
This item was submitted to [Loughborough's Research Repository](#) by the author.
Items in Figshare are protected by copyright, with all rights reserved, unless otherwise indicated.

Systematic design, generation, and application of synthetic datasets for flow cytometry

PLEASE CITE THE PUBLISHED VERSION

<https://doi.org/10.5731/pdajpst.2021.012659>

PUBLISHER

Parenteral Drug Association, Inc.

VERSION

AM (Accepted Manuscript)

PUBLISHER STATEMENT

Reproduced with kind permission of the publisher

LICENCE

CC BY-NC-ND 4.0

REPOSITORY RECORD

Cheung, Melissa, Jonathan J Campbell, Robert J. Thomas, Julian Braybrook, and Jon Petzing. 2022.
"Systematic Design, Generation, and Application of Synthetic Datasets for Flow Cytometry". Loughborough University. <https://hdl.handle.net/2134/18420812.v1>.

Systematic design, generation, and application of synthetic datasets for flow cytometry

Melissa Cheung, Jonathan J Campbell, Robert J Thomas, et al.

PDA Journal of Pharmaceutical Science and Technology **2022**,
Access the most recent version at doi:[10.5731/pdajpst.2021.012659](https://doi.org/10.5731/pdajpst.2021.012659)

1

2

3

4 Systematic design, generation, and application of synthetic datasets for flow
5 cytometry

6

7 Melissa Cheung^{1*}, Jonathan J. Campbell², Robert J. Thomas¹, Julian Braybrook² and Jon
8 Petzing¹

9

10 ¹ Centre for Biological Engineering, Loughborough University, Loughborough,
11 Leicestershire, UK

12 ² National Measurement Laboratory, LGC, Teddington, UK.

13

14 * Corresponding author

15 Email: M.Cheung@lboro.ac.uk

16

17 **Abstract**

18 Application of synthetic datasets in training and validation of analysis tools have led to
19 improvements in many decision-making tasks in a range of domains from computer vision to
20 digital pathology. Synthetic datasets overcome the constraints of real-world datasets, namely
21 difficulties in collection and labelling, expense, time and privacy concerns.

22 In flow cytometry, real cell-based datasets are limited by properties such as size, number of
23 parameters, distance between cell populations and distributions, and are often focused on a
24 narrow range of disease or cell types. Researchers in some cases have designed these desired
25 properties into synthetic datasets, however operators have implemented them in inconsistent
26 approaches and there is a scarcity of publicly available, high-quality synthetic datasets.

27 In this research, we propose a method to systematically design and generate flow cytometry
28 synthetic datasets with highly controlled characteristics. We demonstrate the generation of
29 two-cluster synthetic datasets with specific degrees of separation between cell populations,
30 and of non-normal distributions with increasing levels of skewness and orientations of skew
31 pairs. We apply our synthetic datasets to test the performance of a popular automated cell
32 populations identification software, SPADE3, and define the region where the software
33 performance decreases as the clusters get closer together.

34 Application of the synthetic skewed dataset suggests the software is capable of processing
35 non-normal data. We calculate the classification accuracy of SPADE3 with robustness not
36 achievable with real-world datasets. Our approach aims to advance research towards
37 generation of high-quality synthetic flow cytometry datasets, and to increase their awareness
38 among the community.

39 The synthetic datasets can be utilised in benchmarking studies that critically evaluate cell
40 population identification tools and help illustrate potential digital platform inconsistencies.

41 These datasets have the potential to improve cell characterisation workflows that integrate
42 automated analysis in clinical diagnostics and cell therapy manufacturing.

43

44 Keywords: flow cytometry, synthetic datasets, clusters, separation, skew, accuracy,
45 repeatability.

46

47 **1 Introduction**

48 An important challenge in manufacturing of emerging therapies is the need to develop and
49 satisfy appropriate regulatory standards. To support advanced analytics, toolsets are required
50 to establish quality assurance, such as data reference sets for data analysis, presentation, and
51 interpretation (1).

52 Synthetic datasets are datasets generated by computer simulation, rather than collected
53 through real world observations or experiments. These datasets are often created from
54 mathematical models that approximate aspects of real-world data. Synthetic datasets can be
55 referred to as ‘simulated’, ‘artificial’, ‘mock’, ‘toy’, and more colloquially, ‘dummy data’.

56 In unsupervised machine learning, well known ‘toy’ datasets used to compare different
57 clustering algorithms include clusters in the shape of two rings, crescent moons, spirals, and
58 data with no structures (2). Sophisticated synthetic datasets include urban street images
59 applied to object detection for autonomous driving (3,4), and household objects images for
60 object detection in the field of robotic manipulation (5). In medical imaging fields, synthetic
61 datasets generated based on real images (e.g. magnetic resonance imaging (MRI),
62 mammography, and whole-slide histopathology datasets) have demonstrated utility within
63 computer-aided detection or computer-aided diagnosis systems, and may be useful for
64 educational purposes and in quality control (6–8).

65 Further strategies to generate artificial datasets have made use of data augmentation methods
66 (9). Similarly, in flow cytometry, cell subsets from real samples can be selected and
67 computationally mixed in a copy-and-paste strategy with other real or synthetic cell
68 populations, to create augmented and reprocessed ‘semi-synthetic’ datasets (10).

69

70 Real data are generally required during the development of computational analysis tools to
71 provide means for training and validation, as well as potential decision-making reasons. An
72 analysis tool implies any method that performs detection, recognition, identification,
73 classification, tracking, prediction, or any other function that enables subsequent decision
74 making. Real data also play an important role in benchmarking studies that evaluate the
75 performance of these tools.

76 The key advantage of real data is that the information contains true characteristics of a
77 biological system. Albeit biological systems are very large and have complex signalling
78 pathway interactions between multiple cellular and molecular components. Real data have
79 various limitations that necessitates the creation of synthetic datasets to overcome them.
80 Being mathematical models, synthetic datasets can reduce the complexity of real data to gain
81 insight into how they are processed by computational analysis tools that have hidden ‘black
82 box’ algorithms.

83 A disadvantage of synthetic datasets is the issue of how accurately they represent the real
84 data. The purpose of synthetic datasets is to simplify real data, which requires assumptions to
85 be made and boundary conditions to be set around it, hence it is impossible to produce a
86 faithful replication of the real data. Whilst it may not be necessary to capture all the
87 complexities and features of real data in the design of synthetic ones in order to establish their
88 utility and credibility, there is potentially an expectation from users that higher complexity
89 equates to higher quality, leading to a lack of acceptance of synthetic data amongst the
90 biomedical community.

91 Often, real datasets with predetermined criteria are difficult to collect because of limited
92 availability at certain conditions and time periods. Synthetic data can be designed to mirror
93 existing real data and further optimise the dataset by including rare cases and those at

94 extreme conditions, thereby enhancing the realistic range of features or parameters.

95 Additionally, a high level of control is potentially achievable with synthetic datasets, where
96 designers can quickly change one factor at a time or build up layers of complexity through
97 controlled addition of factors.

98 Once collected, a major shortcoming of real data is the laborious and time-consuming task of
99 labelling observations with meaningful information (e.g. healthy vs abnormal) performed by
100 experienced personnel. Synthetic datasets can be designed with the labels inherent in the data,
101 side-stepping this task. The desired property of the synthetic data is also known and can be
102 applied in performance assessment of analysis tools.

103 A drawback of large-scale real datasets is that they sometimes take a large amount of time to
104 acquire (particularly true concerning collection of rare events). An equivalent large-scale,
105 complex synthetic dataset may also require a large amount of computing time to generate,
106 however this problem is negated as computers become faster.

107 Further benefits of synthetic data are: the potential lower costs associated with use of a
108 modern computer rather than expensive technical equipment, reagents and raw materials; the
109 reproducibility of computer code; and the absence of personal data which means that the
110 processing of synthetic data does not have the same privacy concerns and legal compliance
111 requirements as that of real data (11).

112 In flow cytometry, synthetic datasets usually aim to mimic the properties of real cell
113 populations. The properties of these randomly generated datasets range from simple two-
114 dimensional datasets with four clusters (12), to up to 30 populations in 35 dimensions (13).
115 The statistical distributions of synthetic clusters vary from normal (Gaussian), to non-normal
116 generated from mixtures of several Gaussians, and skewed (14–16). Simulated background
117 noise also features (17,18). Prior synthetic work approaches, however, have not explored

118 other possible characteristics specifically such as distance between clusters (both standard
119 and rare), which is modelled in real data through the comparison of median fluorescence
120 intensities between a stained and an unstained population in terms of population widths or
121 standard deviations, in order to estimate the relative brightness of a fluorophore (19,20).
122 Moreover, in a somewhat fragmented space, there is reason to apply systematic design on
123 existing properties (such as the skewness of clusters) to optimise the coverage of
124 characteristics.

