

# Restricting datasets to classifiable samples augments discovery of immune disease markers

**Gunther Glehr**

UKR <https://orcid.org/0000-0002-1495-9162>

**Paloma Riquelme**

University Hospital Regensburg

**Katharina Kronenberg**

University Hospital Regensburg

**Robert Lohmayer**

Leibniz Institute for Immunotherapy

**Victor Lopez-Madrone**

Institute de Neurosciences des systèmes

**Michael Kapinsky**

Beckman Coulter GmbH

**Hans Schlitt**

University of Regensburg <https://orcid.org/0000-0002-3874-0296>

**Edward Geissler**

University Hospital Regensburg

**Rainer Spang**

University of Regensburg

**Sebastian Haferkamp**

University Hospital Regensburg

**James Hutchinson** (✉ [james.hutchinson@klinik.uni-regensburg.de](mailto:james.hutchinson@klinik.uni-regensburg.de))

UKR <https://orcid.org/0000-0002-1199-4210>

---

## Article

### Keywords:

**Posted Date:** June 6th, 2023

**DOI:** <https://doi.org/10.21203/rs.3.rs-2921819/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

**Additional Declarations:** There is **NO** Competing Interest.

---

2

3 **Restricting datasets to classifiable samples augments**

4 **discovery of immune disease markers**

5

6 Gunther Glehr,<sup>1</sup> Paloma Riquelme,<sup>1</sup> Katharina Kronenberg,<sup>1</sup> Robert Lohmayer,<sup>2</sup>

7 Víctor J. López-Madrona,<sup>3</sup> Michael Kapinsky,<sup>4</sup> Hans J. Schlitt,<sup>1</sup> Edward K. Geissler,<sup>1</sup>

8 Rainer Spang,<sup>5</sup> Sebastian Haferkamp<sup>6</sup> and James A. Hutchinson<sup>1†</sup>

9

- 10 1. Department of Surgery, University Hospital Regensburg, Regensburg, Germany
- 11 2. Algorithmic Bioinformatics Research Group, Leibniz Institute for Immunotherapy,
- 12 Regensburg, Germany
- 13 3. Aix Marseille Univ, INSERM, Institute of Systems Neuroscience, Marseille, France
- 14 4. Beckman Coulter Life Sciences GmbH, Krefeld, Germany
- 15 5. Department of Statistical Bioinformatics, University of Regensburg, Regensburg,
- 16 Germany
- 17 6. Department of Dermatology, University Hospital Regensburg, Regensburg, Germany
- 18

19 † *Corresponding author:* Prof. Dr. James A. Hutchinson

20 Department of Surgery

21 University Hospital Regensburg

22 Franz-Josef-Strauß-Allee-11

23 93053 Regensburg, Germany

24 Tel: +49-941-944-6831

25 Email: [james.hutchinson@ukr.de](mailto:james.hutchinson@ukr.de)

26 **Running title:** Dataset restriction aids immune disease marker discovery.

27 **Grant support:** This work was supported by the Bristol-Myers Squibb (BMS) Immune  
28 Oncology Foundation (Award FA19-009) and the Coronavirus Research Fund of the Bavarian  
29 State Ministry for Science and Art.

30 **Authorship:** GG and JAH designed and performed the project, and wrote the manuscript. KK  
31 and PR performed experiments and analyzed data. RL, VJLM and RS gave expert  
32 computational and statistical advice. MK gave expert advice about flow cytometry. HJS and  
33 EKG gave critical feedback. SH provided clinical samples and information, and gave expert  
34 Dermatological Oncology opinion.

35 **Non-standard abbreviations:** Area under the ROC curve (AUC); cytomegalovirus (CMV);  
36 correct classification rate (CCR); double negatives (DN); double positives (DP); Dirichlet  
37 distribution with K categories and its concentration parameters  $\alpha$ :  $Dir(\alpha_1, \alpha_2, \dots, \alpha_K)$ ; false  
38 positive rate (FPR); glutamic oxaloacetic transaminase (GOT); glutamic pyruvic transaminase  
39 (GPT); immune checkpoint inhibitor (ICI); immune-related adverse events (irAE); Kullback-  
40 Leibler divergence (KL-divergence); ROC curve with a long tail towards bottom-left (left-  
41 skewed ROC curve); samples below restriction (marker<sup>LOW</sup> samples); samples above  
42 restriction (marker<sup>HIGH</sup> samples); median fluorescence intensity (MFI); Normal distribution with  
43 mean  $\mu$  and variance  $\sigma^2$  ( $\mathcal{N}(\mu, \sigma^2)$ ); negative predictive value (NPV); positive predictive value  
44 (PPV); partial area under the curve (pAUC); Receiver operating characteristic (ROC); ROC  
45 curve with a long tail towards top-right (right-skewed ROC curve); restricted AUC (rAUC);  
46 restricted standardized AUC (rzAUC); rzAUC of marker<sup>HIGH</sup> samples ( $rzAUC_{HIGH}$ ); rzAUC of  
47 marker<sup>LOW</sup> samples ( $rzAUC_{LOW}$ ); true positive rate (TPR); effector memory T cells ( $T_{EM}$ );  
48 effector memory T cells re-expressing CD45RA ( $T_{EMRA}$ ); central memory T cells ( $T_{CM}$ ); naïve T  
49 cells ( $T_{naïve}$ ); true positive proportion asymmetric (TPP-asymmetric); true negative proportion  
50 asymmetric (TNP-asymmetric)

52 **ABSTRACT**

53 Immunological diseases are typically heterogeneous in clinical presentation, severity and  
54 response to therapy. Markers of immune diseases often reflect this variability, especially  
55 compared to their regulated behavior in health. This leads to a common, unarticulated problem  
56 that frustrates marker discovery and interpretation: Unequal variance of immune disease marker  
57 expression between patient classes necessarily limits a marker's informative range. To solve  
58 this problem, we introduce dataset restriction, a procedure that splits datasets into classifiable  
59 and unclassifiable samples. Applied to synthetic flow cytometry data, restriction identified  
60 markers that were otherwise disregarded. In advanced melanoma, restriction found new  
61 markers of immune-related adverse event risk after immunotherapy and enabled multivariate  
62 models that accurately predicted immunotherapy-related hepatitis. Hence, dataset restriction  
63 augments discovery of immune disease markers, increases predictive certainty for classifiable  
64 samples and improves multivariate models incorporating markers with a limited informative  
65 range. This principle can be directly extended to any classification task.

## 66 INTRODUCTION

67 The immune system detects pathological challenges with exquisite sensitivity and specificity,  
68 which enables it to mount appropriate protective responses<sup>1</sup>. Widely distributed immune cell  
69 subsets are responsible for sensing pathogens, tissue injury and cellular stress through diverse  
70 receptor systems<sup>2-4</sup>. These disease-related signals are then amplified through humoral and  
71 cellular cascades that stimulate the migration, expansion and activation of particular effector  
72 cell populations<sup>5</sup>. By capturing information about the precise nature of the immune response,  
73 we can draw inferences about the triggering event, allowing us to develop diagnostic or  
74 prognostic models to guide personalized treatment decisions<sup>6</sup>.

75 Flow cytometry is a sophisticated, fast and relatively inexpensive method for analyzing the  
76 properties of single cells from a uniform cell suspension<sup>7</sup>. In clinical practice, flow cytometry  
77 is commonly used to profile leucocyte subset distribution in patient blood samples, especially  
78 in the context of hematological malignancies and infectious diseases<sup>8</sup>. Modern cytometers  
79 simultaneously collect data about expression of multiple proteins in single cells, while also  
80 allowing us to interrogate many millions of cells from a single sample<sup>9</sup>. This enables accurate  
81 identification of narrowly defined cell subsets, including rare populations, as well as broadly  
82 surveying many leucocyte subsets<sup>10</sup>. This rich information is captured as a data matrix for each  
83 sample with an unordered number of rows corresponding to cells and a defined number of cell-  
84 associated features as columns<sup>11</sup>.

85 Applications of flow cytometry in clinical diagnostics are growing rapidly<sup>12</sup>. Of special interest,  
86 recent reports claim that immunophenotyping of peripheral blood leucocytes can be used to  
87 predict immune-related adverse events (irAE) following immune checkpoint inhibitor (ICI)-  
88 therapy<sup>13-16</sup>. Combined treatment with anti-PD-1 (Nivolumab) and anti-CTLA-4 antibody

89 (Ipilimumab) is now first-line therapy for many patients with unresectable metastatic  
90 melanoma<sup>17</sup>. Its effectiveness is remarkable in terms of clinical response rates, progression-free  
91 survival and overall survival; however, immune-mediated complications, such as colitis or  
92 hepatitis, present a significant clinical concern<sup>18</sup>. Life-threatening reactions are uncommon<sup>19</sup>,  
93 but they often require interruption or discontinuation of immunotherapy, and introduction of  
94 glucocorticoids or non-steroidal immunosuppressants<sup>20</sup>. Clinically applicable, robust markers  
95 to guide irAE prevention or treatment strategies in patients would be extremely useful<sup>21</sup>.

96 Extracting reliable predictive information from flow cytometry measurements is difficult  
97 because disease-related changes are often small compared to typical biological and technical  
98 variations<sup>22</sup>. This is especially true when investigating systemic changes in peripheral blood  
99 samples for signals that reflect localized disease<sup>23</sup>. Consequently, we often rely upon  
100 computational methods to perceive small and multivariate, but consistent changes between  
101 patient samples<sup>24</sup>. Most current approaches entail identification of cell populations with  
102 clustering methods like FlowSOM<sup>25</sup>, extracting sample-wise cell frequencies from each cluster  
103 and then comparing between samples to identify significantly differentially represented cell  
104 subsets<sup>26</sup>. Alternatively, some methods identify disease related changes at a single-cell level<sup>27</sup>.

105 Compared to the tightly regulated homeostasis of health, immunological diseases are inherently  
106 more variable<sup>28</sup>. Generally speaking, it follows that immune disease-related markers are more  
107 variably expressed in disease than health<sup>29,30</sup>. As we show, this fundamental biological insight  
108 is important because overlapping marker expression with unequal variance between patient  
109 classes necessarily implies a range of marker values with no discriminatory potential. This  
110 problem is exaggerated when marker distributions with unequal variance substantially overlap  
111 between two patient classes, such as health and disease. Critically, we often find that disease-  
112 related differences in immunological markers are small in relative and absolute terms<sup>31</sup>. This

113 inconvenient and unintuitive property, which is typical of markers measured by flow cytometry,  
114 masks informative markers in discovery studies and limits their clinical utility<sup>32</sup>.

115 In this report, we examine the problem of finding and interpreting disease markers with a  
116 restricted range of informative values from an immunologist's perspective. To do this, we must  
117 first disambiguate some key terms with different meanings for immunologists and computer  
118 scientists. Properties of single cells measured by flow cytometry, such as cell lineage-associated  
119 surface antigen expression, will be called "features." We reserve "marker" to mean a sample-  
120 related quantity, such as cell subset frequency, that signifies information relevant to sample  
121 classification, hence diagnoses. The distribution of marker values within a set of patient samples  
122 is described by its probability density function, or simply "density." Throughout this article, we  
123 present plots of densities that compare marker expression in patient subgroups; crucially, these  
124 should not be mistaken for histograms showing feature expression within samples.

125 We provide a computational method to optimally restrict markers to their informative range,  
126 which makes them easier to discover and interpret. The power of dataset restriction is  
127 demonstrated through its application to flow cytometry markers; in particular, T cell subset  
128 frequencies. For each marker, we calculate a restricted standardized AUC (rzAUC) for every  
129 marker value by splitting the sample set into marker<sup>HIGH</sup> and marker<sup>LOW</sup> parts. We define the  
130 optimal restriction according to the maximum absolute rzAUC of either the marker<sup>HIGH</sup> or  
131 marker<sup>LOW</sup> part. We then assign a permutation p-value to the optimal rzAUC. Finally, we  
132 leverage the adapted range of all restricted markers in a multivariate (random forest) model by  
133 forcing decision tree cuts within each informative range.

134 In essence, restriction identifies the informative range of a marker, which allows us to segregate  
135 datasets into classifiable and unclassifiable samples. Importantly, using information about the  
136 informative range of markers leads to superior multivariate models. We qualify our method



137 using realistically simulated flow cytometry data, then apply it to real T cell subset analyses to  
138 discover new markers of irAE risk in patients receiving immunotherapy for advanced  
139 melanoma. Using a restricted dataset, we were able to train and prospectively validate a  
140 multivariate model to predict immunotherapy-related hepatitis, which failed when using  
141 unrestricted data. Our computational methods can be directly applied to other types of data, not  
142 limited to transcriptomic, (epi-)genomic, proteomic, metabolomics or clinical information.

143 **RESULTS**144 **Two-class distributions resulting in skewed ROC curves**

145 We begin by showing how the distribution of a discriminatory marker that differs in its  
146 expression between diseased (patients) and unaffected (controls) individuals results in skewed  
147 receiver operating characteristic (ROC) curves. ROC curves relate the true positive rate (TPR)  
148 and false positive rate (FPR) for a disease marker at every datapoint in a two-class classification  
149 problem. Area under the ROC curve (AUC) is often used as a measure of the discriminatory  
150 capacity of a disease marker<sup>33</sup>. Throughout this report, we illustrate distributions of marker  
151 expression within classes by plotting probability densities. Densities are normalized to 1  
152 within each class, so the appearance of these plots is independent of class size (see  
153 Supplementary Note 1). In the following sections, we consider hypothetical markers whose  
154 expression is normally distributed  $\mathcal{N}(\mu, \sigma^2)$  with mean  $\mu$  and variance  $\sigma^2$ .

155

156 Perfectly discriminatory markers result in concave ROC curves with an AUC = 1 (Fig 1a). For  
157 imperfect markers, where there is overlap between the distributions of a disease marker  
158 expression in patient and control populations, provided that variance is equal in both classes,  
159 the ROC curve is symmetric about the anti-diagonal with  $1 > \text{AUC} > 0.5$ . In the hypothetical  
160 example, marker expression is normally distributed with equal variances in the patient  $\mathcal{N}(6, 1)$   
161 and control  $\mathcal{N}(5, 1)$  populations, but mean expression is higher in patients (Fig 1b). Entirely  
162 uninformative markers result in straight diagonal ROC curves with an AUC = 0.5 (Fig 1c).

163

164 Interpreting the area under a ROC curve is more complicated when comparing overlapping  
165 marker distributions with unequal variances, which result in ROC curves skewed around the  
166 anti-diagonal. Our first hypothetical example of a skewed ROC curve shows that normally

167 distributed, overlapping marker distributions with a higher mean and variance in the patient  
168 population compared to controls leads to a right-skewed ROC curve that crosses the diagonal  
169 in a region corresponding to low marker expression values (Fig 1d). Indeed, it is generally true  
170 that normally distributed populations with different variances result in non-concave ROC  
171 curves that cross the diagonal<sup>34</sup>. To illustrate this point, we simulated 200 samples by drawing  
172 random values from Normal distributions to show how varying the mean and variance of marker  
173 expression in patient and control distributions affects the shape and AUC of ROC curves  
174 (Supplementary Video 1 and 2). In the context of clinical diagnostics, markers of immune  
175 diseases usually reflect a change between tightly-regulated homeostasis in health and a  
176 disturbed, higher-variability condition in disease. Coupled with the fact that disease-associated  
177 changes in cell subset frequencies in blood are typically small, it is perhaps unsurprising disease  
178 markers measured by flow cytometry frequently result in skewed ROC curves<sup>13</sup>. In support of  
179 this assertion, we present a real-world example of a right-skewed ROC curve with a low AUC  
180 (Fig 1e). Specifically, this example shows that erythrocyte counts were elevated in baseline  
181 blood samples from patients metastatic melanoma who responded to combined Ipilimumab plus  
182 Nivolumab (Ipi-Nivo) therapy compared to non-responders.

183

184 Left-skewed ROC curves arise when the negative population has a lower mean, but higher  
185 variance than the positive population (Fig 1f). We find a real-world example in the previously  
186 unreported association between CD8<sup>+</sup>  $\gamma\delta$  T cells and hepatitis risk after combined Ipi-Nivo  
187 therapy (Fig 1g and Supplementary Fig 2). In this case, the higher variance of the control  
188 population might be due to technical imprecision in quantifying a rare cell population, since the  
189 absolute number of CD8<sup>+</sup>  $\gamma\delta$  T cells in blood was only  $25.6 \pm 19.3$  c/nl.

190

191 We next considered the case of a phenotypically heterogeneous positive population, which  
192 could reflect multiple aetiologies leading to a common clinical presentation, different stages of  
193 a disease that culminate in a common presentation or different treatment responses. In such  
194 scenarios, we expect a bimodal distribution of a disease marker in the positive population that  
195 leads to a skewed ROC curve (Fig 1h). We previously reported the identification of a subset of  
196 patients with advanced melanoma who developed hepatitis after Ipi-Nivo therapy, which was  
197 reliably predicted by CMV-associated expansion of CD4<sup>+</sup> T<sub>EM</sub> cells prior to immunotherapy<sup>35</sup>.  
198 In our melanoma dataset, we show that using CD4<sup>+</sup> T<sub>EM</sub> frequencies to predict hepatitis after  
199 immune checkpoint inhibitor (ICI) therapy leads to a right-skewed ROC curve (Fig 1i). We  
200 know from previous work that baseline CD4<sup>+</sup> T<sub>EM</sub> expansion is only a useful marker of hepatitis  
201 risk in CMV-infected patients, who constituted just 47.3% of our study cohort; therefore, this  
202 is a biologically validated example of alternative immunopathologies contributing to a common  
203 pathological presentation that impacts marker performance.

204

205 These three hypothetical distributions, and their real-world counterparts, demonstrate an  
206 important concept in immune marker discovery – namely, that a disease marker may be highly  
207 informative over a restricted range of measured values, but will consistently misclassify  
208 samples with marker values outside that range. By extension, using AUC across the entire ROC  
209 curve to assess predictive performance leads us to disregard potentially informative markers.  
210 Clearly, we need a method of finding such markers and defining their valid ranges.

