

GMD Vignette:

Generic histogram construction, generic distance measure,
cluster analysis with evaluation and visualization

Xiaobei Zhao* Albin Sandelin*

February 6, 2012

(Version: 0.3.1.1)

Abstract

The purpose of this **GMD** Vignette is to show how to get started with the R package **GMD**. *GMD* denotes **Generalized Minimum Distance between distributions**, which is a true metric that measures the distance between the shapes of any two discrete numerical distributions (*e.g.* histograms).

The vignette includes a brief introduction, an example to illustrate the concepts and the implementation of **GMD** and case studies that were carried out using classical data sets (*e.g.* *iris*) and high-throughput sequencing data (*e.g.* *ChIP-seq*) from biology experiments. The appendix on page 15 contains an overview of package functionality, and examples using primary functions in histogram construction (the **ghist** function), how to measure distance between distributions (the **gdist** function), cluster analysis with evaluation (the “**elbow**” method in the **css** function) and visualization (the **heatmap.3** function).

Keywords: histogram, distance, metric, non-parametric, cluster analysis, hierarchical clustering, sum-of-squares, heatmap.3

Contents

1	Introduction	2
2	Minimal Example: “Hello, GMD!”	3
3	An example to understand GMD	4
3.1	Histogram: construction and visualization	4
3.2	Histogram: distance measure and alignment	7

*Bioinformatics Centre, University of Copenhagen, Department of Biology and Biotech Research and Innovation Centre, Ole Maaløes Vej 5, DK-2200, Copenhagen, Denmark. Email: xiaobei@binf.ku.dk ; albin@binf.ku.dk

4 Case study	9
4.1 CAGE: measuring the dissimilarities among TSSDs	9
4.2 ChIP-seq: measuring the similarities among histone modification patterns	12
A Functionality	15
A.1 An overview	15
A.2 ghist: Generic construction and visualization of histograms	16
A.2.1 Examples using simulated data	16
A.2.2 Examples using iris data	18
A.2.3 Examples using nottem data	20
A.3 gdist: Generic construction and visualization of distances	22
A.4 css: Clustering Sum-of-Squares and the ‘elbow’ plot	24
A.5 heatmap.3: Visualization in cluster analysis, with evaluation	27
A.5.1 Examples using mtcars data	27
A.5.2 Examples using ruspini data	31
B Data	33
B.1 GMD dataset overview	33
B.2 CAGE data: cage and cagel	33
B.3 ChIP-seq data: chipseq_mES and chipseq_hCD4T	34

1 Introduction

Similar to the Earth Mover’s distance, **Generalized Minimum distance of Distributions (GMD)** (based on MDPA - Minimum Difference of Pair Assignment [3]) is a true metric of the similarity between the shapes of two histograms¹. Considering two normalized histograms A and B , GMD measures their similarity by counting the necessary “shifts” of elements between the bins that have to be performed to transform distribution A into distribution B .

The **GMD** package provides classes and methods for computing **GMD** in R [5]. The algorithm has been implemented in C to interface with R for efficient computation. The package also includes downstream cluster analysis in function **css** (A.4 on page 24) that use a pairwise distance matrix to make partitions given variant criteria, including the “elbow” rule as discussed in [7] or desired number of clusters. In addition, the function **heatmap.3** (A.5 on page 27) integrates the visualization of the hierarchical clustering in dendrogram, the distance measure

¹In statistics (and many other fields), histogram refers to a graphical representation of category frequencies in the data. Here, we use this term in a more mathematical sense, defined as a function that counts categorical data or a result returned by such a function.

in heatmap and graphical representations of summary statistics of the resulting clusters or the overall partition. For more flexibility, the function `heatmap.3` can also accept plug-in functions defined by end-users for custom summary statistics.

The motivation to write this package was born with the project [7] on characterizing Transcription Start Site (TSS) landscapes using high-throughput sequencing data, where a non-parametric distance measure was developed to assess the similarity among distributions of high-throughput sequencing reads from biological experiments. However, it is possible to use the method for any empirical distributions of categorical data.

The package is available on CRAN. The source code is available at <http://cran.r-project.org/web/packages/GMD/> under GPL license.

2 Minimal Example: “Hello, GMD!”

```
hello-GMD.R
1  ##' Check GMD's sanity and start up
2  ##' @name hello-GMD
3
4  ## GMD at CRAN, for source code download and installation
5  ## http://cran.r-project.org/web/packages/GMD/index.html
6
7  ## load GMD
8  library(GMD)
9
10 ## version of GMD and description
11 packageVersion("GMD")
12 packageDescription("GMD")
13
14 ## view GMD vignette
15 vignette("GMD-vignette",package="GMD")
16
17 ## list the available data sets in GMD
18 data(package="GMD")
19
20 ## list all the objects in the GMD
21 ls("package:GMD")
22
23 ## help info on GMD
24 help(package="GMD")
25
26 ## run a demo
27 demo("GMD-demo")
28
29 ## cite GMD in publications
30 citation(package="GMD")
31
```

Figure 1: Source code of “hello-GMD.R”.

hello-GMD.R (fig. 1) is a minimal example to load and check of that your GMD installation works. It also

includes code for viewing the package information and this “Vignette”, checking data sets provided by `GMD`, starting a demo and listing the citation of `GMD`.

3 An example to understand GMD

This example, based on simulated data, is designed to illustrate the concepts and the implementation of `GMD` by stepping through the computations in detail.

