# **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>

Torang et al. Page 2 of 29

Correspondence:
avid.klinke@mail.wvu.edu
Department of Chemical and
ciomedical Engineering, West
Girginia University, 1306 Evansdale
Dr., 26506 Morgantown, WV, USA
ull list of author information is
vailable 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; In silico Cytometry

2

# Background

- 4 Host immune response is a coordinated complex system, consisting of different spe-
- 5 cialized innate and adaptive immune cells that vary dynamically and in different
- 6 anatomical locations. As shown in Fig. 1, innate immune cells comprise myeloid
- $_{7}$  cells, which include eosinophils, neutrophils, basophils, monocytes, and mast cells.
- 8 Adaptive immune cells are mainly B lymphocytes and T lymphocytes that specif-
- 9 ically recognize different antigens [1]. Linking innate with adaptive immunity are

Torang et al. Page 3 of 29

Natural Killer cells and antigen presenting cells, like macrophages and dendritic cells. Traditionally, unique cell markers have been used to characterize different immune cell subsets from heterogeneous cell mixtures using flow cytometry [2, 3, 4]. However, flow cytometry measures on the order of 10 parameters simultaneously 13 and relies on prior knowledge for selecting relevant molecular markers, which could provide a biased view of the immune state within a sample [5]. Recent advances in technology, like mass cytometry or multispectral imaging, have expanded the number of molecular markers, but the number of markers used for discriminating 17 among cell types within a sample remains on the order of  $10^{1.5}$ . In the recent years, quantifying tumor immune contexture using bulk transcriptomics or single-cell RNA sequencing data (scRNA-seq) has piqued the interest of the scientific community [6, 7, 8, 9, 10]. Advances in transcriptomics technology, like RNA sequencing, provide a much higher dimensional view of which genes are 22 expressed in different immune cells (i.e., on the order of 10<sup>3</sup>) [11]. Conceptually, inferring cell types from data using an expanded number of biologically relevant genes becomes more tolerant to non-specific noise and non-biological differences among samples and platforms. In practice, cell types can be identified using gene signatures, which are defined as sets of genes linked to common downstream functions or inductive networks that are co-regulated [12, 13], using approaches such as Gene Set Enrichment Analysis (GSEA) [12]. However, as microarray data can inflate detecting low abundance and noisy transcripts and scRNA-seq data can have a lower 30 depth of sequencing, opportunities for refining methods to quantify the immune 31 contexture using gene signatures still remain. 32 Leveraging transcriptomics data to identify immune cell types presents analytic 33 challenges for extracting informative gene signatures hidden within a high dimensional transcriptomics space that is characterized by low sample numbers with noisy and missing values. Typically, the number of cell samples is in the range of hundreds or less, while the number of profiled genes is in the tens of thousands [14]. Yet, only a few number of genes are relevant for discriminating among immune cell subsets. Datasets with a large number of noisy and irrelevant genes decrease the accuracy and computing efficiency of machine learning algorithms, especially when the number of samples are very limited. Hence, feature selection algorithms may be used to reduce the number of redundant genes [15]. Using feature selection Torang et al. Page 4 of 29

methods enable developing gene signatures in different biomedical fields of study [16]. There are many proposed feature selection methods that can select gene sets that enable classifying samples with high accuracy. In recent years, regularization methods have became more popular, which efficiently select features [17] and also 46 control for overfitting [18]. As a machine learning tool, logistic regression is considered to be a powerful discriminative method [18]. However, logistic regression alone is not applicable for high-dimensional cell classification problems [19]. On the other hand, hybrid methods, like regularized logistic regression, have been successfully applied to high-dimensional problems [20]. Regularized logistic regression selects a small set of genes with the strongest effects on the cost function [17]. A regularized logistic regression can be also be applied with different regularization terms. The most popular regularized terms are LASSO, Ridge [21], and elastic-net [22], which 54 impose the l1 norm, l2 norm, and linear combination of l1 norm and l2 norm regularization, respectively, to the cost function. It has been shown that, specifically in very high dimensional problems, elastic-net outperforms LASSO and Ridge [17, 22]. In this study, we focused on two-step regularized logistic regression techniques to develop immune cell signatures and immune cell and T helper cell classifiers using 59 RNA-seq data for the cells highlighted in bold in Fig. 1. The first step of the process included a pre-filtering phase to select the optimal number of genes and implemented an elastic-net model as a regularization method for gene selection in generating the classifiers. The pre-filtering step reduced computational cost and increased final 63 accuracy by selecting the most discriminative and relevant set of genes. Finally, we illustrate the value of the approach in annotating gene expression profiles obtained from single-cell RNA sequencing. The second step generated gene signatures for individual cell types using selected genes from first step and implemented a binary 67 regularized logistic regression for each cell type against all other samples.

# Results

We developed classifiers for subsets of immune cells and T helper cells separately with two main goals. First, we aimed to annotate RNA-seq data obtained from an enriched cell population with information as to the immune cell identity. Second, we developed gene signatures for different immune cells that could be used to quantify the prevalence from RNA-seq data obtained from a heterogeneous cell population.

Torang et al. Page 5 of 29

Prior to developing the classifiers, the data was pre-processed to remove genes that
have low level of expression for most of samples (details can be found in Methods section) and normalized to increase the homogeneity in samples from different
studies and to decrease dependency of expression estimates to transcript length
and GC-content. Genes retained that had missing values for some of the samples
were assigned a value of -1. Next, regularized logistic regression (elastic-net) was
performed and the optimal number of genes and their coefficients were determined.

# Generating and validating an immune cell classifier

In developing the immune cell classifier, we determined the optimal number of genes in the classifier by varying the lambda value used in the regularized logistic regression of the training samples and assessing performance. To quantify the performance using different lambdas, a dataset was generated by combining True-Negative samples, which were created using a bootstrapping approach that randomly resampled associated genes and their corresponding value from the testing datasets to create a synthetic dataset of similar size and complexity, with the original testing data, which were untouched during training and provided True-Positive samples. The accuracy of predicting the True-Positive samples were used to generate Receiver Operating Characteristic (ROC) curves (Fig. 2a). Performance using each lambda was quantified as the Area Under the ROC Curve (AUC).

The optimal lambda for immune cell classifier was the smallest value (i.e., highest number of genes) that maximized the AUC. Functionally, this lambda value
represents the trade-off between retaining the highest number of informative genes
(i.e., classifier signal) for developing the gene signature in the second step, while
not adding non-informative genes (i.e., classifier noise). Consequently, we selected
a lambda value of 1e-4 (452 genes) for the immune cell classifier, where the selected
genes and their coefficients are shown in Table S1.

To explore correlations between the weights of selected genes with their expression level, we generated heatmaps shown in Fig. 2, panels b and c. A high level of gene expression is reflected as a larger positive coefficient in a classifier model, while low or absent expression results in a negative coefficient. This is interpreted as, for example, if gene A is not in cell type 1, the presence of this gene in a sample decreases the probability for that sample to be cell type 1. For instance, E-cadherin (CDH1)

Torang et al. Page 6 of 29

was not detected in almost all monocyte samples and thus has a negative coefficient.