211

### 212 **Dataset restriction is a new method to find disease markers**

213 Disease markers that give rise to skewed ROC curves perform well in a subset of samples,  
214 which may belong to either the positive or negative class, but are only informative over a certain

215 range. This leads us to the idea that particular *samples* may be classifiable or unclassifiable  
216 according to any given disease marker. Here, we present and implement a method of marker  
217 discovery that relies upon restricting training datasets to classifiable samples<sup>36</sup>. In the given  
218 example, we compared the distributions of 2500 positive and 2500 negative simulated samples,  
219 in which 20% of positive and 2% of negative samples were drawn from a normal distribution  
220  $\mathcal{N}(9, 1)$  and all other samples were drawn from  $\mathcal{N}(6, 1)$  (Fig 2a). This resulted in a right-  
221 skewed ROC curve for the complete dataset (Fig 2b). We first generated two ROC curves for  
222 every possible “restriction” of the dataset – explicitly, one for samples above the restriction  
223 ( $\text{marker}^{\text{HIGH}}$  samples, orange; Fig 2c, d-f) and one for samples beneath ( $\text{marker}^{\text{LOW}}$  samples,  
224 blue; Fig 2c, g-i).  $\text{marker}^{\text{HIGH}}$  samples generally correspond to the bottom-left part of the  
225 complete ROC curve (Fig 2d). Considering the densities of only  $\text{marker}^{\text{HIGH}}$  samples (Fig 2e),  
226 the restricted ROC curve had a superior “restricted” AUC (rAUC) of 0.692 (Fig 2f).  $\text{marker}^{\text{LOW}}$   
227 samples generally correspond to the top-right part of the complete ROC curve, but here their  
228 densities overlapped substantially; therefore, the restricted ROC curve was close to diagonal  
229 (Fig 2g-i). Notably, restricted densities are not the same as those in Fig 2a but are instead re-  
230 calculated on either  $\text{marker}^{\text{HIGH}}$  or  $\text{marker}^{\text{LOW}}$  samples.

231  
232 Standardizing each rAUC according to sample size gave the restricted standardized AUC  
233 (rzAUC). The maximum absolute value of rzAUC defined the optimal restriction value (Fig  
234 2c). In our example, rzAUC was maximal at FPR = 0.258, which corresponded to an optimal  
235 marker restriction value of 6.8. Consequently,  $\text{marker}^{\text{HIGH}}$  samples should be kept and  
236  $\text{marker}^{\text{LOW}}$  samples should be discarded – that is to say,  $\text{marker}^{\text{HIGH}}$  samples are classifiable,  
237 whereas  $\text{marker}^{\text{LOW}}$  samples are unclassifiable. Supplementary Video 3 helps to visualize the  
238 rAUC for varying restrictions of the dataset.

239 In other situations, the positive class may have higher or lower marker values, potentially  
240 leading to an  $AUC < 0.5$  and accordingly a negative  $rzAUC$ . In Supplementary Fig 3, we show  
241 that regardless which class is labelled positive or negative, our method identifies the same  
242 restriction value. In such cases,  $marker^{HIGH}$  and  $marker^{LOW}$   $rzAUC$ s are mirrored, meaning the  
243 restriction at the optimal absolute  $rzAUC$  remains identical. Critically, regardless of marker  
244 distribution, because areas under ROC curves are independent of class size, it follows that  
245 restriction values are also independent of class size<sup>37</sup>.

246

### 247 **Restriction identifies classifiable samples in simulated datasets**

248 To test our computational approach, we next applied it to our four preceding examples by  
249 simulating 100 samples from each class. In the first example, the negative class  $\mathcal{N}(5, 1)$  and  
250 positive class  $\mathcal{N}(6, 1)$  gave rise to a symmetrical ROC curve with a maximum  $rzAUC$   
251 corresponding to  $FPR=1$ ; consequently, the optimally informative dataset contained all samples  
252 (Fig 3a). In the second example, the negative class  $\mathcal{N}(5, 1)$  and positive class  $\mathcal{N}(6, 2)$  produced  
253 a right-skewed ROC curve because the variances were unequal (Fig 3b). We see that low marker  
254 values led to a consistent misclassification, indicated by the ROC curve crossing the diagonal.  
255 The maximum  $rzAUC$  of 5.8 for  $marker^{HIGH}$  samples indicated that samples with a marker value  
256  $< 4$  must be discarded. In the third example, the negative class  $\mathcal{N}(5, 2)$  and positive class  
257  $\mathcal{N}(6, 1)$  produced a left-skewed ROC curve (Fig 3c). Here, high marker values led to consistent  
258 misclassification; therefore, the ROC curve deviated below the diagonal. The maximum  $rzAUC$   
259 of 5.8 for  $marker^{LOW}$  samples indicated that samples with a marker value  $> 7$  must be discarded.  
260 In the fourth example, we compared 100 samples from the negative class  $\mathcal{N}(5, 1)$  and a bimodal  
261 positive class consisting of 90 samples from the same distribution  $\mathcal{N}(5, 1)$ , plus 10 samples  
262 from a distribution  $\mathcal{N}(9, 1)$  with a higher mean (Fig 3d). The resulting right-skewed ROC curve

263 reflected the fact that our simulated marker was only informative for higher sample values.  
264 Accordingly, the optimal rzAUC of 2.4 for marker<sup>HIGH</sup> samples restricted our dataset to samples  
265 with a marker value  $\geq 6.2$ . Evidently, our method is able to optimally restrict cleanly simulated  
266 patient populations, such that we retain only classifiable samples.

267

### 268 **Synthesizing realistic flow cytometry datasets**

269 Realistic synthetic data can be extremely valuable in machine-learning; for instance, for  
270 validating new analytical methods, calculating experimental sample sizes or data augmentation.  
271 Because no generative model already existed, we developed a new algorithm to create synthetic  
272 flow cytometry datasets, which differed from the preceding simulated examples in several key  
273 respects – specifically, they comprise multiple covarying markers, incorporate a realistic level  
274 of noise, and were adjusted in biologically meaningful ways. Our web-based interactive gating  
275 tree allows readers to synthesize their own flow cytometry data (Supplementary Website 1).

276

277 In order to validate our restriction method, we needed a way of imitating disease-related  
278 differences between groups of samples. In the method described above, any effect that changes  
279 the proportion of cells in any gate(s) equates to changing the Dirichlet distribution parameters.  
280 In the given example, the originally estimated mean proportions are projected onto the gating  
281 tree and corresponding Dirichlet distribution for three example leafs A, B and L. Here, the mean  
282 proportion of CD8<sup>+</sup> T<sub>EMRA</sub> cells is 7.17 % (Fig 4a). Now, instead of determining the number of  
283 cells in each leaf gate according to the originally estimated distribution, we generate synthetic  
284 cells from a modified Dirichlet distribution in which the mean proportion of CD8<sup>+</sup> T<sub>EMRA</sub> cells  
285 was arbitrarily changed to 33.23 % (Fig 4b). Using our method, changing the proportion of cells  
286 in any gate leads to changes in the proportion of cells in all other gates, which we represent by

287 the different intensities of red shading in the gating trees and the different Dirichlet distribution  
288 for the three example leafs A, B and L. Three examples of gating generated with a mean  
289 proportion  $CD8^+ T_{EMRA}$  cells = 33.23 % are provided (Supplementary Fig 4).

290

### 291 **Applying restriction to realistic synthesized flow cytometry datasets**

292 We next applied our restriction method to synthetic flow cytometry datasets that incorporated  
293 estimated technical and biological noise typical of real-world measurements. Specifically, we  
294 generated synthetic samples that gave rise to marker distributions similar to the preceding  
295 simulated examples. Artificial disease-associations were introduced by changing the frequency  
296 of  $CD4^+ T_{EM}$  cells, which had a baseline mean proportion of 7.7 % among healthy donors. We  
297 subsequently extracted  $CD4^+ T_{EM}$  cell frequencies relative to  $CD3^+$  T cells from all samples by  
298 applying our standard gating strategy, then applied our restriction method. Similar to Figure 3,  
299 we simulated marker values from normal distributions. We then generated synthetic flow  
300 cytometry datasets by setting the  $CD4^+ T_{EM}$  cell Dirichlet parameter to each simulated marker  
301 value.

302

303 In the first example, the negative class  $\mathcal{N}(7.7, 1)$  and positive class  $\mathcal{N}(10.7, 1)$  gave rise to a  
304 symmetrical ROC curve (Fig 5a). As expected, the results were much noisier than those shown  
305 in Figure 3; nevertheless, the maximum  $rzAUC = 7.2$  corresponded to  $FPR = 1$ , so the optimally  
306 informative dataset contained all samples. In the second example, the negative class  
307  $\mathcal{N}(7.7, 1)$  and positive class  $\mathcal{N}(8.7, 3)$  gave rise to a right-skewed ROC curve (Fig 5b). The  
308 maximum  $rzAUC = 4.3$  led us to retain marker<sup>HIGH</sup> samples with  $\geq 5.94\%$   $CD4^+ T_{EM}$  cells. In  
309 the third example, the negative class  $\mathcal{N}(7.7, 3)$  and positive class  $\mathcal{N}(8.7, 1)$  gave rise to a left-  
310 skewed ROC curve (Fig 5c). The maximum  $rzAUC = 4.6$  led to a restriction of the dataset to



311 marker<sup>LOW</sup> samples with  $< 8.37\%$  CD4<sup>+</sup> T<sub>EM</sub>. In the fourth example, we compared the negative  
312 class  $\mathcal{N}(7.7, 1)$  and a bimodal positive class comprising 80 samples showing no effect  
313  $\mathcal{N}(7.7, 1)$  plus 20 samples from a distribution  $\mathcal{N}(16.7, 1)$  with a higher mean (Fig 5d). The  
314 resulting right-skewed ROC curve with a maximum rzAUC = 3.7 led us to keep marker<sup>HIGH</sup>  
315 samples with  $\geq 8.24\%$  CD4<sup>+</sup> T<sub>EM</sub>. Hence, our method is able to appropriately restrict  
316 realistically synthesized flow cytometry datasets for symmetric or skewed ROC curves, such  
317 that we retain only classifiable samples.

318

### 319 **Restriction method improves findability in realistic synthesized datasets**

320 As explained above, introducing an artificial disease association into realistically synthesized  
321 flow cytometry datasets by adjusting the frequency of one cell population (in this case, CD4<sup>+</sup>  
322 T<sub>EM</sub> cells) leads to changes in all other nodes in our gating tree. We next asked whether our  
323 restriction method could also improve the discoverability of these covariant markers in the  
324 synthesized datasets presented above. Of note, rzAUC allows us to compare the discriminatory  
325 performance of different markers within one dataset; however, rzAUC values are not  
326 comparable between datasets, including independent training, validation and test datasets  
327 (Supplementary Fig 5). The AUC is equivalent to the Mann-Whitney U-statistic<sup>33</sup> and we can  
328 extend this equivalence to the rzAUC; however, this is not helpful in assigning significance  
329 values because optimizing for highest rzAUC introduces a bias. Instead, we must calculate  
330 permutation p-values<sup>38</sup>. For each of our four realistic synthesized examples, we calculated  
331 permutation p-values using the unrestricted sample set and the optimally restricted sample set  
332 for every gated cell population. Figure 6 shows these p-values as scatter plots in which the  
333 green-shading demarcates unrestricted p-values  $> 0.05$  and optimally restricted p-values  $< 0.05$   
334 – that is, markers identified as significant using our restriction method, but missed without it.

335 In our example of a symmetric ROC curve, we found that CD4<sup>+</sup> T<sub>EM</sub> cells and 3 subordinate  
336 populations were significant discriminators in both the unrestricted and restricted datasets (Fig  
337 6a). Two further populations were significant only in the restricted dataset. In the second  
338 example, which resulted in a right-skewed ROC curve, we found that two populations (i.e.  
339 CD4<sup>+</sup> T<sub>EM</sub> cells and CD27<sup>+</sup> CD28<sup>+</sup> CD57<sup>-</sup> CD4<sup>+</sup> T<sub>EM</sub> cells) had an optimal restriction  
340 permutation p-value = 0, whereas the corresponding unrestricted permutation p-value was >  
341 0.05 (Fig 6b). 8 other subsets were significant discriminators after restriction, but were not  
342 significant in the unrestricted sample-set. In the third example, which resulted in a left-skewed  
343 ROC curve, we found that CD4<sup>+</sup> T<sub>EM</sub> cells had an optimal restriction permutation p-value = 0,  
344 but were not significant in the unrestricted dataset (Fig 6c). 11 other subsets were significant  
345 discriminators after restriction, but not in the unrestricted sample-set. In the fourth example,  
346 CD4<sup>+</sup> T<sub>EM</sub> cells had an optimal restriction permutation p-value = 0.007, but were not significant  
347 in the unrestricted dataset (Fig 6d). Hence, dataset restriction enables discovery of disease  
348 markers which would otherwise be disregarded in synthesized flow cytometry datasets.

349

### 350 **Dataset restriction discovers valid irAE markers**

351 Having qualified our restriction method using synthesized datasets, we next applied it to real  
352 clinical data. In previous work, we investigated pre-treatment peripheral blood samples from  
353 110 patients with advanced melanoma who received Ipi-Nivo therapy<sup>13</sup>. Using conventional  
354 methods, we found no significant marker after correcting for multiple comparison. Here, we  
355 asked whether our restriction method could reveal any novel markers of hepatitis or colitis risk  
356 in the same dataset. No markers of colitis survived correction for multiple comparison (Fig 6e  
357 and Supplementary Fig 6). However, in predicting hepatitis, our restriction method returned 7  
358 significant markers with an unrestricted permutation p-value > 0.05 (Fig 6f). After correction

359 for multiple testing, 4 of these 7 hepatitis markers remained significant with an FDR < 0.05. By  
360 contrast, no marker identified from the unrestricted dataset returned a significant permutation  
361 p-value after correction for multiple testing. Thus, our restriction method returned significant  
362 disease-associated markers, which were not found using the unrestricted dataset.

363

364 Using our novel restriction method, we identified CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cell frequency  
365 relative to CD4<sup>+</sup> in blood as a marker of hepatitis risk after dataset restriction. To illustrate the  
366 potential utility of restricted markers, we compared the performance of CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup>  
367 T<sub>EM</sub> frequency as a marker of hepatitis risk in our unrestricted and restricted datasets (Fig 7).  
368 The discriminatory cut-off for patient classification, defined by the Youden index, was the same  
369 for both the restricted and unrestricted datasets, such that samples with more than 9.56% of  
370 CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> relative to CD4<sup>+</sup> are predicted hepatitis positive. Accordingly, using  
371 the unrestricted dataset, CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub>(%) correctly predicted incidence of hepatitis  
372 in 74 of 110 patients. The unrestricted cell frequency had a sensitivity (TPR) of 45.8% and a  
373 specificity (true negative rate, TNR) of 83.9%. The positive predictive value (PPV) was 68.8%  
374 and the negative predictive value (NPV) was 66.7%. Our restriction method usually implies  
375 that some cases should be considered unclassifiable based upon its marker values. In this  
376 example, 58 of 110 patients were deemed unclassifiable. Incidence of hepatitis was correctly  
377 predicted in 40 of 52 classifiable samples. The restricted cell frequency had a sensitivity of  
378 91.7% and a specificity of 64.3%. The positive predictive value was 68.8% and the negative  
379 predictive value was 90% .

380 **Multivariate analysis of restricted data predicts hepatitis**

381 Although our restriction method leads to discarding many samples as unclassifiable by any  
382 particular marker, we found that different markers define noncongruent sets of classifiable  
383 samples (Fig 8a). This led us to investigate whether using restricted datasets could improve the  
384 predictive performance of multivariate models. First, we built a random forest model<sup>39</sup> using  
385 all 84 markers from the unrestricted training dataset of 110 patients (Fig 8b). When this model  
386 was applied to a fully independent, prospective validation set of 30 patients, the resulting  
387 predictions were inaccurate (CCR=56.7% vs. 54.8% under the no-information model).

388

389 By contrast, we observed a significant improvement in predictive performance using the  
390 restricted dataset to train our random forest. To leverage information from our restriction  
391 method, we assigned a value of -1 to restricted samples across all 84 markers. When this model  
392 was applied to the validation set, the resulting predictions were significant (Fisher's Exact p-  
393 value = 0.026) and had a correct classification rate of 73.3%. 12 of 16 predictions of hepatitis  
394 were correct (PPV=75%) and 10 of 14 negative predictions were correct (NPV=71.4%). Hence,  
395 in principle, dataset restriction can improve the training and performance of multivariate  
396 predictive models.

397 **DISCUSSION**

398 Immunological diseases are often heterogeneous in clinical presentation and severity, reflecting  
399 the variability of their underlying immunopathologies<sup>15</sup>. It follows, we argue, that immune  
400 disease-associated markers typically exhibit greater variance among diseased patients than  
401 unaffected individuals. This general proposition was broadly corroborated by our real-world  
402 examples of patient groups who were prone to immunotherapy-related complications. Unequal  
403 variance in marker distribution between patient classes affects our ability to identify markers  
404 with discriminatory capacity over a restricted range of marker values. To solve this biological  
405 problem, we introduced dataset restriction as a marker discovery tool. In artificial and real-  
406 world examples, dataset restriction enabled us to find discriminatory markers that were  
407 disregarded by conventional measures of marker performance. Moreover, we showed that  
408 dataset restriction improves performance of multivariate predictive models. Our work  
409 formalizes a new way of evaluating diagnostic results – namely, that certain markers can only  
410 be usefully interpreted over a restricted range of values, and that samples with values outside  
411 this range should be considered as unclassifiable.

412  
413 Flow cytometry is a powerful method for interrogating the phenotype of many single cells  
414 within a heterogeneous mixture. This technique allowed us to estimate the relative numbers of  
415 accurately defined leucocyte subsets in peripheral blood samples, including T cell subsets,  
416 which are direct targets of Ipilimumab (anti-CTLA-4) and Nivolumab (anti-PD-1) therapy<sup>40</sup>.  
417 Although flow cytometry generates rich and immunologically interpretable data, it has two key  
418 limitations – namely, that blood leucocyte frequencies vary within a narrow dynamic range,  
419 and that higher order cell feature combinations may define rare cell subsets<sup>41,42</sup>. Small disease-  
420 related changes in markers are problematic because substantially overlapping marker

421 distributions with unequal variance lead to exaggerated skewness of ROC curves. Rare cell  
422 subsets are problematic because our estimates of their frequency are less reliable<sup>43</sup>. Crucially,  
423 dataset restriction helps to overcome the special difficulties of correctly interpreting flow  
424 cytometry data by limiting marker values to a range in which the signal-to-noise ratio is  
425 increased relative to the full range. Consequently, we reduce the likelihood of false positive or  
426 false negative classification at the cost of discarding some samples as unclassifiable.

427

428 We created an R-package, called *restrictedROC*<sup>36</sup>, that calculates restricted standardized AUC  
429 scores. The *rzAUC* is returned together with a value that delimits the marker's informative  
430 range. This builds upon earlier ideas about partial AUCs, which were introduced to account for  
431 imposed restrictions that capped true and false positive rates<sup>44-46</sup>. Imposed restrictions usually  
432 come from domain knowledge; for instance, tests with a high false positive rate are  
433 inappropriate for expensive diagnostic screening applications, whereas tests with a high false  
434 negative rate are inappropriate when a life-saving treatment is available<sup>47</sup>. McClish introduced  
435 a "standardization" for partial AUCs for a given range of false positive rates, such that a  
436 randomly selected positive sample has a higher value than a randomly selected negative sample  
437 conditional upon the negative sample arising from the false positive range<sup>48</sup>. In our method, we  
438 introduced a scaling factor for the two-way partial AUC<sup>47</sup> resulting in the restricted AUC  
439 (*rAUC*). With this scaling factor, the *rAUC* becomes the probability that a randomly selected  
440 positive sample has a higher value than a randomly selected negative sample conditional upon  
441 *both samples* arising from a range spanned by a minimum true positive rate and a maximal false  
442 positive rate. The restricted standardized AUC (*rzAUC*) then takes into account both the *rAUC*  
443 and the number of samples in the marker<sup>HIGH</sup> or marker<sup>LOW</sup> range leveraging the equivalence  
444 between AUC and Mann-Whitney U test<sup>33</sup>.