3.1 Histogram: construction and visualization

```
> ## -----
> ## chunk1: Load library and simulate data
> ## -----
> require("GMD") # load library
> ## create two normally-distributed samples
> ## with unequal means and unequal variances
> set.seed(2012)
> x1 <- rnorm(1000,mean=-5, sd=10)
> x2 <- rnorm(1000,mean=10, sd=5)

> ## -----
> ## chunk2: Construct histograms
> ## -----
> ## create common bins
> n <- 20 # desired number of bins
> breaks <- gbreaks(c(x1,x2),n) # bin boundaries
> ## make two histograms
> v1 <- ghist(x1,breaks=breaks,digits=0)
> v2 <- ghist(x2,breaks=breaks,digits=0)

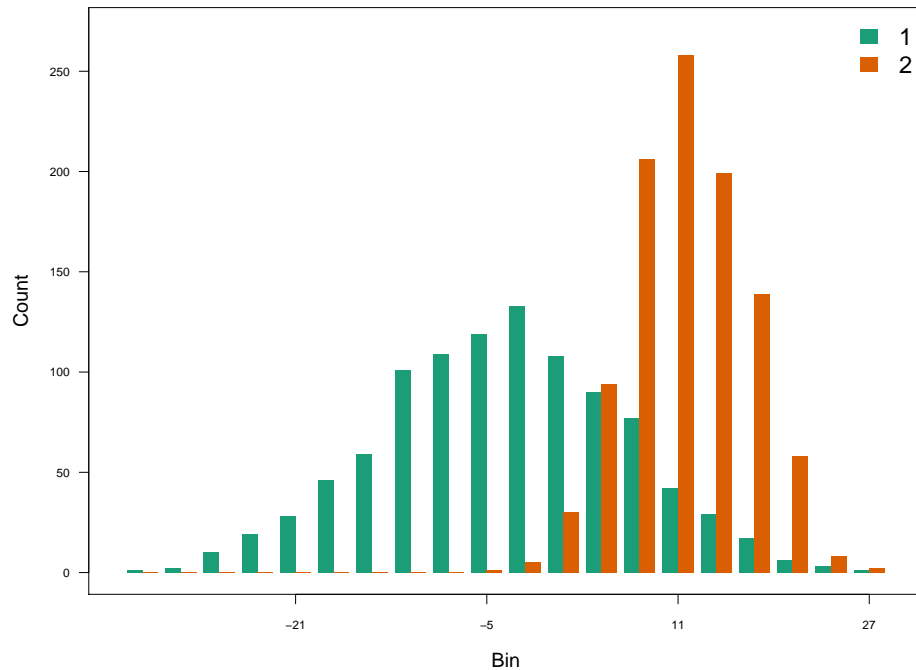
> ## -----
> ## chunk3: Save histograms as multiple-histogram ('mhist') object
> ## -----
> x <- list(v1,v2)
> mhist.obj <- as.mhist(x)
```

```

> ## -----
> ## chunk4: Visualize an `mhist` object
> ## -----
> ## plot histograms side-by-side
> plot(mhist.obj,mar=c(1.5,1,1,0),main="Histograms of simulated normal distributions")

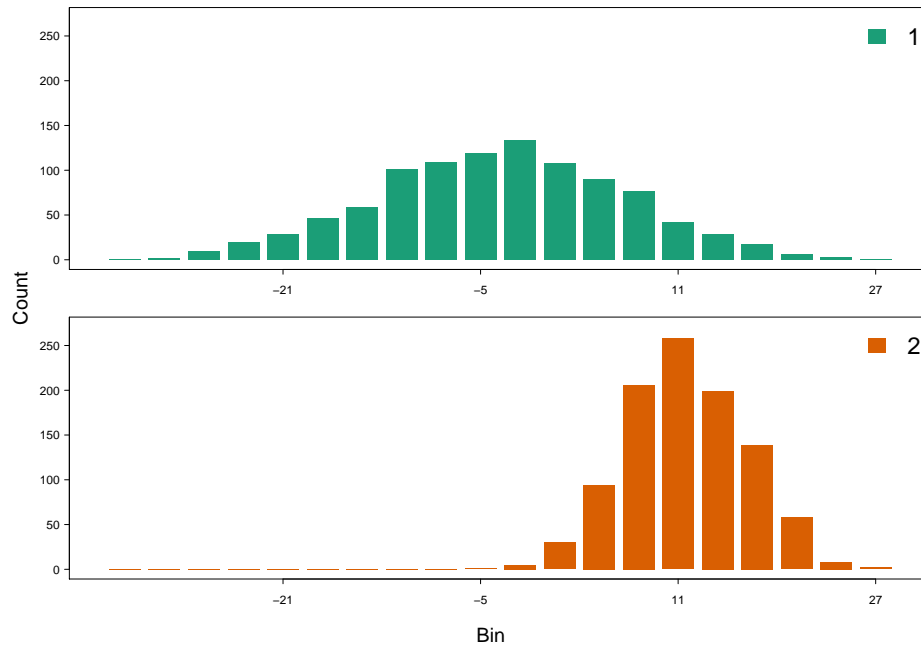
```

Histograms of simulated normal distributions



```
> ## plot histograms as subplots, with corresponding bins aligned
> plot(mhist.obj, beside=FALSE, mar=c(1.5,1,1,0),
+      main="Histograms of simulated normal distributions")
```

Histograms of simulated normal distributions



3.2 Histogram: distance measure and alignment

Here we measure the *GMD* distance between shapes of two histograms with option *sliding* on.

```
> ## -----
> ## chunk5: Measure the pairwise distance between two histograms by GMD
> ## -----
> gmdp.obj <- gmdp(v1,v2,sliding=TRUE)
> print(gmdp.obj)          # print a brief version by default

[1] 1.334

> print(gmdp.obj,mode="detailed") # print a detailed version

GM-Distance: 1.334
Sliding: TRUE
Number of hits: 1
Gap:
      v1      v2
Hit1    5      0

Resolution: 1

> print(gmdp.obj,mode="full")    # print a full version

Distribution of v1:
1 2 10 19 28 46 59 101 109 119 133 108 90 77 42 29 17 6 3 1
(After normalization)
0.001 0.002 0.01 0.019 0.028 0.046 0.059 0.101 0.109 0.119 0.133 0.108 0.09 0.077 0.042 0.029 0.017 0.006 0.003 0.001

Distribution of v2:
0 0 0 0 0 0 0 0 0 1 5 30 94 206 258 199 139 58 8 2
(After normalization)
0 0 0 0 0 0 0 0 0.001 0.005 0.03 0.094 0.206 0.258 0.199 0.139 0.058 0.008 0.002

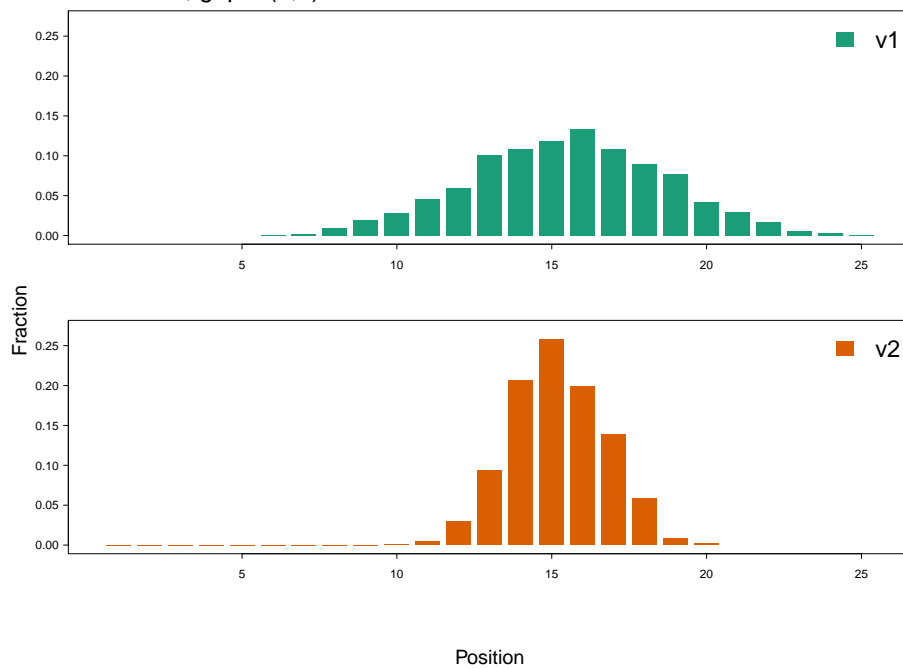
GM-Distance: 1.334
Sliding: TRUE
Number of hits: 1
Gap:
      v1      v2
Hit1    5      0

Resolution: 1
```

