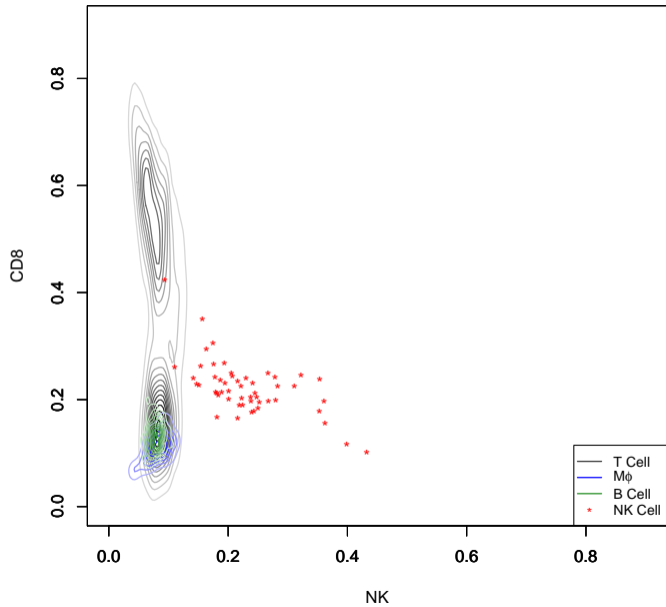
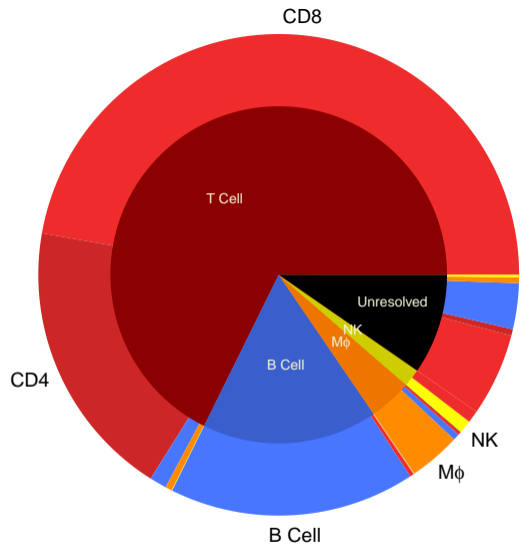


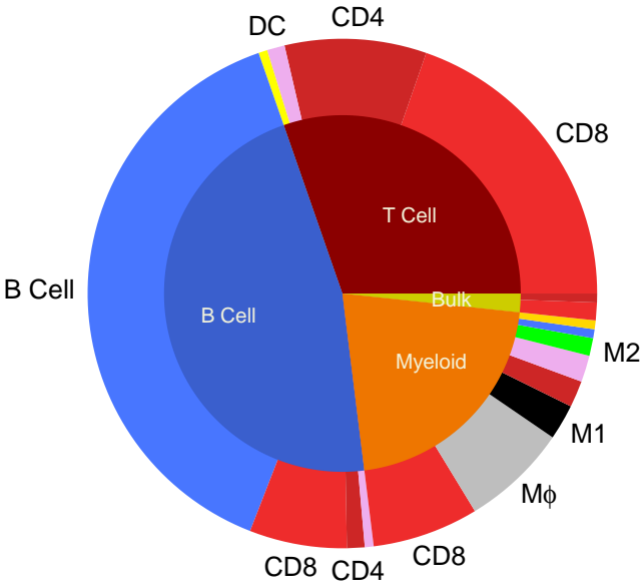


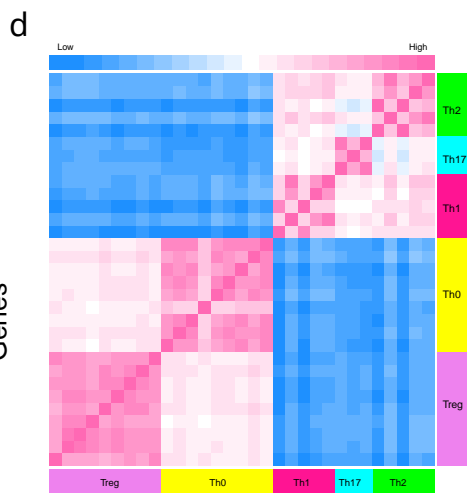