125 Evaluation of developers' own tools using internally generated synthetic datasets is
126 inherently biased, therefore external and independent testing is a prerequisite for software
127 credibility in the clinical and biomanufacturing communities. Benchmarking datasets are
128 used in independent studies to compare software performance, however, existing studies
129 performed have solely relied on experimental data toolsets and have not used synthetic
130 datasets (21,22). This may be related to a limited amount of synthetic datasets available
131 within public flow cytometry repositories for the community to use (23). Software
132 benchmarking studies holds similarities to other quality assurance methods, such as those
133 applied in proficiency testing defined in ISO 13528:2005 (24), and similar statistical methods
134 can be used to evaluate software output performances in an external and independent manner.
135 When using real datasets as the test material, determination of the software performance is
136 achieved through comparison of software results against an estimate of the true value. This
137 value is assigned through a choice between 1) formulation, 2) cellular certified reference
138 materials (of which very few exist for flow cytometry) (25), 3) manually gated analysis from
139 one expert, 4) consensus manual analysis results from a group of experts, or 5) consensus
140 values from participant results. Possible bias from the results of experts or participants
141 reduces the robustness of the test. Synthetic datasets can be used adjacent to certified
142 reference materials with potential benefit.

143

144 In this research we propose the use of synthetic datasets for benchmarking unsupervised
145 learning automated flow cytometry data analysis software of which there are a large array of
146 options available to the data analyst (26). We define a description of the data characteristics
147 of flow cytometry data and demonstrate two methods to generate highly controlled,
148 systematically designed synthetic datasets with different degrees of separation between
149 clusters, and different levels of skew. We illustrate the use of our synthetic datasets using an
150 exemplar software, SPADE3 (27), and present results that allow robust calculations of
151 performance metrics not possible with real cell data. This work starts to explore the role of
152 synthetic datasets as digital reference materials and standards, and the potential regulatory
153 implications as the biomedical and biomanufacturing fields move increasingly towards using
154 automated systems, machine learning and artificial intelligence techniques.

155

156 **2 Materials and Methods**

157 **2.1 Target characteristics for synthetic flow cytometry datasets**

158 We identified certain commonly recognised data characteristics or potential statistical
 159 attributes of flow cytometry data and put forward a strategy to control and modify these
 160 characteristics to create systematic scenarios for testing software (Table I). In this research
 161 we targeted the separation / overlap and the skew properties in our simulation studies because
 162 these had not been addressed in previous work and/or the designs had not been approached in
 163 a systematic way. In order to focus on these properties, non-target characteristics such as the
 164 number of clusters, number of datapoints, and number of dimensions were kept constant, and
 165 noise was excluded in our simulations.

166 **2.2 Description of the Separation Index**

167 The separation index (SI) is used throughout this research to define the distance between
 168 clusters. The SI measures the magnitude of the gap between a pair of clusters based on the
 169 upper and lower percentiles of the two clusters (28). In the one-dimensional example (Figure
 170 1), the SI can be summarised as Eq. 1:

$$SI = \frac{L_2(\alpha/2) - U_1(\alpha/2)}{U_2(\alpha/2) - L_1(\alpha/2)} \quad \text{Eq. 1}$$

171
 172 where $L_i(\alpha/2)$ and $U_i(\alpha/2)$ are the sample lower and upper ($\alpha/2$) quantiles of cluster i . The
 173 interpretation of the SI is relatively straight forward, the range is [-0.999, +0.999] with values
 174 approaching +1 indicating increasing separation, SI of 0 indicating clusters touching, and SI
 175 approaching -1 indicating total overlap. In practice, our working range for the SI was [-0.3,
 176 +0.3]. These limits were defined because at a SI of +0.3 clusters were already very well
 177 separated, and at a SI of -0.3 clusters appeared well overlapped or merged.

178

179 **2.3 Hardware and software**

180 Dataset generation and analysis was run on a 64-bit Windows 10 operating system with a
181 3.00 GHz processor and 64 GB of RAM. Computational tools used are listed in Table II.
182 Throughout this paper, we use regular type to refer to software or computing environments,
183 *italics* for packages, and `monospace` font to designate functions.

184 **2.4 Synthetic datasets**

185 The concept of creating artificial, computer-generated flow cytometry datasets is essentially
186 random number generation, with numbers typically drawn from a normal distribution. Other
187 probability distributions are available e.g. binomial, exponential, Poisson, Student's t, etc. If
188 flow cytometry data are considered as mixtures of subpopulations of a heterogenous sample,
189 then the generation of synthetic data is a process of creating a mixture of random clusters.

190

191 **2.4.1 Separation dataset generation**

192 We designed a library of two-cluster synthetic datasets in two dimensions with 1,000
193 datapoints per cluster as an exemplar size of cell populations in real flow cytometry data,
194 with different degrees of separation between neighbouring clusters ranging from well-
195 separated to merged. The datasets were prepared using the R *clusterGeneration* package, with
196 SI values ranging from [-0.3, +0.3] at 0.05 intervals. Nine random normally distributed
197 cluster replicates were generated at each SI value. Covariance matrices were randomly
198 generated from eigenvalues between 1 and 5 to give a variability in the diameter and shape of
199 clusters that is similar to those seen in real flow cytometry data. These parameters produced
200 clusters with known separation, but which were random in their elliptical shape attribute.
201 Datasets were converted to FCS3.1 format using the R package *flowCore*.

202

203 **2.4.2 Skew dataset generation**

204 We designed a library of two-cluster synthetic datasets in two dimensions with 1,000
205 datapoints per cluster, with different levels of skew and skew-direction pairs. Single skew
206 clusters were prepared with the package *sn*, of which the α parameter regulates asymmetry
207 (29). Likewise random cluster replicates were generated at each skew direction (left and
208 right) along the x-axis in addition to each level of skew input α values between 2.5 to 10, at
209 intervals of 2.5. Applying the skewing α parameter causes the diameter of the elliptical
210 cluster to reduce along the x-axis. To compensate for this, clusters were elongated to obtain a
211 pre-skew diameter using the package *rescale*. The skewness of the clusters before and after
212 rescaling were identical (measured using the package *psych*) determined by the asymmetry
213 around the mean remaining unchanged (Figure 2A). Two clusters were joined together, and
214 one cluster shifted further away from the other through vector arithmetic operations in R
215 (Figure 2B). The distance between two clusters was measured with the *clusterGeneration*
216 package, datasets with a SI value between -0.25 and -0.15 were selected for further
217 processing. Files were converted to FCS3.1 standard using *flowCore* and visualised within
218 FlowJo software.

219

220 **2.5 Real datasets**

221 All material was obtained with the approval of and in accordance with the respective Ethics
222 Committees of Loughborough University and LGC, and under jurisdiction of the Human
223 Tissue Authority.

224 **2.5.1 Real cell PBMC dataset 1**

225 Fresh whole blood from healthy donors (Cambridge Bioscience, UK) was processed using
226 Ficoll-Paque (Fisher) to isolate the buffy coat layer containing peripheral blood mononuclear

227 cells (PBMCs). Cells were single-stained separately with CD4-PerCP-Vio700, CD45RO-
228 APC-Vio770, and CCR7-VioBlue (all from Miltenyi Biotech). Data were acquired using BD
229 FACSCantoII cytometer equipped with 3 lasers (405nm/ 30mW, 488nm/ 20mW, 633nm/
230 17mW). 100,000 cell events were collected.

231 **2.5.2 Real cell PBMC dataset 2**

232 PBMCs (LGC, UK) were stained with CD3-BB515, CD4-BB700, CD45RA-BV786 (all from
233 BD Biosciences), and live/dead fixable aqua dead cell stain (Invitrogen). Data were acquired
234 using a BD LSRFortessa cell analyser equipped with four lasers (355nm /20mW, 405nm/
235 50mW, 488nm/ 50mW, 640nm/ 40mW). 200,000 cell events were collected. Single-stained
236 beads and fluorescence-minus-one controls were used to calculate compensation.

237

238 **2.6 SPADE3 analysis of synthetic datasets**

239 SPADE3 was run within Matlab R2019a. Each FCS file was run separately. User input
240 parameters that were selected were: overlapping markers used for SPADE tree = CH1, CH2;
241 ignore compensation; no transformation. All other settings were left as default values (local
242 density neighbourhood size = 5, local density approximation factor = 1.5, maximum
243 allowable cells = 50,000, outlier density = 1, target density = 20,000 cells, algorithm = K-
244 means, number of desired clusters = 100).