Now, let's have a look at the alignment by *GMD*, with a distance 1.334 and a “shift” of 5 in the 1st distribution. It is important to note that the specific features (the values in this case) of the original bins in the histograms are ignored with *sliding* on. To keep original bin-to-bin correspondence, please set *sliding* to **FALSE** (see examples in [section 4.2](#) on [page 12](#)).

```
> ## -----
> ## chunk6: Show alignment
> ## -----
> plot(gmdp.obj,beside=FALSE)
```

Optimal alignment between distributions (with sliding)
GMD=1.334, gap=c(5,0)



4 Case study

4.1 CAGE: measuring the dissimilarities among TSSDs

Studies have demonstrated that the spatial distributions of read-based sequencing data from different platforms often indicate functional properties and expression profiles (reviewed in [6] and [8]). Analyzing the distributions of DNA reads is therefore often meaningful. To do this systematically, a measure of similarity between distributions is necessary. Such measures should ideally be true metrics, have few parameters as possible, be computationally efficient and also make biological sense to end-users. Case studies were made in section 4.1 and 4.2 to demonstrate the applications of GMD using distributions of CAGE and ChIP-seq reads.

In this section we demonstrate how GMD is applied to measure the dissimilarities among TSSDs, histograms of transcription start site (TSS) that are made of CAGE tags, with option *sliding* on. The spatial properties of TSSDs vary widely between promoters and have biological implications in both regulation and function. The raw data were produced by CAGE and downloaded from FANTOM3 ([2]) and CAGE sequence reads were preprocessed as did in [7]. The following codes `case-cage.R` (fig. 2) are sufficient to perform both pairwise GMD calculation by function `gmdp` and to construct a GMD distance matrix by function `gmdm`. A handful of options are available for control and flexibility, particularly, the option `sliding` is enabled by default to allow partial alignment.

```
case-cage.R
1  require("GMD") # load library
2  data(cage)     # load data
3
4  ## measure pairwise distance
5  x <- gmdp(cage[["Pfkfb3 (T02R00AEC2D8)"]],cage[["Csfl (T03R0672174D)"]])
6  print(x)      # print a brief version by default
7  print(x, mode="full") # print a full version by default
8
9  ## show alignment
10 plot(x,labels=c("Pfkfb3","Csfl"),beside=FALSE)
11
12 ## show another alignment
13 plot(gmdp(cage[["Hig1 (T09R0743763C)"]],cage[["Cd72 (T04R028B8BC9)"]]),
14      labels=c("Hig1 (T09R0743763C)","Cd72 (T04R028B8BC9)"),
15      beside=FALSE)
16
17 ## construct a distance matrix and visualize it
18 short.labels <- gsub("(.)" "\\(.", "\\1",names(cage)) # get short labels
19 x <- gmdm(cage[1:6],labels=short.labels[1:6])
20 plot(x)
```

Figure 2: Source code of “case-cage.R”.

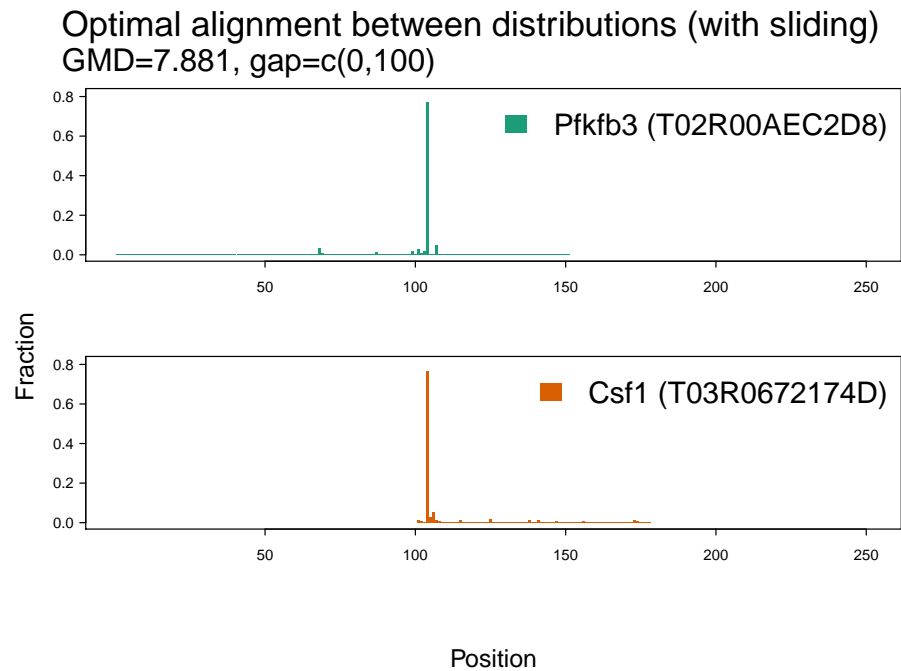


Figure 3: Graphical output 1 of source code “case-cage.R”.

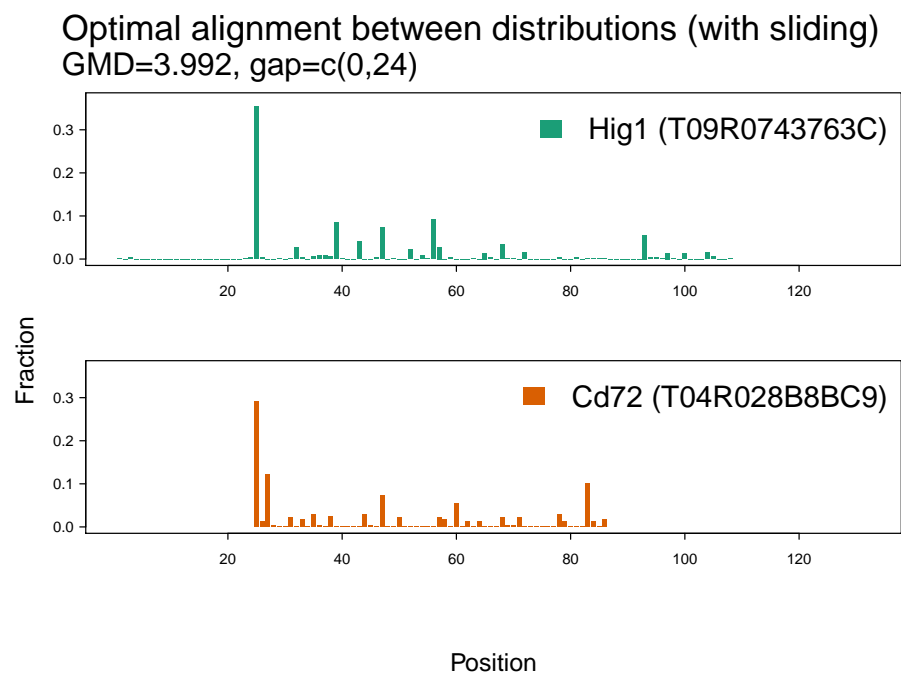


Figure 4: Graphical output 2 of source code “case-cage.R”.