445 We further developed our method to determine the optimal range of marker values that correctly  
446 classifies samples. Specifically, we optimize a restriction that either includes samples with  
447 higher marker values ( $\text{marker}^{\text{HIGH}}$ ) or lower marker values ( $\text{marker}^{\text{LOW}}$ ) and has the highest  
448 possible absolute rzAUC. The rzAUC can be directly compared within one dataset, but depends  
449 on the total number of samples. By calculating permutation p-values<sup>38</sup> for the rzAUC, we  
450 remove this dependence and attribute significance values.

451  
452 There are alternative ways of describing the geometric symmetry of ROC curves apart from  
453 graphical skewness. Left-skewed ROC curves are also described as True Negative Proportion  
454 (TNP)-asymmetric and right-skewed ROC curves as True Positive Proportion (TPP)-  
455 asymmetric. These asymmetries can be defined by Kullback-Leibler divergences<sup>49</sup> (KL-  
456 divergences). Therefore, KL-divergence could be used to assess whether restriction should be  
457 applied to a given marker; however, in the case of symmetric ROC curves, our restriction keeps  
458 all samples, so such preselection of markers is unnecessary. Importantly, excluding samples to  
459 minimize KL-divergence is not the equivalent of dataset restriction.

460  
461 In principle, dataset restriction can be applied to optimize any marker range. However,  
462 following from our immunological rationale, restricting the upper or lower range is especially  
463 applicable in clinical diagnostics. For completeness of our discussion, we can imagine a marker  
464 with both uninformative  $\text{marker}^{\text{HIGH}}$  and  $\text{marker}^{\text{LOW}}$  values (ie. where only mid-range values  
465 are informative) that might only be discovered by applying our restriction method twice in  
466 succession. In theory, restriction could be iteratively applied to a dataset until no more samples  
467 are identified as unclassifiable, but the practical value of multiply restricting datasets is unclear.  
468

469 In order to validate our restriction method, we developed a novel method for synthesizing  
470 realistic flow cytometry data with class-related effects. Because no generative method  
471 previously existed, our approach represents a significant contribution to cytometry analysis,  
472 particularly for benchmarking of diagnostic flow cytometry algorithms, sample size  
473 calculations or data augmentation. Our method uses an expert-given hierarchical gating  
474 strategy, where the proportions of cells per gate are described with a Dirichlet distribution.  
475 Within each terminal (leaf) gate, the cells are described using a normal distribution. Thus, we  
476 effectively created a Gaussian mixture distribution with the number of components defined by  
477 the number of terminal gates. In cytometry, (Gaussian) mixture models are an established  
478 method for unsupervised cell population identification<sup>50,51</sup>. In principle, these earlier  
479 approaches could be used to generate new cells from estimated distributions, although their  
480 focus was labeling existing cells rather than creating artificial ones. Our use of a hierarchical  
481 gating strategy and a Gaussian mixture model allows for the creation of complex data  
482 distributions. In future, for certain applications, further adaptations of our approach, such as  
483 multivariate skew t-distributions, could be used to improve the accuracy of simulated data<sup>52</sup>.

484

485 Restricting markers to an informative range of values is important because it improves  
486 classification performance. We emphasize that classification cut-offs and restriction values are  
487 different concepts. Classification cutoffs, such as the Youden index<sup>53</sup>, divide a sample set into  
488 predicted positive and predicted negative classes. By contrast, restriction divides a sample set  
489 into classifiable and unclassifiable samples. In the context of individualized patient care, it  
490 might seem unproductive to label samples as unclassifiable. On the contrary, we argue that the  
491 clinical utility of a predictive marker improves if its certainty is high, even if it only works in  
492 a small subset of patients. Consider a disease-related marker giving a right-skewed ROC curve:



493 Conventional approaches return a reliable positive classification and an unreliable negative  
494 classification; in contrast, our restriction method returns a reliable positive classification, a  
495 reliable negative classification and a set of unclassifiable samples, which do not necessarily  
496 have the most negative values. Of note, the discriminatory cut-off determined by the Youden  
497 index is often the same after restriction but may change in some cases. When interpreting a  
498 single marker, our restriction method improves either the positive or the negative predictive  
499 value, so improves certainty of our predictions.

500

501 Our method may concern some clinicians, who will legitimately ask about unclassifiable  
502 patients<sup>54</sup>. Here, we provide an answer by building an informative and prospectively validated  
503 random forest model after replacing all restricted values with a constant outside the informative  
504 range. Consequently, we force each tree of the random forest to select a cut-off within the  
505 informative range or a cut-off between the classifiable and unclassifiable regions. More  
506 sophisticated methods may be developed in future, but our experimentally validated random  
507 forest is a proof-of-principle that differently restricted markers can be usefully combined in  
508 multivariate models.

509

510 To demonstrate the potential clinical utility of dataset restriction, we applied our method to the  
511 clinically significant problem of immune-related adverse events following combined  
512 immunotherapy. In univariate analyses, dataset restriction identified new markers associated  
513 with ICI-related hepatitis, including CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cells, that were not returned by  
514 conventional methods. Of clinical importance, dataset restriction increased NPV without  
515 compromising PPV. Combining many restricted markers into a random forest model generated  
516 an informative model, whereas training on unrestricted data from the same set of 110 samples

517 returned no valid models. To validate our predictive model, we assessed its performance in an  
518 independent, prospectively collected set of 30 samples, where it returned significant  
519 predictions, which were superior to the performance of any single marker alone. Beyond the  
520 scope of this article, such multivariate models could be extended to include markers from  
521 multiple flow cytometry panels or other patient-related information, such as age, sex or clinical  
522 chemistry results.

523

524 Clinical manifestations of immune disease are often heterogeneous. This is certainly true of  
525 irAE after immunotherapy, which vary greatly in severity, time-of-onset, clinical features and  
526 response to treatment<sup>28</sup>. Further, there is increasing evidence that multiple immune aetiologies  
527 lead to common clinical presentations, such as colitis<sup>55</sup>, myositis<sup>56</sup> or hepatitis<sup>57</sup>. This  
528 heterogeneity connotes individual genetic predisposition<sup>58,59</sup>, environmental factors<sup>60,61</sup> and  
529 past immunological challenges<sup>35</sup>. In particular, we now recognize the contribution of previous  
530 viral infections in preconditioning towards adverse reactions. An unanticipated consequence of  
531 dataset restriction is that disease markers with a bimodal distribution in the positive class, such  
532 as might arise from multiple aetiologies, are findable. Excitingly, combining features from a  
533 restricted dataset into multivariate models should, in principle, enable predictions about  
534 diseases with multiple aetiotypes – a situation where conventional methods are unsuitable.  
535 Extending this idea of dataset restriction as a way of classifying samples with intraclass  
536 heterogeneity to unsupervised methods, such as PCA or clustering, could aid discovery of  
537 previously unknown patient subsets.

538

539 In summary, clinical markers that can only be interpreted over a restricted range are inherently  
540 likely in immune diseases. Where classical methods fail, dataset restriction solves the problem

541 of discovering and interpreting such markers. Our approach is not limited to prospective data,  
542 but can also be used retrospectively to find new markers or improve existing ones. Dataset  
543 restriction was developed here to analyze flow cytometry data; however, it is directly applicable  
544 to any sample classification problem. In immunological research, this could include  
545 transcriptomic, proteomic, (epi-)genomic, metabolomic or imaging data. We hope others will  
546 apply our method to existing datasets, perhaps leading to valuable new markers or novel  
547 biological insights.

548 **MATERIALS AND METHODS**

549 **Collection of clinical information**

550 Data from three sources were used in this study: (1) a training set (n=48) from a cohort of  
551 healthy humans used to develop our flow cytometry data simulations; (2) a previously reported  
552 training set (n=110) from patients with advanced melanoma used for marker discovery<sup>13</sup>; and  
553 (3) a new prospective validation set (n=30) from patients with advanced melanoma. Whole  
554 blood was collected from healthy thrombocyte donors with approval from the Ethics Committee  
555 of the University of Regensburg (approval 22-2780-01). All donors gave full, written consent  
556 to sample and data collection. Clinical samples for the marker training and validation sets were  
557 collected within a single-center, non-interventional study<sup>62</sup>, which was conducted in accordance  
558 with the Declaration of Helsinki and all applicable German and European laws and ethical  
559 standards. This observational study was authorized by the Ethics Committee of the University  
560 of Regensburg (approval 16-101-0125) and registered with clinicaltrials.gov (NCT04158544).  
561 Blood samples were obtained from patients with Stage III/IV melanoma under the care of the  
562 Department of Dermatology at University Hospital Regensburg (UKR). Eligible patients were  
563 consecutively recruited without stratification or matching. All participants gave full, informed  
564 written consent. For the training set, the first reported case was recruited in OCT-2016 and the  
565 last reported case was recruited in JUN-2021. For the prospective validation set, the first  
566 reported case was recruited in JUN-2021 and the last reported case was recruited in JAN-2023  
567 (Supplementary Table 1). All study participants received standard-of-care treatment according  
568 to local guidelines. Specifically, patients with unresectable metastatic disease who received  
569 first- or second-line checkpoint inhibitor therapy were initially treated with Nivolumab ( $\alpha$ PD-  
570 1; 1 mg/kg; Bristol-Myers Squibb) plus Ipilimumab ( $\alpha$ CTLA-4; 3 mg/kg; Bristol-Myers

571 Squibb) for up to four cycles at 3 week intervals. Thereafter, patients received 480 mg  
572 Nivolumab monotherapy at 4 week intervals.

573

574 **Diagnosis and grading of clinical outcomes.**

575 All irAE were evaluated by an expert Dermatological Oncologist. ICI-related hepatitis was  
576 diagnosed when: (i) GOT, GPT,  $\gamma$ -GT or total bilirubin substantially deviated from pretreatment  
577 values; (ii) this change was not attributable to other causes, such as co-medication or viral  
578 disease; and (iii) liver injury was sufficiently severe that ICI therapy was suspended or stopped,  
579 or immunosuppression was started. Colitis was diagnosed when increased stool frequency or  
580 loose consistency, accompanied by abdominal discomfort led to suspension or cessation of ICI  
581 therapy and introduction of immunosuppressive treatment. Clinical responses were assessed  
582 using the Response Evaluation Criteria in Solid Tumors (RECIST 1.1)<sup>63</sup>. Patients with  
583 progressive disease were categorized as non-responders, whereas those with complete or partial  
584 responses, and those with stable disease, were categorized as responders.

585

586 **Flow cytometry.**

587 Step-by-step protocols for preparing and analyzing clinical samples by flow cytometry can be  
588 accessed through Nature Protocol Exchange<sup>64</sup>. Briefly, blood was collected into EDTA-  
589 vacutainers by peripheral venepuncture then delivered to the responsible lab at ambient  
590 temperature. Samples were stored at 4°C for up to 4 h before processing. Whole blood samples  
591 were stained using the DURAClone IM T Cell Subsets Tube (Beckman Coulter, B53328). Data  
592 were collected using a Navios<sup>TM</sup> cytometer running Cytometry List Mode Data Acquisition and  
593 Analysis Software version 1.3 (Beckman Coulter). An experienced operator performed blinded

594 analyses following a conventional workflow that entailed sample-wise recompensation, arcsinh  
595 transformation and rescaling before applying a uniform gating strategy (Supplementary Fig 4).

596

### 597 **Restriction method**

598 We propose a method for finding markers with high performance in subsets of samples that  
599 involves: 1) “restricting” samples into marker<sup>HIGH</sup> and marker<sup>LOW</sup> sets for every unique marker  
600 value; 2) calculating the corresponding restricted receiver operating characteristic(ROC) curve;  
601 3) calculating the area under the restricted ROC curve; 4) adjusting the restricted AUC (rAUC)  
602 for sample size; 5) selecting the optimal restriction level, 6) calculating permutation p-values;  
603 and 7) reporting performance and significance. This algorithm is implemented as an R package  
604 called *restrictedROC*<sup>36</sup>.

605

606 *To define our nomenclature, we first introduce ROC curve analysis.* Let a cut-off  $c \in \mathbb{R}$ , a  
607 continuous marker  $Y \in \mathbb{R}$  and a grouping of samples into *diseased* (positive,  $D = 1$ ) and *non-*  
608 *diseased* (negative,  $D = 0$ ). A sample can be classified as diseased if  $Y \geq c$  and into non-  
609 diseased if  $Y < c$ . The true positive rate (*TPR*) and false positive rate (*FPR*) at cut-off  $c$  are  
610 defined as  $TPR(c) = P[Y \geq c | D = 1] = P[Y_D \geq c]$  and  $FPR(c) = P[Y \geq c | D = 0] =$   
611  $P[Y_{\bar{D}} \geq c]$ . The ROC curve relates the *TPR* and *FPR* for all possible cut-offs  $c$ , including  
612  $\{\infty, -\infty\}$  (*nb.* compare with Supplementary Fig 7). We can write the value of the ROC curve  
613 at any false positive rate  $t \in (0,1)$  as  $ROC(t) = TPR(FPR^{-1}(t))$ . For empirical survival  
614 functions  $S_D$  and  $S_{\bar{D}}$  we write  $ROC(t) = S_D(S_{\bar{D}}^{-1}(t))$  by substituting *TPR* and *FPR*. The area  
615 under the ROC curve (AUC) is then defined as

$$AUC = \int_0^1 ROC(t) dt . \quad (1)$$

616

617 Consequently, a perfectly discriminating marker with higher values corresponding to the  
 618 positive class translates to a perfect ROC curve with  $AUC = 1$ . An uninformative marker has  
 619 an AUC of 0.5, corresponding to  $ROC(t) = t \forall t \in (0, 1)$ . A perfectly discriminating marker  
 620 but with higher values corresponding to the negative class has an AUC of 0. From a  
 621 probabilistic point of view, the AUC equals the probability that the marker value of a random  
 622 positive sample will be higher than that of a random negative sample:  $AUC = P[Y_D > Y_{\bar{D}}]$   
 623 <sup>33,65,66</sup>. The derivation is given in Supplementary Note 2.

624  
 625 *Next, we introduce the concept of restricted ROC curves.* Our “restriction” is a marker value  
 626 that splits the samples into marker<sup>HIGH</sup> and marker<sup>LOW</sup> sets. For both sets, we separately  
 627 calculate “restricted” ROC curves and their corresponding restricted AUC (rAUC). See  
 628 supplement for the full derivation. In Supplementary Note 3, we prove that calculating rAUC  
 629 is identical to scaling a partial AUC (pAUC). Therefore, before we describe our computational  
 630 method, we consider the (two-way) pAUC<sup>47,67</sup>. The partial AUC (pAUC) is defined as the AUC  
 631 up to a certain false positive rate. Its probabilistic correspondence has been shown<sup>45,66</sup>:

$$pAUC(t_0) = \int_0^{t_0} ROC(t) dt = P[Y_D > Y_{\bar{D}} | Y_{\bar{D}} > S_{\bar{D}}^{-1}(t_0)] \cdot t_0 \quad (2)$$

632 The pAUC was recently extended to two-way partial AUCs<sup>47</sup>. Here, the area is calculated  
 633 between  $S_D(t) \geq 1 - \alpha$  and  $S_{\bar{D}}(t) \leq \beta$ . This area, shown in Supplementary Fig 8 as shaded  
 634 area **A**, can be written as

$$AUC_{\alpha}^{\beta} = \int_{S_{\bar{D}}(S_{\bar{D}}^{-1}(1-\alpha))}^{\beta} ROC(t) dt - (1 - \alpha) \left( \beta - S_{\bar{D}}(S_D^{-1}(1 - \alpha)) \right) \quad (3)$$

$$= P[Y_D > Y_{\bar{D}}, Y_D \leq S_D^{-1}(1 - \alpha), Y_{\bar{D}} \geq S_{\bar{D}}^{-1}(\beta)] \quad (4)$$

635 Our restriction method uses two special cases of  $AUC_{\alpha}^{\beta}$ , shown in Supplementary Fig 9:

636 1) The left part of the area under the curve up to a false positive rate  $\beta$ , which is identical  
 637 to the pAUC described earlier

$$AUC_{high}(\beta) = AUC_{\alpha=1}^{\beta} = \tag{5}$$

$$= \int_0^{\beta} ROC(t) dt \tag{6}$$

$$= P[Y_D > Y_{\bar{D}}, Y_{\bar{D}} > S_{\bar{D}}^{-1}(\beta)] \tag{7}$$

638

639 2) The right part of the area under the curve with at least a true positive rate of  $1 - \alpha$

$$AUC_{low}(\beta) = AUC_{\alpha}^{\beta=1} = \tag{8}$$

$$= \int_{S_{\bar{D}}(S_{\bar{D}}^{-1}(1-\alpha))}^1 ROC(t) dt (1 - \alpha) \left(1 - S_{\bar{D}}(S_{\bar{D}}^{-1}(1 - \alpha))\right) \tag{9}$$

$$= P[Y_D > Y_{\bar{D}}, Y_D \leq S_D^{-1}(1 - \alpha)] \tag{10}$$

640 Partial AUCs consider only a specific part of the original ROC curve, therefore the  
 641 interpretation of perfect ( $AUC = 1$ ) or non-informative ( $AUC = 0.5$ ) becomes invalid. For  
 642 pAUC, the following standardization was proposed to restore this interpretation<sup>48</sup>

$$standardized \text{ pAUC} = \frac{1}{2} \left(1 + \frac{\text{pAUC} - \min}{\max - \min}\right) \tag{11}$$

643 where, min is the pAUC given an non-informative marker ( $\min = \frac{\beta^2}{2}$ ), and max is the pAUC  
 644 given a perfect marker ( $\max = \beta$ ) up to an false positive rate of  $\beta$ .