245

246 **2.7 Statistics and performance metrics**

247 Methods used for statistical analysis included the mean, standard deviation, coefficient of
248 variation, and metrics derived from the confusion matrix as shown in Table III.

249 The performances of software runs were calculated using Eq. 2 and Eq. 3:

$$\text{Absolute difference to reference count} = |A - B| \quad \text{Eq. 2}$$

$$\text{Population percentage difference (\%)} = \frac{|A - B|}{\text{Total population}} \times 100 \quad \text{Eq. 3}$$

250 where A is the software output count, and B is the reference value defined here as the known
 251 count of cluster 1 (1,000 events) which was designed inherently in the dataset.

252 As in binary classification (30), here a true positive (a ‘hit’) is defined as the correct SPADE3
 253 assignment of a target cell to its reference target population set during cluster generation.

254 Events in cluster 1 of the synthetic datasets were arbitrarily selected as the ‘target’ cases.

255 SPADE3 assignment of a non-target cell to its non-target population is a true negative.

256 Misclassification of a non-target cell to a target population is a false positive, and

257 misclassification of a target cell to a non-target population is a false negative (a ‘miss’). The

258 evaluation metrics calculated from the confusion matrix include the accuracy, precision,

259 recall and F1 measure, and are defined in Eq. 4 to Eq. 7 (31). We compared the individual

260 cell assignments to a cluster predicted from SPADE3 with the reference cell assignments,

261 using the R package *caret*.

262

$$\text{Accuracy} = \frac{\text{True positive} + \text{True negative}}{\text{Total population}} \quad \text{Eq. 4}$$

263

$$\text{Precision} = \frac{\text{True positive}}{\text{True positive} + \text{False positive}} \quad \text{Eq. 5}$$

264

$$\text{Recall} = \frac{\text{True positive}}{\text{True positive} + \text{False negative}} \quad \text{Eq. 6}$$

265

$$F1 = 2 \times \frac{\textit{Precision} \times \textit{Recall}}{\textit{Precision} + \textit{Recall}} \quad \text{Eq. 7}$$

266

267 **3 Discussion**

268 In flow cytometry, data typically contain cell populations that are positive or negative for a
269 marker of interest. The distance between the positive and negative cell populations is
270 variable, ranging from well resolved to merged. Multiple factors affect this separation,
271 including but not limited to biological attributes (the level of marker expression, affinity and
272 avidity of antibody binding, the number of antibody bound per cell) and assay variables
273 (antibody panel design and concentrations used in the staining process, fluorophore
274 brightness and dye stability, sensitivity and resolution of the detectors). This separation
275 directly impacts the accuracy and precision of manual data analysis, with significantly higher
276 technical variation seen in poorly resolved populations compared with clearly defined cell
277 populations in human peripheral blood (32).

278

279 **3.1 Distance between clusters**

280 To simulate flow cytometry data with different distances between a positive and a negative
281 population we generated a normally distributed synthetic two-cluster dataset with different
282 degrees of separation between clusters. For comparison, we measured the SI between two
283 clusters in a real cell dataset. The PBMCs dataset 1 contained negatively and positively
284 stained populations in each fluorescent channel. These subpopulations were separated using
285 the automated cell population identification software SPADE3 (27). Then the magnitude of
286 the gap between pairs of real cell clusters along each channel was measured using the
287 `sepIndex` function in the *clusterGeneration* package. We observed similar SI values
288 between the real-world positive and negatively stained cell populations and the synthetic
289 clusters, within the range of -0.3 and +0.2 (Figure 3). These results show our method defined
290 here for generating synthetic flow cytometry datasets is able to successfully simulate the
291 distance parameters seen between clusters in a real example of flow cytometry data.

292 Individual cluster statistics can potentially vary by a small amount, however at this stage of
293 the research this was not a focus of this study.

294

295 **3.2 Clusters with non-normal distributions**

296 The synthetic datasets generated in Section 3.1 contain clusters following a normal
297 distribution, visualised as symmetrical bell curves for univariate data or symmetrical circles
298 and ellipses in scatterplots for multivariate data. However, real flow cytometry data consist of
299 cell population clusters that follow a normal distribution as well as those that display non-
300 normal distributions. The exact distribution along a marker channel is difficult to predict and
301 may depend on the state of the cell along a differentiation pathway. For example, a stable
302 haematopoietic stem cell population may display a normal distribution of CD34+ expression
303 that transitions to a non-normal distribution during cell differentiation as CD34 expression
304 decreases (33).

305 Non-normal data are characterised by asymmetry around the sample mean. These cell
306 populations can display positive (right) skew or negative (left) skew. The skewness can be
307 estimated using the adjusted Fisher-Pearson coefficient of skewness (p value) (34), where a
308 normal distribution has a skewness value of $p = 0$, a positive skewness value indicates a tail
309 pointing to the right, and a negative skewness value indicates a tail pointing to the left. The
310 further away the value is from 0, the greater the skew and typically the longer the tail.

311 There are different strategies to generate synthetic flow cytometry datasets with skewed cell
312 populations. One method is to create multiple Gaussian distributions that can then be merged
313 together to form an overall distribution with the desired skew. This strategy has been used
314 previously to create synthetic data with non-convex shapes (15,16). This method may require
315 many rounds of trial and error. To avoid this shortcoming, here we developed and tested a

316 different method using the *sn* R package to generate random clusters with multivariate skew-
317 normal distributions (29).

318 The multivariate skew normal distribution extends the class of normal distribution (defined
319 through a mean vector and covariance matrix) by the addition of a skew parameter.

320 A comparison of the skewed clusters generated from computer simulation against real cell
321 populations from the PBMCs dataset 2 demonstrates that the synthetic and real cases are
322 comparable (Figure 4). This result shows that the simulated data we generated is a realistic
323 model of both positive and negative skew observed in real flow cytometry cell populations,
324 and therefore has biological relevance. Thus, the synthetic dataset can be reliably used to gain
325 understanding on how automated software responds to skewed flow cytometry data, with the
326 additional benefit of the ability to systematically control the strength of the skew as well as
327 the absolute cell number.

328 The strategy we devised to create a dataset with multiple skewed clusters was to generate
329 individual clusters in parts then combine them together to form one whole dataset. The gap
330 between the clusters can be controlled by shifting one cluster closer or further away to the
331 other through vector arithmetic operations, with the SI being measured after the clusters were
332 combined. With skewed clusters, a new level of complexity is introduced compared to
333 normally distributed clusters, because assuming the skew is introduced only in one parameter,
334 each cluster can be left-skewed or right-skewed. Thus, the possible permutations of pairs of
335 skewed clusters in a two-cluster dataset increases from one to three. In this paper we refer to
336 these combinations as: head-to-head, head-to-tail (this is the same orientation as tail-to-head),
337 and tail-to-tail (Figure 5).

338

339 **3.3 Application of synthetic datasets to a cell population identification software**

340 To demonstrate the efficacy of the synthetic datasets, they were passed through an exemplar
341 software, SPADE, in order to illustrate how synthetic data can reveal limitations of
342 automated software, and can provide a deeper understanding of the inner workings of ‘black
343 box’ algorithms in a way that real cell datasets are unable to.

344 SPADE, named for spanning-tree progression analysis for density-normalised events, is a
345 widely applied software package that uses automated down-sampling, clustering and
346 minimum spanning tree construction to aid analysis of high-dimensional flow cytometry data
347 (35). There are two versions of SPADE with different algorithms. The original SPADE1
348 applies a stochastic down-sampling algorithm paired with an agglomerative hierarchical
349 clustering algorithm that produces different outputs when run on the same data. This specific
350 issue of reproducibility in SPADE1 was subsequently resolved in SPADE3 by removing the
351 stochastic algorithms and replacing them with deterministic ones (27). In addition, a tree
352 partitioning function was introduced to assist interpretation of the outputs.

353 Here, we used the Matlab implementation of SPADE3 to process our synthetic datasets using
354 default parameters (as described in section 2.6). We used the auto tree partitioning tool to
355 split the spanning tree into two populations, then compared the population number to our
356 known reference value, which was 1,000 cells per cluster, or 50% of total cells events, for
357 both the separation and skew datasets.

358 For the separation dataset, the absolute difference in cell count of each cluster between the
359 software output and the reference value was calculated for each SI condition. The results
360 show that the accuracy and precision of SPADE3 decreased as the SI decreased from +0.3 to
361 -0.3, with performance deteriorating noticeably at a SI value of -0.2 and below (Figure 6).