Conversely, other genes are only expressed in certain cell types, which results in a 108 high positive coefficient. For instance, CYP27B1, INHBA, IDO1, NUPR1, and UBD are only expressed by M1 macrophages and thus have high positive coefficients. 110 The differential expression among cell types suggests that the set of genes in-111 cluded in the classifier model may also be a good starting point for developing gene signatures, which is highlighted in Fig. 2d. Here, we focused on the expression of 113 the 452 genes included in the classifier model and the correlations between samples clustered based on cell types. The off-diagonal entries in the correlation matrix are 115 colored by euclidean distance with the color indicating similarity or dissimilarity 116 using pink and blue, respectively. Color bars along the axes also highlight the cell 117 types for the corresponding RNA-seq samples. As expected, RNA-seq samples from 118 the same cell type were highly similar. More interestingly, correlation between dif-119 ferent cell types can also be seen, like high similarity between CD4+ and CD8+ T cell samples, CD8+ T cell and NK cell samples, and monocyte and dendritic 121 cell samples. Collectively, these heatmaps illustrate that the selected genes are a highly condensed but are still a representative set of genes that include the main 123 characteristics of the immune cell types. It is also notable to compare the clustering result of cell types based on their coefficients in the classifier shown in Fig. 2b with 125 similarity matrix in Fig. 2d. Since in the classifier coefficients are forcing the model 126 to separate biologically close cell types (like CD4+ T cell and CD8+ T cell), the 127 clustering results suggest that the coefficient vectors are equally dissimilar (Fig. 128 2b). However, in the case of their expression values, their similarity remains (Fig. 129 2d).

Evaluating the Immune Cell classifier using scRNA-seq datasets

To evaluate the proposed classifier in immune cell classification, two publicly accessible datasets generated by scRNA-seq technology were used [23, 24]. The first dataset
included malignant, immune, stromal and endothelial cells from 15 melanoma tissue
samples [23]. We focused on the immune cell samples, which included 2761 annotated samples of T cells, B cells, Mphi and NK cells, and 294 unresolved samples.
The immune cells in this study were recovered by flow cytometry by gating on
CD45 positive cells. Annotations were on the basis of expressed marker genes while

Torang et al. Page 7 of 29

unresolved samples were from the CD45-gate and classified as non-malignant based on inferred copy number variation (CNV) patterns (i.e., CNV score < 0.04).

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

The only cell type with low similarity between our classifier predictions and prior 151 annotations was NK cells, where we classified almost half of samples annotated 152 previously as NK cells as CD8+ T cell. Discriminating between these two cell types 153 is challenging as they share many of the genes related to cytotoxic effector function 154 and can also be subclassified into subsets, like CD56bright and CD56dim NK subsets 155 [25]. To explore this discrepancy, we compared all annotated samples based on their 156 CD8 score and NK score provided by the classifier, as shown in Fig. 3b. Although 157 the number of NK cell samples are relatively low, it seems that the NK samples 158 consist of two groups of samples: one with a higher likelihood of being a NK cell 159 and a second with almost equal likelihood for being either CD8+ T cell or NK cell. 160 We applied principal component analysis (PCA) to identify genes associated with 161 this difference and used Enrichr for gene set enrichment [26, 27]. Using gene sets 162 associated with the Human Gene Atlas, the queried gene set was enriched for genes 163 associated with CD56 NK cells, CD4+ T cell and CD8+ T cell. Collectively, the 164 results suggests that the group of cells with similar score for NK and CD8 in the 165 classifier model are Natural Killer T cells. 166

We also analyzed a second dataset that included 317 epithelial breast cancer cells, 175 immune cells and 23 non-carcinoma stromal cells, from 11 patients diagnosed with breast cancer [24]. We only considered samples annotated previously as immune cells, which were annotated as T cells, B cells, and myeloid samples by clustering the gene expression signatures using non-negative factorization. The

Torang et al. Page 8 of 29

scRNA-seq samples were similarly pre-processed and analyzed using the proposed classifier, with the results shown in Fig. 4. The inner pie chart shows the prior cell 173 annotations reported by [24] and the outer chart shows the corresponding predicted cell annotation by our proposed classifier. Considering T cells as either CD4+ T 175 cell or CD8+ T cell, 94.4% of reported T cells are predicted as the same cell type and other 5.6% is predicted to be DC or NK cells. However, for reported B cells 177 and myeloid cells, we predicted relatively high portion of samples to be T cells ( 178 15.7% of B cells and 40% of myeloid cells). The rest of the myeloid samples were 179 predicted to be macrophages or dendritic cells. Collectively, our proposed classifier agreed with many of the prior cell annotations and annotated many of the samples 181 that were previously unresolved.

# Developing a classifier for T Helper cell subsets

To further apply this methodology to transcriptomic data, a separate classifier for distinguishing among T helper cells was developed using a similar approach to 185 the immune cell classifier. We explored different values of the regression parameter lambda to find the optimal number of genes for this new dataset and visualized 187 the performance of different lambdas by generating True-Negative samples using a bootstrapping approach whereby synthetic datasets were created by randomly 189 resampling testing datasets. Original testing data that were completely untouched 190 during training were used as True-Positive samples. The resulting True-Negative 191 and True-Positive samples were used to generate ROC curves (Fig. 5a) and the 192 AUC was used to score each lambda value. Generally, the lambda values for T 193 helper cell classifier represents the trade-off between retaining genes and keeping the AUC high. However, there appeared to be an inflection point at a lambda value 195 of 0.05 whereby adding additional genes, by increasing lambda, reduced the AUC. Consequently, we selected a lambda value equal to 0.05 (72 genes) for the T helper 197 classifier. The selected genes and their coefficients are listed in Table S1. The gene list was refined subsequently by developing a gene signature. 199

Similar to the immune cell classifier, the coefficients of the selected genes for the T
helper cell classifier correlated with their expression levels, as seen by comparing the
heatmaps shown in Fig. 5, panels b and c. For instance, FUT7 has been expressed in
almost all T helper cell samples except for iTreg that result in a negative coefficient

Torang et al. Page 9 of 29

for this cell type. In addition, there are sets of genes for each cell type that have large coefficients only for certain T helper cell subsets, like ALPK1, TBX21, IL12RB2, IFNG, RNF157 for Th1 that have low expression in other cells. As illustrated in 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. 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 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 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 cell types. 215 Finally by comparing the results of immune cell classifier with that of the T helper 216

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

### 222 Application of the Classifiers

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 immune contexture in melanoma and explored differences in immune contex-230 ture that associate with immune checkpoint response. Of note, some patients with 231 melanoma respond to ICIs durably but many others show resistance [28]. Specifi-232 cally, 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 Torang et al. Page 10 of 29

of cell type in each tumor sample as it was more informative and meaningful than using absolute cell numbers. It is notable that untreated and NoInfo samples likely include both ICI-resistant and ICI-sensitive tumors.