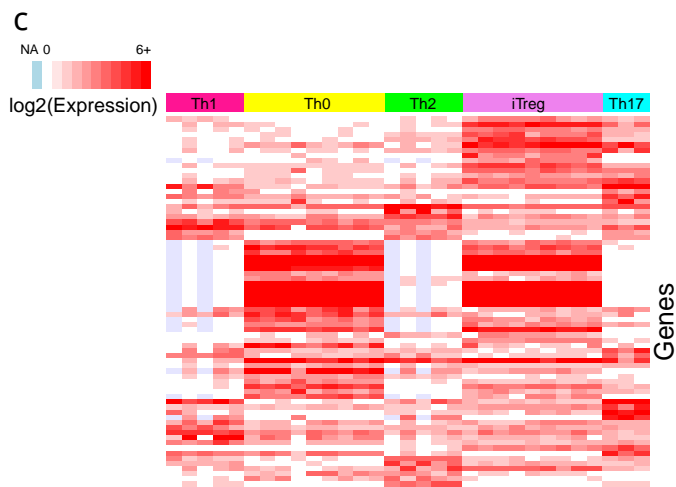
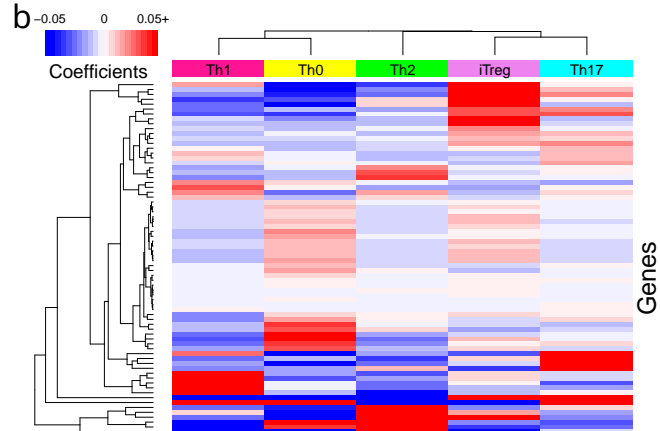
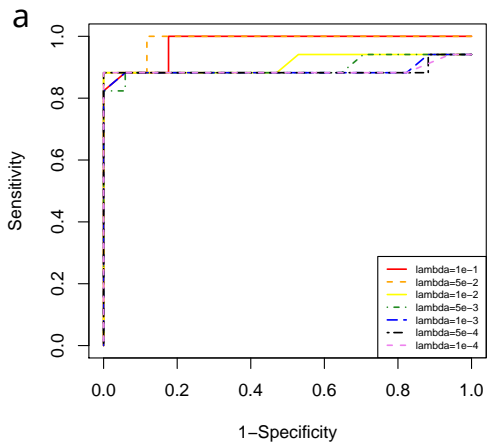