362 These results were to be expected, because defining the boundary between one cluster and
363 another becomes progressively more difficult as clusters get closer together.

364 The benefit of applying the synthetic datasets to test software such as SPADE3 was the
365 ability to quantify for the first time the SI value where the software began to lose
366 performance. The high level of control in designing the gap between cell populations within
367 the synthetic datasets would have been very difficult to achieve with real cell data.
368 Furthermore, since the absolute counts and frequencies of each cell population was known in
369 the synthetic dataset, the evaluation of the software was based on robust absolute traceable
370 figures, and did not rely on comparison with a manually gated reference subpopulation count,
371 which has already been shown to be operator dependent and potentially biased (36,37).

372 In the design of the skew dataset, a constant SI value of -0.2 between clusters was chosen
373 because it fell in the critical region where the SPADE3 software began to deteriorate. The
374 skew dataset was processed through SPADE3, then the difference in cell population
375 percentage of the cluster between the SPADE3 output and the reference value was calculated
376 for each skew condition and cluster pair orientation. We found that, for each cluster pair
377 orientation, increasing the level of skewness in the clusters had no effect on the accuracy and
378 precision of SPADE3. However, at each level of skewness, SPADE3 was able to partition the
379 two clusters with improved performance when the orientation was tail-to-tail, followed by
380 head-to-tail and finally head-to-head (Figure 7).

381

382 This pattern of performance appeared to correlate with the density of points between the two
383 clusters. The skew dataset was planned with skewness and skew orientation among the
384 variable design factors, and the separation between clusters as constant factors. The
385 systematic way this skew dataset was designed allowed for the pattern of behaviour of

386 SPADE3 to become apparent. This finding suggests the SPADE3 algorithm is well suited to
387 analysis of skewed data albeit with a performance bias and sensitivity depending on skewed
388 cluster orientation. This may not be the case for other algorithms that use different clustering
389 techniques, in particular those that use a Gaussian mixture model-based clustering approach.
390 Further work to investigate this in a software comparison study is warranted.

391

392 **3.4 Assessment of software performance based on synthetic data**

393 One of the benefits of synthetic data is that, as well as ‘true’ population counts and
394 frequencies, an estimate of the true membership of a cell to its cluster is known *a priori*. This
395 is not the case with real cell data, where membership of a cell to a population is estimated by
396 an analyst performing manual gating. Here, we demonstrate the evaluation of robust
397 performance metrics of SPADE3 runs on synthetic datasets using confusion matrix analysis.
398 Each cell event in the synthetic dataset was pre-assigned a cluster membership on generation.
399 These cluster memberships were withheld for the SPADE3 analyses. After running the
400 datasets through SPADE3, the software predictions of cluster memberships for all 2,000 cell
401 events were compared with the reference cluster memberships using the R *caret* package.
402 Events in cluster 1 were arbitrarily assigned as positive cases.

403 The results from the SPADE3 analysis of the synthetic separation dataset (Table IV) showed
404 a classification accuracy (Eq. 4) greater than 90% with SI values of -0.1 or greater. Accuracy
405 fell to 86% and 43% at SI values of -0.2 and -0.3 respectively (Figure 8A).

406 The same pattern appeared with precision, also called positive predictive value (Eq. 5) with
407 values greater than 90% at SI values of -0.1 or larger, then falling to 81% and 50% at SI
408 values of -0.2 and -0.3 respectively (Figure 8B).

409 The recall metric that measures the rate of true positives identified (Eq. 6) gave perfect scores
410 of 100% at SI of +0.3, scored greater than 90% at SI values of -0.2 or greater, and
411 deteriorated to 46% with SI of -0.3 (Figure 8C).

412 The F1 score (Eq. 7) was calculated from the precision and accuracy to give the overall
413 accuracy of SPADE3. The F1 score remained above 90% for SI values greater than -0.1, then
414 fell to 85% at a SI value of -0.2 and reduced further to 47% at a SI value of -0.3 (Figure 8D).

415 The results from the classification analysis reinforce the finding that SPADE3 performs
416 strongly when clusters are well-separated, but the accuracy decreases as clusters approach a
417 SI value of -0.2 and falls to below 50% when processing overlapping clusters with SI values
418 of -0.3. Our application of this synthetic flow cytometry dataset in this instance has helped to
419 illustrate good performance characteristics and ranges, but also limitations of SPADE3 with
420 respect to cluster separation both with normal and skewed cluster probability distributions.

421

422 **4 Conclusion**

423

424 In this article we have introduced a systematic method of designing and generating synthetic
425 flow cytometry datasets, with specific focus on control of the distance between clusters and
426 the probability distributions of events within clusters. We applied our computer-generated
427 flow cytometry datasets to an automated data analysis software, SPADE3, and have shown
428 that the synthetic datasets are capable of critically assessing the quality of the software
429 outputs and hence the software performance. In addition, we have given an example of
430 quantifying performance assessment using synthetic datasets that is robust compared with
431 using real-world datasets.

432

433 The systematic approach we have implemented to produce flow cytometry datasets is straight
434 forward to execute computationally, but would be complicated to achieve experimentally due
435 to uncontrollable external sources of variation within real cell datasets, thus synthetic datasets
436 here overcome the limitations of acquiring real datasets. The synthetic datasets have the same
437 range of data properties as their biological equivalents and can serve as credible substitutes
438 for real flow cytometry datasets for the testing of automated cell population identification
439 software.

440

441 It is noted here that in most cases algorithms that underpin flow cytometry analysis software
442 have been previously published, but it is (at times) inherently opaque implementation of
443 algorithms in executable code that whilst allows understanding of inputs and outputs, does
444 not allow a full understanding of the data transfer functions. Application of our synthetic
445 datasets to an automated cell population identification software such as SPADE3 can

446 therefore help users understand how the underlying ‘black box’ algorithm works. Here we
447 have identified the regions where SPADE3 begins to lose performance, specifically where
448 two clusters are located at a SI value of -0.2 or less. Our results suggest that SPADE3 is not
449 specifically affected by the probability distributions of data, but is more sensitive to the
450 relative density of data points between two clusters. Findings such as these can provide
451 guidance to users on software selection when having to contend with large array of potential
452 software solutions (26) (i.e. in this exemplar whether SPADE3 would be an appropriate tool
453 of choice for automated analysis of real data containing heavy overlapping of clusters), and
454 then help to understand the validity of their automated analysis outputs.

455

456 A further benefit of synthetic datasets was apparent when we assessed the performance of
457 SPADE3 using the metrics calculated from the confusion matrix. This classification analysis
458 relied on comparison of software predictions with known ‘true’ conditions hence an absolute
459 analysis. In the synthetic dataset, an estimate of the true assignment of cells to clusters was
460 designed into the data, making this analysis relatively rapid and robust. With real data, the
461 assignment of cells to subpopulations must first be manually determined (often with potential
462 difficulty and error), then the dataset labelled, before a confusion matrix can be calculated.
463 These additional steps are time-consuming and prone to error. The variability observed in
464 manually gated datasets means either the analysis from a single expert must be taken as the
465 best estimate of the ‘true value’, or a pooled manual analysis from a group of experts is used.
466 The first option risks bias, and the second is dependent on the precision and accuracy of the
467 group. In both instances, it is difficult for the final analysis to be as robust compared with
468 synthetic datasets.

469 There are also a few disadvantages with synthetic datasets. The datasets are built on
470 assumptions of real data. These approximations are based on mathematical models and have
471 limits. Although the aim is to create synthetic datasets that are as realistic as possible, there
472 may be features missing as a function of boundary conditions and design assumptions.

473

474 Further investigations on flow cytometry synthetic datasets will follow two main directions.
475 The first is the extension of the work to generate datasets with controls on other flow
476 cytometry data properties identified in Table I, and more complex datasets with multiple
477 controlled factors. We are currently developing and optimising synthetic datasets with rare
478 cell populations with and without skewed distributions, and noise characteristics. The second
479 area is the comparison of software performance across multiple platforms when challenged
480 with synthetic datasets targeting flow cytometry automated data analysis software that
481 employ various clustering algorithms such as K-means, hierarchical, partition, density-based,
482 model-based, spectral clustering and self-organising maps.

483

484 Besides being a benchmarking tool for software developers, possible further applications of
485 synthetic datasets include their use as educational and training tools for manual gating, as part
486 of external quality assessment (EQA) and proficiency testing schemes. In addition, as cell
487 identification and quantification in medical diagnostics and cell therapy/ regenerative
488 medicines manufacturing fields move increasingly towards automated machine learning and
489 artificial intelligence techniques, it is likely that synthetic datasets will have important
490 regulatory applications as digital reference materials and standards, as well as potential
491 regulatory implications.