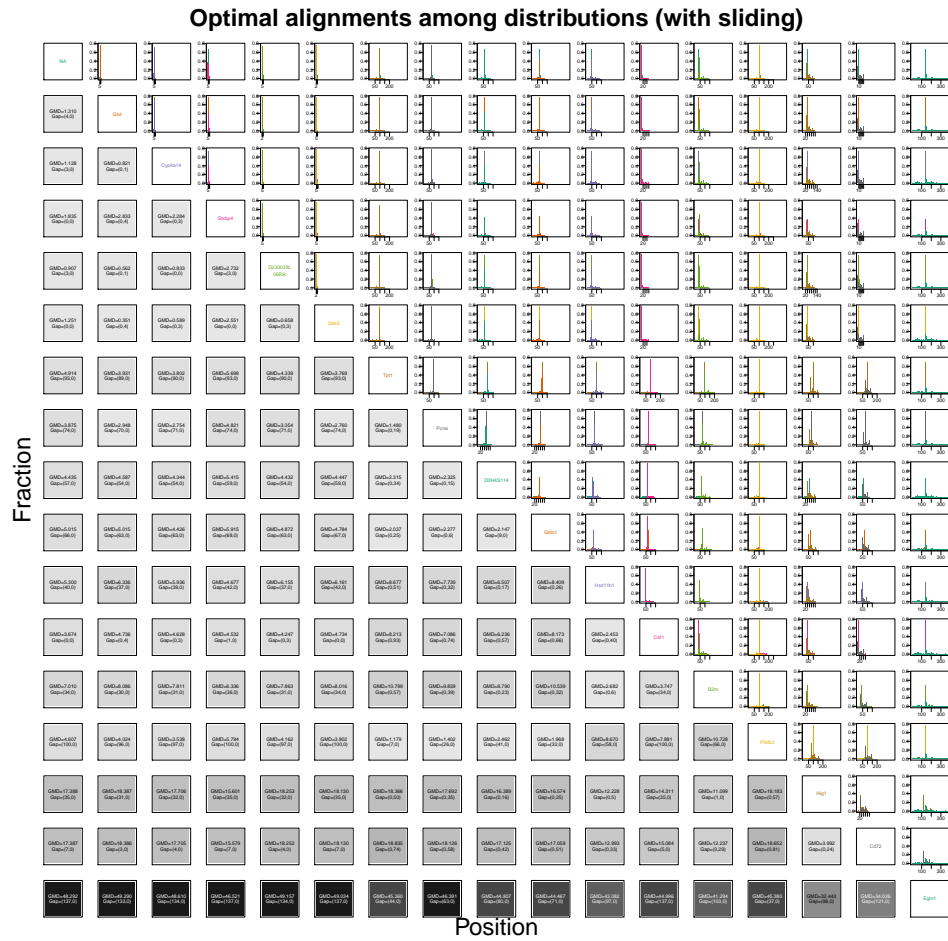


Figure 5: Graphical output 1 of source code “case-cage.R”.

4.2 ChIP-seq: measuring the similarities among histone modification patterns

In this section we demonstrate how **GMD** is applied to measure the dissimilarities between histone modifications represented by ChIP-seq reads. Distinctive patterns of chromatin modifications around the TSS are associated with transcription regulation and expression variation of genes. Comparing the chromatin modification profiles (originally produced by [1] and [4], and preprocessed by [7]), the **sliding** option is disabled for fixed alignments at the TSSs and the flanking regions. The GMD measure indicates how well profiles are co-related to each other. In addition, the downstream cluster analysis is visualized with function **heatmap.3** that use GMD distance matrix to generate clustering dendrograms and make partitions given variant criteria, including the ‘‘**elbow**’’ rule (discussed in [7]) or desired number of clusters.

```

1  require("GMD")      # load library
2  data(chipseq_mES)    # load data
3  data(chipseq_hCD4T) # load data
4
5  ## pairwise distance and alignment based on GMD metric
6  plot(gmdm(chipseq_mES,sliding=FALSE))
7
8  ## clustering on spatial distributions of histone modifications
9  x <- gmdm(chipseq_hCD4T,sliding=FALSE,resolution=10)
10 heatmap.3(x,revC=TRUE)
11
12 ## Determine the number of clusters by "Elbow" criterion
13 main <- "Heatmap of ChIP-seq data (human CD4+ T cells)"
14 dist.obj <- gmdm2dist(x)
15 css.multi.obj <- css.hclust(dist.obj,hclust(dist.obj))
16 elbow.obj <- elbow.batch(css.multi.obj,ev.thres=0.90,inc.thres=0.05)
17 heatmap.3(dist.obj, main=main, revC=TRUE, kr=elbow.obj$k, kc=elbow.obj$k)
18
19 ## more strict threshold
20 elbow.obj <- elbow.batch(css.multi.obj,ev.thres=0.75,inc.thres=0.1)
21 heatmap.3(dist.obj, main=main, revC=TRUE, kr=elbow.obj$k, kc=elbow.obj$k)
22
23 ## side plots
24 normalizeVector <- function(v){v/sum(v)} # a function to normalize a vector
25 dev.new(width=12,height=8)
26
27 ## summary of row clusters
28 expr1 <- list(quote(op <- par(mar = par("mar")*2)),
29               quote(plot(mhist.summary(as.mhist(i.x)),if.plot.new=FALSE)),
30               quote(par(op))
31             )
32
33 ## summary of row clustering
34 expr2 <- list(quote(tmp.clusters <- cutree(hclust(dist.row),k=kr)),
35               quote(tmp.css <- css(dist.row,tmp.clusters)),
36               quote(print(tmp.css)),
37               quote(tmp.wev <- tmp.css$wss/tmp.css$tss),
38               quote(names(tmp.wev) <- as.character(unique(tmp.clusters))),
39               quote(tmp.wev <- tmp.wev[order(unique(tmp.clusters))]),
40               quote(barplot(tmp.wev,main="Cluster Explained Variance", xlab="Cluster",
41                             ylab="EV",col="white",border="black",
42                             ylim=c(0,max(tmp.wev)*1.1),cex.main=1)))
43 expr3 <- list(quote(op <- par(mar = par("mar")*2)),
44               quote(plot.elbow(css.multi.obj,elbow.obj,if.plot.new=FALSE)),
45               quote(par(op))
46             )
47
48 heatmap.3(dist.obj, main=main, cex.main=1.25, revC=TRUE, kr=elbow.obj$k, kc=elbow.obj$k,
49           keysizes=1,mapsizes=4.5,row.data=lapply(chipseq_hCD4T,normalizeVector),
50           plot.row.clusters=TRUE,plot.row.clusters.list=list(expr1),
51           plot.row.clustering=TRUE,plot.row.clustering.list=list(expr2,expr3))
52

```

Figure 6: Source code of “case-chipseq.R”.