645 In contrast, our restriction method applies the following two scaling factors to any two-way  
 646 partial  $AUC_{\alpha}^{\beta}$

$$rAUC_{\alpha}^{\beta} := AUC_{\alpha}^{\beta} \cdot \frac{1}{\beta - S_{\bar{D}}(S_{\bar{D}}^{-1}(1 - \alpha))} \cdot \frac{1}{S_D(S_{\bar{D}}^{-1}(\beta)) - (1 - \alpha)} \tag{12}$$

647 Effectively, these two scaling factors rescale the area spanned through  $\alpha$  and  $\beta$  to 1.  
 648 Importantly, this is equivalent to calculating rAUC considering only samples with  $S_{\bar{D}}^{-1}(\beta) < t$   
 649  $< S_{\bar{D}}^{-1}(1 - \alpha)$ . This has a probabilistic interpretation of



$$\text{rAUC}_\alpha^\beta = P[Y_D > Y_{\bar{D}} \mid S_{\bar{D}}^{-1}(\beta) \leq Y \leq S_{\bar{D}}^{-1}(1 - \alpha)] \quad (13)$$

650 Here, the  $\text{rAUC}_\alpha^\beta$  is defined in terms of maximum false positive rate  $1 - \alpha$  and minimum true  
 651 positive rate  $\beta$ . Alternatively, we introduce a “restriction”  $r \in \mathbb{R}$  which splits the data into  
 652 marker<sup>HIGH</sup> and marker<sup>LOW</sup> sets where  $\alpha := 1 - S_D(r)$  and  $\beta := S_{\bar{D}}(r)$ . With this, our two  
 653 special cases become

$$\text{rAUC}_{high}(r) = \text{rAUC}_{\alpha=1}^{\beta=S_{\bar{D}}(r)} = \text{AUC}_{high}(S_{\bar{D}}(r)) \cdot \frac{1}{S_{\bar{D}}(r)} \cdot \frac{1}{S_D(r)} \quad (14)$$

$$\text{rAUC}_{low}(r) = \text{rAUC}_{\alpha=1-S_D(r)}^{\beta=1} = \text{AUC}_{low}(1 - S_D(r)) \cdot \frac{1}{1-S_{\bar{D}}(r)} \cdot \frac{1}{1-S_D(r)} \quad (15)$$

654 This is equivalent to keeping marker<sup>HIGH</sup> samples with values  $> r$  ( $\text{rAUC}_{high}$ ) or to keeping  
 655 marker<sup>LOW</sup> samples with values  $\leq r$  ( $\text{rAUC}_{low}$ ), then calculating AUC on the restricted dataset.  
 656 Supplementary Video 3 uses a hypothetical dataset to visualize the rAUC and show the visual  
 657 equivalence of our scaling factor and when restricting the dataset.

658 More extreme restrictions result in fewer samples, so our estimates of  $\text{rAUC}(r)$  become  
 659 increasingly unreliable; therefore, we adjust  $\text{rAUC}(r)$  for sample size after restriction. Here, we  
 660 leverage the equality of the AUC to the Mann-Whitney U test<sup>33</sup> in order to calculate the  
 661 restricted standardized AUC ( $\text{rzAUC}_X$ ) for  $X$  either marker<sup>HIGH</sup> and marker<sup>LOW</sup> sets by  
 662 calculating the test statistic

$$\text{rzAUC}_X(r) = \frac{\text{rAUC}_X(r) - 0.5}{\sqrt{\text{var}_{H_0}(\text{rAUC}_X(r))}} \quad (16)$$

663 where  $\text{var}_{H_0}(\text{rAUC}_X(r))$  is the variance under the null hypothesis  $H_0$  that positive and negative  
 664 samples are independent and identically distributed. This demands no assumption of normality.

665 Then this  $\text{var}_{H_0}$  is<sup>68,69</sup>

$$\text{var}_{H_0}(\text{rAUC}_X(r)) = \frac{m + n + 1}{12mn} \quad (17)$$

666 where  $m$  is the number of positive samples and  $n$  is the number of negative samples with marker  
667 values higher ( $\text{rAUC}_{high}$ ) or lower or equal ( $\text{rAUC}_{low}$ ) than the restriction  $r$ . With this  
668 adjustment, a higher number of samples reduces variance, hence  $\text{rAUC}_x$  becomes more  
669 reliable. For a visual example, see Supplementary Fig 10 where the  $\text{rAUC}$  and  $\text{rAUC}$  are shown  
670 for all possible restrictions in terms of the false positive rate. The  $\text{rAUC}_x$  can be negative if  
671 the corresponding  $\text{rAUC}_x$  is below 0.5, decreases with fewer samples and increases in absolute  
672 value the further  $\text{rAUC}_x$  is from 0.5.

673

674 *After calculating the  $\text{rAUC}$ , we next identify the optimal restriction, which is defined as the*  
675 *highest absolute value of  $\text{rAUC}_{high}$  or  $\text{rAUC}_{low}$ . Including more samples would result in a*  
676 *smaller  $\text{rAUC}_x$  and therefore smaller  $\text{rAUC}_x$ . Excluding more samples would result in an equal*  
677 *or higher  $\text{rAUC}_x$  but also a higher variance and therefore also a smaller  $\text{rAUC}_x$ . With this*  
678 *restriction we include some and potentially, but not necessarily, exclude other samples in the*  
679 *calculation of the  $\text{rAUC}$ . We describe the excluded samples as “unclassifiable” and remove*  
680 *them from further calculation of usual performance measures like accuracy, specificity or*  
681 *sensitivity.*

682

683 *Finally, we calculate permutation p-values for the unrestricted AUC and  $\text{rAUC}$ . After*  
684 *obtaining the unrestricted AUC for an unrestricted dataset or the  $\text{rAUC}_x$  for an optimized*  
685 *subset of samples, we need to assign a p-value using permutation tests. This a non-parametric*  
686 *way to determine statistical significance based upon a null hypothesis that class labels assigned*  
687 *to samples are interchangeable<sup>38</sup>. Following this approach, we first calculate unrestricted AUC,*  
688  *$\text{rAUC}_{high}$  and  $\text{rAUC}_{low}$  using the correct labels. Then we permute the labels 10,000 times*  
689 *before recalculating unrestricted AUC,  $\text{rAUC}_{high}$  and  $\text{rAUC}_{low}$ . The unrestricted permutation*

690 p-value is then the number of times the permuted unrestricted AUC is above the original  
691 unrestricted AUC. Likewise, the restricted permutation p-value is then the number of times either  
692  $\text{rzAUC}_{high}$  or  $\text{rzAUC}_{low}$  are absolutely higher than the optimal  $\text{rzAUC}_X$ .

693

#### 694 **Multivariate restriction analysis**

695 Our restriction method identifies only a part of the samples as classifiable and cannot make  
696 predictions for the unclassifiables. This potentially excludes many samples, so constrains  
697 predictive power. To circumvent this problem, we replace the marker values of unclassifiable  
698 samples with a clearly distinct value (-1) and then apply a random forest. With this substitution,  
699 we can predict all given samples, regardless if they are unclassifiable by some markers. In our  
700 melanoma dataset, per sample we first downsampled 10,000 CD3<sup>+</sup> T cells. We then restricted  
701 our set of features to 84 gates where at least 10% of 110 training samples contained more than  
702 10 counts. Then we calculated relative proportion to either CD4<sup>+</sup> CD8<sup>-</sup> or CD4<sup>-</sup> CD8<sup>+</sup> T cells.  
703 We also used CD4<sup>+</sup>CD8<sup>+</sup> (double positive), CD4<sup>-</sup> CD8<sup>-</sup> (double negative), CD4<sup>+</sup> CD8<sup>-</sup> and CD4<sup>-</sup>  
704 CD8<sup>+</sup> T cell counts, which are relative to the fixed parent gate of 10,000 CD3<sup>+</sup> T cells.

705

706 For our unrestricted, classical multivariate approach, we used the proportions and counts of all  
707 110 previously published training samples. We then trained a random forest<sup>39</sup> model using the  
708 H2O R library<sup>70</sup> with 1000 trees and the number of bins of 100, a random manual seed for  
709 reproducibility of the results the remaining default parameters. Explicitly, a maximum depth of  
710 20, a minimum number of samples in a node of 1, logloss stopping metric, the number of  
711 randomly sampled candidate markers as floor of the square root of 84 (9), a sample rate of  
712 0.632, minimum split improvement of  $10^{-5}$  and an automatic histogram type. Finally, we applied  
713 the random forest on a prospective cohort of n=30 patients.

714

715 For our restricted multivariate approach, we performed a marker-wise restriction to samples,  
716 then replaced all unclassifiable marker values with -1. (We chose this value because all  
717 classifiable values are strictly positive as they represent either proportions of CD4<sup>+</sup> or CD8<sup>+</sup> T  
718 cells, or absolute T cell counts.) This substitution forces each tree in the random forest to select  
719 discriminatory cutoffs within the range of informative marker values. We then trained a random  
720 forest model with the same settings as for the unrestricted multivariate approach. We finally  
721 applied the restriction values obtained from the training set to the prospective validation set,  
722 replaced the unclassifiable marker values with -1 and applied the random forest on the  
723 prospective cohort.

724

### 725 **Synthesizing realistic flow cytometry data**

726 Our method to synthesize realistic flow cytometry data is accessible as python<sup>71</sup> package  
727 *NBNode* via github <https://github.com/ggrrlab/NBNode><sup>72</sup>. The process of hierarchically gating  
728 cells and simulating data with any given effect in any cell population involves five steps. In the  
729 following, **bold** letters or arrows above the letter ( $\vec{a}$ ) denote vectors, regular letters single  
730 values.

731

732 *In the first step, we applied a uniform manual gating* to 48 human peripheral blood samples  
733 stained with the DURAClone IM T Cell Subsets Tube (Beckman Coulter GmbH). Data were  
734 preprocessed by manually recompensating the samples, removing the TIME feature, and asinh  
735 transforming all cell features  $x$

$$\text{asinh}_{\text{cofactor}}(x) = \text{asinh}\left(\frac{x}{\text{cofactor}}\right) \quad (18)$$

736 with the following cofactors: FS INT: 1, FS TOF: 1, SS INT: 1, CD45RA FITC: 1000, CCR7  
 737 PE: 2000, CD28 ECD: 2000, PD1 PC5.5: 800, CD27 PC7: 3000, CD4 APC: 4000, CD8 AF700:  
 738 10000, CD3 AA750: 500, CD57 PB: 2000, CD45 KrO: 20. Because the channel-wise median  
 739 fluorescence intensity (MFI) varied between samples, this alone was not sufficient to apply the  
 740 exact same gating to all samples. Therefore, we performed a sample-wise rescaling  
 741 (Supplementary Fig 11 and Supplementary Video 4). For every cell feature  $x$ , we identified the  
 742 positive and negative population of all cells and found the corresponding  $MFI_x^+$  and  $MFI_x^-$ .  
 743 Using these, the rescaling min-max standardizes all cells per sample,

$$rescale(x) := \frac{x - MFI_x^-}{MFI_x^+ - MFI_x^-} \quad (19)$$

744 leading to a rescaled  $MFI_{rescale(x)}^+$  of 1 and a rescaled  $MFI_{rescale(x)}^-$  of 0.

745 We then applied a standard gating strategy, which is shown schematically (Supplementary Fig  
 746 12a) and explicitly for a real-world sample (Supplementary Fig 4). This hierarchical gating of  
 747 biaxial scatter plots is effectively a decision tree with 98 “leaf” gates (Supplementary Fig 12a).  
 748 Each leaf gate corresponds to a terminal gating node and all supraordinate nodes are  
 749 “intermediate” gates. Importantly, every cell must fall into one, and only one, of the subordinate  
 750 98 leaf gates.

751  
 752 *In the second step, we model the proportion of cells in each leaf gate after uniformly gating all*  
 753 *cells from all samples. Specifically, we describe the proportion of cells in each gate according*  
 754 *to a Dirichlet distribution  $Dir(\alpha)$  (Supplementary Fig 12b,c). The Dirichlet distribution is a*  
 755 *suitable choice after its mass is only on non-negative compositions that sum up to one.*  
 756 *Following Minka *et al.*<sup>73</sup>, let  $\mathbf{p} \in (0, 1)^K$  be one random vector of proportions such that*  
 757  *$\sum_k^K p_k = 1$  for  $k \in \{1, \dots, K\}$  for  $K$  cell populations. In our case, all cells of a sample fall into*  
 758 *one and only one of the 98 terminal gates. Therefore, the sum of the cell percentages in each*

759 terminal gate adds up to 100%. The probability density under the Dirichlet model with a  
 760 parameter vector  $\boldsymbol{\alpha} \in \mathbb{R}_{>0}^K$  is defined as

$$p(\mathbf{p}) \sim \text{Dir}(\alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k p_k^{\alpha_k - 1} \quad (20)$$

$$\text{with } p_k > 0 \text{ and } \sum_k p_k = 1$$

761 More intuitively, the  $\boldsymbol{\alpha}$  parameters can be split into mean proportions per cell population and a  
 762 precision:

$$\mathbf{m} = E[\mathbf{p}] = \frac{\boldsymbol{\alpha}}{\sum_k \alpha_k} \quad (\text{mean vector}) \quad (21)$$

$$s = \sum_k \alpha_k \quad (\text{precision}) \quad (22)$$

763 Hence, a useful explanation of the parameters is that the higher the precision, the more localized  
 764 the probability becomes around the means.  $\alpha_x > \alpha_y$  indicates that, on average, the proportion  
 765 of cell population  $x$  is higher than the proportion of cell population  $y$ . If  $0 < \alpha_k < 1$ , the  
 766 distribution is effectively pushed away from the corresponding cell population. See  
 767 Supplementary Fig 13 and Supplementary Table 2 for examples of the Dirichlet distribution  
 768 with  $K=3$  and different parametrizations of  $\boldsymbol{\alpha}$ . Plots were built using the R-package<sup>74</sup>. We set a  
 769 maximum likelihood estimation of the distribution parameters  $\boldsymbol{\alpha}$ <sup>75</sup> based on  $N_{\text{samples}} = 48$   
 770 measured cell population proportions  $\mathbf{p}^{(i)}$  for  $i \in \{1, \dots, N_{\text{samples}}\}$ . In some cell populations  
 771 and samples there were no cells so the proportion became zero. Because the estimation cannot  
 772 handle proportions equal to zero, we added a pseudo-proportion to all proportions and  
 773 normalized to 1 before applying maximum likelihood estimation. With this, the zero-adjusted  
 774 proportion  $\mathbf{p}_k^{(i)''}$  of sample  $i$  and cell population  $k$  becomes

$$\mathbf{p}_k^{(i)''} = \frac{\mathbf{p}_k^{(i)'}}{\sum_k^K \mathbf{p}_k^{(i)'}} \quad (23)$$

$$\text{with } \mathbf{p}_k^{(i)'} = \mathbf{p}_k^{(i)} + 0.001 \cdot \min(\text{all proportions}) \quad (24)$$

775 We end up with a Dirichlet distribution with estimates for the parameter  $\hat{\alpha}$

$$Dir(\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_K) \quad (25)$$

776

777 *In the third step, we build a gating hierarchy using the estimated  $\alpha$  parameters* corresponding  
 778 to the leaf nodes. We used the estimated Dirichlet parameters and manual gating structure to  
 779 create a probabilistic representation of the gating hierarchy. In this structure, all cells fall into  
 780 one and only one gate. To calculate intermediate nodes, we sum the estimated  $\alpha$  parameters  
 781 according to the manual gating tree, starting from the bottom and working to the top. Given a  
 782 Dirichlet distributed variable with  $K$  cell populations

$$p(\mathbf{p}) = (p_1, \dots, p_K) \sim Dir(\alpha_1, \alpha_2, \dots, \alpha_K) \quad (26)$$

783 the sum of any two cell populations is again Dirichlet distributed

$$(p_1, \dots, p_i + p_j, \dots, p_K) \sim Dir(\alpha_1, \dots, \alpha_i + \alpha_j, \dots, \alpha_K) \quad (27)$$

784 Therefore, every intermediate or leaf node is described by a Dirichlet distribution. Intuitively,  
 785 all cells of any gate must fall in one of the subsequent gates and can, therefore, reflect a Dirichlet  
 786 distribution. To visualize proportions corresponding to these parameters, the decision tree was  
 787 shaded in red, such that deeper red indicates a higher proportion of cells in that gate  
 788 (Supplementary Fig 12a).

789

790 *In the fourth step, we fit a cell feature distribution using cells from all samples per leaf gate*  
 791 (Supplementary Fig 12d,e). The Dirichlet distributions only describe the number of cells in  
 792 every gate – that is, a vector of  $K$  cell population proportions  $\mathbf{p} \in (0, 1)^K$ . However, a flow  
 793 cytometry measurement results in a  $\mathbb{R}^{n \times m}$  matrix with  $n$  cells and  $m$  cell features where every

794 cell comes from a specific cell population. Each such cell population is defined by the  $m$   
795 continuous cell feature values. Accordingly, we model the cells for each leaf node  $l$  by a  
796 multivariate normal distribution  $\mathcal{N}(\vec{\mu}_l, \Sigma_l)$  with mean  $\vec{\mu}_l \in \mathbb{R}^m$  and covariance matrix  $\Sigma_l \in$   
797  $\mathbb{R}^{m \times m}$ . In the illustrated example, we show the parameters of one gate's normal distribution  
798 with the centers of the ellipsoids  $\vec{\mu}_l$  and the shaded areas  $\vec{\mu}_l \pm \sigma$  (Supplementary Fig 12e). We  
799 estimated the normal distributions using all cells from  $n=48$  samples. For populations with  $< 2$   
800 cells, a covariance matrix was not calculable, so such populations were removed.

801

802 *In the fifth step, we use the estimated cell population and cell feature distributions to generate*  
803 *realistic flow cytometry datasets.* We use the estimated parameters of the Dirichlet distribution  
804 and the normal distributions of each leaf node to generate cells. As shown in Supplementary  
805 Fig 14, this simulation involves: (a) drawing a vector  $\mathbf{p} \in \mathbb{R}^K$  from the estimated Dirichlet  
806 distribution  $Dir(\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_K)$ , which represents the proportion of cells in each leaf node; (b)  
807 calculating the number of synthetic cells per leaf node using the expected number of cells for  
808 the sample (e.g. 10,000 cells); and (c) Finally, drawing the required number of synthetic cells  
809 from the normal distribution of each corresponding leaf node for each sample. By repeating this  
810 process for each sample, we generate a synthetic dataset that reflects the underlying population  
811 of cells. We visualize our complete decision tree as an interactive online tool  
812 (<https://vissim.gunthergl.com/>) (Supplementary Website 1).