492

493 **Acknowledgements**

494 The authors would like to thank Lydia Beeken and Dr Ben Diffey from the Centre for
495 Biological Engineering, Loughborough University, and Dr Shiqiu Xiong at LGC, for their
496 help with acquiring real flow cytometry datasets.

497

498 **Conflict of Interest Declaration**

499 The authors declare that they have no competing interests.

500

501 **References**

- 502 (1) Phillips, W.; Medcalf, N.; Dalgarno, K.; Makatoris, H.; Sharples, S.; Srai, J.; Hourd,
503 P.; Kapletia, D. *Redistributed Manufacturing in Healthcare: Creating New Value*
504 *through Disruptive Innovation*; 2019.
- 505 (2) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.;
506 Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.;
507 Courneau, D.; Brucher, M.; Perrot, M.; Duchesnay, É. Scikit-Learn: Machine
508 Learning in Python. *J. Mach. Learn. Res.* **2011**, *12* (Oct), 2825–2830.
- 509 (3) Ros, G.; Sellart, L.; Materzynska, J.; Vazquez, D.; Lopez, A. M. The SYNTHIA
510 Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban
511 Scenes. *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* **2016**, 3234–
512 3243. <https://doi.org/10.1109/CVPR.2016.352>.
- 513 (4) Wrenninge, M.; Unger, J. Synscapes: A Photorealistic Synthetic Dataset for Street
514 Scene Parsing. *arXiv Preprint*. **2018**. <https://arxiv.org/abs/1810.08705> (accessed 2020-
515 11-20)
- 516 (5) Tremblay, J.; To, T.; Birchfield, S. Falling Things: A Synthetic Dataset for 3D Object
517 Detection and Pose Estimation. *IEEE Comput. Soc. Conf. Comput. Vis. Pattern*
518 *Recognit. Work.* **2018**, 2038–2041. <https://doi.org/10.1109/CVPRW.2018.00275>.
- 519 (6) Hagiwara, A.; Warntjes, M.; Hori, M.; Andica, C.; Nakazawa, M.; Kumamaru, K. K.;
520 Abe, O.; Aoki, S. SyMRI of the Brain: Rapid Quantification of Relaxation Rates and
521 Proton Density, with Synthetic MRI, Automatic Brain Segmentation, and Myelin
522 Measurement. *Invest. Radiol.* **2017**, *52* (10), 647–657.
523 <https://doi.org/10.1097/RLI.0000000000000365>.
- 524 (7) Ratanaprasatporn, L.; Chikarmane, S. A.; Giess, C. S. Strengths and Weaknesses of
525 Synthetic Mammography in Screening. *Radiographics* **2017**, *37* (7), 1913–1927.
526 <https://doi.org/10.1148/rg.2017170032>.
- 527 (8) Niazi, M. K. K.; Parwani, A. V.; Gurcan, M. N. Digital Pathology and Artificial
528 Intelligence. *Lancet Oncol.* **2019**, *20* (5), e253–e261. <https://doi.org/10.1016/S1470->
529 [2045\(19\)30154-8](https://doi.org/10.1016/S1470-2045(19)30154-8).
- 530 (9) Perez, L.; Wang, J. The Effectiveness of Data Augmentation in Image Classification
531 Using Deep Learning. *arXiv Preprint*. **2017**. <https://arxiv.org/abs/1712.04621>
532 (accessed 2021-01-05).
- 533 (10) Arvaniti, E.; Claassen, M. Sensitive Detection of Rare Disease-Associated Cell
534 Subsets via Representation Learning. *Nat. Commun.* **2017**, *8* (1), 1–10.
535 <https://doi.org/10.1038/ncomms14825>.
- 536 (11) The European Parliament and the Council of the European Union. Regulation (EU)
537 2016/679 of the European Parliament and of the Council of 27 April 2016 on the
538 Protection of Natural Persons with Regard to the Processing of Personal Data and on
539 the Free Movement of Such Data, and Repealing Directive 95/46/EC. *Off. J. Eur.*
540 *Union* **2016**, 1–88.
- 541 (12) Sugar, I. P.; Sealfon, S. C. Misty Mountain Clustering: Application to Fast
542 Unsupervised Flow Cytometry Gating. *BMC Bioinformatics* **2010**, *11* (1), 1–14.

- 543 <https://doi.org/10.1186/1471-2105-11-502>.
- 544 (13) Samusik, N.; Good, Z.; Spitzer, M. H.; Davis, K. L.; Nolan, G. P. Automated Mapping
545 of Phenotype Space with Single-Cell Data. *Nat. Methods* **2016**, *13* (6), 493–496.
546 <https://doi.org/10.1038/nmeth.3863>.
- 547 (14) Naim, I.; Datta, S.; Rebhahn, J.; Cavanaugh, J. S.; Mosmann, T. R.; Sharma, G.
548 SWIFT-Scalable Clustering for Automated Identification of Rare Cell Populations in
549 Large, High-Dimensional Flow Cytometry Datasets, Part 1: Algorithm Design. *Cytom.*
550 *Part A* **2014**, *85* (5), 408–421. <https://doi.org/10.1002/cyto.a.22446>.
- 551 (15) Pyne, S.; Hu, X.; Wang, K.; Rossin, E.; Lin, T.-I.; Maier, L. M.; Baecher-Allan, C.;
552 McLachlan, G. J.; Tamayo, P.; Hafler, D. A.; De Jager, P. L.; Mesirov, J. P.
553 Automated High-Dimensional Flow Cytometric Data Analysis. *Proc. Natl. Acad. Sci.*
554 **2009**, *106* (21), 8519–8524. <https://doi.org/10.1073/pnas.0903028106>.
- 555 (16) Ge, Y.; Sealfon, S. C. Flowpeaks: A Fast Unsupervised Clustering for Flow Cytometry
556 Data via K-Means and Density Peak Finding. *Bioinformatics* **2012**, *28* (15), 2052–
557 2058. <https://doi.org/10.1093/bioinformatics/bts300>.
- 558 (17) Zare, H.; Shooshtari, P.; Gupta, A.; Brinkman, R. R. Data Reduction for Spectral
559 Clustering to Analyze High Throughput Flow Cytometry Data. *BMC Bioinformatics*
560 **2010**, *11* (1), 1–16. <https://doi.org/10.1186/1471-2105-11-403>.
- 561 (18) Bendall, S. C.; Davis, K. L.; Amir, E. A. D.; Tadmor, M. D.; Simonds, E. F.; Chen, T.
562 J.; Shenfeld, D. K.; Nolan, G. P.; Pe’Er, D. Single-Cell Trajectory Detection Uncovers
563 Progression and Regulatory Coordination in Human B Cell Development. *Cell* **2014**,
564 *157* (3), 714–725. <https://doi.org/10.1016/j.cell.2014.04.005>.
- 565 (19) Bigos, M. Separation Index: An Easy-to-Use Metric for Evaluation of Different
566 Configurations on the Same Flow Cytometer. *Curr. Protoc. Cytom.* **2007**, *40* (1), 1–21.
567 <https://doi.org/10.1002/0471142956.cyo121s40>.
- 568 (20) Telford, W. G.; Babin, S. A.; Khorev, S. V.; Rowe, S. H. Green Fiber Lasers: An
569 Alternative to Traditional DPSS Green Lasers for Flow Cytometry. *Cytom. Part A*
570 **2009**, *75* (12), 1031–1039. <https://doi.org/10.1002/cyto.a.20790>.
- 571 (21) Aghaeepour, N.; Finak, G.; Hoos, H.; Mosmann, T. R.; Brinkman, R.; Gottardo, R.;
572 Scheuermann, R. H.; Gottardo, R.; Scheuermann, R. H. Critical Assessment of
573 Automated Flow Cytometry Data Analysis Techniques. *Nat. Methods* **2013**, *10* (3),
574 228–238. <https://doi.org/10.1038/nmeth.2365>.
- 575 (22) Weber, L. M.; Robinson, M. D. Comparison of Clustering Methods for High-
576 Dimensional Single-Cell Flow and Mass Cytometry Data. *Cytom. Part A* **2016**, *89*
577 (12), 1084–1096. <https://doi.org/10.1002/cyto.a.23030>.
- 578 (23) Spidlen, J.; Breuer, K.; Rosenberg, C.; Kotecha, N.; Brinkman, R. R. FlowRepository:
579 A Resource of Annotated Flow Cytometry Datasets Associated with Peer-Reviewed
580 Publications. *Cytom. Part A* **2012**, *81* (9), 727–731.
581 <https://doi.org/10.1002/cyto.a.22106>.
- 582 (24) International Organization for Standardization. *ISO 13528:2005 Statistical Methods*
583 *for Use in Proficiency Testing by Interlaboratory Comparison*; Geneva, 2005.
- 584 (25) Wang, L.; Abbasi, F.; Ornatsky, O.; Cole, K. D.; Misakian, M.; Gaigalas, A. K.; He,
585 H. J.; Marti, G. E.; Tanner, S.; Stebbings, R. Human CD4 + Lymphocytes for Antigen