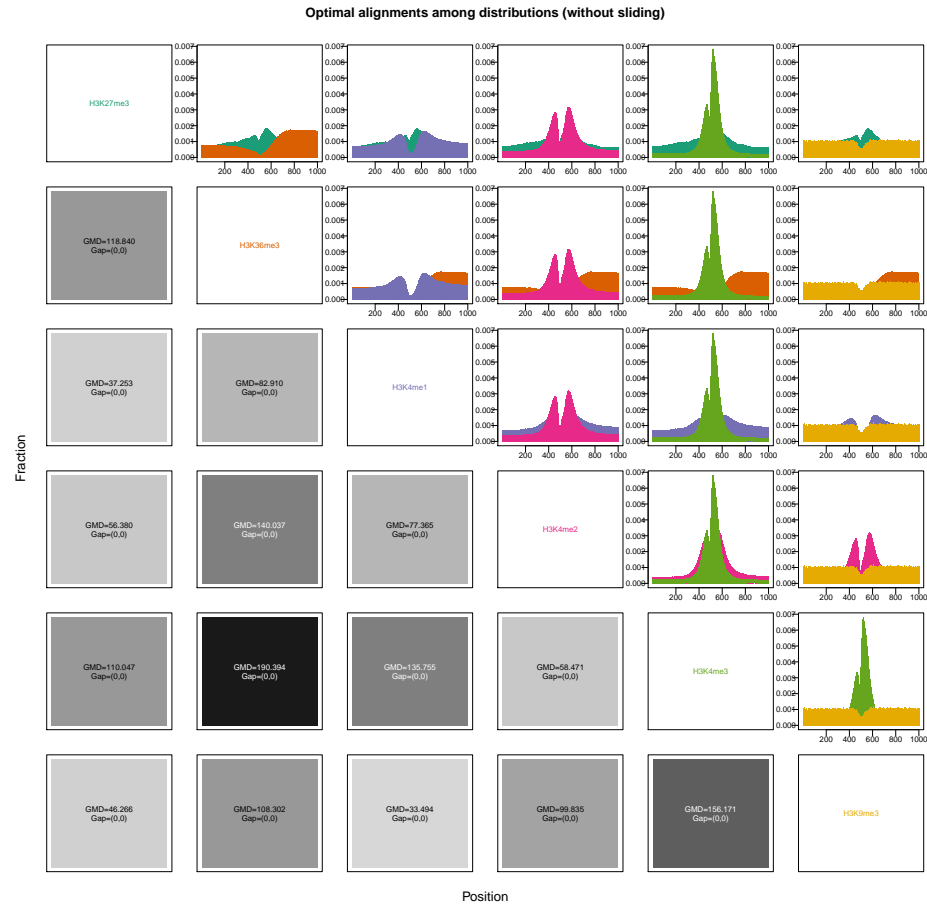


Figure 7: Graphical output 1 of source code “case-chipseq.R”.

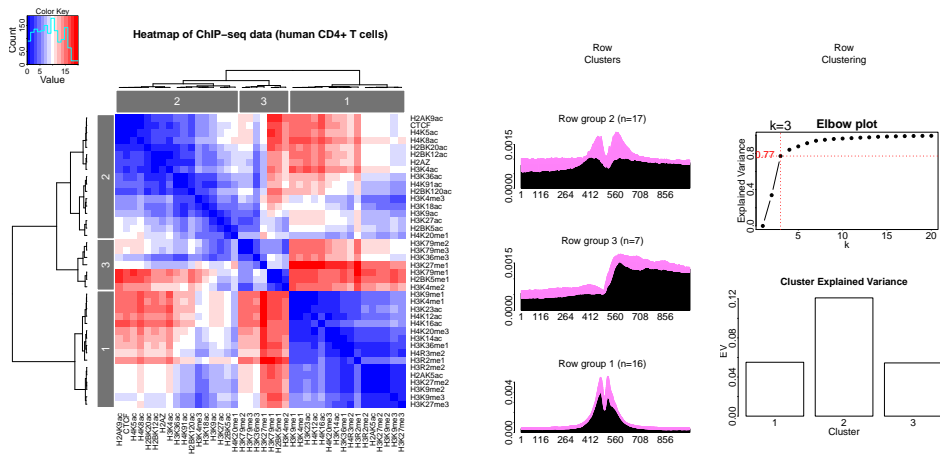


Figure 8: Graphical output 1 of source code “case-chipseq.R”.

A Functionality

A.1 An overview

Table 1. Functions of the **GMD** R package

Function	Description
<code>ghist</code>	Generalized Histogram Computation and Visualization
<code>gdist</code>	Generalized Distance Matrix Computation
<code>css</code>	Computing Clustering Sum-of-Squares and evaluating the clustering by the ‘ <code>elbow</code> ’ method
<code>heatmap.3</code>	Enhanced Heatmap Representation with Dendrogram and Partition
<code>gmdp</code>	Computation of GMD on a pair of histograms
<code>gmdm</code>	Computation of GMD Matrix on a set of histograms

A.2 ghist: Generic construction and visualization of histograms

A.2.1 Examples using simulated data

example-ghist.R (fig. 9) is an example on how to construct a histogram object from raw data and make a visualization based on this.

```
example-ghist.R
1  ## load library
2  require("GMD")
3
4  ## create two normally-distributed samples
5  ## with unequal means and unequal variances
6  set.seed(2012)
7  v1 <- rnorm(1000,mean=-5, sd=10)
8  v2 <- rnorm(1000,mean=10, sd=5)
9
10 ## create common bins
11 n <- 20 # desired number of bins
12 breaks <- gbreaks(c(v1,v2),n) # bin boundaries
13 x <-
14   list(ghist(v1,breaks=breaks,digits=0),
15         ghist(v2,breaks=breaks,digits=0))
16 mhist.obj <- as.mhist(x)
17
18 ## plot histograms side-by-side
19 plot(mhist.obj,mar=c(1.5,1,1,0),main="Histograms of simulated normal distributions")
20
21 ## plot histograms as subplots,
22 ## with corresponding bins aligned
23 plot(mhist.obj,beside=FALSE,mar=c(1.5,1,1,0),
24       main="Histograms of simulated normal distributions")
```

Figure 9: Source code of “example-ghist.R”.

Histograms of simulated normal distributions

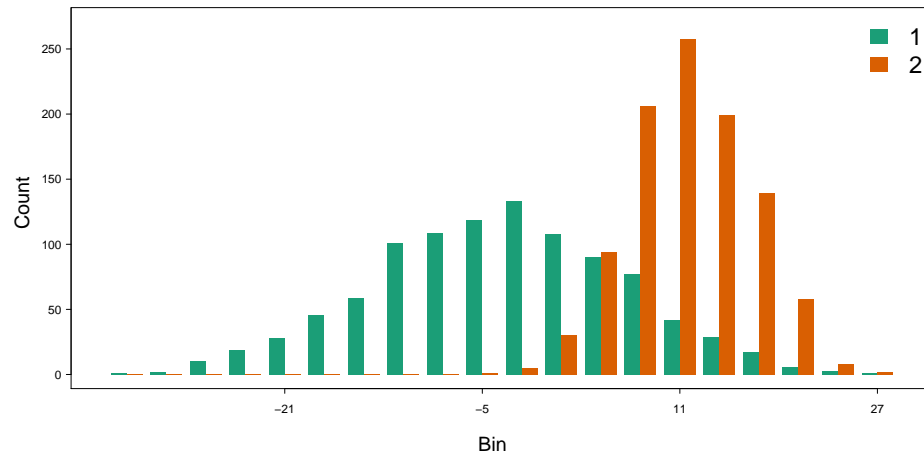


Figure 10: Graphical output 1 of source code “example-ghist.R”.