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

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 eigenvectors, the low values of PC2 correspond to a low prevalence of M $\phi$  or high percentage 255 of B cells. In addition, based on the first principal component (PC1), resistant sam-256 ples have either the lowest values of PC1 (Mel74, Mel75, Mel58, Mel 78), which 257 correspond to higher than average prevalence of CD8+ T cells, or the highest val-258 ues of PC1 (Mel60, Mel72, Mel94), which show a higher than average prevalence of 259 B cells. 260

In hierarchical clustering, the optimal number of clusters was selected based on calculation of different cluster indices using the NbClust R package [30] which mainly
identified two or three clusters as the optimal number. In considering three groupings of the hierarchical clustering results shown in 6d, seven out of eight ICI-resistant
samples clustered in first two clusters while the third cluster mainly contained untreated samples. The comparison of results from PCA and clustering analyses shows
that the first cluster contained samples with extreme low value of PC1 which itself
divided into two groups; one with extreme low value of PC2 and the other with

Torang et al. Page 11 of 29

higher amount of PC2. The second cluster located in highest amount of PC1 and lowest amount of PC2. All remained samples were clustered as third group, which were predominantly untreated samples. The difference in clustering suggests dissimilarities between ICI-resistant and untreated samples and the possibility of having ICI-sensitive tumors in untreated samples.

### 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, called in silico cytometry [31]. Lever-278 aging the classifier results, we generated corresponding gene signatures using binary 279 elastic-net logistic regression. Specifically, classifier genes with non-zero coefficients were used as initial features of the models, which were then regressed to the same 281 training and testing datasets as used for developing the classifiers. Lambda values were selected for each immune and T helper cell subset based on similar method of 283 lambda selection for classifiers and their values and corresponding AUC are shown 284 in Table S2. Finally, all generated signatures are summarized in Table S3. 285

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

In general, the gene sets generated from the logistic regression model performed well with far fewer requisite genes in the testing set, a desirable result for a gene set intended to be used for immunophenotyping. In Fig. 8, the results of the benchmarking are shown separated by comparative gene set. Both the CIBERSORT and Single-Cell derived gene sets contain an average of 64 and 135 genes, respectively, while the logistic regression gene set contains an average of just 19. The new logistic regression gene set performed comparably to the existing contemporary gene sets

Torang et al. Page 12 of 29

and far exceeded the performance of the manually curated gene set used previously
[6]. The benchmarking results indicate that the logistic regression gene sets are an
improvement in efficacy over compact gene sets, such as those that are manually
annotated or hand-picked. Meanwhile, the logistic regression gene sets also demonstrate an optimization of broader gene sets that contain too many genes for deep
specificity when used in further analysis. The inclusion of too many genes in a set
can dilute the real data across a constant level of noise, while including too few lacks
the power to draw conclusions with high confidence. The logistic regression gene
sets demonstrate a balance of these two issues through its highly refined selection
of genes that can be fine-tuned using its lambda parameter.

# Discussion

311

Recent developments in RNA sequencing enable a high fidelity view of the tran-312 scriptomic landscape associated with host immune response. Despite considerable progress in parsing this landscape using gene signatures, gaps remain in developing 314 unbiased signatures for individual immune cell types from healthy donors using high dimensional RNA-seq data. Here, we developed two classifiers - one for immune cell 316 subsets and one for T helper cell subsets - using elastic-net logistic regression with cross validation. The features of these classifiers were used as a starting point for 318 generating gene signatures that captured with fifteen binary elastic-net logistic re-319 gression models the most relevant gene sets to distinguish among different immune 320 cell types without including too much noise.

Gene signatures in previous studies have been developed and used mainly as a base 322 for deconvoluting the tumor microenvironment to find the presence of immune cells 323 from bulk RNA measures. Therefore, as the first step, determining cell-specific gene 324 signatures critically influences the results of deconvolution methods [32]. Newman 325 et al. defined gene signatures for immune cells using two-sided unequal variances t-test as base matrix for CIBERSORT [8]. In another study, Li et al. in devel-327 oping TIMER, generated gene signatures for six immune cell types with selecting genes with expression levels that have a negative correlation with tumor purity [9]. 329 More recently, Racle et al. developed a deconvolution tool based on RNA-seq data 330 (EPIC) by pre-selecting genes based on ranking by fold change and then selected 331 genes by manually curating and comparing the expression levels in blood and tuTorang et al. Page 13 of 29

mor microenvironment [10]. Finally, quanTIseq (the most recently developed tool for deconvolution) was developed for RNA-seq data based on the gene signatures generated by quantizing the expression levels into different bins and selecting high quantized genes for each cell type that have low or medium expression in other cell types [7]. Although all methods obtained high accuracy based on their developed signatures, a more rigorous and unbiased gene signature developed by RNA-seq data and precise feature selection methods can further improve the accuracy and validate the process for downstream analyses.

In addition, to identify cell types based on their transcriptome, clustering tech-341 niques have been used in many studies [33, 34]. However, there are high variability levels of gene expression even in samples from the same cell type. Moreover, tran-343 scriptomics data has high dimensions (tens of thousands) and this is too complicated for clustering techniques as only few number of genes are discriminative. To over-345 come these problems some studies used supervised machine learning methods like Support Vector Machine (SVM) [35, 36]. However, to the best of our knowledge, this 347 paper is the first to apply two-step regularized logistic regression on RNA-seq transcriptomic of immune cells. This method increases the chance to capture the most 349 discriminative set of genes for each cell type based on the power of an elastic-net [22]. In addition, using a two-step elastic net logistic regression enabled eliminating 351 the most irrelevant genes while keeping the highest number of possible significant 352 genes in the first step and more deeply selecting among them in the second step to 353 generate robust gene signatures for immune cells. 354

Moreover, contemporary methods have only considered a limited number of im-355 mune cell types, and specifically T helper subsets as individual cell types have been 356 neglected [23, 29, 24] in comprehensive studies. Therefore, the other novel aspect 357 of this study is the separation of models for immune cells and T helper cells and 358 development of gene signatures for a large number of immune cell types (fifteen different immune cell types) including different T helper cell subsets. The ability to 360 identify a greater number of immune cell types enables studying immune system in different diseases in more depth. As we used publicly available RNA-seq datasets for 362 immune cells and T helper cells, we acknowledge that our developed classifiers and 363 gene signatures may be still constrained by the limited number of samples specifically for T helper cells. As more data describing the transcriptome of immune cells Torang et al. Page 14 of 29

will become accessible, one can update the classifiers and gene signatures. Despite
the limited number of samples used in the approach, the developed classifiers can
even be applied to completely untouched and large datasets [23, 24] that have been
generated using scRNA-Seq technology which creates noisier data.