- 586 Quantification: Characterization Using Conventional Flow Cytometry and Mass
587 Cytometry. *Cytom. Part A* **2012**, *81 A* (7), 567–575.
588 <https://doi.org/10.1002/cyto.a.22060>.
- 589 (26) Cheung, M.; Campbell, J. J.; Whitby, L.; Thomas, R. J.; Braybrook, J.; Petzing, J. N.
590 Current Trends in Flow Cytometry Automated Data Analysis Software. *Cytom. Part A*
591 **2021**, *99* (10), 1007–1021. <https://doi.org/10.1002/cyto.a.24320>.
- 592 (27) Qiu, P. Toward Deterministic and Semiautomated SPADE Analysis. *Cytom. Part A*
593 **2017**, *91* (3), 281–289. <https://doi.org/10.1002/cyto.a.23068>.
- 594 (28) Qiu, W.; Joe, H. Separation Index and Partial Membership for Clustering. *Comput.*
595 *Stat. Data Anal.* **2006**, *50* (3), 585–603. <https://doi.org/10.1016/j.csda.2004.09.009>.
- 596 (29) Azzalini, A.; Capitanio, A. Statistical Applications of the Multivariate Skew Normal
597 Distribution. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **1999**, *61* (3), 579–602.
598 <https://doi.org/10.1111/1467-9868.00194>.
- 599 (30) Fleiss, J. L.; Levin, B.; Paik, M. C. *Statistical Methods for Rates and Proportions*, 3rd
600 ed.; John Wiley & Sons: Hoboken, NJ, USA, 2003.
- 601 (31) Tharwat, A. Classification Assessment Methods. *Appl. Comput. Informatics* **2021**, *17*
602 (1), 168–192. <https://doi.org/10.1016/j.aci.2018.08.003>.
- 603 (32) Burel, J. G.; Qian, Y.; Lindestam Arlehamn, C.; Weiskopf, D.; Zapardiel-Gonzalo, J.;
604 Taplitz, R.; Gilman, R. H.; Saito, M.; de Silva, A. D.; Vijayanand, P.; Scheuermann, R.
605 H.; Sette, A.; Peters, B. An Integrated Workflow to Assess Technical and Biological
606 Variability of Cell Population Frequencies in Human Peripheral Blood by Flow
607 Cytometry. *J. Immunol.* **2017**, *198* (4), 1748–1758.
608 <https://doi.org/10.4049/jimmunol.1601750>.
- 609 (33) Salati, S.; Zini, R.; Bianchi, E.; Testa, A.; Mavilio, F.; Manfredini, R.; Ferrari, S. Role
610 of CD34 Antigen in Myeloid Differentiation of Human Hematopoietic Progenitor
611 Cells. *Stem Cells* **2008**, *26* (4), 950–959. <https://doi.org/10.1634/stemcells.2007-0597>.
- 612 (34) Joanes, D. N.; Gill, C. A. Comparing Measures of Sample Skewness and Kurtosis. *J.*
613 *R. Stat. Soc. Ser. D (The Stat.)* **1998**, *47* (1), 183–189. <https://doi.org/10.1111/1467-9884.00122>.
- 615 (35) Qiu, P.; Simonds, E. F.; Bendall, S. C.; Gibbs, K. D.; Bruggner, R. V.; Linderman, M.
616 D.; Sachs, K.; Nolan, G. P.; Plevritis, S. K. Extracting a Cellular Hierarchy from High-
617 Dimensional Cytometry Data with SPADE. *Nat. Biotechnol.* **2011**, *29* (10), 886–893.
618 <https://doi.org/10.1038/nbt.1991>.
- 619 (36) Grant, R.; Coopman, K.; Medcalf, N.; Silva-Gomes, S.; Campbell, J. J.; Kara, B.;
620 Braybrook, J.; Petzing, J. Understanding the Contribution of Operator Measurement
621 Variability within Flow Cytometry Data Analysis for Quality Control of Cell and Gene
622 Therapy Manufacturing. *Measurement* **2020**, *150*, 106998.
623 <https://doi.org/10.1016/j.measurement.2019.106998>.
- 624 (37) Grant, R.; Coopman, K.; Medcalf, N.; Silva-Gomes, S.; Campbell, J. J.; Kara, B.;
625 Braybrook, J.; Petzing, J. N. Quantifying Operator Subjectivity within Flow Cytometry
626 Data Analysis as a Source of Measurement Uncertainty and the Impact of Experience
627 on Results. *PDA J. Pharm. Sci. Technol.* **2021**, *75* (1), 33–47.
628 <https://doi.org/10.5731/pdajpst.2019.011213>.

629

630 **Tables**

631 Table I. Characteristics of flow cytometry datasets

Characteristic	Description
Number of clusters	Number of cell subpopulations in a sample
Number of datapoints	Number of cell events acquired from a sample
Number of dimensions	Number of parameters recorded in the experiment, e.g. forward scatter, side scatter, fluorescent markers
Separation	The gap between negatively and positively stained cell populations
Overlap	Poorly resolved populations that appear merged
Placement	Projection direction of one cluster to another in space
Distribution	The shape of the cell population, as modelled on probability distributions e.g. Gaussian, Student's t, exponential, Chi-squared
Spread	The variance of the cell population
Skew and kurtosis	The level of asymmetry around the mean of the cell cluster
Orientation	The direction of the asymmetry
Elongation	Stretched out populations with long tails
Noise	Events that are excluded from analysis e.g. outliers, dead cells, debris, doublets, false events detected in the region of interest.

632

633

634 Table II. Toolset used in this research for the generation of synthetic datasets, automated cell
635 population identification, and performance evaluation. R packages in italics.

Tool	Version	Purpose in this research
R	3.5.1	Programming
RStudio IDE	1.2	Programming environment
Matlab	R2019a	Environment for SPADE analysis
FlowJo	10.6	Flow cytometry data analysis and visualisation
SPADE	3	Automated analysis of synthetic datasets
<i>caret</i>	6.0-82	Calculate performance metrics, confusion matrix
<i>clusterGeneration</i>	1.3.4	Generate synthetic clusters
<i>flowCore</i>	1.48.1	Manipulate flow cytometry data
<i>psych</i>	1.8.12	Measure skew
<i>scales</i>	1.1.0	Scale functions for visualisation
<i>sn</i>	1.5-3	Build and manipulate probability distributions of the skew-normal family
<i>tidyverse</i>	1.3.0	Data manipulation, analysis and visualisation

636

637

638

639

Table III. Confusion Matrix

		Reference	
		Target	Non-target
Predicted	Target	True positive	False positive
	Non-target	False negative	True negative

640

641

642

Table IV. Performance metrics.

Separation index	Mean Accuracy (%)	Mean Precision (%)	Mean Recall (%)	Mean F1 (%)
-0.3	43.0	49.8	45.8	47.4
-0.2	86.1	80.5	91.6	85.1
-0.1	93.9	93.5	94.4	93.9
0	97.0	97.0	97.2	97.0
0.1	98.6	98.7	98.5	98.6
0.2	99.7	99.8	99.7	99.7
0.3	99.98	99.97	100.00	99.98

643

644

645 **Figure captions**

646 Figure 1. One-dimensional example of the separation index (SI) that measures the magnitude of the
 647 gap between two clusters. Vertical lines indicate the lower and upper quantiles of the clusters. The
 648 difference between U1 and L2 (numerator) is divided by the difference between L1 and U2
 649 (denominator) to calculate the SI value. This method is robust against outliers in between the two
 650 clusters that may affect the SI. Figure adapted from (28).

651 Figure 2. Workflow for skew dataset generation. Top panels show scatterplots, bottom panels show
 652 density estimates. A) A cluster with a normal distribution is generated, then skew is added through the
 653 alpha parameter in the R package *sn*, then the cluster is rescaled. B) Two clusters are combined, then
 654 the distances between them can be varied through vector arithmetic operations.

