

# Getting started with flowStats

F. Hahne, N. Gopalakrishnan

May 1, 2024

## Abstract

`flowStats` is a collection of algorithms for the statistical analysis of flow cytometry data. So far, the focus is on automated gating and normalization.

## 1 Introduction

Since `flowStats` is more a collection of algorithms, writing a coherent Vignette is somewhat difficult. Instead, we will present a hypothetical data analysis process that also makes heavy use of the functionality provided by `flowWorkspace`, primarily through `GatingSets`.

We start by loading the `ITN` data set

```
> library(flowCore)
> library(flowStats)
> library(flowWorkspace)
> library(ggcyto)
> data(ITN)
```

The data was acquired from blood samples by 3 groups of patients, each group containing 5 samples. Each `flowFrame` includes, in addition to FSC and SSC, 5 fluorescence parameters: CD3, CD4, CD8, CD69 and HLADR.

First we need to transform all the fluorescence channels. This is a good point to start using a `GatingSet` object to keep track of our progress.

```
> require(scales)
> gs <- GatingSet(ITN)
> trans.func <- asinh
> inv.func <- sinh
> trans.obj <- trans_new("myAsinh", trans.func, inv.func)
> transList <- transformerList(colnames(ITN)[3:7], trans.obj)
> gs <- transform(gs, transList)
```

In an initial analysis step we first want to identify and subset all T-cells. This can be achieved by gating in the CD3 and SSC dimensions. However, there are several other sub-populations and we need to either specify our selection further or segment the individual sub-populations. One solution for the latter approach is to use the mixture modelling infrastructure provided by the `flowClust` package. However, since we are only interested in one single sub-population, the T-cell, it is much faster and easier to use

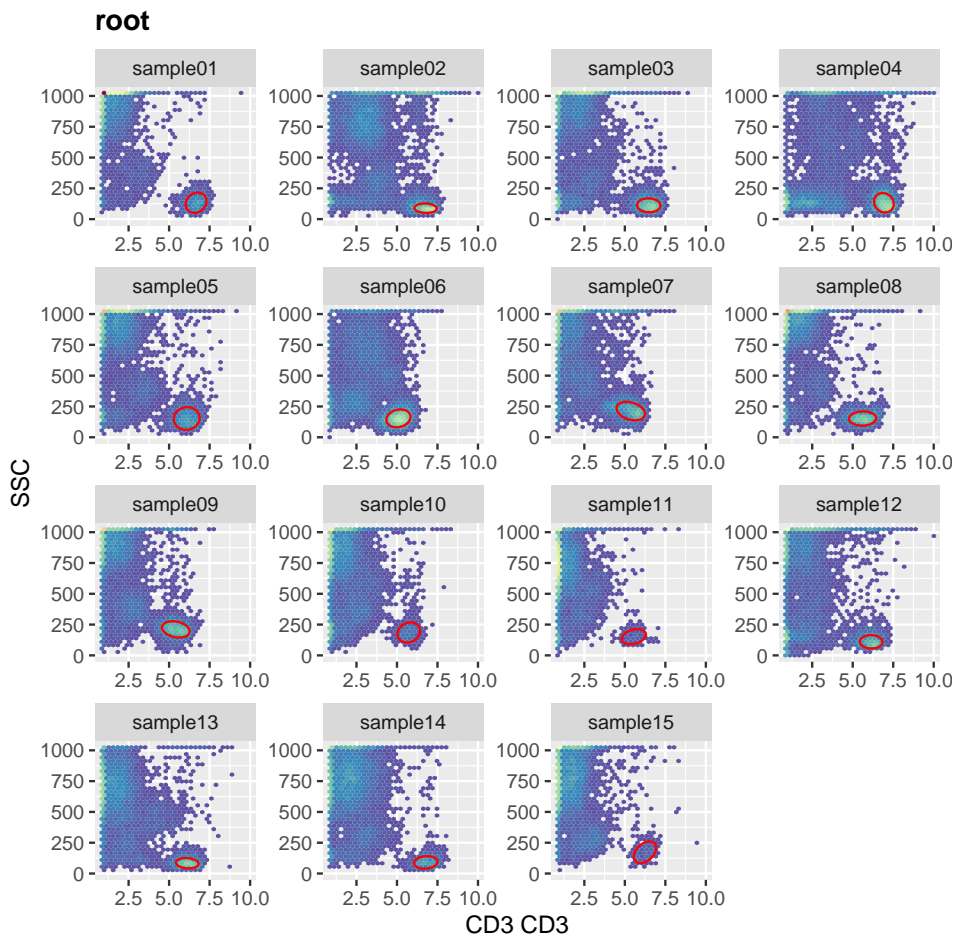
the `lymphGate` function in the `flowStats` package. The idea here is to first do a rough preselection in the two-dimensional projection of the data based on expert knowledge or prior experience and subsequently to fit a `norm2Filter` to this subset. The function also allows us to derive the pre-selection through back-gating: we know that CD4 positive cells are a subset of T-cells, so by estimating CD4 positive cells first we can get a rough idea on where to find the T-cells in the CD3 SSC projection.

```
> lg <- lymphGate(gs_cyto_data(gs), channels=c("CD3", "SSC"),
+               preselection="CD4", filterId="TCells",
+               scale=2.5)
> gs_pop_add(gs, lg)

[1] 2

> recompute(gs)

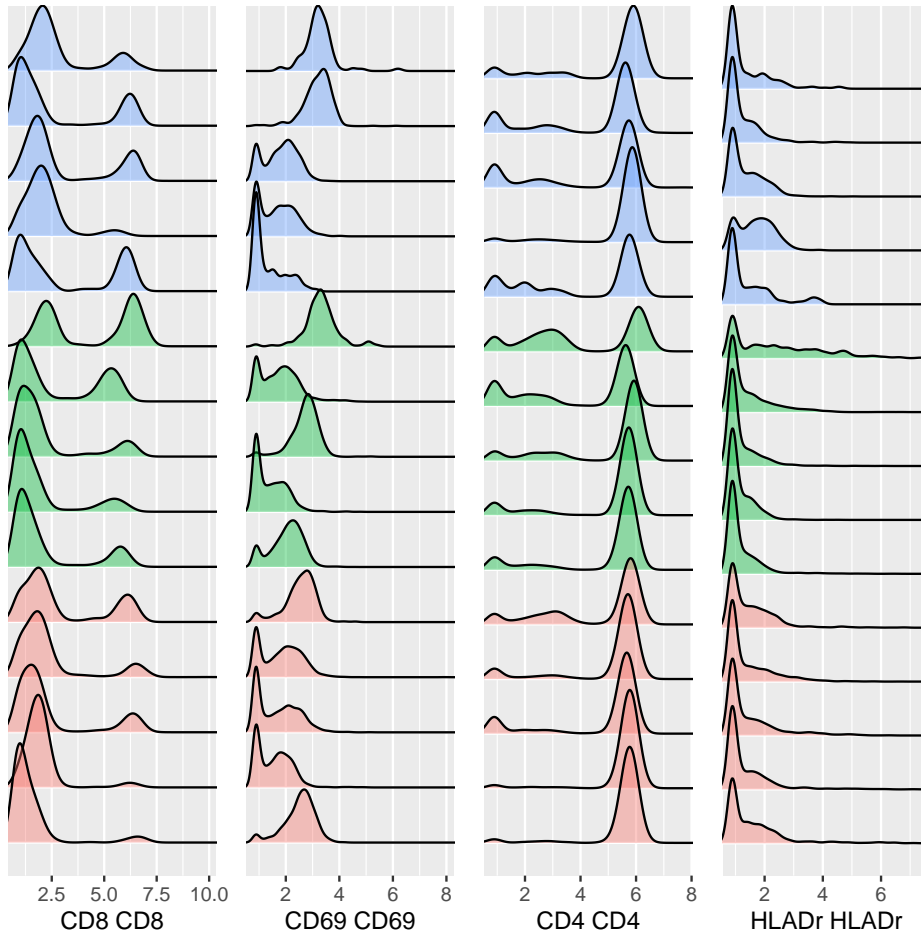
> ggcyto(gs, aes(x = CD3, y = SSC)) + geom_hex(bins = 32) + geom_gate("TCells")
```