813

#### 814 **Realistic imitation of disease associated effects**

815 We can now introduce any given effect in any given cell population and obtain cells from a  
816 realistic synthetic sample. For that, we change the underlying Dirichlet distribution and then  
817 sample from the existing normal distributions as before. To change the proportion of cell



818 population  $x$ , we have to change its parameter  $\alpha_x$ . However, simply changing  $\alpha_x$ , e.g. by a  
 819 factor  $f \in \mathbb{R}_{>0}$  ( $\alpha'_x := f \cdot \alpha_x$ ) also changes the precision and therefore the effective change of  
 820 the population proportion is different than multiplying with  $f$

$$E[p'_x] = \frac{\alpha'_x}{s'} = \frac{f \cdot \alpha_x}{s - \alpha_x + f \cdot \alpha_x} = \frac{f \cdot \alpha_x}{s + (1 - f) \cdot \alpha_x} \neq f \frac{\alpha_x}{s} \quad (28)$$

821 Therefore, we calculate the new  $\alpha_x$  by the share of the expected target proportion in the total  
 822 old precision and the remaining precision is shared across all other nodes

$$\alpha'_x := target_{\%} \cdot s \quad (29)$$

$$\alpha'_{not\ x} := (1 - target_{\%}) \cdot s = \sum_{k \in K \setminus \{x\}} \alpha_k \quad (30)$$

823 where *not*  $x$  corresponds to all nodes which are not the changed node  $x$  **nor** subordinate nodes.  
 824 After a single synthetic cell comes from a specific leaf node distribution, we still have to express  
 825 the changed intermediate node  $x$  by its leaf nodes. After parameter  $\alpha_k$  of any node is the sum  
 826 of all leaf node parameters  $\alpha_l$  below node  $k$ , we calculate the new leaf node parameter  $\alpha'_l$  as  
 827 the old  $\alpha_l$  multiplied with the ratio of the new and old changed node above

$$\alpha'_l = \alpha_l \frac{\alpha'_x}{\alpha_x} \text{ or } \alpha'_l = \alpha_l \frac{\alpha'_{not\ x}}{\alpha_{not\ x}} \quad (31)$$

828  
 829 This finally leads us to a change in the expected proportion of the target population  $x$ .

830  
 831  
 832

833 **ACKNOWLEDGMENTS**

834 The authors thank The Bristol Myers Squibb Foundation for Immuno-Oncology (Award FA-  
835 19-009) and the Bavarian State Ministry for Science and Art's Coronavirus Research  
836 Programme for funding this work. We thank Beckman Coulter Life Sciences GmbH for its  
837 continuing support of our research. We are grateful to Ian V. Hutchinson (Providence St John's  
838 Cancer Institute, CA, USA) for proof-reading our manuscript. This work was only possible with  
839 Erika Ostermeier's outstanding technical support.

840 **DATA AVAILABILITY STATEMENT**

841 The authors declare that all data supporting the findings of this study are available within the  
842 paper, its supplementary information files and downloadable files deposited at figshare (private  
843 reviewer link: <https://figshare.com/s/27719825e2afdb5f5a14>).

844 **CODE AVAILABILITY STATEMENT**

845 The authors declare that all computer code supporting the findings of this study are available  
846 as supplementary information files and downloadable files deposited at figshare (private  
847 reviewer link: <https://figshare.com/s/27719825e2afdb5f5a14>). The python package NBNNode is  
848 accessible at <https://github.com/ggrlab/NBNNode>. The R package restrictedROC is accessible at  
849 <https://github.com/ggrlab/restrictedROC>.

850 **DISCLOSURES**

851 M.K. is a Beckman Coulter Life Sciences associate. S.H. has received consulting fees and  
852 speaker's honoraria from BMS and Merck Sharp & Dohme (MSD). J.A.H. has received  
853 consulting fees and speaker's honoraria from Mallinckrodt Pharmaceuticals. Other authors have  
854 no conflict-of-interest to declare.

## 855 REFERENCES

- 856 1 Medzhitov, R. The spectrum of inflammatory responses. *Science* **374**, 1070-1075  
857 (2021).
- 858 2 Bartok, E. & Hartmann, G. Immune Sensing Mechanisms that Discriminate Self from  
859 Altered Self and Foreign Nucleic Acids. *Immunity* **53**, 54-77 (2020).
- 860 3 Du, H. *et al.* Tuning immunity through tissue mechanotransduction. *Nat Rev Immunol*  
861 **23**, 174-188 (2023).
- 862 4 Rumpret, M., von Richthofen, H. J., Peperzak, V. & Meyaard, L. Inhibitory pattern  
863 recognition receptors. *J Exp Med* **219** (2022).
- 864 5 Deets, K. A. & Vance, R. E. Inflammasomes and adaptive immune responses. *Nat*  
865 *Immunol* **22**, 412-422 (2021).
- 866 6 Willis, J. C. & Lord, G. M. Immune biomarkers: the promises and pitfalls of  
867 personalized medicine. *Nat Rev Immunol* **15**, 323-329 (2015).
- 868 7 Scheffold, A. & Kern, F. Recent developments in flow cytometry. *J Clin Immunol* **20**,  
869 400-407 (2000).
- 870 8 Cossarizza, A. *et al.* Guidelines for the use of flow cytometry and cell sorting in  
871 immunological studies (third edition). *Eur J Immunol* **51**, 2708-3145 (2021).
- 872 9 Liechti, T. *et al.* An updated guide for the perplexed: cytometry in the high-dimensional  
873 era. *Nat Immunol* **22**, 1190-1197 (2021).
- 874 10 Maecker, H. T., McCoy, J. P. & Nussenblatt, R. Standardizing immunophenotyping for  
875 the Human Immunology Project. *Nat Rev Immunol* **12**, 191-200 (2012).
- 876 11 Spidlen, J. *et al.* Data File Standard for Flow Cytometry, Version FCS 3.2. *Cytometry*  
877 *A* **99**, 100-102 (2021).
- 878 12 Liechti, T. *et al.* Immune phenotypes that are associated with subsequent COVID-19  
879 severity inferred from post-recovery samples. *Nat Commun* **13**, 7255 (2022).
- 880 13 Glehr, G. *et al.* External validation of biomarkers for immune-related adverse events  
881 after immune checkpoint inhibition. *Front Immunol* **13**, 1011040 (2022).
- 882 14 Das, R. *et al.* Early B cell changes predict autoimmunity following combination immune  
883 checkpoint blockade. *J Clin Invest* **128**, 715-720 (2018).
- 884 15 Bukhari, S. *et al.* Single-cell RNA sequencing reveals distinct T cell populations in  
885 immune-related adverse events of checkpoint inhibitors. *Cell Rep Med* **4**, 100868  
886 (2023).
- 887 16 Lozano, A. X. *et al.* T cell characteristics associated with toxicity to immune checkpoint  
888 blockade in patients with melanoma. *Nat Med* **28**, 353-362 (2022).

- 889 17 Livingstone, E. *et al.* Adjuvant nivolumab plus ipilimumab or nivolumab alone versus  
890 placebo in patients with resected stage IV melanoma with no evidence of disease  
891 (IMMUNED): final results of a randomised, double-blind, phase 2 trial. *Lancet* **400**,  
892 1117-1129 (2022).
- 893 18 Schneider, B. J. *et al.* Management of Immune-Related Adverse Events in Patients  
894 Treated With Immune Checkpoint Inhibitor Therapy: ASCO Guideline Update. *J Clin*  
895 *Oncol* **39**, 4073-4126 (2021).
- 896 19 Wang, D. Y. *et al.* Fatal Toxic Effects Associated With Immune Checkpoint Inhibitors:  
897 A Systematic Review and Meta-analysis. *JAMA Oncol* **4**, 1721-1728 (2018).
- 898 20 Conroy, M. & Naidoo, J. Immune-related adverse events and the balancing act of  
899 immunotherapy. *Nat Commun* **13**, 392 (2022).
- 900 21 Esfahani, K. *et al.* Moving towards personalized treatments of immune-related adverse  
901 events. *Nat Rev Clin Oncol* **17**, 504-515 (2020).
- 902 22 Ganesan, S. & Mehnert, J. Biomarkers for Response to Immune Checkpoint Blockade.  
903 *Annual Review of Cancer Biology* **4**, 331-351 (2020).
- 904 23 Maecker, H. T. *et al.* New tools for classification and monitoring of autoimmune  
905 diseases. *Nat Rev Rheumatol* **8**, 317-328 (2012).
- 906 24 Fox, A. *et al.* Cyto-Feature Engineering: A Pipeline for Flow Cytometry Analysis to  
907 Uncover Immune Populations and Associations with Disease. *Sci Rep* **10**, 7651 (2020).
- 908 25 Van Gassen, S. *et al.* FlowSOM: Using self-organizing maps for visualization and  
909 interpretation of cytometry data. *Cytometry A* **87**, 636-645 (2015).
- 910 26 Brummelman, J. *et al.* Development, application and computational analysis of high-  
911 dimensional fluorescent antibody panels for single-cell flow cytometry. *Nat Protoc* **14**,  
912 1946-1969 (2019).
- 913 27 Hu, Z., Tang, A., Singh, J., Bhattacharya, S. & Butte, A. J. A robust and interpretable  
914 end-to-end deep learning model for cytometry data. *Proc Natl Acad Sci U S A* **117**,  
915 21373-21380 (2020).
- 916 28 Martins, F. *et al.* Adverse effects of immune-checkpoint inhibitors: epidemiology,  
917 management and surveillance. *Nat Rev Clin Oncol* **16**, 563-580 (2019).
- 918 29 McKean, W. B., Moser, J. C., Rimm, D. & Hu-Lieskovan, S. Biomarkers in Precision  
919 Cancer Immunotherapy: Promise and Challenges. *Am Soc Clin Oncol Educ Book* **40**,  
920 e275-e291 (2020).
- 921 30 Lakshmikanth, T. *et al.* Human Immune System Variation during 1 Year. *Cell Rep* **32**,  
922 107923 (2020).
- 923 31 Harrington, C. *et al.* Noninvasive biomarkers for the diagnosis and management of  
924 autoimmune hepatitis. *Hepatology* **76**, 1862-1879 (2022).
-

- 925 32 Schilling, H. L. *et al.* Development of a Flow Cytometry Assay to Predict Immune  
926 Checkpoint Blockade-Related Complications. *Front Immunol* **12**, 765644 (2021).
- 927 33 Hanley, J. A. & McNeil, B. J. The meaning and use of the area under a receiver operating  
928 characteristic (ROC) curve. *Radiology* **143**, 29-36 (1982).
- 929 34 Gneiting, T. & Vogel, P. Receiver operating characteristic (ROC) curves: equivalences,  
930 beta model, and minimum distance estimation. *Machine Learning* **111**, 2147-2159  
931 (2022).
- 932 35 Hutchinson, J. A. *et al.* Virus-specific memory T cell responses unmasked by immune  
933 checkpoint blockade cause hepatitis. *Nat Commun* **12**, 1439 (2021).
- 934 36 Glehr, G. Restriction method for marker discovery in R. R package version 2.3.4.  
935 <https://github.com/ggrrlab/restrictedROC> (2023).
- 936 37 Fawcett, T. An introduction to ROC analysis. *Pattern Recognition Letters* **27**, 861-874  
937 (2006).
- 938 38 Good, P. *Permutation Tests*. 2 edn, (Springer New York, 2000).
- 939 39 Breiman, L. Random Forests. *Machine Learning* **45**, 5-32 (2001).
- 940 40 Burke, K. P., Grebinoski, S., Sharpe, A. H. & Vignali, D. A. A. Understanding adverse  
941 events of immunotherapy: A mechanistic perspective. *J Exp Med* **218** (2021).
- 942 41 Kverneland, A. H. *et al.* Age and gender leucocytes variances and references values  
943 generated using the standardized ONE-Study protocol. *Cytometry A* **89**, 543-564 (2016).
- 944 42 Arvaniti, E. & Claassen, M. Sensitive detection of rare disease-associated cell subsets  
945 via representation learning. *Nat Commun* **8**, 14825 (2017).
- 946 43 Roederer, M. How many events is enough? Are you positive? *Cytometry A* **73**, 384-385  
947 (2008).
- 948 44 Robin, X. *et al.* pROC: an open-source package for R and S+ to analyze and compare  
949 ROC curves. *BMC Bioinformatics* **12**, 77 (2011).
- 950 45 Dodd, L. E. & Pepe, M. S. Partial AUC estimation and regression. *Biometrics* **59**, 614-  
951 623 (2003).
- 952 46 Parodi, S. *et al.* Restricted ROC curves are useful tools to evaluate the performance of  
953 tumour markers. *Stat Methods Med Res* **25**, 294-314 (2016).
- 954 47 Yang, H., Lu, K., Lyu, X. & Hu, F. Two-way partial AUC and its properties. *Stat*  
955 *Methods Med Res* **28**, 184-195 (2019).
- 956 48 McClish, D. K. Analyzing a portion of the ROC curve. *Med Decis Making* **9**, 190-195  
957 (1989).
-

- 
- 958 49 Bhattacharya, B. & Hughes, G. Symmetry of receiver operating characteristic curves  
959 and Kullback–Leibler divergences between the signal and noise populations. *Journal of*  
960 *Mathematical Psychology* **55**, 365-367 (2011).
- 961 50 Cron, A. *et al.* Hierarchical modeling for rare event detection and cell subset alignment  
962 across flow cytometry samples. *PLoS Comput Biol* **9**, e1003130 (2013).
- 963 51 Johnsson, K., Wallin, J. & Fontes, M. BayesFlow: latent modeling of flow cytometry  
964 cell populations. *BMC Bioinformatics* **17**, 25 (2016).
- 965 52 Boris, P. H., Chariff, A., Raphael, G., François, C. & Rodolphe, T. Sequential Dirichlet  
966 process mixtures of multivariate skew  $t$ -distributions for model-based clustering of  
967 flow cytometry data. *The Annals of Applied Statistics* **13**, 638-660 (2019).
- 968 53 Youden, W. J. Index for rating diagnostic tests. *Cancer* **3**, 32-35 (1950).
- 969 54 Cassel, C. K. & Guest, J. A. Choosing wisely: helping physicians and patients make  
970 smart decisions about their care. *JAMA* **307**, 1801-1802 (2012).
- 971 55 Perez-Ruiz, E. *et al.* Prophylactic TNF blockade uncouples efficacy and toxicity in dual  
972 CTLA-4 and PD-1 immunotherapy. *Nature* **569**, 428-432 (2019).
- 973 56 Pinal-Fernandez, I. *et al.* Transcriptomic profiling reveals distinct subsets of immune  
974 checkpoint inhibitor induced myositis. *Ann Rheum Dis* (2023).
- 975 57 De Martin, E. *et al.* Characterization of liver injury induced by cancer immunotherapy  
976 using immune checkpoint inhibitors. *J Hepatol* **68**, 1181-1190 (2018).
- 977 58 Groha, S. *et al.* Germline variants associated with toxicity to immune checkpoint  
978 blockade. *Nat Med* **28**, 2584-2591 (2022).
- 979 59 Khan, Z. *et al.* Genetic variation associated with thyroid autoimmunity shapes the  
980 systemic immune response to PD-1 checkpoint blockade. *Nat Commun* **12**, 3355 (2021).
- 981 60 McQuade, J. L. *et al.* Association of Body Mass Index With the Safety Profile of  
982 Nivolumab With or Without Ipilimumab. *JAMA Oncol* **9**, 102-111 (2023).
- 983 61 McCulloch, J. A. *et al.* Intestinal microbiota signatures of clinical response and  
984 immune-related adverse events in melanoma patients treated with anti-PD-1. *Nat Med*  
985 **28**, 545-556 (2022).
- 986 62 McShane, L. M. *et al.* Reporting recommendations for tumor marker prognostic studies  
987 (REMARK). *J Natl Cancer Inst* **97**, 1180-1184 (2005).
- 988 63 Eisenhauer, E. A. *et al.* New response evaluation criteria in solid tumours: revised  
989 RECIST guideline (version 1.1). *Eur J Cancer* **45**, 228-247 (2009).
- 990 64 Kronenberg, K., Riquelme, P. & Hutchinson, J. A. Standard protocols for immune  
991 profiling of peripheral blood leucocyte subsets by flow cytometry using DuraClone IM  
992 reagents. *Protocol Exchange* (2019).
-

- 993 65 Bamber, D. The area above the ordinal dominance graph and the area below the receiver  
994 operating characteristic graph. *Journal of Mathematical Psychology* **12**, 387-415  
995 (1975).
- 996 66 Sullivan Pepe, M. *The Statistical Evaluation of Medical Tests for Classification and*  
997 *Prediction*. (Oxford University Press, 2004).
- 998 67 Yang, Z. *et al.* Optimizing Two-way Partial AUC with an End-to-end Framework. *IEEE*  
999 *Trans Pattern Anal Mach Intell* **PP** (2022).
- 1000 68 Mann, H. B. & Whitney, D. R. On a Test of Whether one of Two Random Variables is  
1001 Stochastically Larger than the Other. *The Annals of Mathematical Statistics* **18**, 50-60  
1002 (1947).
- 1003 69 Rosner, B. & Glynn, R. J. Power and sample size estimation for the Wilcoxon rank sum  
1004 test with application to comparisons of C statistics from alternative prediction models.  
1005 *Biometrics* **65**, 188-197 (2009).
- 1006 70 LeDell, E. & Poirier, S. Scalable Automatic Machine Learning. *7th ICML Workshop on*  
1007 *Automated Machine Learning (AutoML)* (2020).
- 1008 71 Van Rossum, G. & Drake, F. L. *Python 3 Reference Manual*. (CreateSpace, 2009).
- 1009 72 Glehr, G. Non-binary node and flow cytometry simulation python library. Python  
1010 package version 1.0.0. <https://github.com/ggrrlab/NBNode>. (2023).
- 1011 73 Minka, T. Estimating a Dirichlet distribution. (2003).
- 1012 74 Kahle, D. dirichlet: A light weight package for the (d/r)dirichlet functions for the  
1013 Dirichlet distribution in R. R package version 0.0.999.  
1014 <https://github.com/dkahle/dirichlet>. (2023).
- 1015 75 Suh, E. J. Dirichlet MLE python library. Python package version v0.9.0.  
1016 <https://github.com/ericsuh/dirichlet>. (2023).  
1017

1018 **FIGURE LEGENDS**

1019 **FIGURE 1.** *Two-class distributions resulting in asymmetric ROC curves.* We present  
1020 hypothetical and real-world examples of marker distributions in two classes, which are intended  
1021 to represent sets of patients with different clinical outcomes. The distribution of values from  
1022 positive (i.e. diseased) class are coloured green and values from negative (i.e. control) class are  
1023 coloured red; the overlapping density areas are coloured in purple. For each example, we  
1024 present the corresponding ROC curve. **(a)** A hypothetical example of perfectly discriminatory  
1025 marker with negative  $\mathcal{N}(5, 1)$  and positive  $\mathcal{N}(15, 1)$  populations that give rise to a symmetrical  
1026 ROC curve. The area under the ROC curve (AUC) is 1.0. **(b)** A hypothetical example of  
1027 substantially overlapping marker distributions in the negative  $\mathcal{N}(5, 1)$  and positive  
1028  $\mathcal{N}(6, 1)$  populations that gives rise to a symmetrical ROC curve with AUC=0.76. **(c)** A  
1029 hypothetical example of a non-informative marker distribution with negative  $\mathcal{N}(5, 1)$  and  
1030 positive  $\mathcal{N}(5, 1)$  populations that gives rise to a diagonal ROC curve with AUC = 0.5. **(d)** A  
1031 hypothetical example of substantially overlapping marker distributions in the negative  $\mathcal{N}(5, 1)$   
1032 and positive  $\mathcal{N}(6, 2)$  populations with unequal variance that gives rise to right-skewed ROC  
1033 curve with AUC = 0.67. **(e)** A real-world example of substantially overlapping distributions of  
1034 absolute erythrocyte counts with unequal variance in patients with metastatic melanoma who  
1035 responded (n = 61) or did not respond (n = 44) to combined Ipilimumab and Nivolumab (Ipi-  
1036 Nivo) therapy. We observe a right-skewed ROC curve with AUC = 0.62. **(f)** A hypothetical  
1037 example of substantially overlapping marker distributions in the negative  $\mathcal{N}(5, 2)$  and positive  
1038  $\mathcal{N}(6, 1)$  populations with unequal variance that gives rise to a left-skewed ROC curve with  
1039 AUC = 0.67. **(g)** A real-world example of substantially overlapping distributions of CD8<sup>+</sup>  $\gamma\delta$  T  
1040 cells with unequal variance in patients with metastatic melanoma who did (n = 22) or did not  
1041 (n = 42) develop treatment-related hepatitis after Ipi-Nivo therapy. We observe a left-skewed



1042 ROC curve with AUC = 0.69. **(h)** A hypothetical example of substantially overlapping marker  
1043 distributions in the normally distributed negative  $\mathcal{N}(5, 1)$  and bimodally distributed positive  
1044 populations. In this example, the positive population comprises 10 % cases with elevated  
1045 marker expression  $\mathcal{N}(9, 1)$  and 90 % cases with unaltered marker expression  $\mathcal{N}(5, 1)$ .  
1046 Phenotypic heterogeneity in the diseased cases gives rise to a right-skewed ROC curve with  
1047 AUC = 0.55. **(i)** A real-world example of a phenotypically heterogeneous set of patients with  
1048 metastatic melanoma who did (n = 48) or did not (n=62) develop treatment-related hepatitis  
1049 after Ipi-Nivo therapy. As previously described, a subset of these patients exhibited a baseline  
1050 expansion of CD4<sup>+</sup> T<sub>EM</sub> cells, which was likely driven by subclinical cytomegalovirus (CMV)  
1051 reactivation. Consequently, CD4<sup>+</sup> T<sub>EM</sub> cell frequency before therapy is a weakly discriminatory  
1052 marker of hepatitis-risk that gives rise to a right-skewed ROC curve with AUC = 0.64.