Histograms of simulated normal distributions

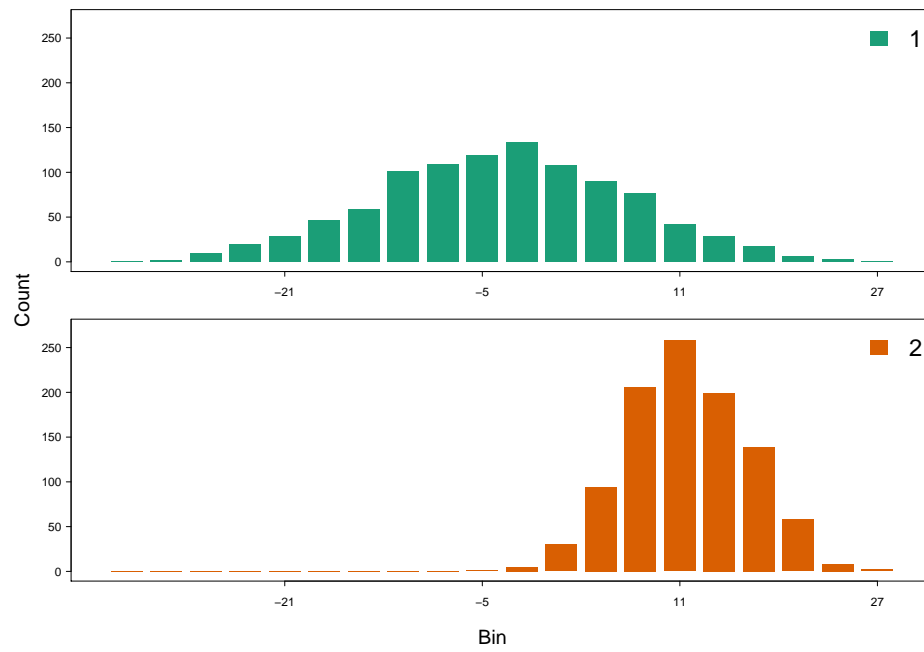


Figure 11: Graphical output 1 of source code “example-ghist.R”.

A.2.2 Examples using iris data

case-iris.R (fig. 12) is a study on how to obtain and visualize histograms, using Fisher's *iris* data set.

```
case-iris.R
1  ## load library
2  require("GMD")
3
4  ## load data
5  data(iris)
6
7  ## create common bins
8  n <- 30                                # the number of bins
9  breaks <- gbreaks(iris[, "Sepal.Length"], n) # the boundary of bins
10
11 ## create a list of histograms
12 Sepal.Length <-
13   list(setosa=ghist(iris[iris$Species=="setosa", "Sepal.Length"], breaks=breaks),
14         versicolor=ghist(iris[iris$Species=="versicolor", "Sepal.Length"], breaks=breaks),
15         virginica=ghist(iris[iris$Species=="virginica", "Sepal.Length"], breaks=breaks)
16   )
17
18 ## convert to a `hist' object
19 x <- as.mhist(Sepal.Length)
20
21 ## get bin-wise summary statistics
22 summary(x)
23
24 ## visualize the histograms
25 plot(x, beside=FALSE,
26       main="Histogram of Sepal.Length of iris", xlab="Sepal.Length (cm)")
```

Figure 12: Source code of “case-iris.R”.

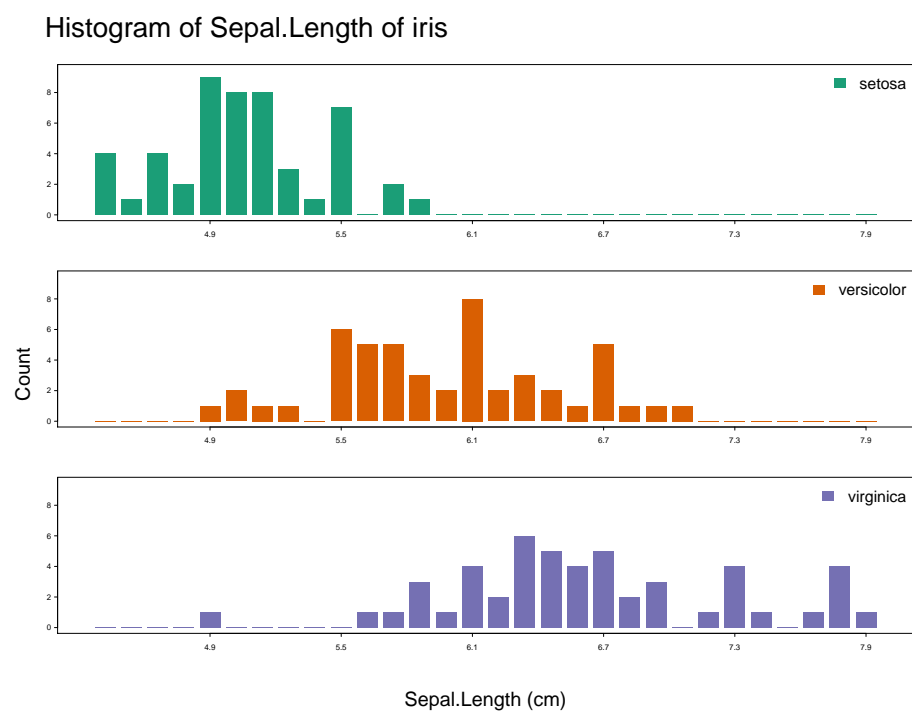


Figure 13: Graphical output 1 of source code “case-iris.R”.

A.2.3 Examples using nottem data

case-nottem.R (fig. 14) is a study on how to draw histograms side-by-side and to compute and visualize a bin-wise summary plot, using air temperature data at Nottingham Castle.

```
case-nottem.R
1  ## load library
2  require("GMD")
3
4  ## load data
5  data(nottem)
6
7  class(nottem)      # a time-series (ts) object
8  x <- ts2df(nottem) # convert ts to data.frame
9  mhist1 <- as.mhist(x[1:3,])
10
11 ## plot multiple discrete distributions side-by-side
12 plot(mhist1,xlab="Month",ylab="Degree Fahrenheit",
13      main="Air temperatures at Nottingham Castle")
14
15 ## make summary statistics for each bin
16 mhist2 <- as.mhist(x)
17 ms <- mhist.summary(mhist2)
18 print(ms)
19
20 ## plot bin-wise summary statistics with
21 ## confidence intervals over the bars
22 plot(ms, main="Mean air temperatures at Nottingham Castle (1920-1939)",
23      xlab="Month", ylab="Degree Fahrenheit")
```

Figure 14: Source code of “case-nottem.R”.