# 370 Conclusions

Here, we developed an immune cell classifier and classifier for T helper cell subsets 371 along with gene signatures to distinguish among fifteen different immune cell types. 372 Elastic-net logistic regression was used to generate classifiers with 10-fold cross-373 validation after normalizing and filtering two separate RNA-seq datasets that were 374 generated using defined homogeneous cell populations. Subsequently, we generated gene signatures using a second step of binary regularized logistic regression applied 376 to the RNA-seq data using previously selected classifier genes. As an external val-377 idation, the resulting classifiers accurately identified the type of immune cells in 378 scRNA-seq datasets. Our classifiers and gene signatures can be considered for dif-379 ferent downstream applications. First, the classifiers may be used to detect the type 380 of immune cells in under explored bulk tissue samples profiled using RNA-seq and 381 to verify the identity of immune cells annotated with low confidence. Second, the 382 gene signatures could be used to study tumor micro-environments and the interdependence of immune response with cancer cell phenotypes, which is emerging to 384 be an important clinical question. 385

# 386 Methods

B7 Data Acquisition

RNA-seq datasets for 15 different immune cell types including T helper cells, were obtained from ten different studies [37, 38, 39, 40, 41, 42, 43, 44, 45, 46], which were publicly accessible via the *Gene Expression Omnibus* [47]. The list of samples is provided as Supplementary Table S1. The cell types were divided into two groups: immune cells that include B cells, CD4+ and CD8+ T cells, monocytes (Mono), neutrophils (Neu), natural killer (NK) cells, dendritic cells (DC), macrophage (M $\phi$ ), classically (M1) and alternatively (M2) activated macrophages, and the T helper cells that include Th1, Th2, Th17, Th0, and Regulatory T cells (Treg). The goal was to train the gene selection model on immune cell types, and CD4+ T cell subsets (T helper cells), separately. If these two groups of cells are analyzed together, many

Torang et al. Page 15 of 29

of the genes that potentially could be used to discriminate among T helper cell

subsets might be eliminated as they overlap with genes associated with CD4+ T cells.

In short, a total of 233 samples were downloaded and divided into two sets of 185 and 48 samples, for immune cells and T helper cells, respectively. Moreover, immune cell samples were further divided into 108 training and 77 testing samples.

Training and testing numbers for T helper samples were 31 and 17, respectively. Training and testing data include samples from all studies. For a verification dataset, scRNA-seq data derived from CD45+ cell samples obtained from breast cancer

[24] and melanoma [23] were used with GEO accession numbers of GSE75688 and

#### 409 Data Normalization

408

GSE72056, respectively.

The expression estimates provided by the individual studies were used, regardless 410 of the underlying experimental and data processing methods (Table S1). For devel-411 oping individual gene signatures and cell classification models, we did not use raw 412 data due to sample heterogeneity such as different experimental methods and data processing techniques used by different studies as well as differences across biolog-414 ical sources. Rather, we applied a multistep normalization process before training models. To eliminate obvious insignificant genes from our data, for immune cell sam-416 ples, genes with expression values higher than or equal to five counts, in at least 417 five samples were kept, otherwise, they were eliminated from the study. However, 418 for T helper samples, due to fewer number of samples, four samples with values 419 higher than or equal to five counts were enough to be considered in the study. After 420 first step of filtering, the main normalization step was used to decrease dependency 421 of expression estimates to transcript length and GC-content[48, 49]. For all four 422 sets of samples, including training and testing samples for immune cells and for T helper cells, expression estimates were normalized separately by applying within La-424 neNormalization and betweenLaneNormalization functions from EDASeq package [50] in the R programming language (R 3.5.3), to remove GC-content biases and 426 between-lane differences in count distributions [50]. After normalization, the second 427 step of filtration, which was similar to the first step, was applied to eliminate genes 428 with insignificant expression.

Torang et al. Page 16 of 29

# 30 Missing Values

In contrast to previous studies that only considered intersection genes [51] and to avoid deleting discriminative genes, we kept genes with high expression as much as 432 possible. However, for most of genes, values for some samples were not reported. Hence, to deal with these missing values, we used an imputation method [52] and 434 instead of mean imputation we set a dummy constant since mean imputation in this case is not meaningful and can increase error. Specifically, we generated a training 436 set for each group of cell types, by duplicating the original training set 100 times and randomly eliminating ten percent of expression values. We next set -1 for all 438 these missing values (both original missing values and those we eliminated) as a dummy constant because all values are positive and it is easier for the system to 440 identify these values as noise. This approach makes the system learn to neglect a specific value (-1) and treat it like noise, instead of learning it as a feature of the 442 samples.

#### 444 Classifier Training and Testing

Considering the few number of training samples in comparison with the high dimensions (15453 genes in immune cell samples and 9146 genes in the T helper 446 samples) and to avoid both over fitting the model and adding noise to the prediction model, we used regularization with logistic regression to decrease the total 448 number of genes and select the most discriminative set of genes. To perform gene selection, we trained a lasso-ridge logistic regression (elastic-net) model, which au-450 tomatically sets the coefficients of a large number of genes to zero and prunes 451 the number of genes as features of the classifier. We cross-validated the model by 452 implementing cv.glmnet function with nfold=10 from glmnet package [21] in R pro-453 gramming language, using training sets for both groups of cell types. We normalized 454 the gene expression values using a log2 transform over training sets to decrease the range of values that can affect the performance of the model (log2(counts+1)). In 456 order to find the optimal number of genes, we tried seven different lambdas and tested the results over the testing samples (cv.qlmnet(family="multinomial", al-458 pha=0.93, thresh=1e-07, lambda=c(0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001), 459 type.multinomial="grouped", nfolds=10)). To select the optimal value for lambda, 460 True-Negative samples were generated using a bootstrapping approach that ranTorang et al. Page 17 of 29

domly samples testing datasets to create a synthetic dataset with similar size and complexity but without underlying biological correlation, then we generated ROC curves and considered original testing datasets as True-Positive samples.

# Developing Gene Signatures

Genes selected by the classifier models were used as initial sets to build gene signa-466 tures. In this case, we trained a new binary elastic-net model for each cell type by considering a certain cell type as one class and all other cell types as another class. 468 The training and testing samples used to build gene signatures were the training 460 and testing samples used in developing the classifiers with the difference being that 470 they only contained the selected genes. Similar steps including dealing with missing 471 values, applying log2 and visualization by ROC to select optimal number of genes 472 were applied for each cell type. This two-step gene selection approach has the ad-473 vantage that it eliminates a large number of undiscriminating genes at the first and 474 finally select few number of genes for each cell type. 475

#### 476 Benchmarking

Fisher exact testing was used for each gene set to characterize true and systematically scrambled data as a measure of performance of the gene set as a means of 478 distinguishing between cell subtypes. In order to establish negative control values for determining specificity, a bootstrapping approach was used [53], where data was 480 scrambled by randomly resampling with replacement expression values by gene as 481 well as by patient to create a synthetic dataset with a similar size and complexity 482 of the original dataset. The threshold for expression binarization for Fisher exact 483 testing was selected based on gene expression histograms of the data to separate 484 the measured expression from background noise levels, with 2.48 being used as 485 the threshold (after  $\log 2$  normalization). One-thousand iterations  $(N_{boot})$  were pro-486 cessed and compiled in order to produce ROC curves with 95% confidence intervals shaded about the averaged ROC curve for each gene set's performance. A boot-488 strapping approach for generating a negative control sample is appropriate when a sufficiently large bootstrap sample (i.e.,  $N_{boot} \geq 1000$ ) and the original dataset 490 is sufficiently diverse (i.e.,  $N_{data} \geq 30$ ) [54]. The tested gene sets were the logistic 491 regression gene set, the CIBERSORT gene set [8], the single cell gene set [29], and 492 the manually curated gene set that had been used previously [6].