1053 **FIGURE 2.** *Method to optimally restrict datasets to classifiable samples.* We present a  
1054 simulated example of marker distributions in two classes, which are intended to represent sets  
1055 of patients with different clinical outcomes. **(a)** The distribution of values from positive (i.e.  
1056 diseased,  $n=2500$ ) class are coloured green and values from negative (i.e. control,  $n=2500$ ) class  
1057 are coloured red; the overlapping density areas are coloured in purple. In this example, 20 % of  
1058 positive samples and 2 % of negative samples were drawn from a population with elevated  
1059 marker expression  $\mathcal{N}(9, 1)$ . All other samples were drawn from a population with unaltered  
1060 marker expression  $\mathcal{N}(6, 1)$ . The optimal restriction of this dataset lies at a marker value of 6.8,  
1061 which is marked with a red line. Restriction of the dataset defines two subsets of samples –  
1062 explicitly, marker<sup>HIGH</sup> (orange) and marker<sup>LOW</sup> (blue) samples. **(b)** A complete ROC curve  
1063 marked at the optimal restriction point (red lines) that corresponds to  $FPR = 0.258$ . Restricting  
1064 the parts of the ROC curve corresponding to marker<sup>HIGH</sup> or marker<sup>LOW</sup> samples gives us  
1065 restricted ROC curves for which restricted AUCs (rAUCs) can be calculated. This is equivalent  
1066 to rescaling the respective part of the ROC curve. **(c)** Adjusting the rAUC for the number of  
1067 samples delimited by the restriction gives the restricted standardized AUC (rzAUC). Hence, we  
1068 can plot rzAUC for marker<sup>HIGH</sup> and marker<sup>LOW</sup> samples at all possible restriction values. The  
1069 optimal restriction value is defined as the maximum absolute rzAUC for either the marker<sup>HIGH</sup>  
1070 or marker<sup>LOW</sup> samples. **(d)** A complete ROC curve to illustrate the delimitation of marker<sup>HIGH</sup>  
1071 values (orange rectangle) according to the optimal restriction. **(e)** Densities of the negative and  
1072 positive classes after restriction to marker<sup>HIGH</sup> values. **(f)** ROC curve constructed from  
1073 marker<sup>HIGH</sup> samples. Intuitively, we see that recalculating the ROC curve using only marker<sup>HIGH</sup>  
1074 samples is equivalent to rescaling the partial ROC curve. **(g)** A complete ROC curve to illustrate  
1075 the delimitation of marker<sup>LOW</sup> values (blue rectangle) according to the optimal restriction. **(h)**

1076 Densities of the negative and positive classes after restriction to marker<sup>LOW</sup> values. (i) ROC  
1077 curve constructed from marker<sup>LOW</sup> samples.

1078 **FIGURE 3.** *Optimal restriction of two-class distributions results in asymmetric ROC curves.*  
1079 We present four simulated examples of marker distributions in two classes, which are intended  
1080 to represent sets of patients with different clinical outcomes. The distribution of values from  
1081 the positive (i.e. diseased) class are coloured green and values from the negative (i.e. control)  
1082 class are coloured red; the overlapping density areas are coloured in purple. For each example,  
1083 we present the following: a plot of positive and negative class densities; the complete ROC  
1084 curve; a plot of marker values against FPR; a plot of rzAUC calculated for marker<sup>HIGH</sup> (orange)  
1085 and marker<sup>LOW</sup> (blue) samples at all FPR values. In each plot, red lines indicate the optimal  
1086 restriction as a marker value or FPR value. **(a)** A simulated example of a symmetric ROC curve  
1087 from 100 negative  $\mathcal{N}(5, 1)$  and 100 positive  $\mathcal{N}(6, 1)$  samples. **(b)** A simulated example of a  
1088 right-skewed ROC curve from 100 negative  $\mathcal{N}(5, 1)$  and 100 positive  $\mathcal{N}(6, 2)$  samples. **(c)** A  
1089 simulated example of a left-skewed ROC curve from 100 negative  $\mathcal{N}(5, 2)$  and 100 positive  
1090  $\mathcal{N}(6, 1)$  samples. **(d)** Results for a right-skewed ROC curve from 100 negative  
1091  $\mathcal{N}(5, 1)$  samples and 100 positive samples from a bimodally distributed positive population. In  
1092 this example, the positive population comprises 10 % cases with elevated marker expression  
1093  $\mathcal{N}(9, 1)$  and 90 % cases with unaltered marker expression  $\mathcal{N}(5, 1)$ .

1094 **Figure 4.** *Realistically synthesizing flow cytometry data from two class distributions.* Various  
1095 applications of our synthetic flow cytometry data depend upon generating samples with  
1096 differences in cell subset distributions. Here, we provide an example of increasing the  
1097 proportion of CD8<sup>+</sup> T<sub>EMRA</sub> cells in synthetic samples from a baseline value of 7.17 % in the  
1098 negative class to an altered value of 33.23 % in the positive class. The intensity of red shading  
1099 in the gating trees illustrates this change in CD8<sup>+</sup> T<sub>EMRA</sub> cells and contingent changes in other  
1100 populations. **(a)** Gating tree with 7.17% CD8<sup>+</sup> T<sub>EMRA</sub> cells, its corresponding Dirichlet  
1101 distribution parameter ( $\alpha = 12.34$ ) and the density for three example gates: A, B and L. The  
1102 cell count table for three samples drawn from this distribution is shown. **(b)** Gating tree with  
1103 33.23% CD8<sup>+</sup> T<sub>EMRA</sub> cells, its corresponding Dirichlet distribution parameter ( $\alpha = 57.18$ ) and  
1104 the density for three example gates: A, B and L. The cell count table for three samples drawn  
1105 from this distribution is shown. Of special note, percentages of cells in all other gates also  
1106 changed according to the Dirichlet distribution, leading to changes in simulated cell counts  
1107 across all leaf gates.  
1108

1109 **FIGURE 5.** *Restriction of synthetheized flow cytometry datasets.* We present examples of  
1110 marker distributions in two classes, which are intended to represent sets of patients with  
1111 different clinical outcomes. The distribution of values from positive (i.e. diseased) class are  
1112 coloured green and values from negative (i.e. control) class are coloured red; overlapping  
1113 density areas are coloured in purple. For each example, we present the following: a plot of  
1114 positive and negative class densities; the complete ROC curve; a plot of marker values against  
1115 FPR; a plot of rAUC calculated for marker<sup>HIGH</sup> (orange) and marker<sup>LOW</sup> (blue) samples at all  
1116 FPR values. In each plot, red lines indicate the optimal restriction as a marker value or FPR  
1117 value. **(a)** A synthetic example of a symmetrical ROC curve from 100 negative  $\mathcal{N}(7.7, 1)$  and  
1118 100 positive  $\mathcal{N}(10.7, 1)$  samples. **(b)** A synthetic example of a right-skewed ROC curve from  
1119 from 100 negative  $\mathcal{N}(7.7, 1)$  and 100 positive  $\mathcal{N}(8.7, 3)$  samples. **(c)** A synthetic example of  
1120 a left-skewed ROC curve from 100 negative  $\mathcal{N}(7.7, 3)$  and 100 positive  $\mathcal{N}(8.7, 1)$  samples.  
1121 **(d)** Results for a synthetic right-skewed ROC curve from 100 negative  $\mathcal{N}(7.7, 1)$  samples and  
1122 100 positive samples from a bimodally distributed positive population. In this example, the  
1123 positive population comprises 20 % cases with elevated marker expression  $\mathcal{N}(16.7, 1)$  and 80  
1124 % cases with unaltered marker expression  $\mathcal{N}(7.7, 1)$ .

1125 **FIGURE 6.** *Restriction aids marker discovery in synthetic and real-world flow cytometry*  
1126 *datasets.* Using synthetic and real-world datasets, we demonstrate that dataset restriction  
1127 augments discovery of novel immune markers with discriminatory capacity over a limited range  
1128 of values. In each plot, the x-axis shows permutation p-values for the AUC of complete ROC  
1129 curves for every gated cell population; the y-axis shows permutation p-values for the AUC of  
1130 optimally restricted ROC curves. Points within the green-shaded rectangles represent cell  
1131 subsets for which p-values derived from unrestricted data are not significant ( $p \geq 0.05$ ), but p-  
1132 values derived from optimally restricted data are significant ( $p < 0.05$ ). Named subsets falling  
1133 in the green-shaded area are disease-related markers discovered exclusively through dataset  
1134 restriction. **(a)** Permutation p-values from synthetic samples in which a disease-related effect  
1135 was introduced into  $CD4^+ T_{EM}$  resulting in a symmetric ROC curve. 100 samples in the negative  
1136  $\mathcal{N}(7.7, 1)$  class and 100 samples in the positive  $\mathcal{N}(10.7, 1)$  class were generated. **(b)**  
1137 Permutation p-values from synthetic samples in which a disease-related effect was introduced  
1138 into  $CD4^+ T_{EM}$  resulting in a right-skewed ROC curve. 100 samples in the negative  
1139  $\mathcal{N}(7.7, 1)$  class and 100 samples in the positive  $\mathcal{N}(8.7, 3)$  class were generated. **(c)**  
1140 Permutation p-values from synthetic samples in which a disease-related effect was introduced  
1141 into  $CD4^+ T_{EM}$  resulting in a left-skewed ROC curve. 100 samples in the negative  
1142  $\mathcal{N}(7.7, 3)$  class and 100 samples in the positive  $\mathcal{N}(8.7, 1)$  class were generated. **(d)**  
1143 Permutation p-values from synthetic samples in which a disease-related effect was introduced  
1144 into  $CD4^+ T_{EM}$  resulting in a right-skewed ROC curve. 100 samples in the negative  
1145  $\mathcal{N}(7.7, 1)$  class and 100 samples from a bimodally distributed positive class were generated. In  
1146 this example, the positive population comprises 20 % cases with elevated marker expression  
1147  $\mathcal{N}(16.7, 1)$  and 80% cases with unaltered marker expression  $\mathcal{N}(7.7, 1)$ . **(e)** Permutation p-  
1148 values from a training set of real-world clinical flow cytometry samples (n=110). 84 markers

1149 were selected where  $\geq 10$  % of samples had more than 10 counts. Dataset restriction reveals 4  
1150 previously undescribed markers of treatment-related colitis risk in metastatic melanoma  
1151 patients receiving Ipi-Nivo therapy. **(f)** Permutation p-values from a training set of real-world  
1152 clinical flow cytometry samples (n=110). 84 markers were selected where  $\geq 10$  % of samples  
1153 had more than 10 counts. Dataset restriction reveals 7 previously undescribed markers of  
1154 treatment-related hepatitis risk in metastatic melanoma patients receiving Ipi-Nivo therapy.

1155



1156 **FIGURE 7.** *Clinical interpretation of restricted markers in predicting disease.* Our method of  
1157 dataset restriction leads to counterintuitive clinical interpretations of marker values. This is  
1158 illustrated by our discovery of CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cells as a univariate marker of hepatitis  
1159 risk after immunotherapy. Here, we illustrate the conventional evaluation of marker  
1160 performance across all samples with evaluation of marker performance in a restricted dataset.  
1161 **(a)** Densities of CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cells in all samples from patients with metastatic  
1162 melanoma who developed hepatitis (n = 48) or did not (n = 62) after starting Ipi-Nivo therapy.  
1163 **(b)** Following the classical approach of determining a classification cut-off for CD27<sup>+</sup> CD28<sup>+</sup>  
1164 CD4<sup>+</sup> T<sub>EM</sub> frequency relative to CD4<sup>+</sup> T cells using the Youden Index, we predict hepatitis if >  
1165 9.62% and then assess the correct classification rate (CCR), negative predictive value (NPV),  
1166 positive predictive value (PPV), sensitivity (or true positive rate, TPR) and specificity (or true  
1167 negative rate, TNR) for all samples. **(c)** Our restriction method is predicated on there being a  
1168 range of values over which a marker provides no discriminatory information. Optimally  
1169 restricting CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cell values leads us to discard 58 of 110 samples as  
1170 “unclassifiable.” For the remaining 42 samples where CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> frequency  
1171 relative to CD4<sup>+</sup> T cells > 7.62% , we determine a classification cut-off using the Youden Index,  
1172 again predicting hepatitis if > 9.62%. Accordingly, we obtain a confusion table with CCR =  
1173 76.9%, specificity = 64.3% sensitivity = 91.7%, PPV = 68.8% and NPV = 90% across the  
1174 classifiable samples.

1175

1176 **Figure 8.** *Significant markers according to classical and our analysis predicting hepatitis. (a)*  
1177 Heatmap showing significant markers of hepatitis risk after Ipi-Nivo therapy based upon  
1178 permutation p-values for unrestricted and restricted AUC. Only permutation p-values for  
1179 markers discovered using our restriction method remained significant after correction for  
1180 multiple testing, indicated by marker names in red text. Each further column reflects one  
1181 sample, each row a feature. The samples are grouped into patients who did (green) or did not  
1182 (red) develop treatment-related hepatitis, shown in the very first row. The main matrix consists  
1183 of three values: Those excluded according by restriction (white); those included and predicted  
1184 positive (dark green); and those included and predicted negative (dark red). Columns were  
1185 clustered, rows in increasing order according to the number of excluded samples. **(b)** Random  
1186 forest predictions and performances on the prospective validation cohort (n=30) trained on  
1187 unrestricted marker values from the 110 training samples. **(c)** Random forest model predictions  
1188 and performances on the prospective validation cohort (n=30) trained on restricted marker  
1189 values from the 110 training samples. In this case, restricted marker values were replaced with  
1190 -1 before training the random forest.

# Figures

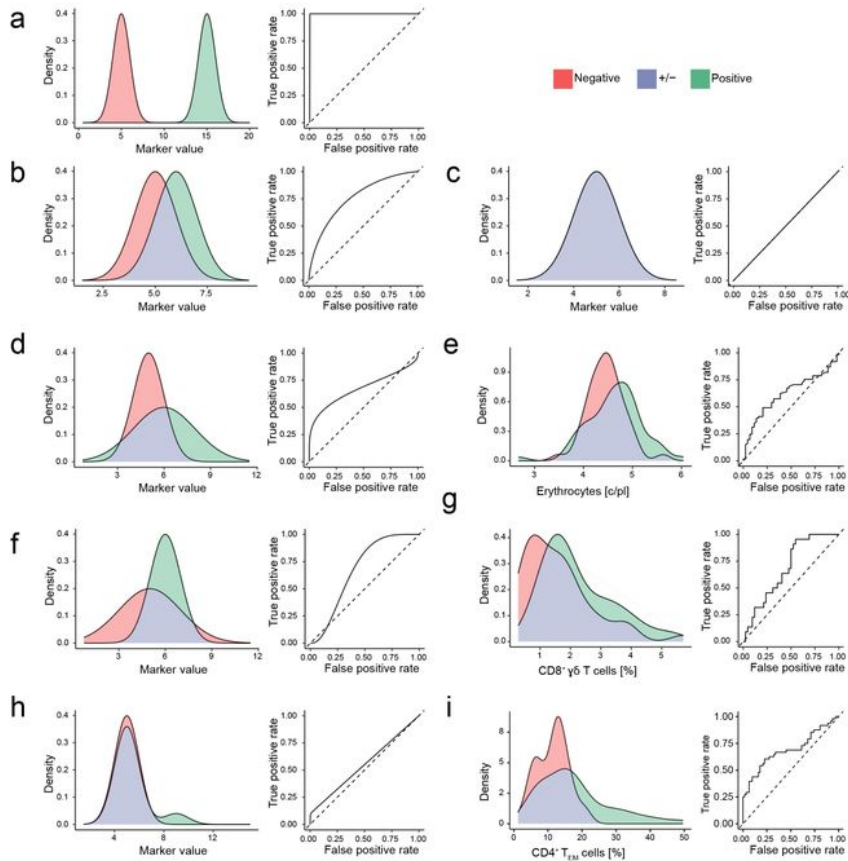


Figure 1

**FIGURE 1.** Two-class distributions resulting in asymmetric ROC curves. We present hypothetical and real-world examples of marker distributions in two classes, which are intended to represent sets of patients with different clinical outcomes. The distribution of values from positive (i.e. diseased) class are colored green and values from negative (i.e. control) class are colored red; the overlapping density areas are colored in purple. For each example, we present the corresponding ROC curve. (a) A hypothetical example of perfectly discriminative marker with negative  $N(5, 1)$  and positive  $N(15, 1)$  populations that give rise to a symmetrical ROC curve. The area under the ROC curve (AUC) is 1.0. (b) A hypothetical example of substantially overlapping marker distributions in the negative  $N(5, 1)$  and positive  $N(6, 1)$  populations that give rise to a symmetrical ROC curve with AUC=0.76. (c) A hypothetical example of a non-informative marker distribution with negative  $N(5, 1)$  and positive  $N(5, 1)$  populations that give rise to a diagonal ROC curve with AUC = 0.5. (d) A hypothetical example of substantially overlapping marker distributions in the negative  $N(5, 1)$  and positive  $N(6, 2)$  populations with unequal variance that give rise to right-skewed ROC curve with AUC = 0.67. (e) A real-world example of substantially overlapping distributions of absolute erythrocyte counts with unequal variance in patients with metastatic sarcoma who responded ( $n = 61$ ) or did not respond ( $n = 44$ ) to combined lipiramicin and Nivolumab (Ipi-Nivo) therapy. We observe a right-skewed ROC curve with AUC = 0.62. (f) A hypothetical example of substantially overlapping marker distributions in the negative  $N(5, 2)$  and positive  $N(6, 1)$  populations with unequal variance that give rise to a left-skewed ROC curve with AUC = 0.67. (g) A real-world example of substantially overlapping distribution of CD8<sup>+</sup> γδ T cells with unequal variance in patients with metastatic sarcoma who did ( $n = 22$ ) or did not ( $n = 42$ ) develop treatment-related hepatitis after Ipi-Nivo therapy. We observe a left-skewed ROC curve with AUC = 0.69. (h) A hypothetical example of substantially overlapping marker distributions in the normally distributed negative  $N(5, 1)$  and broadly distributed positive population. In this example, the positive population comprises 10 % cases with elevated marker expression  $N(9, 1)$  and 90 % cases with unaltered marker expression  $N(5, 1)$ . Phenotypic heterogeneity in the diseased cases gives rise to a right-skewed ROC curve with AUC = 0.55. (i) A real-world example of a phenotypically heterogeneous set of patients with metastatic sarcoma who did ( $n = 48$ ) or did not ( $n=62$ ) develop treatment-related hepatitis after Ipi-Nivo therapy. As previously described, a subset of these patients exhibited a basophilic expression of CD4<sup>+</sup> T<sub>H17</sub> cells, which was likely driven by subclinical cytomegalovirus (CMV) reactivation. Consequently, CD4<sup>+</sup> T<sub>H17</sub> cell frequency before therapy is a weakly discriminatory marker of hepatitis-risk that gives rise to a right-skewed ROC curve with AUC = 0.64.