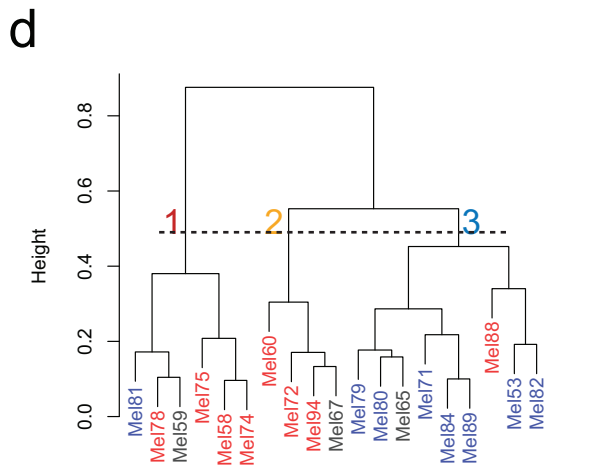
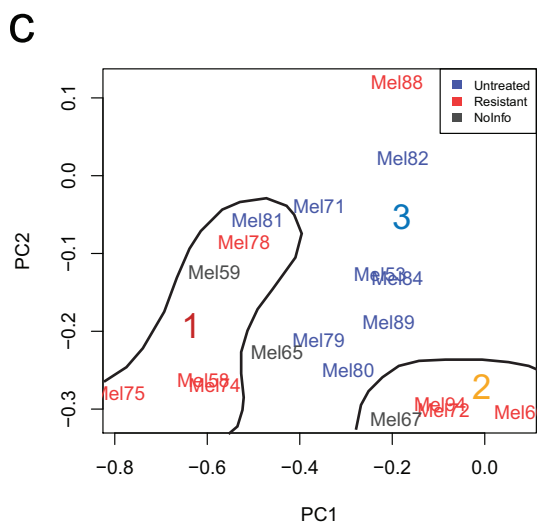
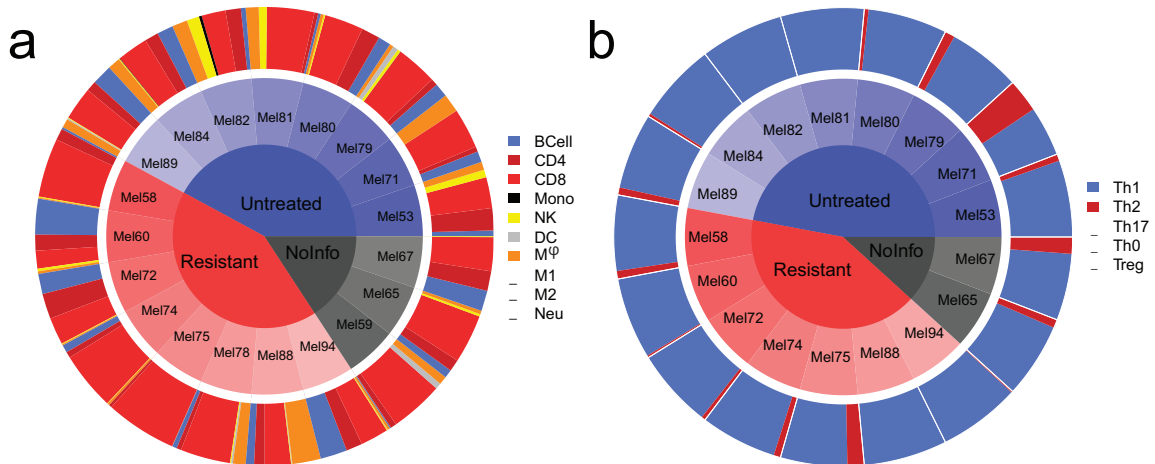




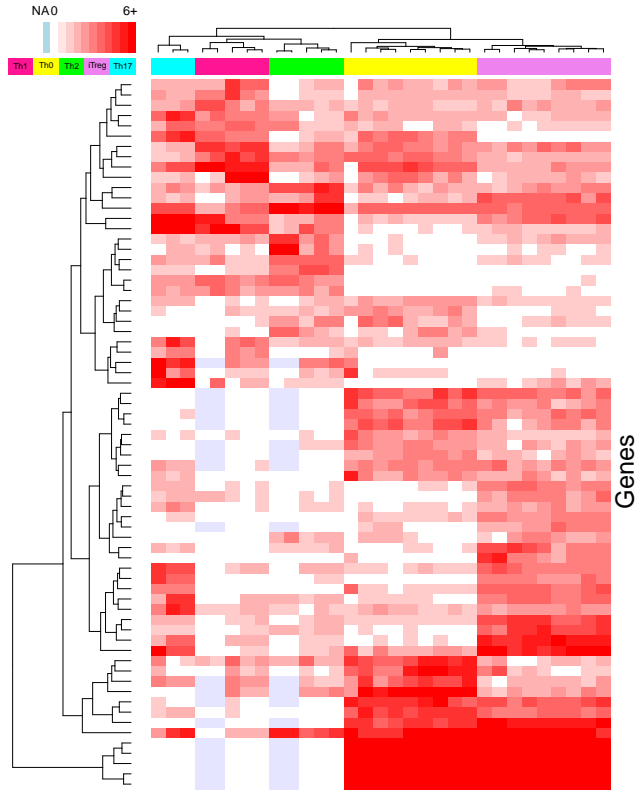