In the next step we want to separate T-helper and NK cells using the CD4 and CD8 stains. A convenient way of doing this is to apply a `quadGate`, assuming that both CD4 and CD8 are binary markers

(cells are either positive or negative for CD4 and CD8). Often investigators use negative samples to derive a split point between the positive and negative populations, and apply this constant gate on all their samples. This will only work if there are no unforeseen shifts in the fluorescence intensities between samples which are purely caused by technical variation rather than biological phenotype. Let's take a look at this variation for the T-cell subset and all 4 remaining fluorescence channels:

```
> require(ggribes)
> require(gridExtra)
> pars <- colnames(gs)[c(3,4,5,7)]
> data <- gs_pop_get_data(gs, "TCells")
> plots <- lapply(pars, function(par)
+   as.ggplot(ggcyto(data, aes(x = !!par, fill = GroupID)) +
+     geom_density_ridges(mapping = aes(y = name), alpha = 0.4) +
+     theme(axis.title.y = element_blank(),
+           axis.text.y = element_blank(),
+           axis.ticks.y = element_blank())) +
+   facet_null()
+ ))
> do.call("grid.arrange", c(plots, nrow=1))
```



Indeed the data, especially for CD4 and CD8, don't align well. At this point we could decide to compute the *quadGates* for each sample separately. Alternatively, we can try to normalize the data and then compute a common gate. The `warpSet` function can be used to normalize data according to a set of landmarks, which essentially are the peaks or high-density areas in the density estimates shown before. The ideas here are simple:

- High density areas represent particular sub-types of cells.
- Markers are binary. Cells are either positive or negative for a particular marker.
- Peaks should align if the above statements are true.

The algorithm in `warpSet` performs the following steps:

1. Identify landmarks for each parameter using a `curv1Filter`
2. Estimate the most likely total number ( $k$ ) of landmarks
3. Perform k-means clustering to classify landmarks
4. Estimate warping functions for each sample and parameter that best align the landmarks, given the underlying data. This step uses functionality from the `fda` package.

## 5. Transform the data using the warping functions.

The algorithm should be robust to missing peaks in some of the samples, however the classification in step 3 becomes harder since it is not clear which cell population it represents.