Torang et al. Page 18 of 29

# List of abbreviations

495 ROC: receiver-operator curves

scRNA-seq: single-cell RNA-seq

497 AUC: area under the ROC curve

498 CNV: copy number variation

499 PCA: principal component analysis

500 ICI: immune checkpoint inhibitor

501 SVM: support vector machine

502

#### Declarations

#### 604 Ethics approval and consent to participate

- The results described in this manuscript consist of secondary analyses of existing data and was determined by the
- 506 West Virginia University IRB to qualify for an exemption from human subject research under U.S. HHS regulations
- 507 45 CFR 46.101(b)(4).

#### 508 Consent for publication

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

#### 510 Availability of data and material

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

#### 515 Competing interests

The authors declare that they have no competing interests.

#### 517 Funding

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

# 522 Authors' contributions

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

### 525 Acknowledgements

#### 526 Author details

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

#### 530 References

- 1. Carmona, S.J., Teichmann, S.A., Ferreira, L., Macaulay, I.C., Stubbington, M.J., Cvejic, A., Gfeller, D.:
- Single-cell transcriptome analysis of fish immune cells provides insight into the evolution of vertebrate immune
- cell types. Genome research, 207704 (2017)

Torang et al. Page 19 of 29

- 2. Bendall, S.C., Simonds, E.F., Qiu, P., El-ad, D.A., Krutzik, P.O., Finck, R., Bruggner, R.V., Melamed, R.,
- Trejo, A., Ornatsky, O.I., *et al.*: Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum. Science **332**(6030), 687–696 (2011)
- 3. Shay, T., Kang, J.: Immunological genome project and systems immunology. Trends in immunology **34**(12), 602–609 (2013)
- 4. Kinter, A.L., Hennessey, M., Bell, A., Kern, S., Lin, Y., Daucher, M., Planta, M., McGlaughlin, M., Jackson,
- 8., Ziegler, S.F., et al.: Cd25+ cd4+ regulatory t cells from the peripheral blood of asymptomatic hiv-infected
- individuals regulate cd4+ and cd8+ hiv-specific t cell immune responses in vitro and are associated with
- favorable clinical markers of disease status. Journal of Experimental Medicine 200(3), 331–343 (2004)
- 5. Vegh, P., Haniffa, M.: The impact of single-cell rna sequencing on understanding the functional organization of the immune system. Briefings in functional genomics (2018)
- 6. Kaiser, J.L., Bland, C.L., Klinke, D.J.: Identifying causal networks linking cancer processes and anti-tumor immunity using bayesian network inference and metagene constructs. Biotechnology progress **32**(2), 470–479 (2016)
- 7. Finotello, F., Mayer, C., Plattner, C., Laschober, G., Rieder, D., Hackl, H., Krogsdam, A., Posch, W.,
  Wilflingseder, D., Sopper, S., et al.: quantiseq: quantifying immune contexture of human tumors. bioRxiv,
  223180 (2017)
- 8. Newman, A.M., Liu, C.L., Green, M.R., Gentles, A.J., Feng, W., Xu, Y., Hoang, C.D., Diehn, M., Alizadeh,
  A.A.: Robust enumeration of cell subsets from tissue expression profiles. Nature methods 12(5), 453 (2015)
- 9. Li, B., Severson, E., Pignon, J.-C., Zhao, H., Li, T., Novak, J., Jiang, P., Shen, H., Aster, J.C., Rodig, S., *et al.*: Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. Genome biology **17**(1), 174 (2016)
- 10. Racle, J., de Jonge, K., Baumgaertner, P., Speiser, D.E., Gfeller, D.: Simultaneous enumeration of cancer and immune cell types from bulk tumor gene expression data. Elife **6**, 26476 (2017)
- 11. Kidd, B.A., Peters, L.A., Schadt, E.E., Dudley, J.T.: Unifying immunology with informatics and multiscale biology. Nature immunology **15**(2), 118 (2014)
- 12. Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A.,
- Pomeroy, S.L., Golub, T.R., Lander, E.S., et al.: Gene set enrichment analysis: a knowledge-based approach for
- interpreting genome-wide expression profiles. Proceedings of the National Academy of Sciences 102(43),
- 563 **15545–15550** (2005)
- 13. Liberzon, A., Subramanian, A., Pinchback, R., Thorvaldsdóttir, H., Tamayo, P., Mesirov, J.P.: Molecular
   signatures database (msigdb) 3.0. Bioinformatics 27(12), 1739–1740 (2011)
- I4. Zheng, C.-H., Chong, Y.-W., Wang, H.-Q.: Gene selection using independent variable group analysis for tumor
   classification. Neural Computing and Applications 20(2), 161–170 (2011)
- 568
   15. Wu, M.-Y., Dai, D.-Q., Shi, Y., Yan, H., Zhang, X.-F.: Biomarker identification and cancer classification based
   569
   on microarray data using laplace naive bayes model with mean shrinkage. IEEE/ACM transactions on
   570
   computational biology and bioinformatics 9(6), 1649–1662 (2012)
- 16. Cui, Y., Zheng, C.-H., Yang, J., Sha, W.: Sparse maximum margin discriminant analysis for feature extraction and gene selection on gene expression data. Computers in biology and medicine **43**(7), 933–941 (2013)
- 17. Algamal, Z.Y., Lee, M.H.: Regularized logistic regression with adjusted adaptive elastic net for gene selection in high dimensional cancer classification. Computers in biology and medicine **67**, 136–145 (2015)
- 18. Liang, Y., Liu, C., Luan, X.-Z., Leung, K.-S., Chan, T.-M., Xu, Z.-B., Zhang, H.: Sparse logistic regression with a I 1/2 penalty for gene selection in cancer classification. BMC bioinformatics **14**(1), 198 (2013)
- 19. Bielza, C., Robles, V., Larrañaga, P.: Regularized logistic regression without a penalty term: An application to cancer classification with microarray data. Expert Systems with Applications **38**(5), 5110–5118 (2011)
- 20. Cawley, G.C., Talbot, N.L.: Gene selection in cancer classification using sparse logistic regression with bayesian regularization. Bioinformatics **22**(19), 2348–2355 (2006)
- Friedman, J., Hastie, T., Tibshirani, R.: Regularization paths for generalized linear models via coordinate
   descent. Journal of statistical software 33(1), 1 (2010)
- Zou, H., Hastie, T.: Regularization and variable selection via the elastic net. Journal of the Royal Statistical
   Society: Series B (Statistical Methodology) 67(2), 301–320 (2005)