655 Figure 3. Comparison of synthetic and real clusters with representative separation index values
 656 ranging from -0.3, 0.0 to +0.2, showing overlapping, touching and well-separated clusters,
 657 respectively. Top panel shows synthetic data generated from R package *clusterGeneration*, bottom
 658 panel shows real peripheral blood mononuclear cells (PBMCs) data after automated population
 659 detection and partition with SPADE3 software followed by separation index calculation.

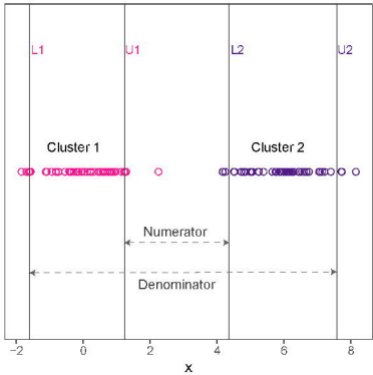
660 Figure 4. Comparison of synthetic (A) and real (B) flow cytometry data with skewed distributions.
 661 Both left-skewed and right-skewed synthetic clusters can be generated that mimic real data.
 662 Asymmetry around the mean is clearly shown with contour plots (top) and histograms (bottom).

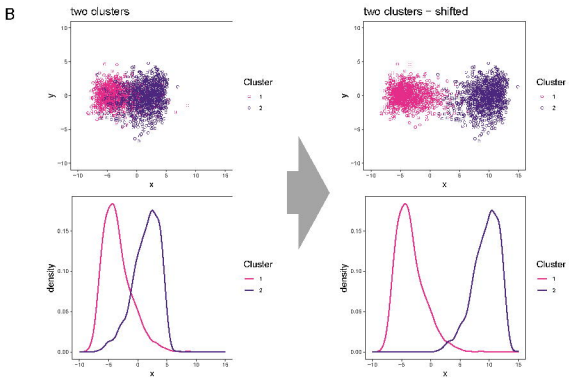
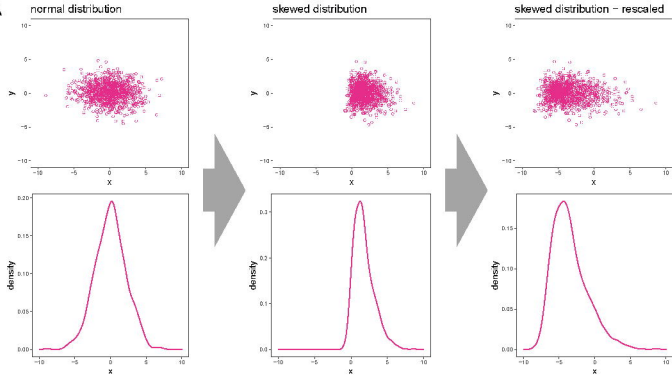
663 Figure 5. Generation of synthetic two-cluster skewed datasets. Three combinations of cluster pairs are
 664 shown here; head-to-head, head-to-tail and tail-to-tail.

665 Figure 6. The synthetic two-cluster separation dataset was run through SPADE3, then the absolute
 666 difference between the SPADE3 cell count to the reference value of cluster 1 was calculated. Result
 667 demonstrates the accuracy and repeatability of SPADE3 deteriorates as the distance between clusters
 668 decreases. Data represents mean $\pm 1SD$.

669 Figure 7. The synthetic two-cluster skewed dataset was run through SPADE3, then the gap between
 670 the SPADE3 cell population percentages and the reference cell population percentages was calculated.
 671 Increasing the cluster skewness in the datasets did not affect SPADE3 performance. However, the
 672 accuracy and repeatability of SPADE3 improved as direction of the cluster pairs changed from head-
 673 to-head to head-to-tail, and then to tail-to-tail. All clusters had separation index of -0.2. Data
 674 represents mean $\pm 1SD$.

675 Figure 8. Performance metrics for SPADE3 analysis of Separation dataset.





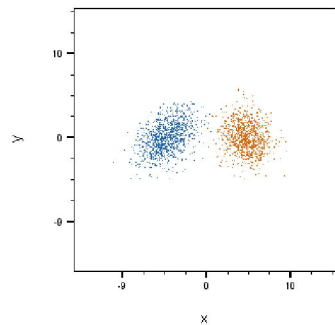
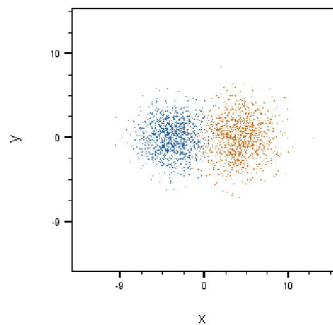
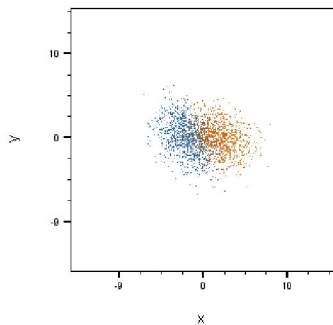
Separation index

-0.3

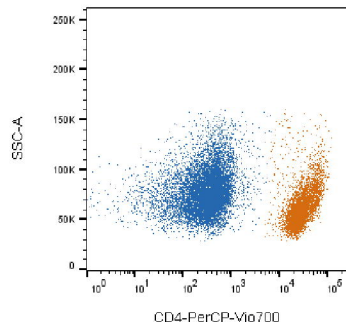
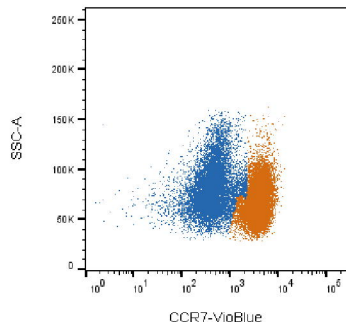
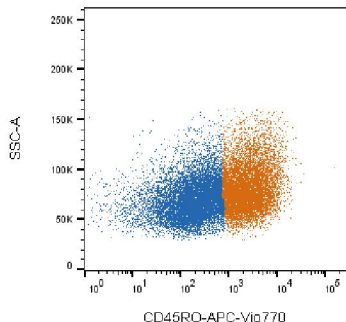
0.0