Figure 1

See image above for figure legend.

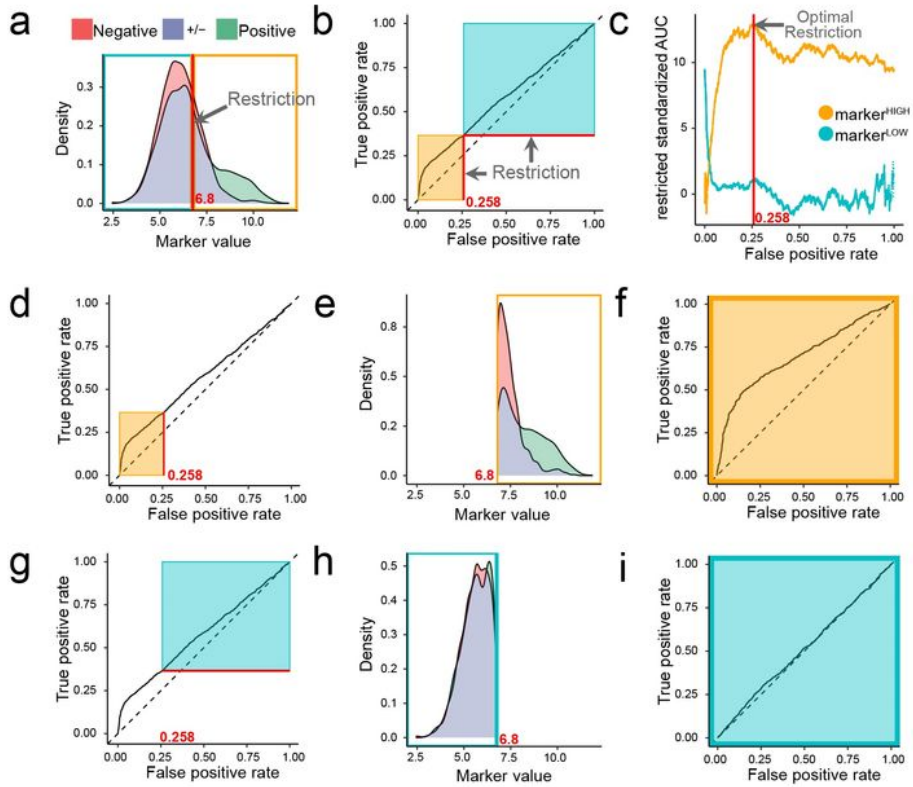


Figure 2

**FIGURE 2.** Method to optimally restrict datasets to classifiable samples. We present a simulated example of marker distributions in two classes, which are intended to represent sets of patients with different clinical outcomes. **(a)** The distribution of values from positive (i.e. diseased,  $n=2500$ ) class are coloured green and values from negative (i.e. control,  $n=2500$ ) class are coloured red; the overlapping density areas are coloured in purple. In this example, 20% of positive samples and 2% of negative samples were drawn from a population with elevated marker expression  $N(9, 1)$ . All other samples were drawn from a population with unaltered marker expression  $N(6, 1)$ . The optimal restriction of this dataset lies at a marker value of 6.8, which is marked with a red line. Restriction of the dataset defines two subsets of samples – explicitly,  $\text{marker}^{\text{HIGH}}$  (orange) and  $\text{marker}^{\text{LOW}}$  (blue) samples. **(b)** A complete ROC curve marked at the optimal restriction point (red lines) that corresponds to  $\text{FPR} = 0.258$ . Restricting the parts of the ROC curve corresponding to  $\text{marker}^{\text{HIGH}}$  or  $\text{marker}^{\text{LOW}}$  samples gives us restricted ROC curves for which restricted AUCs ( $r\text{AUC}$ s) can be calculated. This is equivalent to rescaling the respective part of the ROC curve. **(c)** Adjusting the  $r\text{AUC}$  for the number of samples delimited by the restriction gives the restricted standardized AUC ( $rz\text{AUC}$ ). Hence, we can plot  $rz\text{AUC}$  for  $\text{marker}^{\text{HIGH}}$  and  $\text{marker}^{\text{LOW}}$  samples at all possible restriction values. The optimal restriction value is defined as the maximum absolute  $rz\text{AUC}$  for either the  $\text{marker}^{\text{HIGH}}$  or  $\text{marker}^{\text{LOW}}$  samples. **(d)** A complete ROC curve to illustrate the delimitation of  $\text{marker}^{\text{HIGH}}$  values (orange rectangle) according to the optimal restriction. **(e)** Densities of the negative and positive classes after restriction to  $\text{marker}^{\text{HIGH}}$  values. **(f)** ROC curve constructed from  $\text{marker}^{\text{HIGH}}$  samples. Intuitively, we see that recalculating the ROC curve using only  $\text{marker}^{\text{HIGH}}$  samples is equivalent to rescaling the partial ROC curve. **(g)** A complete ROC curve to illustrate the delimitation of  $\text{marker}^{\text{LOW}}$  values (blue rectangle) according to the optimal restriction. **(h)** Densities of the negative and positive classes after restriction to  $\text{marker}^{\text{LOW}}$  values. **(i)** ROC curve constructed from  $\text{marker}^{\text{LOW}}$  samples.

## Figure 2

See image above for figure legend.

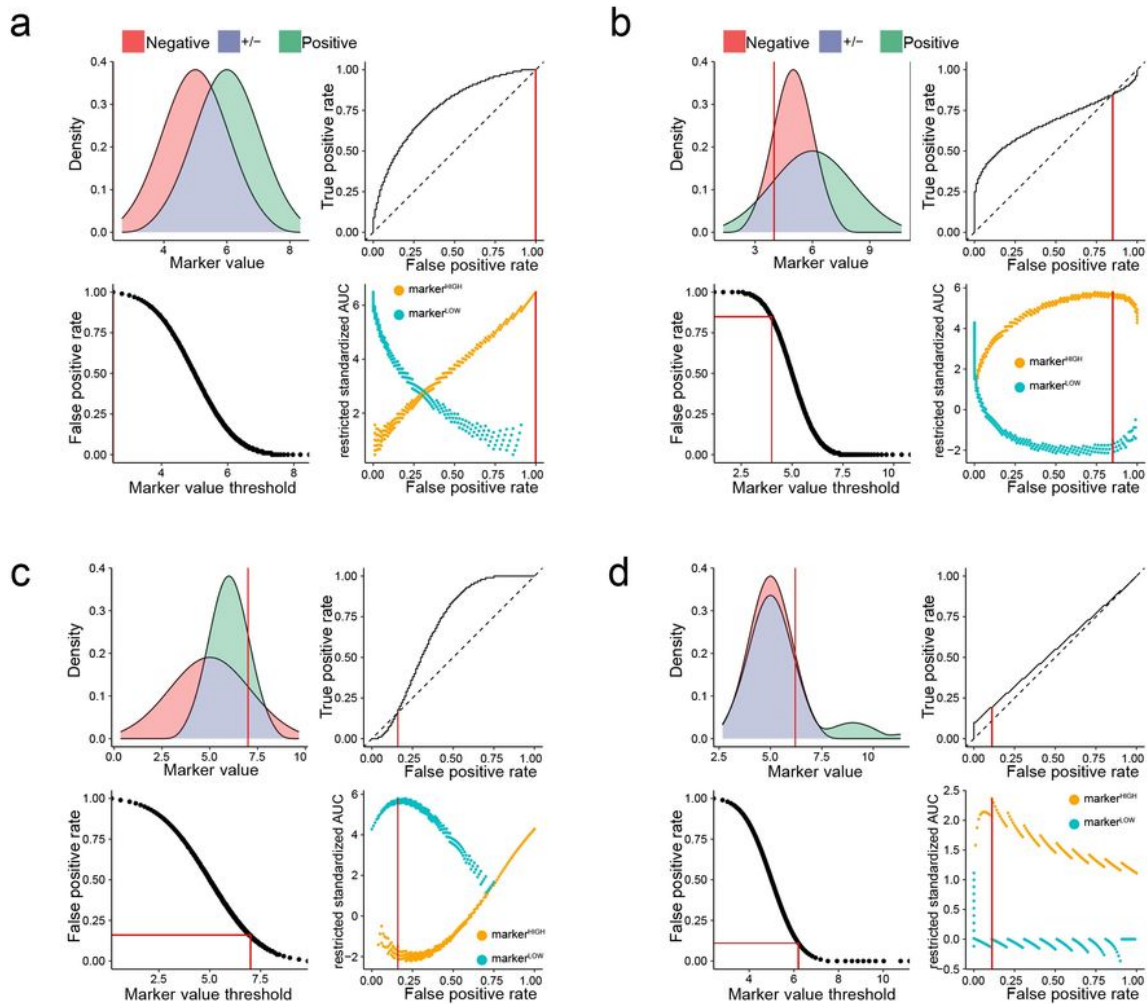


Figure 3

**FIGURE 3.** Optimal restriction of two-class distributions results in asymmetric ROC curves.

We present four simulated examples of marker distributions in two classes, which are intended to represent sets of patients with different clinical outcomes. The distribution of values from the positive (i.e. diseased) class are coloured green and values from the negative (i.e. control) class are coloured red; the overlapping density areas are coloured purple. For each example, we present the following: a plot of positive and negative class densities; the complete ROC curve; a plot of marker values against FPR; a plot of  $rAUC$  calculated for  $\text{marker}^{\text{HIGH}}$  (orange) and  $\text{marker}^{\text{LOW}}$  (blue) samples at all FPR values. In each plot, red lines indicate the optimal restriction as a marker value or FPR value. (a) A simulated example of a symmetric ROC curve from 100 negative  $N(5, 1)$  and 100 positive  $N(6, 1)$  samples. (b) A simulated example of a right-skewed ROC curve from 100 negative  $N(5, 1)$  and 100 positive  $N(6, 2)$  samples. (c) A simulated example of a left-skewed ROC curve from 100 negative  $N(5, 2)$  and 100 positive  $N(6, 1)$  samples. (d) Results for a right-skewed ROC curve from 100 negative  $N(5, 1)$  samples and 100 positive samples from a bimodally distributed positive population. In this example, the positive population comprises 10% cases with elevated marker expression  $N(9, 1)$  and 90% cases with unaltered marker expression  $N(5, 1)$ .

## Figure 3

See image above for figure legend.

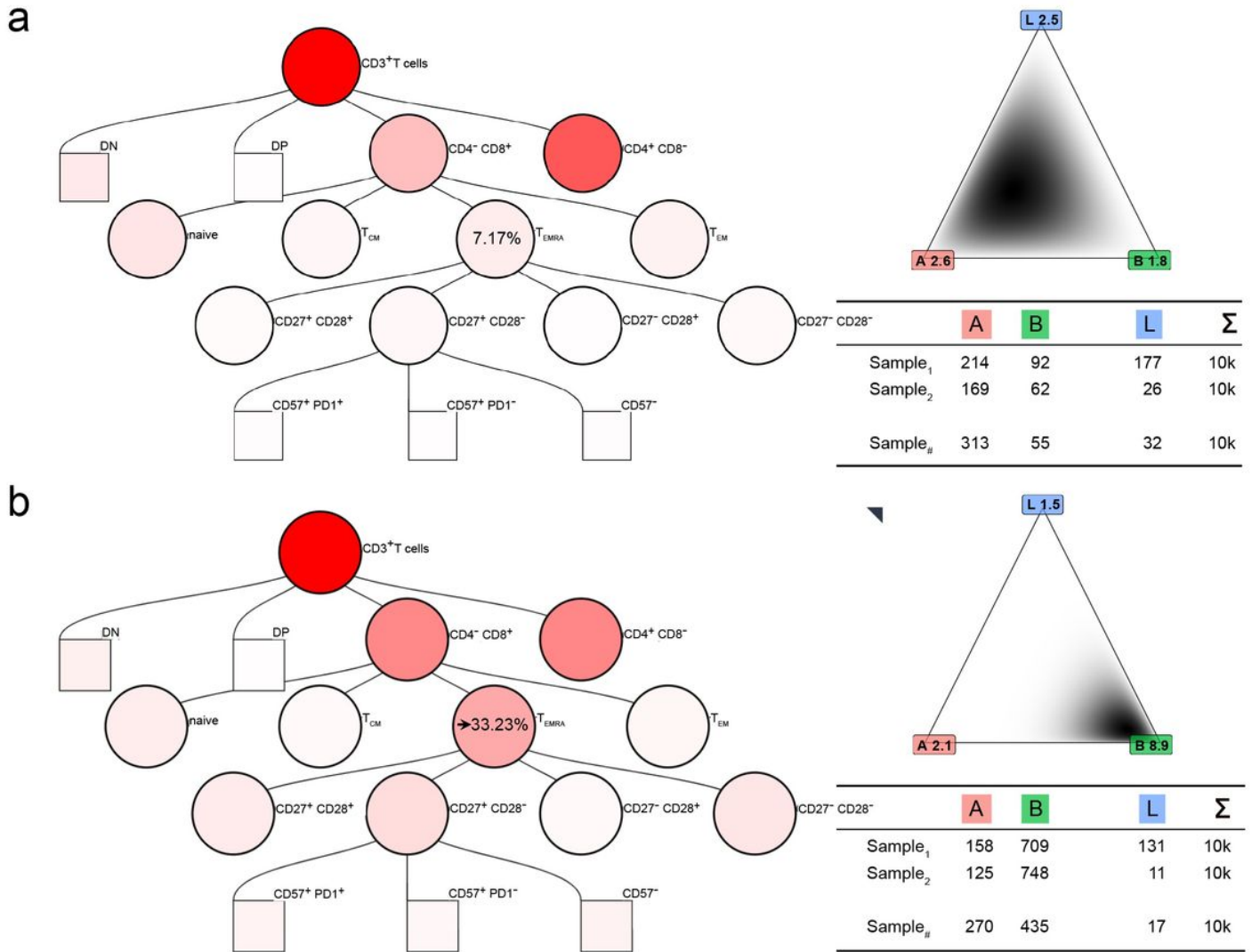


Figure 4

## Figure 4

*Realistically synthesizing flow cytometry data from two class distributions.* Various applications of our synthetic flow cytometry data depend upon generating samples with differences in cell subset distributions. Here, we provide an example of increasing the proportion of CD8<sup>+</sup> T<sub>EMRA</sub> cells in synthetic samples from a baseline value of 7.17 % in the negative class to an altered value of 33.23 % in the positive class. The intensity of red shading in the gating trees illustrates this change in CD8<sup>+</sup> T<sub>EMRA</sub> cells and contingent changes in other populations. **(a)** Gating tree with 7.17% CD8<sup>+</sup> T<sub>EMRA</sub> cells, its corresponding Dirichlet distribution parameter ( $\alpha = 12.34$ ) and the density for three example gates: A, B and L. The cell count table for three samples drawn from this distribution is shown. **(b)** Gating tree with 33.23% CD8<sup>+</sup> T<sub>EMRA</sub> cells, its corresponding Dirichlet distribution parameter ( $\alpha = 57.18$ ) and the density for three example gates: A, B and L. The cell count table for three samples drawn from this distribution is

shown. Of special note, percentages of cells in all other gates also changed according to the Dirichlet distribution, leading to changes in simulated cell counts across all leaf gates.