While `warpSet` can be used directly on a `flowSet`, the `normalize` method was written to integrate with `GatingSet` objects. We must first create the gate whose channels we would like to normalize, which in this case will be done using the `quadrantGate` function to automatically estimate a `quadGate` in the CD4 and CD8 channels.

```
> qgate <- quadrantGate(gs_pop_get_data(gs, "TCells"), stains=c("CD4", "CD8"),
+                       filterId="CD4CD8", sd=3)
> gs_pop_add(gs, qgate, parent = "TCells")

[1] 3 4 5 6

> recompute(gs)
```

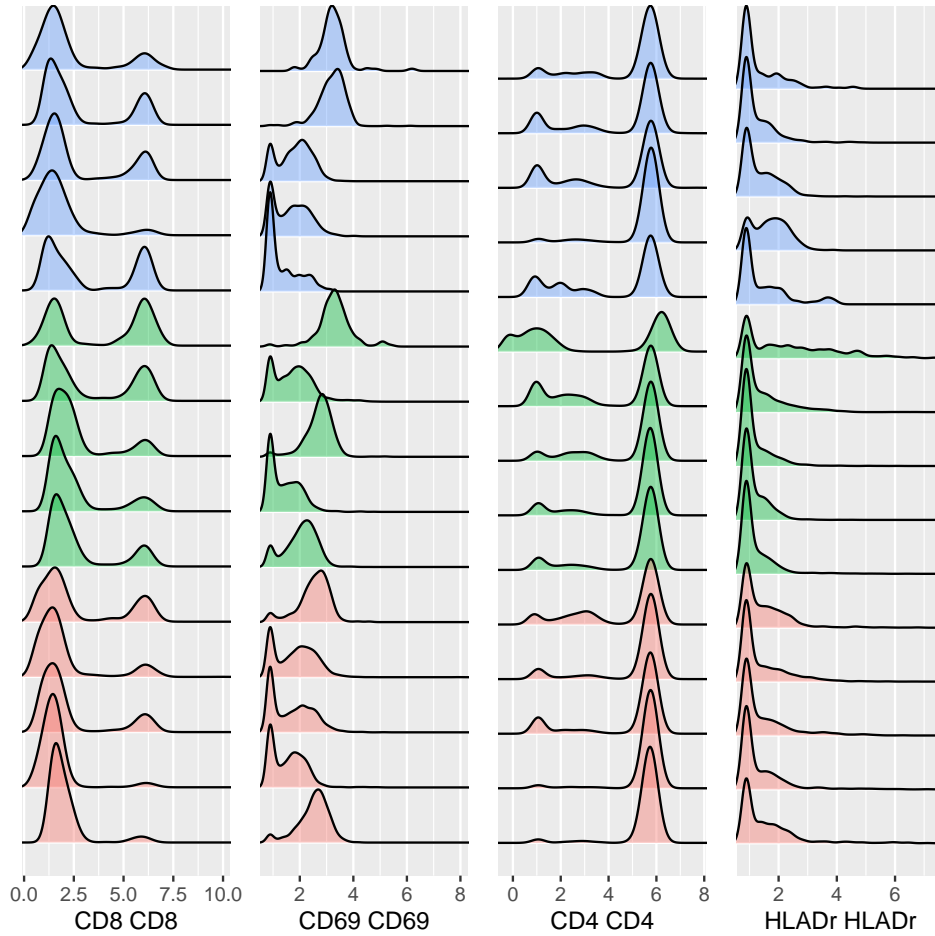
We can now normalize the channels relevant to this gate using `normalize`, so that we can successfully apply a single `quadGate` to all of the samples.

```
> gs <- normalize(gs, populations=c("CD4+CD8+", "CD4+CD8-", "CD4-CD8+", "CD4-CD8-"),
+                 dims=c("CD4", "CD8"), minCountThreshold = 50)

Estimating landmarks for channel CD4 ...
Estimating landmarks for channel CD8 ...
Registering curves for parameter CD4 ...
Registering curves for parameter CD8 ...
```