Torang et al. Page 20 of 29

- Tirosh, I., Izar, B., Prakadan, S.M., Wadsworth, M.H., Treacy, D., Trombetta, J.J., Rotem, A., Rodman, C.,
   Lian, C., Murphy, G., et al.: Dissecting the multicellular ecosystem of metastatic melanoma by single-cell
   rna-sea. Science 352(6282), 189–196 (2016)
- 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, Y.H., *et al.*: Single-cell rna-seq enables comprehensive tumour and immune cell profiling in primary breast
- cancer. Nature communications **8**, 15081 (2017)
- 591 25. Caligiuri, M.A.: Human natural killer cells, Blood 112(3), 461-469 (2008)
- 592 26. Chen, E.Y., Tan, C.M., Kou, Y., Duan, Q., Wang, Z., Meirelles, G.V., Clark, N.R., Ma'ayan, A.: Enrichr: 593 interactive and collaborative html5 gene list enrichment analysis tool. BMC bioinformatics **14**(1), 128 (2013)
- 27. Kuleshov, M.V., Jones, M.R., Rouillard, A.D., Fernandez, N.F., Duan, Q., Wang, Z., Koplev, S., Jenkins, S.L.,
- Jagodnik, K.M., Lachmann, A., et al.: Enrichr: a comprehensive gene set enrichment analysis web server 2016
- update. Nucleic acids research 44(W1), 90–97 (2016)
- 597 28. Sharma, P., Hu-Lieskovan, S., Wargo, J.A., Ribas, A.: Primary, adaptive, and acquired resistance to cancer 598 immunotherapy. Cell **168**(4), 707–723 (2017)
- 599 29. Jerby-Arnon, L., Shah, P., Cuoco, M.S., Rodman, C., Su, M.-J., Melms, J.C., Leeson, R., Kanodia, A., Mei, S.,
- Lin, J.-R., et al.: A cancer cell program promotes t cell exclusion and resistance to checkpoint blockade. Cell 175(4), 984–997 (2018)
- 30. Charrad, M., Ghazzali, N., Boiteau, V., Niknafs, A., Charrad, M.M.: Package 'nbclust'. Journal of statistical
   software 61, 1–36 (2014)
- 604 31. Newman, A.M., Steen, C.B., Liu, C.L., Gentles, A.J., Chaudhuri, A.A., Scherer, F., Khodadoust, M.S.,
- Esfahani, M.S., Luca, B.A., Steiner, D., Diehn, M., Alizadeh, A.A.: Determining cell type abundance and expression from bulk tissues with digital cytometry. Nat. Biotechnol. (2019)
- 32. Finotello, F., Trajanoski, Z.: Quantifying tumor-infiltrating immune cells from transcriptomics data. Cancer
   Immunology, Immunotherapy 67(7), 1031–1040 (2018)
- 33. Xu, C., Su, Z.: Identification of cell types from single-cell transcriptomes using a novel clustering method.
  Bioinformatics 31(12), 1974–1980 (2015)
- 34. Grün, D., Lyubimova, A., Kester, L., Wiebrands, K., Basak, O., Sasaki, N., Clevers, H., van Oudenaarden, A.:
   Single-cell messenger rna sequencing reveals rare intestinal cell types. Nature 525(7568), 251 (2015)
- 35. Hu, Y., Hase, T., Li, H.P., Prabhakar, S., Kitano, H., Ng, S.K., Ghosh, S., Wee, L.J.K.: A machine learning approach for the identification of key markers involved in brain development from single-cell transcriptomic
- data. BMC genomics 17(13), 1025 (2016)
- 36. Yao, F., Zhang, C., Du, W., Liu, C., Xu, Y.: Identification of gene-expression signatures and protein markers for breast cancer grading and staging. PloS one **10**(9), 0138213 (2015)
- 37. Linsley, P.S., Speake, C., Whalen, E., Chaussabel, D.: Copy number loss of the interferon gene cluster in melanomas is linked to reduced t cell infiltrate and poor patient prognosis. PloS one **9**(10), 109760 (2014)
- 620 38. Hoek, K.L., Samir, P., Howard, L.M., Niu, X., Prasad, N., Galassie, A., Liu, Q., Allos, T.M., Floyd, K.A., Guo,
- Y., et al.: A cell-based systems biology assessment of human blood to monitor immune responses after influenza vaccination. PloS one **10**(2). 0118528 (2015)
- 623 39. Beyer, M., Mallmann, M.R., Xue, J., Staratschek-Jox, A., Vorholt, D., Krebs, W., Sommer, D., Sander, J.,
- Mertens, C., Nino-Castro, A., et al.: High-resolution transcriptome of human macrophages. PloS one 7(9),
- 625 45466 (2012)
- 40. Şenbabaoğlu, Y., Gejman, R.S., Winer, A.G., Liu, M., Van Allen, E.M., de Velasco, G., Miao, D., Ostrovnaya,
- 627 I., Drill, E., Luna, A., et al.: Tumor immune microenvironment characterization in clear cell renal cell carcinoma
- identifies prognostic and immunotherapeutically relevant messenger rna signatures. Genome biology **17**(1), 231 (2016)
- 41. Corces, M.R., Buenrostro, J.D., Wu, B., Greenside, P.G., Chan, S.M., Koenig, J.L., Snyder, M.P., Pritchard,
- J.K., Kundaje, A., Greenleaf, W.J., *et al.*: Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. Nature genetics **48**(10), 1193 (2016)
- 42. Kumar, N.A., Cheong, K., Powell, D.R., da Fonseca Pereira, C., Anderson, J., Evans, V.A., Lewin, S.R.,
- 634 Cameron, P.U.: The role of antigen presenting cells in the induction of hiv-1 latency in resting cd4+ t-cells.
- 635 Retrovirology **12**(1), 76 (2015)