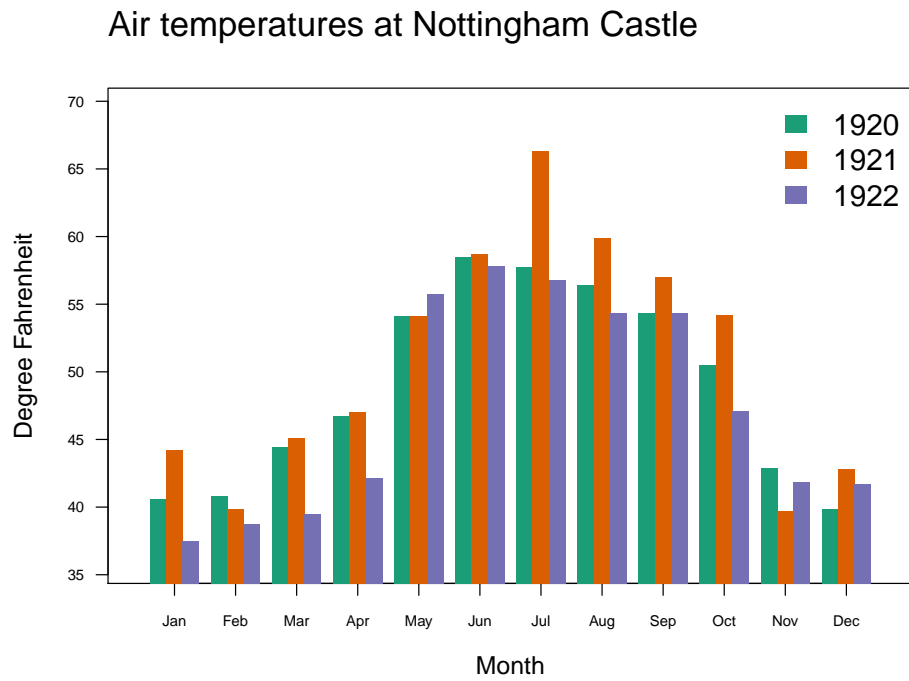


Figure 15: Graphical output 1 of source code “case-nottem.R”.

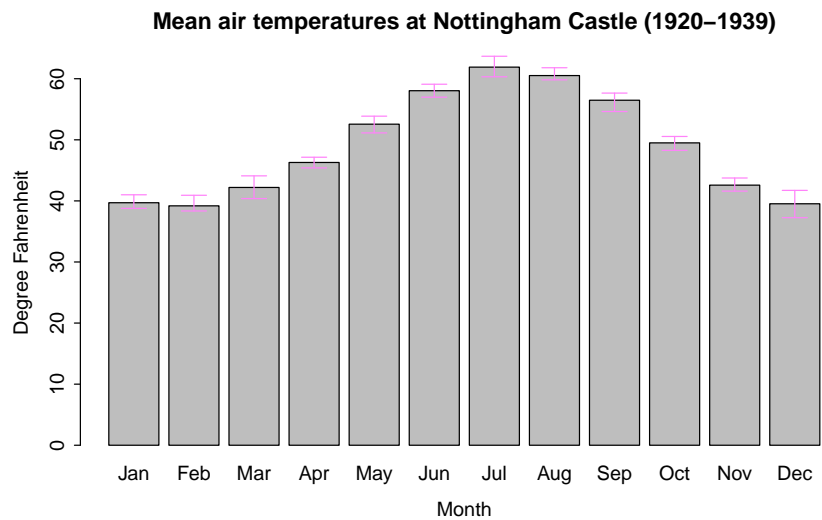


Figure 16: Graphical output 2 of source code “case-nottem.R”.

A.3 gdist: Generic construction and visualization of distances

`example-gdist.R` (fig. 17) is an example on how to measure distances using a user-defined metric, such as *correlation distance* and *GMD*.

```
example-gdist.R
1  ## load library
2  require("GMD")
3  require(cluster)
4
5  ## compute distance using Euclidean metric (default)
6  data(ruspini)
7  x <- gdist(ruspini)
8
9  ## see a dendrogram result by hierarchical clustering
10 dev.new(width=12, height=6)
11 plot(hclust(x),
12      main="Cluster Dendrogram of Ruspini data",
13      xlab="Observations")
14
15 ## convert to a distance matrix
16 m <- as.matrix(x)
17
18 ## convert from a distance matrix
19 d <- as.dist(m)
20 stopifnot(d == x)
21
22 ## Use correlations between variables "as distance"
23 data(USJudgeRatings)
24 dd <- gdist(x=USJudgeRatings,method="correlation.of.variables")
25 dev.new(width=12, height=6)
26 plot(hclust(dd),
27      main="Cluster Dendrogram of USJudgeRatings data",
28      xlab="Variables")
```

Figure 17: Source code of “example-gdist.R”.

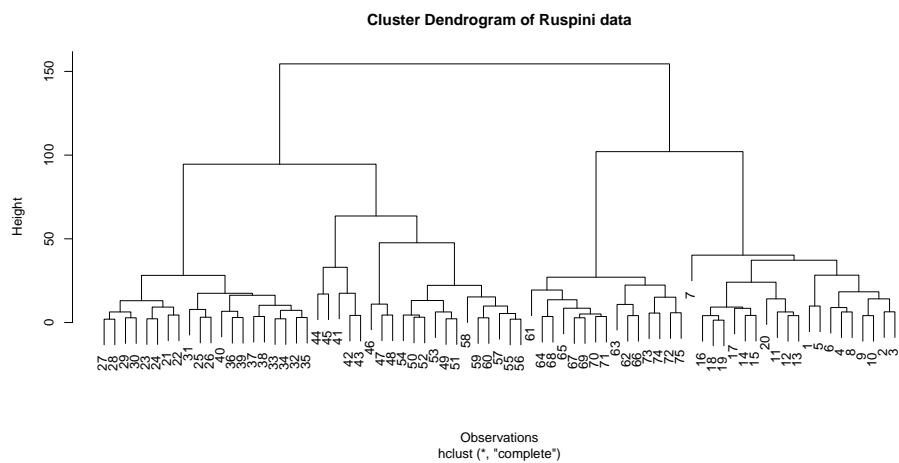


Figure 18: Graphical output 1 of source code “example-gdist.R”.

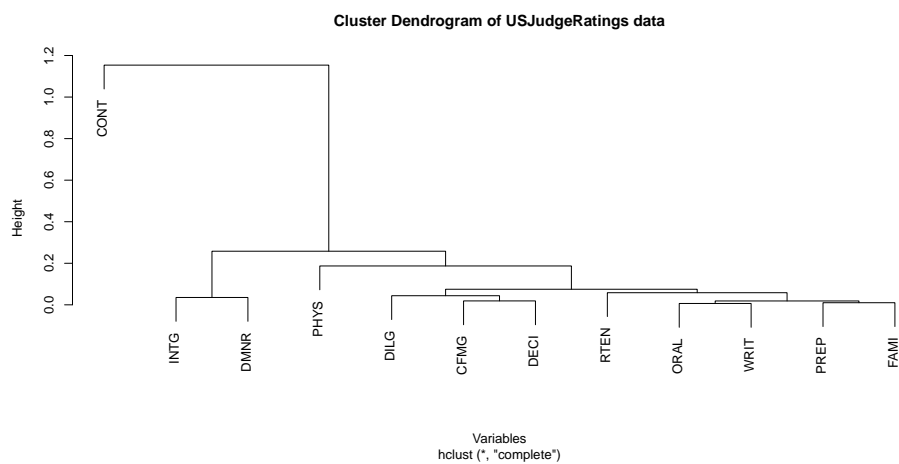


Figure 19: Graphical output 2 of source code “example-gdist.R”.