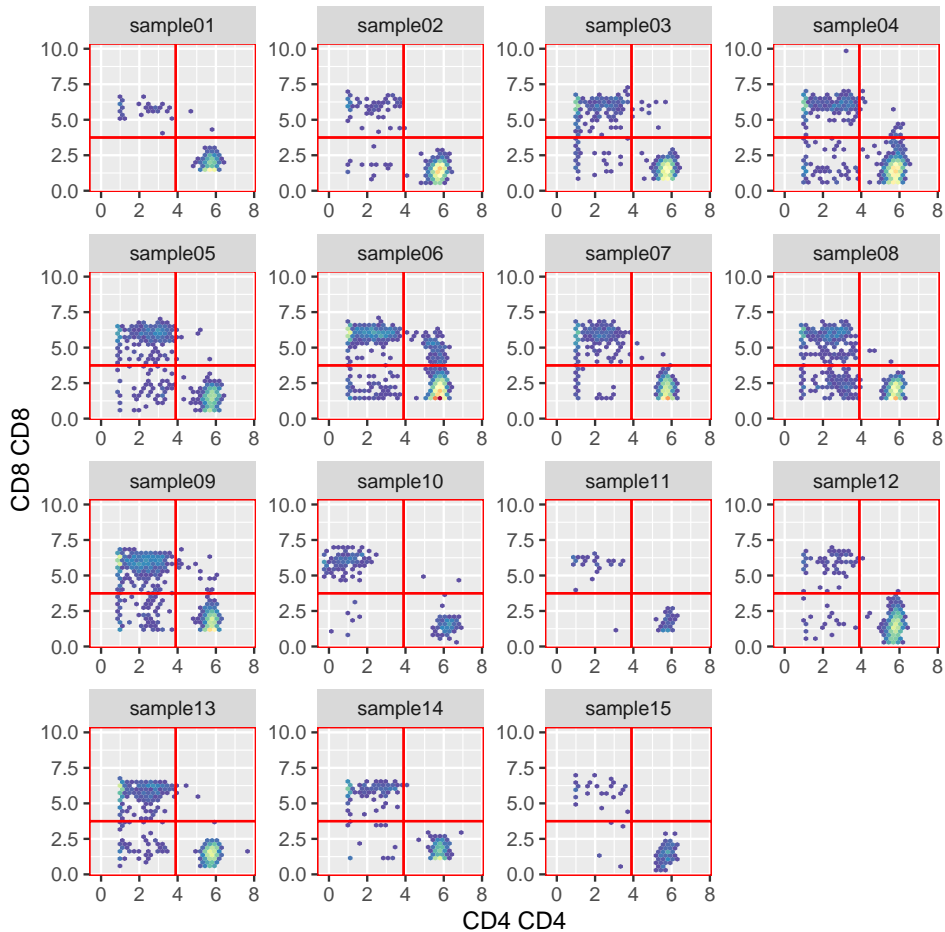
The normalization aligns the peaks in the selected channels across all samples.

```
> data <- gs_pop_get_data(gs, "TCells")
> plots <- lapply(pars, function(par)
+   as.ggplot(ggcyto(data, aes(x = !!par, fill = GroupID)) +
+     geom_density_ridges(mapping = aes(y = name), alpha = 0.4) +
+     theme(axis.title.y = element_blank(),
+           axis.text.y = element_blank(),
+           axis.ticks.y = element_blank()) +
+     facet_null()
+   ))
> do.call("grid.arrange", c(plots, nrow=1))
```



This ensures that the single quadGate is appropriately placed for each sample.

```
> ggcyto(gs_pop_get_data(gs, "TCells"), aes(x=CD4, y=CD8)) +
+   geom_hex(bins=32) +
+   geom_gate(gs_pop_get_gate(gs, "CD4-CD8-")) +
+   geom_gate(gs_pop_get_gate(gs, "CD4-CD8+")) +
+   geom_gate(gs_pop_get_gate(gs, "CD4+CD8-")) +
+   geom_gate(gs_pop_get_gate(gs, "CD4+CD8+"))
```



In a final step we might be interested in finding the proportion of activated T-helper cells by means of the CD69 stain. The `rangeGate` function is helpful in separating positive and negative peaks in 1D.

```
> CD69rg <- rangeGate(gs_cyto_data(gs), stain="CD69",
+                     alpha=0.75, filterId="CD4+CD8-CD69", sd=2.5)
> gs_pop_add(gs, CD69rg, parent="CD4+CD8-")
```

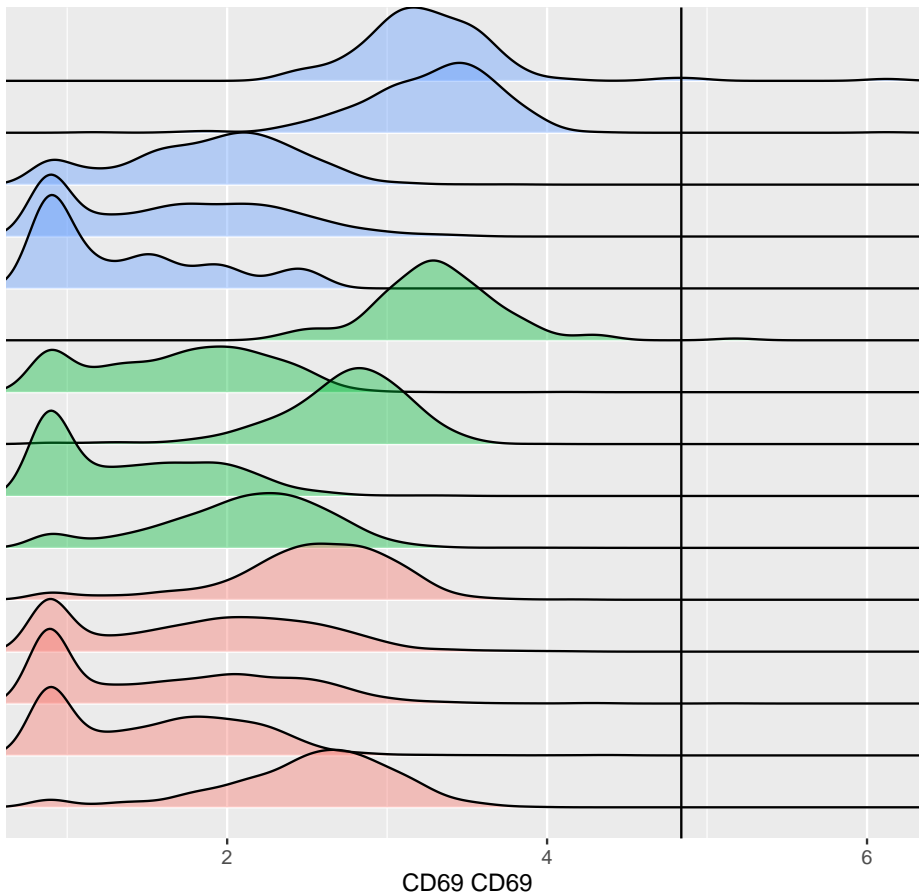
```
[1] 7
```

```
> # gs_pop_add(gs, CD69rg, parent="CD4+CD8-", name = "CD4+CD8-CD69-", negated=TRUE)
> recompute(gs)
```

```
> ggcyto(gs_pop_get_data(gs, "CD4+CD8-"), aes(x=CD69, fill = GroupID)) +
+   geom_density_ridges(mapping = aes(y = name), alpha = 0.4) +
+   theme(axis.title.y = element_blank(),
+         axis.text.y = element_blank(),
+         axis.ticks.y = element_blank()) +
+   geom_vline(xintercept = CD69rg@min) +
```

```
+ facet_null() +
+ ggtitle("CD4+")
```