Torang et al. Page 21 of 29

- 43. Zhang, H., Xue, C., Shah, R., Bermingham, K., Hinkle, C.C., Li, W., Rodrigues, A., Tabita-Martinez, J., Millar,
   J.S., Cuchel, M., et al.: Functional analysis and transcriptomic profiling of ipsc-derived macrophages and their
   application in modeling mendelian disease. Circulation research, 114 (2015)
- 639 44. Kanduri, K., Tripathi, S., Larjo, A., Mannerström, H., Ullah, U., Lund, R., Hawkins, R.D., Ren, B.,
- Lähdesmäki, H., Lahesmaa, R.: Identification of global regulators of t-helper cell lineage specification. Genome medicine **7**(1), 122 (2015)
- 45. Spurlock III, C.F., Tossberg, J.T., Guo, Y., Collier, S.P., Crooke III, P.S., Aune, T.M.: Expression and functions
   of long noncoding rnas during human t helper cell differentiation. Nature communications 6, 6932 (2015)
- 46. Schmidt, A., Marabita, F., Kiani, N.A., Gross, C.C., Johansson, H.J., Éliás, S., Rautio, S., Eriksson, M.,
- Fernandes, S.J., Silberberg, G., et al.: Time-resolved transcriptome and proteome landscape of human
- regulatory t cell (treg) differentiation reveals novel regulators of foxp3. BMC biology 16(1), 47 (2018)
- Edgar, R., Domrachev, M., Lash, A.E.: Gene expression omnibus: Ncbi gene expression and hybridization array
   data repository. Nucleic acids research 30(1), 207–210 (2002)
- 48. Oshlack, A., Wakefield, M.J.: Transcript length bias in rna-seq data confounds systems biology. Biology direct 4(1), 14 (2009)
- 49. Robinson, M.D., Oshlack, A.: A scaling normalization method for differential expression analysis of rna-seq
   data. Genome biology 11(3), 25 (2010)
- 653 50. Risso, D., Schwartz, K., Sherlock, G., Dudoit, S.: Gc-content normalization for rna-seq data. BMC 654 bioinformatics **12**(1), 480 (2011)
- 55. Schwalie, P.C., Ordóñez-Morán, P., Huelsken, J., Deplancke, B.: Cross-tissue identification of somatic stem and
   progenitor cells using a single-cell rna-sequencing derived gene signature. Stem Cells 35(12), 2390–2402 (2017)
- 52. García-Laencina, P.J., Sancho-Gómez, J.-L., Figueiras-Vidal, A.R.: Pattern classification with missing data: a
   review. Neural Computing and Applications 19(2), 263–282 (2010)
- 659 53. Efron, B., Tibshirani, R.: An Introduction to the Bootstrap. Chapman and Hall, London (1993)
- 54. Chernick, M.R.: Bootstrap Methods: A Practitioner's Guide, pp. 150–151. Wiley, New York (1999)

Torang et al. Page 22 of 29

#### 661 Figures

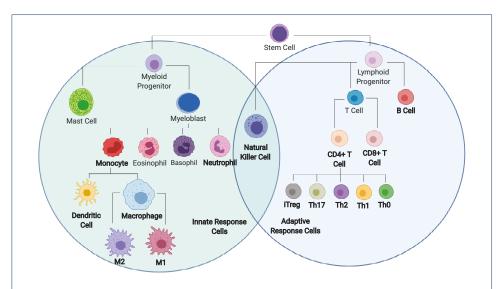


Figure 1 Lineage tree representation of cells of the immune system. Immune cells are derived from hematopoietic stem cells (HSCs). HSCs differentiate into lymphoid and myeloid progenitors that further branch out to the more specific cell types associated with adaptive and innate immunity. This Figure indicates the main immune cell subsets and arrows are to show lineage relationships. Gene signatures were developed in this study for immune cells highlighted in bold.

Torang et al. Page 23 of 29

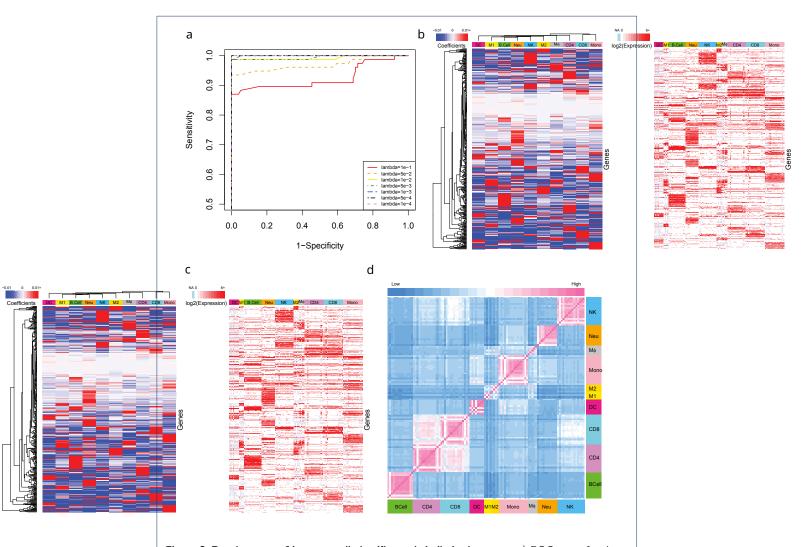


Figure 2 Development of immune cell classifier and similarity heatmap. a) ROC curve for the immune cell classifier was calculated using the indicated lambda values (shown in different colors and line styles) and 10-fold cross validation. The lambda value that maximized the AUC value was used for subsequent calculations. Elastic-net logistic regression was used to discriminate among ten immune cell types, where the value of the non-zero coefficients (panel b), expression levels (panel c), and similarity map (panel d) for the 452 genes included in the classifier are indicated by color bars for each panel . In panel b, blue to red color scheme indicates coefficients ranging from negative to positive values. Ordering of the genes is the same in panels b and c. In panel c, light blue indicates missing values and the intensity of red color (white/red color scale on the top-left) shows the log base 2 expression level. A color bar on top of this panel was used to separate samples of each cell type. Panel d illustrates the similarity between samples calculated using distance matrix based on same 452 genes. Color bars on the left and bottom sides are to separate samples of each cell type and the top color bar (light blue/pink color scale) shows the intensity of similarity or dissimilarity of samples.

Torang et al. Page 24 of 29

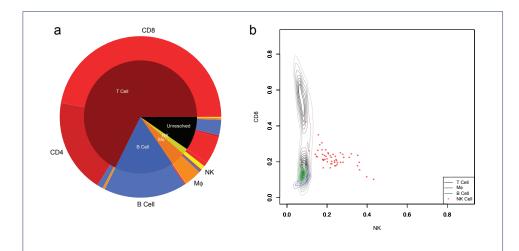


Figure 3 Immune cell annotation prediction based on scRNA-seq data against prior annotations reported in melanoma dataset. a) The inner pie chart summarizes the cell annotations reported by Tirosh et al [23] and includes 298 unannotated CD45-positive non-malignant cells (labeled as Unresolved) isolated from melanoma tissue samples. Unannotated samples were acquired following gating for CD45+ single cells and classified as non-malignant based on inferred copy number variation patterns. Using gene expression values reported for each scRNA-seq sample, a new cell annotation was determined based on the closest match with the alternative cell signatures determined using elastic-net logistic regression, which are summarized in outer pie chart. b) The contour plot for the likelihood of a sample to be either an NK cell or CD8+ T cell based on gene expression stratified by cells previously annotated by [23] to be T cells, macrophages, B cells, or NK cells.