+0.2

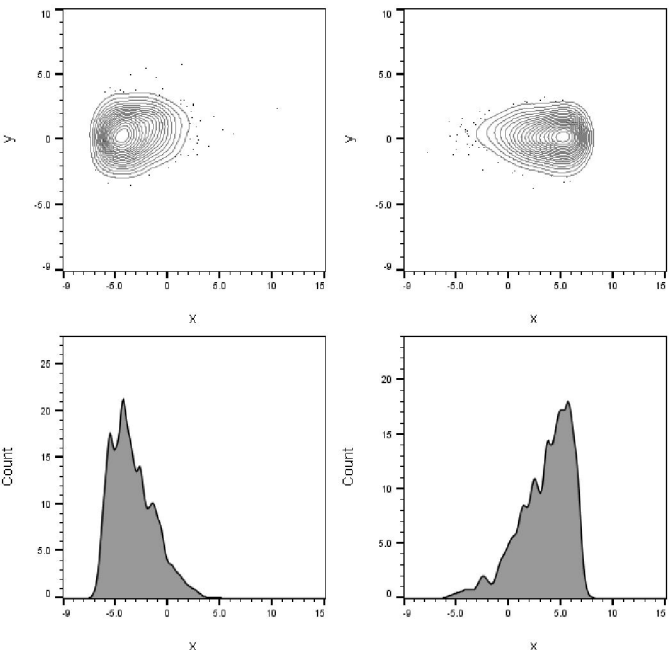
Synthetic data



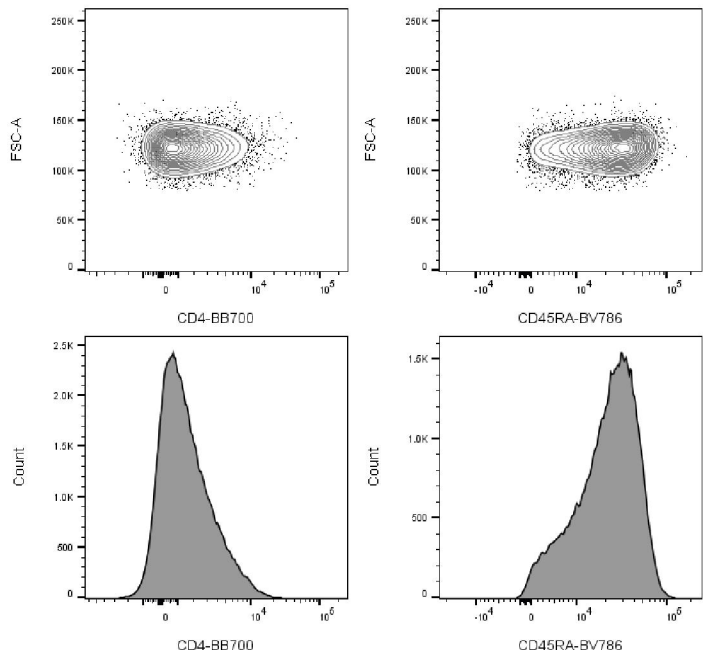
Real data

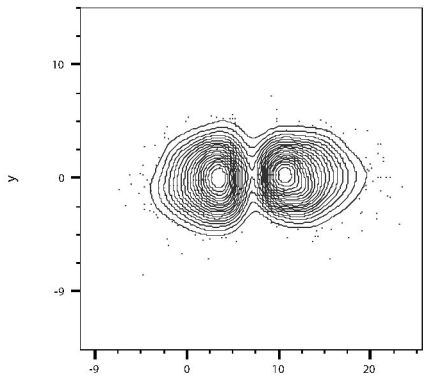


(A) Synthetic one-cluster

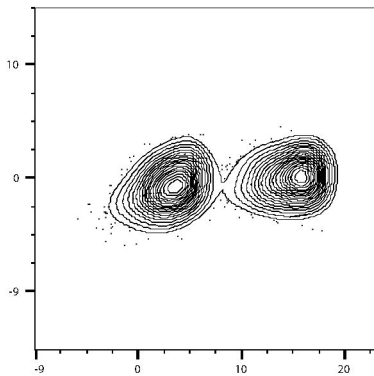


(B) Real one-cluster

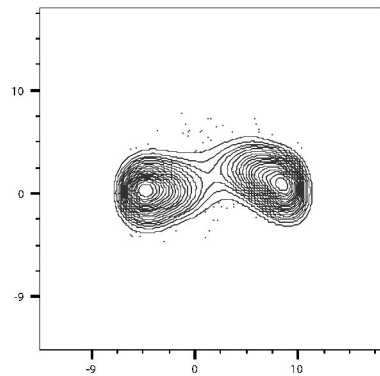




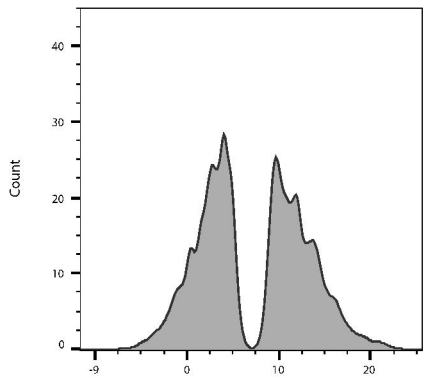
x



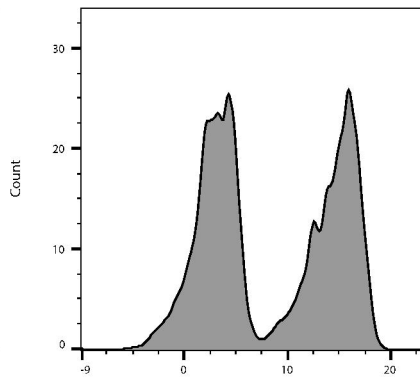
x



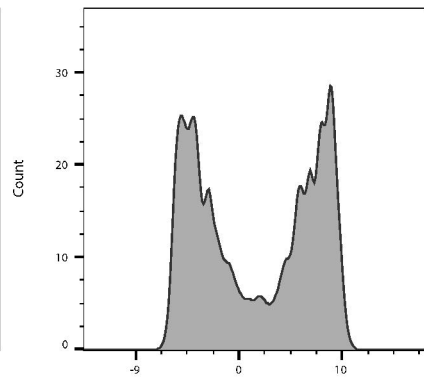
x



x

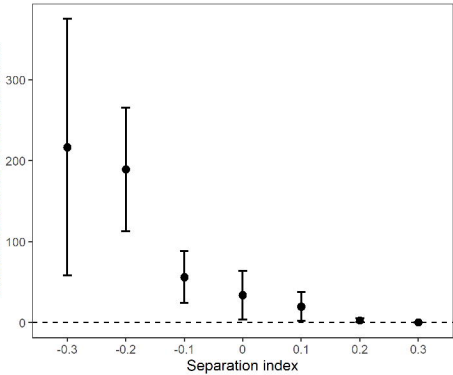


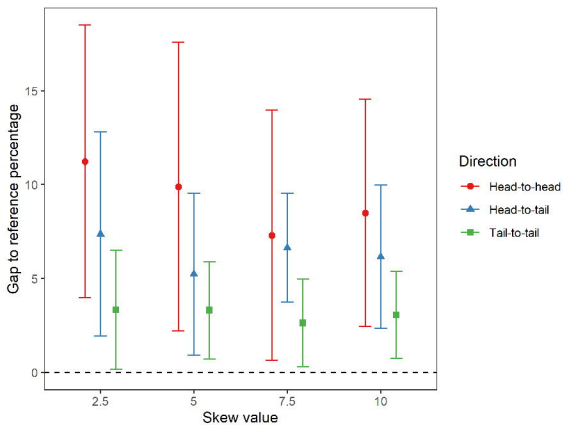
x

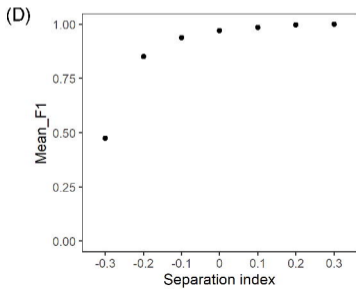
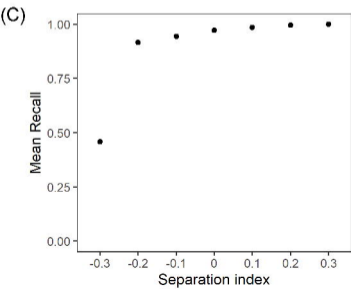
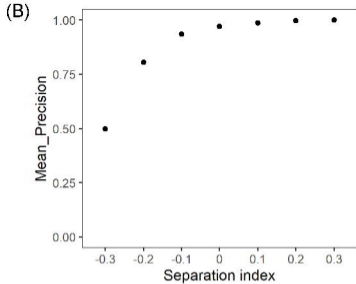
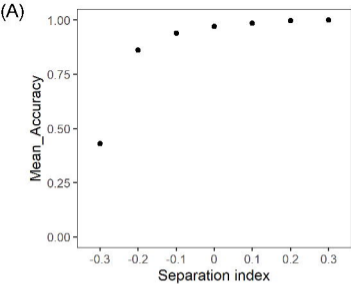


x

Absolute difference to reference count







PDA Journal of Pharmaceutical Science and Technology



An Authorized User of the electronic PDA Journal of Pharmaceutical Science and Technology (the PDA Journal) is a PDA Member in good standing. Authorized Users are permitted to do the following:

- Search and view the content of the PDA Journal
- Download a single article for the individual use of an Authorized User
- Assemble and distribute links that point to the PDA Journal
- Print individual articles from the PDA Journal for the individual use of an Authorized User
- Make a reasonable number of photocopies of a printed article for the individual use of an Authorized User or for the use by or distribution to other Authorized Users

Authorized Users are not permitted to do the following:

- Except as mentioned above, allow anyone other than an Authorized User to use or access the PDA Journal
- Display or otherwise make any information from the PDA Journal available to anyone other than an Authorized User
- Post articles from the PDA Journal on Web sites, either available on the Internet or an Intranet, or in any form of online publications
- Transmit electronically, via e-mail or any other file transfer protocols, any portion of the PDA Journal
- Create a searchable archive of any portion of the PDA Journal
- Use robots or intelligent agents to access, search and/or systematically download any portion of the PDA Journal
- Sell, re-sell, rent, lease, license, sublicense, assign or otherwise transfer the use of the PDA Journal or its content
- Use or copy the PDA Journal for document delivery, fee-for-service use, or bulk reproduction or distribution of materials in any form, or any substantially similar commercial purpose
- Alter, modify, repackage or adapt any portion of the PDA Journal
- Make any edits or derivative works with respect to any portion of the PDA Journal including any text or graphics
- Delete or remove in any form or format, including on a printed article or photocopy, any copyright information or notice contained in the PDA Journal