CD4+



This channel could also benefit from normalization before determining a single rangeGate to apply to all samples.

```
> gs <- normalize(gs, populations=c("CD4+CD8-CD69"),
+               dims=c("CD69"), minCountThreshold = 50)
```

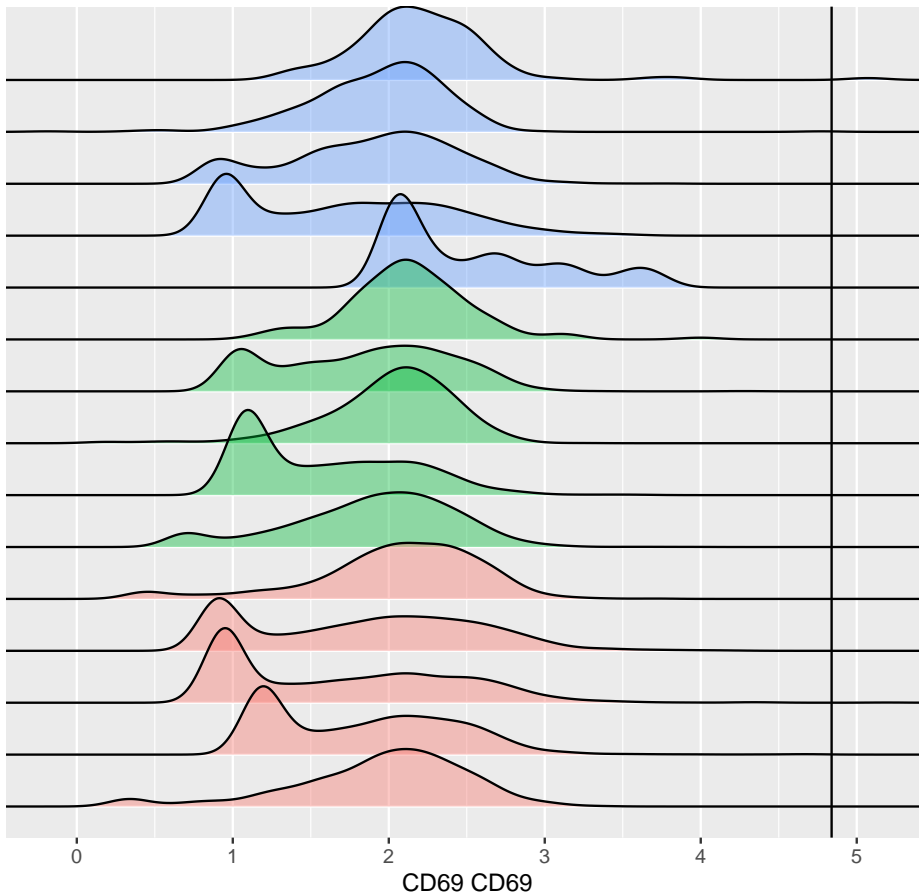
Estimating landmarks for channel CD69 ...

Registering curves for parameter CD69 ...

```
> ggcyto(gs_pop_get_data(gs, "CD4+CD8-"), aes(x=CD69, fill = GroupID)) +
+   geom_density_ridges(mapping = aes(y = name), alpha = 0.4) +
+   theme(axis.title.y = element_blank(),
+         axis.text.y = element_blank(),
+         axis.ticks.y = element_blank()) +
+   geom_vline(xintercept = CD69rg@min) +
+   facet_null() +
+   ggtitle("CD4+")
```



CD4+



## 2 Probability Binning

A probability binning algorithm for quantitating multivariate distribution differences was described by Roederer et al. The algorithm identifies the flow parameter in a flowFrame with the largest variance and divides the events in the flowFrame into two subgroups based on the median of the parameter. This process continues until the number of events in each subgroup is less than a user specified threshold.