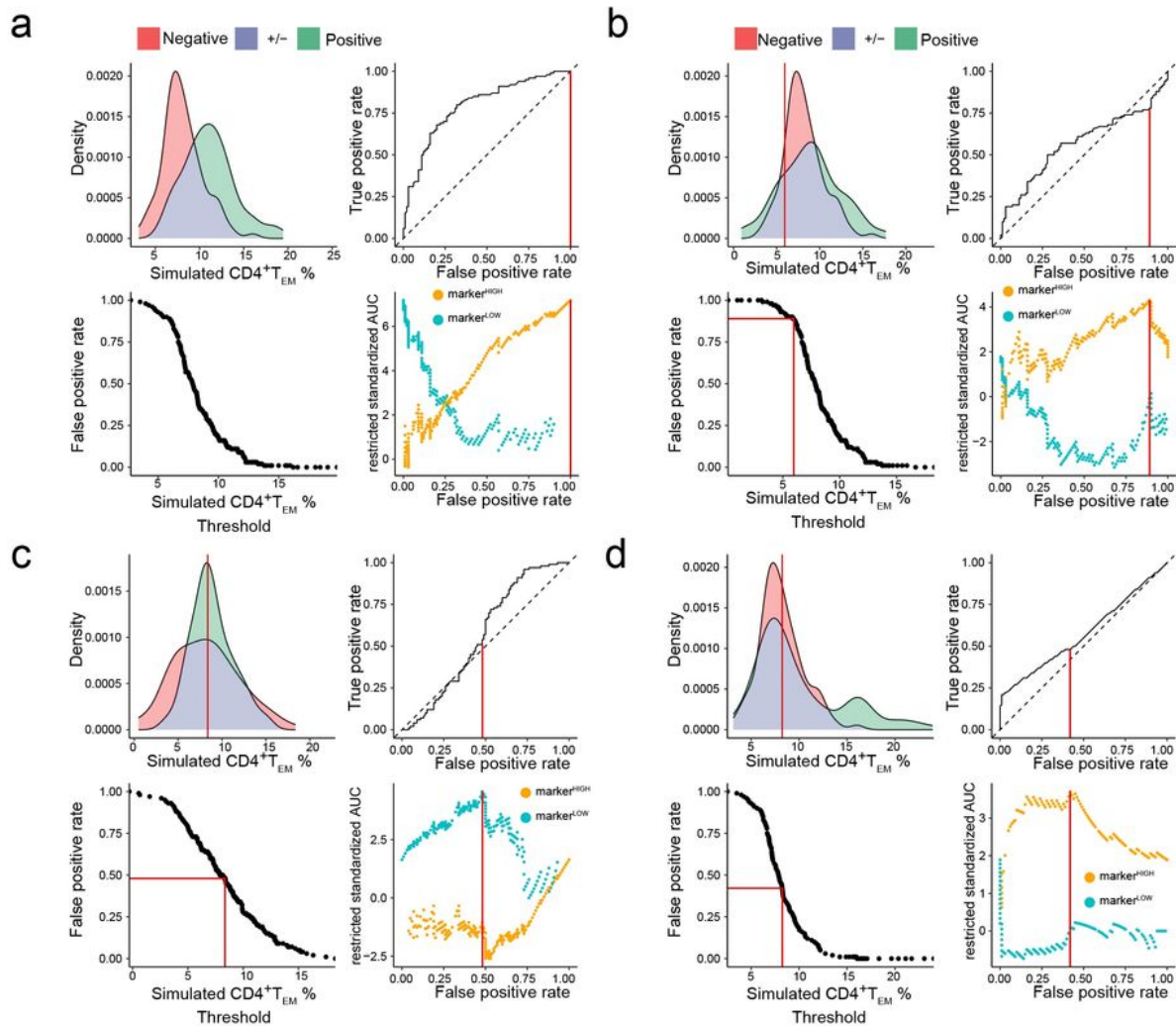


Figure 5

**FIGURE 5. Restriction of synthesized flow cytometry datasets.** We present examples of marker distributions in two classes, which are intended to represent sets of patients with different clinical outcomes. The distribution of values from positive (i.e. diseased) class are coloured green and values from negative (i.e. control) class are coloured red; overlapping density areas are coloured in purple. For each example, we present the following: a plot of positive and negative class densities; the complete ROC curve; a plot of marker values against FPR; a plot of  $rzAUC$  calculated for  $marker^{HIGH}$  (orange) and  $marker^{LOW}$  (blue) samples at all FPR values. In each plot, red lines indicate the optimal restriction as a marker value or FPR value. **(a)** A synthetic example of a symmetrical ROC curve from 100 negative  $N(7.7, 1)$  and 100 positive  $N(10.7, 1)$  samples. **(b)** A synthetic example of a right-skewed ROC curve from 100 negative  $N(7.7, 1)$  and 100 positive  $N(8.7, 3)$  samples. **(c)** A synthetic example of a left-skewed ROC curve from 100 negative  $N(7.7, 3)$  and 100 positive  $N(8.7, 1)$  samples. **(d)** Results for a synthetic right-skewed ROC curve from 100 negative  $N(7.7, 1)$  samples and 100 positive samples from a bimodally distributed positive population. In this example, the positive population comprises 20% cases with elevated marker expression  $N(16.7, 1)$  and 80% cases with unaltered marker expression  $N(7.7, 1)$ .

## Figure 5

See image above for figure legend.

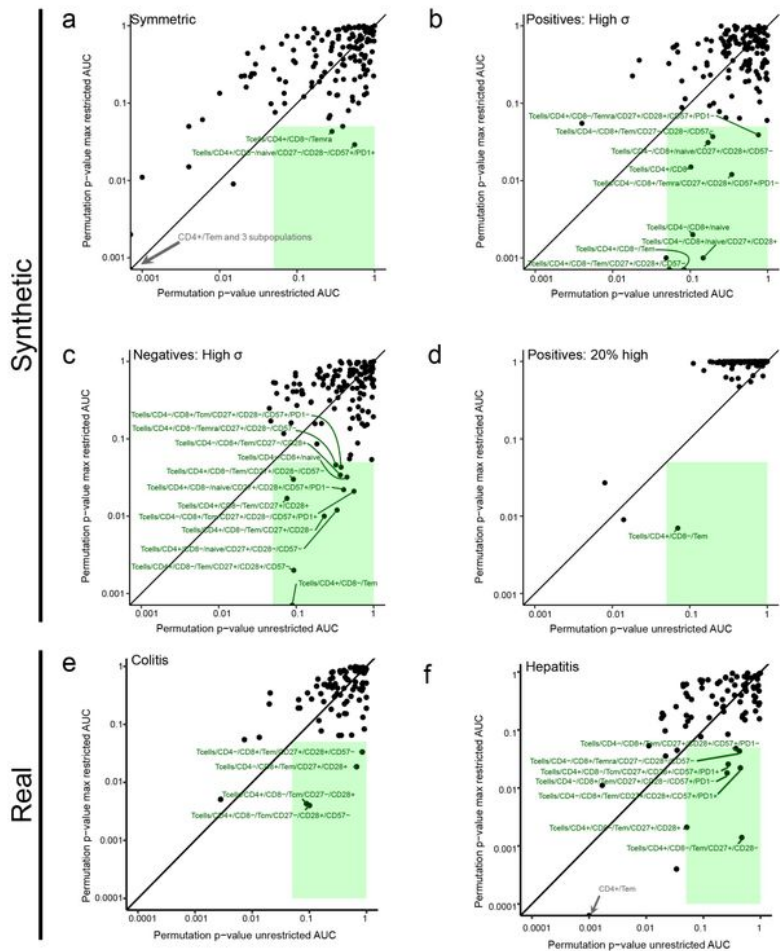


Figure 6

**FIGURE 6. Distinct sub-marker discovery in synthetic and real-world flow cytometry datasets.** Using synthetic and real-world datasets, we demonstrate that distinct restriction ensembles discover different marker sets with discriminatory capacity over a limited range of values. In each plot, the x-axis shows permutation p-values for the AUC of complete BOC curves for every gated cell population, the y-axis shows permutation p-values for the AUC of optimally restricted BOC curves. Points within the green-shaded rectangles represent cell subsets for which p-values derived from unrestricted data are not significant ( $p > 0.05$ ), but p-values derived from optimally restricted data are significant ( $p < 0.05$ ). Normal column filling in the green-shaded area are discovered markers discovered exclusively through distinct restriction. (a) Permutation p-values from synthetic samples in which a disease-related effect was introduced into CD4<sup>+</sup>Tem resulting in a symmetric BOC curve. 100 samples in the negative  $N(7, 1)$  class and 100 samples in the positive  $N(10, 2)$  class were generated. (b) Permutation p-values from synthetic samples in which a disease-related effect was introduced into CD4<sup>+</sup>Tem resulting in a right-skewed BOC curve. 100 samples in the negative  $N(7, 1)$  class and 100 samples in the positive  $N(8, 2)$  class were generated. (c) Permutation p-values from synthetic samples in which a disease-related effect was introduced into CD4<sup>+</sup>Tem resulting in a left-skewed BOC curve. 100 samples in the negative  $N(7, 1)$  class and 100 samples in the positive  $N(8, 2)$  class were generated. (d) Permutation p-values from synthetic samples in which a disease-related effect was introduced into CD4<sup>+</sup>Tem resulting in a bimodally distributed positive class was generated. In this example, the positive population comprises 20% cases with elevated marker expression  $N(16, 1)$  and 80% cases with standard marker expression  $N(7, 1)$ . (e) Permutation p-values from a training set of real-world classical flow cytometry samples ( $n=115$ ). 84 markers were selected where  $\geq 10\%$  of samples had more than 10 counts. Distinct restriction reveals 4 previously undetected markers of treatment-related colitis risk in autoimmune autoimmune patients receiving Ipilimumab therapy. (f) Permutation p-values from a training set of real-world classical flow cytometry samples ( $n=110$ ). 84 markers were selected where  $\geq 10\%$  of samples had more than 10 counts. Distinct restriction reveals 7 previously undetected markers of treatment-related hepatitis risk in autoimmune autoimmune patients receiving Ipilimumab therapy.

Figure 6

See image above for figure legend.



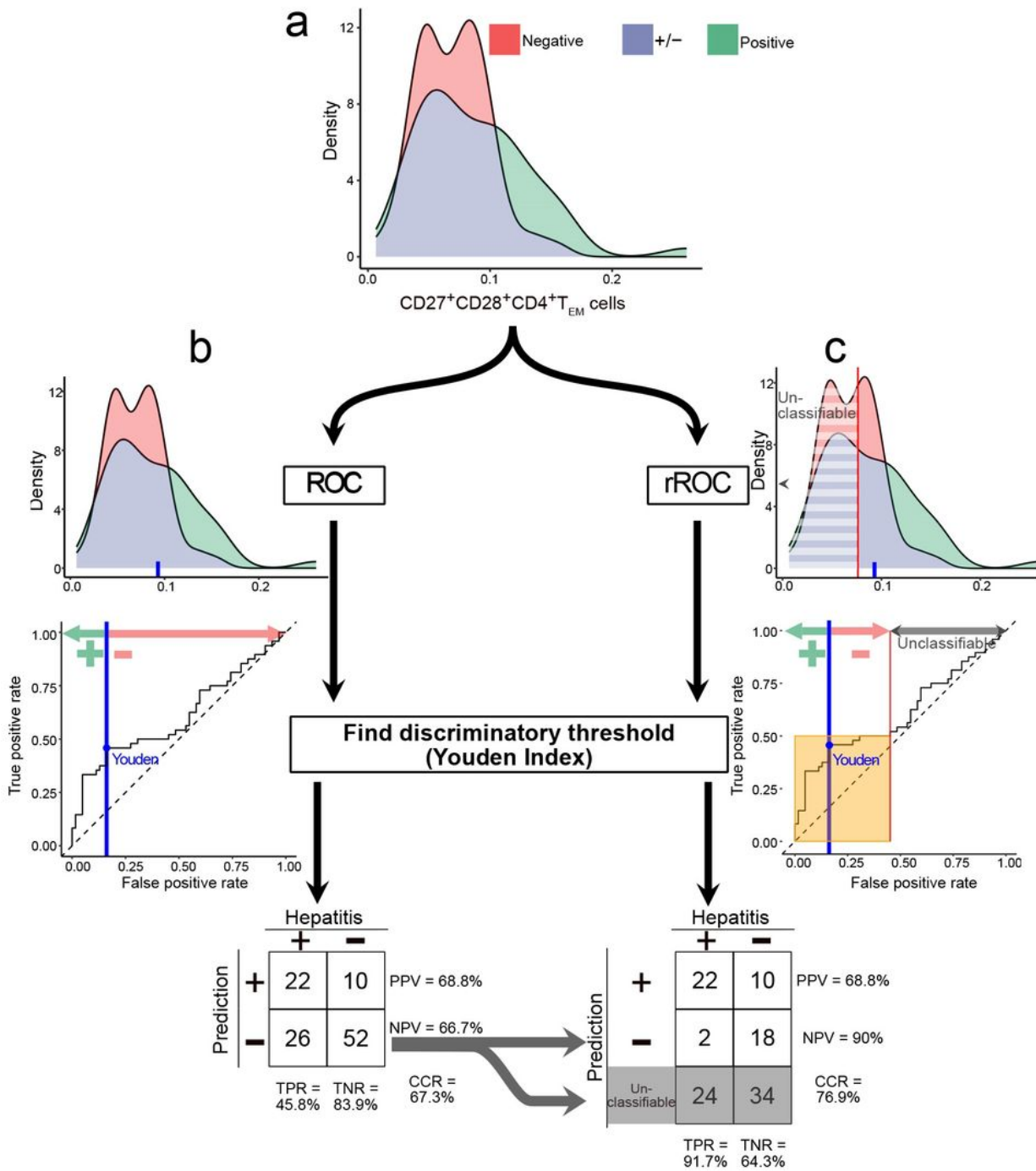


Figure 7

## Figure 7

*Clinical interpretation of restricted markers in predicting disease.* Our method of dataset restriction leads to counterintuitive clinical interpretations of marker values. This is illustrated by our discovery of CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cells as a univariate marker of hepatitis risk after immunotherapy. Here, we illustrate the conventional evaluation of marker performance across all samples with evaluation of marker

performance in a restricted dataset. **(a)** Densities of CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cells in all samples from patients with metastatic melanoma who developed hepatitis (n = 48) or did not (n = 62) after starting Ipi-Nivo therapy. **(b)** Following the classical approach of determining a classification cut-off for CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> frequency relative to CD4<sup>+</sup> T cells using the Youden Index, we predict hepatitis if > 9.62% and then assess the correct classification rate (CCR), negative predictive value (NPV), positive predictive value (PPV), sensitivity (or true positive rate, TPR) and specificity (or true negative rate, TNR) for all samples. **(c)** Our restriction method is predicated on there being a range of values over which a marker provides no discriminatory information. Optimally restricting CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> cell values leads us to discard 58 of 110 samples as “unclassifiable.” For the remaining 42 samples where CD27<sup>+</sup> CD28<sup>+</sup> CD4<sup>+</sup> T<sub>EM</sub> frequency relative to CD4<sup>+</sup> T cells > 7.62%, we determine a classification cut-off using the Youden Index, again predicting hepatitis if > 9.62%. Accordingly, we obtain a confusion table with CCR = 76.9%, specificity = 64.3% sensitivity = 91.7%, PPV = 68.8% and NPV = 90% across the classifiable samples.

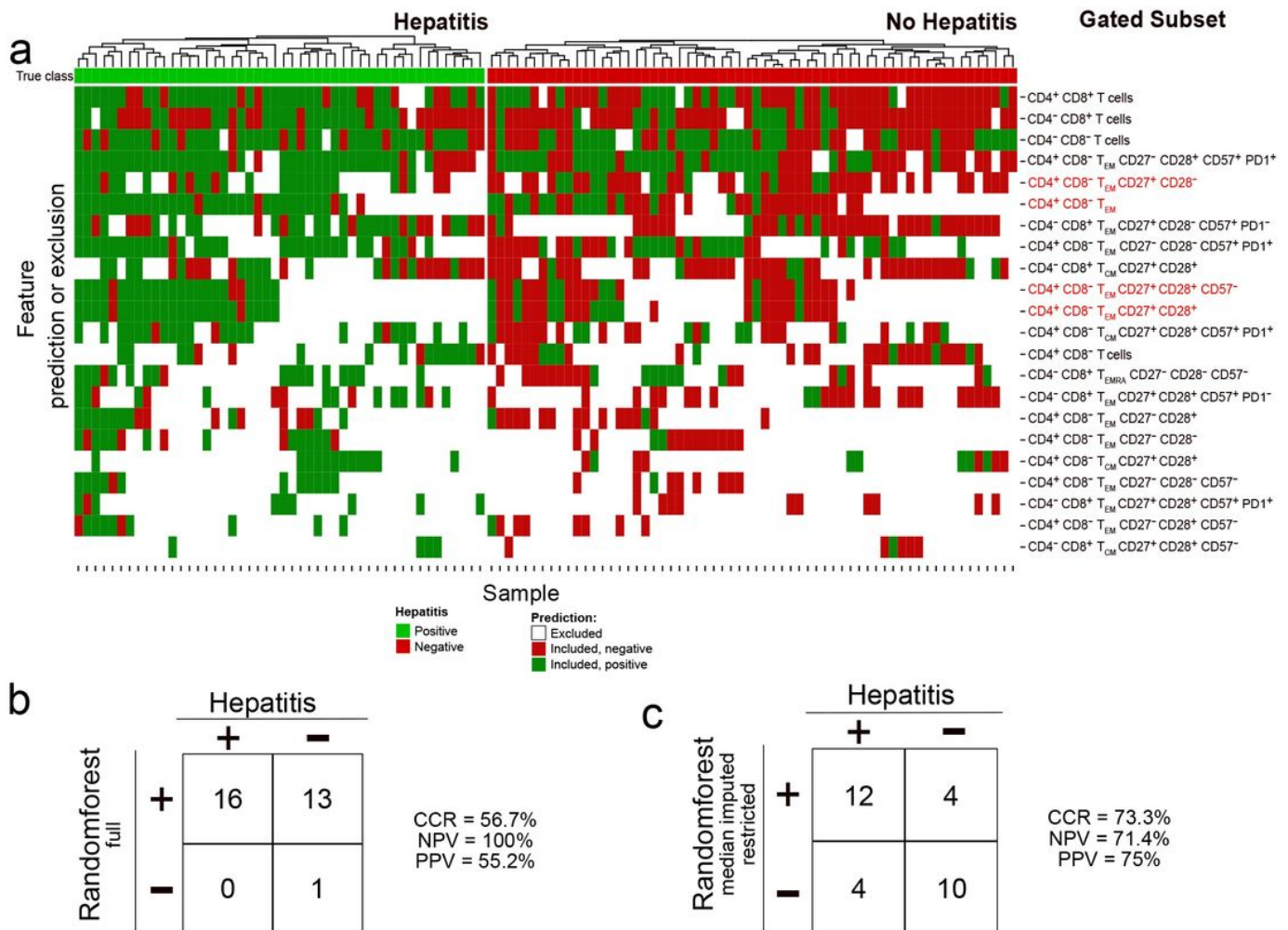


Figure 8

## Figure 8

*Significant markers according to classical and our analysis predicting hepatitis.* **(a)** Heatmap showing significant markers of hepatitis risk after Ipi-Nivo therapy based upon permutation p-values for unrestricted and restricted AUC. Only permutation p-values for markers discovered using our restriction method remained significant after correction for multiple testing, indicated by marker names in red text. Each further column reflects one sample, each row a feature. The samples are grouped into patients who did (green) or did not (red) develop treatment-related hepatitis, shown in the very first row. The main matrix consists of three values: Those excluded according by restriction (white); those included and predicted positive (dark green); and those included and predicted negative (dark red). Columns were clustered, rows in increasing order according to the number of excluded samples. **(b)** Random forest predictions and performances on the prospective validation cohort (n=30) trained on unrestricted marker values from the 110 training samples. **(c)** Random forest model predictions and performances on the prospective validation cohort (n=30) trained on restricted marker values from the 110 training samples. In this case, restricted marker values were replaced with -1 before training the random forest.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Glehr2023ReportingSummaryFINAL.pdf](#)
- [Glehr2023SupplementaryMaterialFINAL.pdf](#)
- [Glehr2023SoftwarePolicyFINAL.pdf](#)
- [Supplementaryvideo1.gif](#)
- [Supplementaryvideo2.gif](#)
- [Supplementaryvideo3.gif](#)
- [Supplementaryvideo4.gif](#)