Torang et al. Page 25 of 29

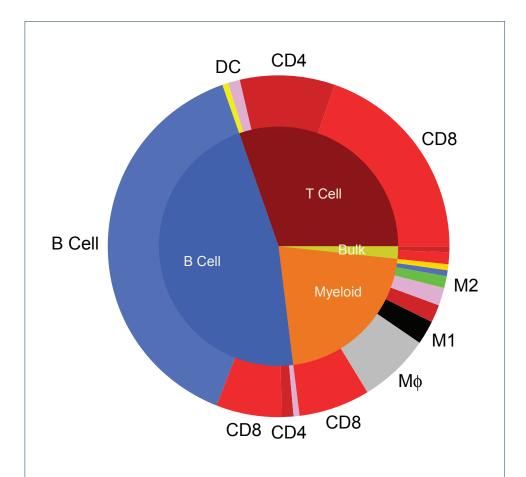


Figure 4 Immune cell annotation prediction against prior annotations reported in breast cancer scRNA-seq dataset. The inner pie chart summarizes the cell annotations reported by Chung et al [24], which annotated scRNA-seq results by clustering by gene ontology terms using likelihood ratio test. Using the gene expression profile reported for each scRNA-seq sample, a new cell annotation was determined based on the closest match with the alternative cell signatures determined using elastic-net logistic regression, which is summarized in the outer pie chart.

Torang et al. Page 26 of 29

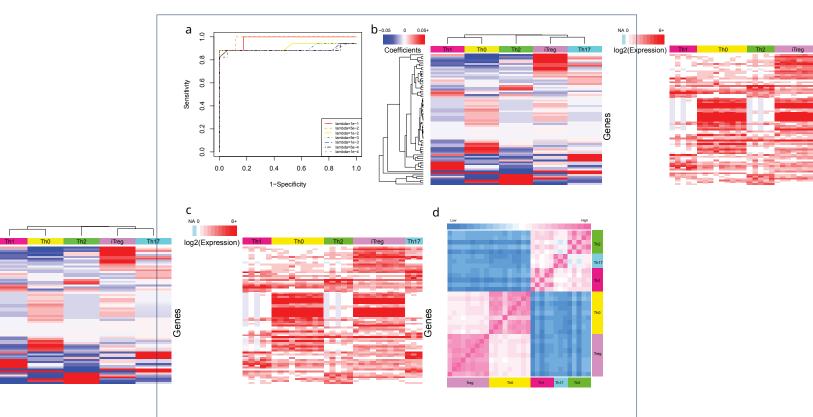


Figure 5 Development of T helper cell classifier and similarity heatmaps a) ROC curve for the T helper cell classifier was calculated using the indicated lambda values (shown in different colors and line styles) and 10-fold cross validation. The lambda value that maximized the AUC value was used for subsequent calculations. Elastic-net logistic regression to discriminate among five T helper cell types, where the value of the non-zero coefficients (panel b), expression levels (panel c), and similarity map (panel d) for the 72 genes included in the classifier are indicated by color bars for each panel. In panel b, blue to red color scheme indicates coefficients ranging from negative to positive values. Ordering of the genes is the same in panels b and c. In panel c, light blue indicates missing values and the intensity of red color (white/red color scale on the top-left) indicates the log base 2 expression level. A color bar on top of this panel was used to separate samples of each cell type. Panel d illustrates the similarity between samples calculated using an euclidean distance matrix based on the same 72 genes, where the color indicates the distance (pink: high similarity/low distance; blue: low similarity/high distance). Color bar on the top/side of the heatmap indicates the cell type of origin.

Torang et al. Page 27 of 29

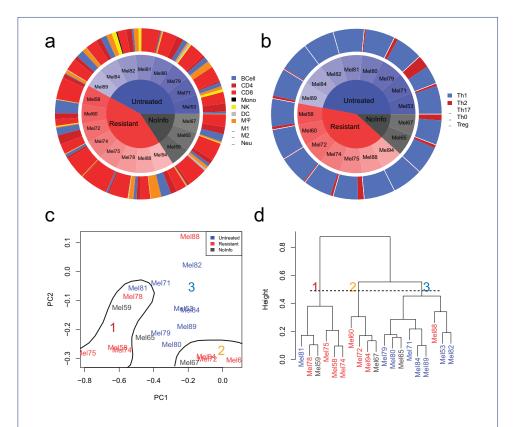


Figure 6 Annotation of scRNA-seq results from melanoma dataset stratified by patient treatment status. Treatment status of patients diagnosed with melanoma was stratified based on their response to ICIs ([23, 29]). a) The distribution in immune cell annotations and b) T helper cell annotations based on scRNA-seq data were separated into samples obtained from ICI-resistant tumors, untreated tumors, and tumors reported in melanoma data without information about treatment status. Distributions are shown based on the percentage of all immune cells measured for each patient. Cell annotations were based on immune cell classifier and T helper cell classifier results. c) PCA analysis was applied to the data obtained from both classifiers and the results for the first and second principal components were plotted. Red, blue, and grey colors indicate resistant, untreated and NoInfo (samples that have no information about their treatment status in the reference works) tumors, respectively. d) Samples were hierarchically clustered based on the percentages of the nine immune cells and five T helper cells and same coloring applied to show tumor types.

Torang et al. Page 28 of 29

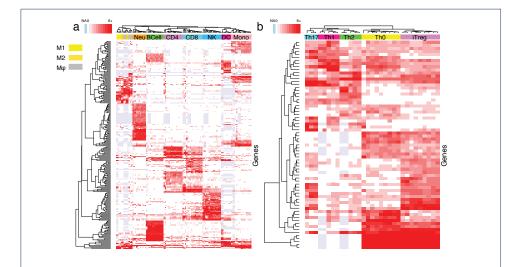


Figure 7 Heatmaps of the expression levels for the final list of genes created by gene signatures. The expression of genes retained in immune cell signatures (panel a) and T helper cell signatures (panel b) were clustered by similarity in expression levels (rows) and by similarity in samples (columns). The color bar at the top indicates the samples cell type. Light blue shows missing values and the intensity of red color (white/red color scale on the top-left color bar) indicates the log base 2 expression level in both panels.

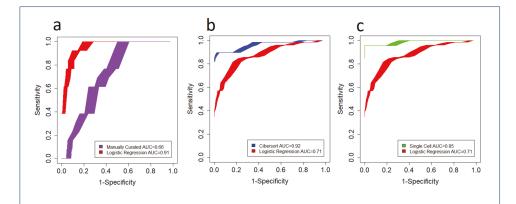


Figure 8 Benchmarking ROC performance curves. ROC curves to illustrate relative performance between logistic regression gene set and the manually curated (Panel A), CIBERSORT (Panel B), and single cell gene sets (Panel C). The logistic regression gene set's performance is shown in red. Shaded regions are 95% confidence intervals about the average ROC curve simulated from 1000 iterations.

Torang et al. Page 29 of 29

- 2 Additional Files
- $_{\rm 663}$   $\,$  Table S1. Coefficients of immune cell classifier and T helper cell classifier
- 664 Coefficients of immune cell classifier were located in the first sheet and coefficients of T helper cells were located in
- the second sheet.
- 666 Table S2. Lambda Selection by AUC Values
- 667 Lambdas with corresponding calculated AUC. The final column shows the selected lambdas
- 668 Table S3. Genes in developed gene signature for immune and T helper cells
- Yellow boxes show genes with negative impact in possibility of being related cell type.
- $\,$  Table S4. Data information used in training models.
- $\,$  The second sheet shows the names that were used in creating the datasets.