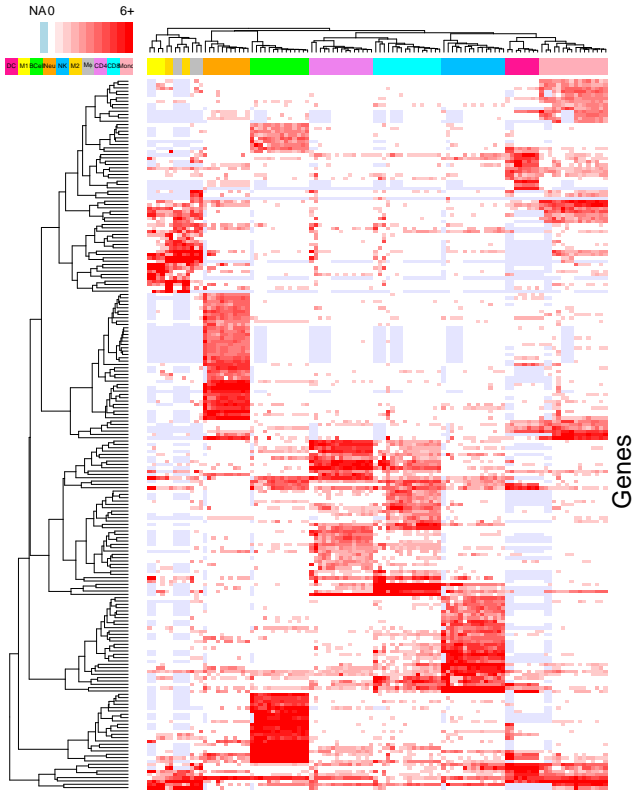


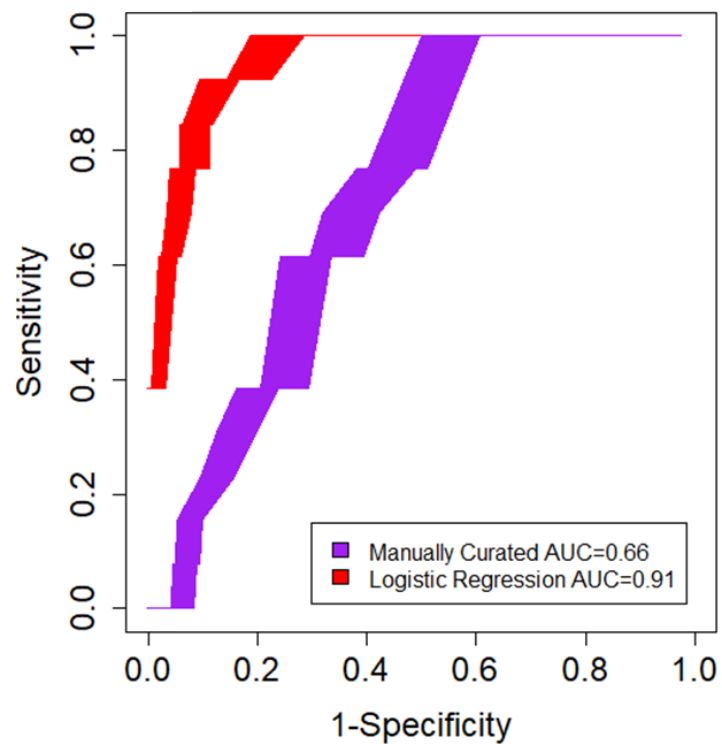
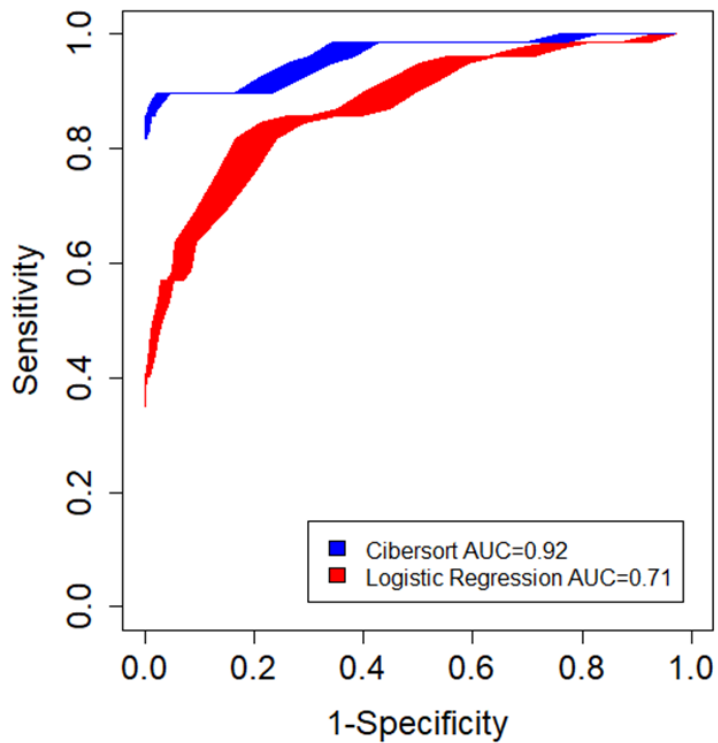










**a****b****c**