For comparison across multiple samples, probability binning algorithm can be applied to a control dataset to obtain the position of bins and the same bins can be applied to the experimental dataset. The number of events in the control and sample bins can then be compared using the Pearson's chi-square test or the probability binning metric defined by Roederer et al.

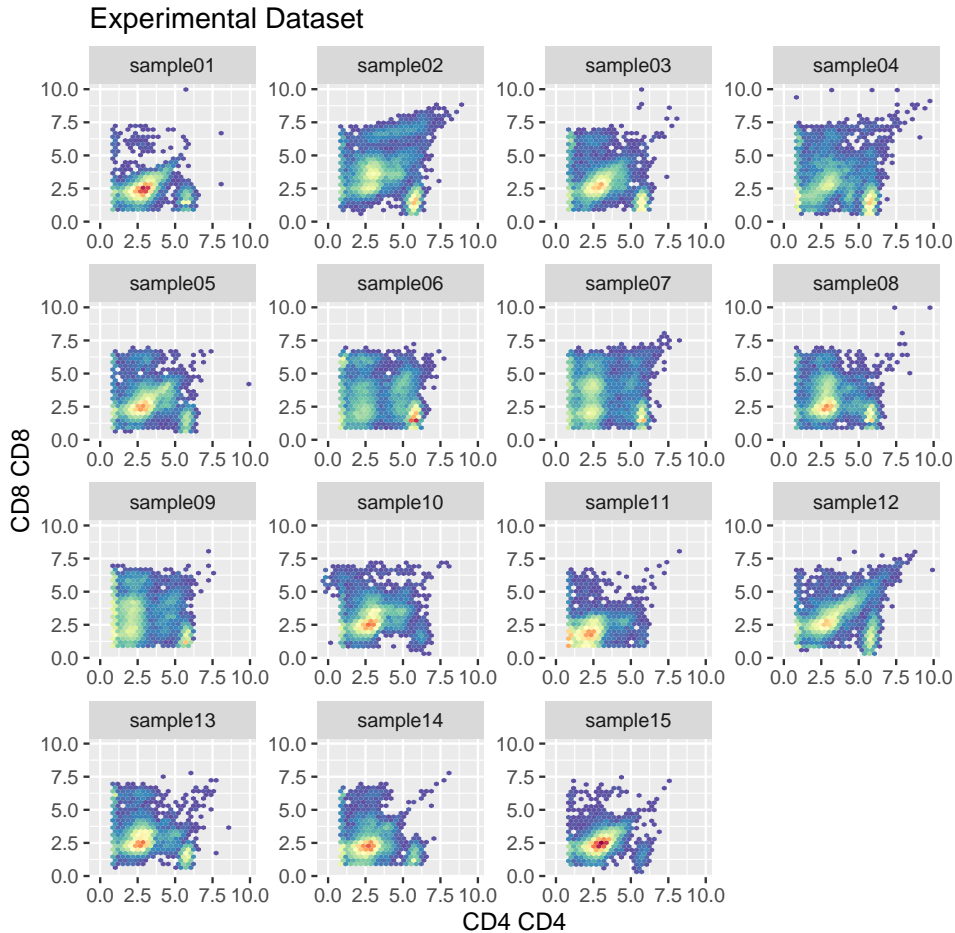
Although probability binning can be applied simultaneously to all parameters in a flowFrame with bins in  $n$  dimensional hyperspace, we proceed with a two dimensional example from our previous discussion involving CD4 and CD8 populations. This helps to simplify the demonstration of the method and interpretation of results.

From the workflow object containing the warped data, we extract our data frame of interest. We try to compare the panels using probability binning to identify patients with CD4, CD8 populations different from a control flowFrame that we create using the data from all the patients.

```
> dat <- gs_cyto_data(gs)
>
```

The dat is visualized below

```
> autoplot(dat, "CD4", "CD8") +
+   ggtitle("Experimental Dataset")
```



The control dataset is created by combining all the flowFrames in the flowSet. The flowFrame is then subsetted after applying a sampleFilter so that the control flowSet created has approximately the same number of events as the other flowSets in our example.

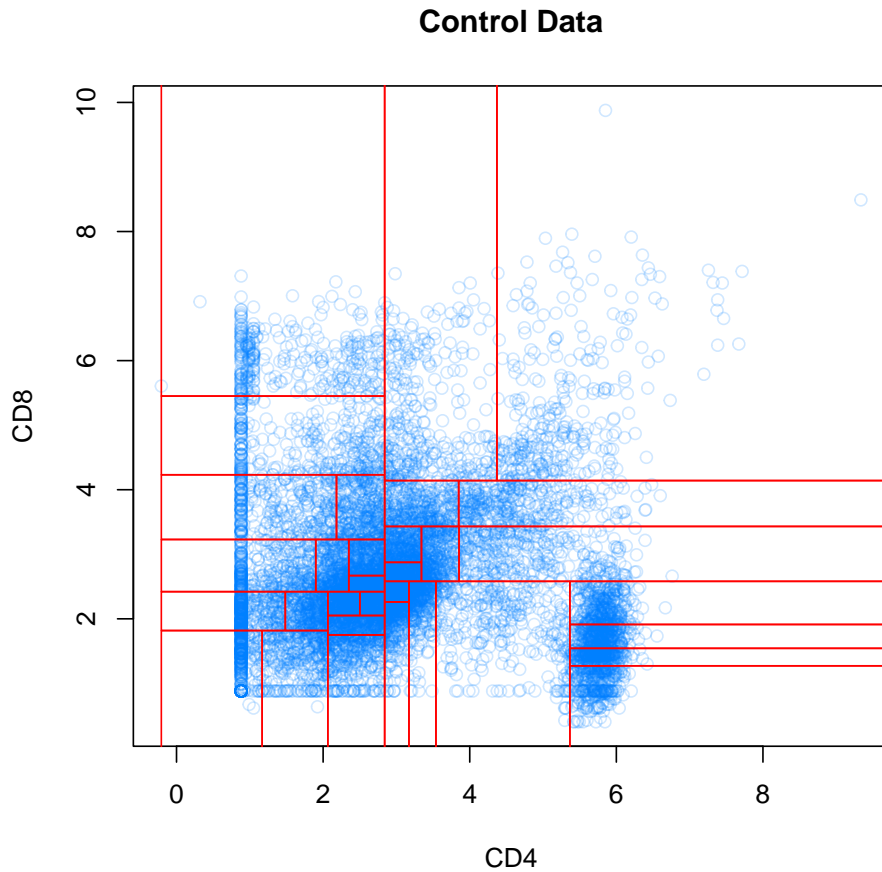
```
> datComb <- as(dat, "flowFrame")
> subCount <- nrow(exprs(datComb))/length(dat)
> sf <- sampleFilter(filterId="mySampleFilter", size=subCount)
> fres <- filter(datComb, sf)
> ctrlData <- Subset(datComb, fres)
> ctrlData <- ctrlData[, -ncol(ctrlData)] ##remove the column name "origi
>
```

The probability binning algorithm can then applied to the control data. The terminating condition for the algorithm is set so that the number of events in each bin is approximately 5 percent of the total number of events in the control data.

```
> minRow=subCount*0.05
> refBins<-proBin(ctrlData,minRow,channels=c("CD4","CD8"))
>
```

The binned control Data can be visualized using the `plotBins` function. Areas in the scatter plot with a large number of data points have a higher density of bins. Each bin also has approximately same number of events.

```
> plotBins(refBins,ctrlData,channels=c("CD4","CD8"),title="Control Data")
>
```



The same bin positions from the control data set are then applied to each flowFrame in our sample Data set.

```
> sampBins <- fsApply(dat,function(x){
+   binByRef(refBins,x)
+ })
```

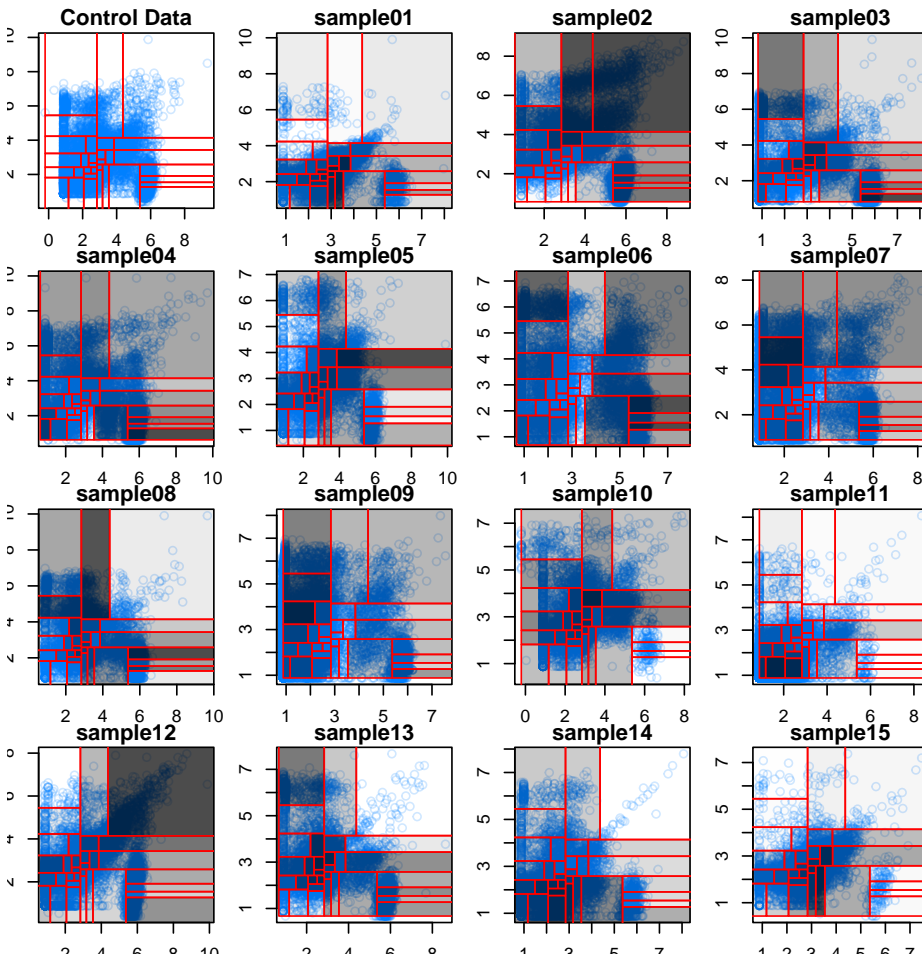
For each patient, the number events in the control and sample bins can be compared using the `calcPearsonChi` or using Roederers probability binning metric.

```
> pearsonStat <- lapply(sampBins, function(x) {
+   calcPearsonChi(refBins, x)
+ })

> sCount <- fsApply(dat, nrow)
> pBStat <- lapply(seq_along(sampBins), function(x) {
+   calcPBChiSquare(refBins, sampBins[[x]], subCount, sCount[x])
+ })
```

For each sample, the results can be visualized using the `plotBins` function. The residuals from Roeders probability binning metric or the Pearsons chi square test can be used to shade bins to highlight bins in each sample that differ the most from the control sample.

```
> par(mfrow=c(4, 4), mar=c(1.5, 1.5, 1.5, 1.5))
> plotBins(refBins, ctrlData, channels=c("CD4", "CD8"), title="Control Data")
> patNames <- sampleNames(dat)
> tm<-lapply(seq_len(length(dat)), function(x) {
+   plotBins(refBins, dat[[x]], channels=c("CD4", "CD8"),
+   title=patNames[x],
+   residuals=pearsonStat[[x]]$residuals[2, ],
+   shadeFactor=0.7)
+ })
+ )
```



The patient with CD4/CD8 populations most different from that of the control group can be identified from the magnitude of Pearson-chi square statistic (or Probability binning statistic).

	chi_Square_Statistic	pBin_Statistic
sample01	2388.29	299.38
sample02	2877.82	361.55
sample03	686.72	83.28
sample04	1784.35	222.68
sample05	594.86	71.61
sample06	3106.44	390.58
sample07	2141.89	268.08
sample08	865.63	106.00
sample09	3001.37	377.24
sample10	1503.84	187.05
sample11	5749.78	726.29
sample12	1129.43	139.50
sample13	1051.87	129.65
sample14	1858.50	232.09
sample15	3697.78	465.68

- R version 4.4.0 RC (2024-04-16 r86468), x86\_64-pc-linux-gnu
- Locale: LC\_CTYPE=en\_US.UTF-8, LC\_NUMERIC=C, LC\_TIME=en\_GB, LC\_COLLATE=C, LC\_MONETARY=en\_US.UTF-8, LC\_MESSAGES=en\_US.UTF-8, LC\_PAPER=en\_US.UTF-8, LC\_NAME=C, LC\_ADDRESS=C, LC\_TELEPHONE=C, LC\_MEASUREMENT=en\_US.UTF-8, LC\_IDENTIFICATION=C
- Time zone: America/New\_York
- TZcode source: system (glibc)
- Running under: Ubuntu 22.04.4 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.20-bioc/R/lib/libRblas.so
- LAPACK: /usr/lib/x86\_64-linux-gnu/lapack/liblapack.so.3.10.0
- Base packages: base, datasets, grDevices, graphics, methods, stats, utils
- Other packages: BH 1.84.0-0, flowCore 2.17.0, flowStats 4.17.0, flowWorkspace 4.17.0, ggcyto 1.33.0, ggplot2 3.5.1, ggridges 0.5.6, gridExtra 2.3, ncdFlow 2.51.0, scales 1.3.0, xtable 1.8-4
- Loaded via a namespace (and not attached): Biobase 2.65.0, BiocGenerics 0.51.0, DEoptimR 1.1-3, IDPmisc 1.1.21, KernSmooth 2.23-22, MASS 7.3-60.2, Matrix 1.7-0, R6 2.5.1, RColorBrewer 1.1-3, RCurl 1.98-1.14, RProtoBufLib 2.17.0, Rcpp 1.0.12, Rgraphviz 2.49.0, S4Vectors 0.43.0, XML 3.99-0.16.1, bitops 1.0-7, cli 3.6.2, clue 0.3-65, cluster 2.1.6, colorspace 2.1-0, compiler 4.4.0, corpcor 1.6.10, cytolib 2.17.0, data.table 1.15.4, deSolve 1.40, deldir 2.0-4, dplyr 1.1.4, fansi 1.0.6, farver 2.1.1, fda 6.1.8, fds 1.8, flowViz 1.69.0, generics 0.1.3, glue 1.7.0, graph 1.83.0, grid 4.4.0, gtable 0.3.5, hdrde 3.4, hexbin 1.28.3, interp 1.1-6, jpeg 0.1-10, ks 1.14.2, labeling 0.4.3, lattice 0.22-6, latticeExtra 0.6-30, lifecycle 1.0.4, magrittr 2.0.3, matrixStats 1.3.0, mclust 6.1.1, mnormt 2.1.1, munsell 0.5.1, mvtnorm 1.2-4, parallel 4.4.0, pcaPP 2.0-4, pillar 1.9.0, pkgconfig 2.0.3, plyr 1.8.9, png 0.1-8, pracma 2.4.4, rainbow 3.8, rlang 1.1.3, robustbase 0.99-2, rrcov 1.7-5, splines 4.4.0, stats4 4.4.0, tibble 3.2.1, tidyselct 1.2.1, tools 4.4.0, utf8 1.2.4, vctrs 0.6.5, withr 3.0.0, zlibbioc 1.51.0