A.4 `css`: Clustering Sum-of-Squares and the ‘‘elbow’’ plot: determining the number of clusters in a data set

A good clustering yields clusters where the total within-cluster sum-of-squares (WSSs) is small (i.e. cluster cohesion, measuring how closely related are objects in a cluster) and the total between-cluster sum-of-squares (BSSs) is high (i.e. cluster separation, measuring how distinct or well-separated one cluster is from the other).

`example-css.R` (fig. 20) is an example on how to make correct choice of k using “**elbow criterion**”. A good k is selected according a) how much of the total variance in the whole data that the clusters can explain, and b) how large gain in explained variance we obtain by using these many clusters compared to one less or one more, the so-called “elbow” criterion.

The optimal choice of k will strike a balance between maximum compression of the data using a single cluster, and maximum accuracy by assigning each data point to its own cluster. More important, an ideal k should also be relevant in terms of what it reveals about the data, which typically cannot be measured by a metric but by a human expert. Here we present a way to measure such performance of a clustering model, using squared Euclidean distances. The evaluation is based on pairwise distance matrix and therefore more generic in a way that doesn’t involve computating the “centers” of the clusters in the raw data, which are often not available or hard to obtain.


```

1  ## load library
2  require("GMD")
3
4  ## simulate data around 12 points in Euclidean space
5  pointv <- data.frame(x=c(1,2,2,4,4,5,5,6,7,8,9,9),y=c(1,2,8,2,4,4,5,9,9,8,1,9))
6  set.seed(2012)
7  mydata <- c()
8  for (i in 1:nrow(pointv)){
9    mydata <- rbind(mydata,cbind(rnorm(10,pointv[i,1],0.1),rnorm(10,pointv[i,2],0.1)))
10 }
11 mydata <- data.frame(mydata); colnames(mydata) <- c("x","y")
12 plot(mydata,type="p",pch=21, main="Simulated data")
13
14 ## determine a "good" k using elbow
15 dist.obj <- dist(mydata[,1:2])
16 hclust.obj <- hclust(dist.obj)
17 css.obj <- css.hclust(dist.obj,hclust.obj)
18 elbow.obj <- elbow.batch(css.obj)
19 print(elbow.obj)
20
21 ## make partition given the "good" k
22 k <- elbow.obj$k; cutree.obj <- cutree(hclust.obj,k=k)
23 mydata$cluster <- cutree.obj
24
25 ## draw a elbow plot and label the data
26 dev.new(width=12, height=6)
27 par(mfcol=c(1,2),mar=c(4,5,3,3),omi=c(0.75,0,0,0))
28 plot(mydata$x,mydata$y,pch=as.character(mydata$cluster),col=mydata$cluster,cex=0.75,
29      main="Clusters of simulated data")
30 plot.elbow(css.obj,elbow.obj,if.plot.new=FALSE)
31
32 ## clustering with more relaxed thresholds (, resulting a smaller "good" k)
33 elbow.obj2 <- elbow.batch(css.obj,ev.thres=0.90,inc.thres=0.05)
34 mydata$cluster2 <- cutree(hclust.obj,k=elbow.obj2$k)
35
36 dev.new(width=12, height=6)
37 par(mfcol=c(1,2), mar=c(4,5,3,3),omi=c(0.75,0,0,0))
38 plot(mydata$x,mydata$y,pch=as.character(mydata$cluster2),col=mydata$cluster2,cex=0.75,
39      main="Clusters of simulated data")
40 plot.elbow(css.obj,elbow.obj2,if.plot.new=FALSE)

```

Figure 20: Source code of “example-css.R”.

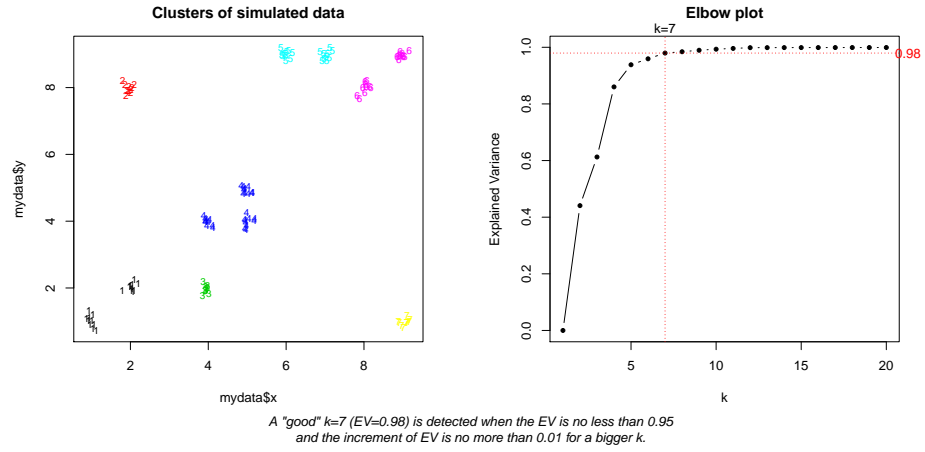


Figure 21: Graphical output 1 of source code “example-css.R”.

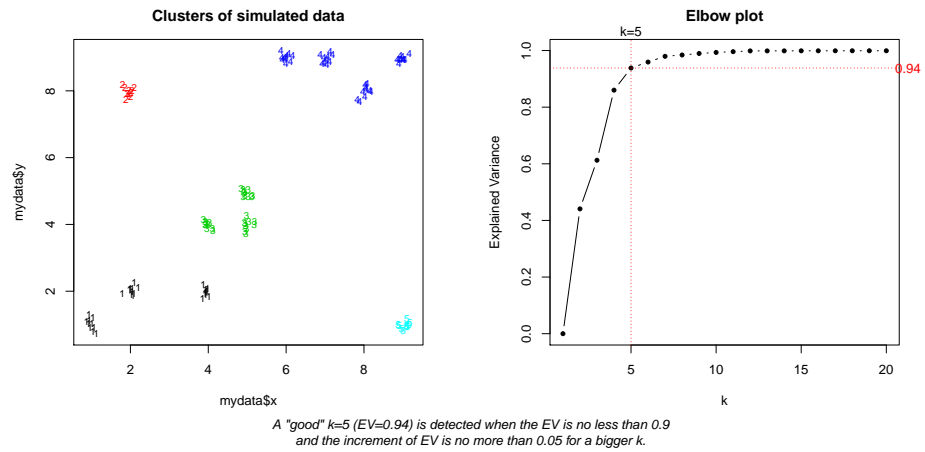


Figure 22: Graphical output 2 of source code “example-css.R”.

A.5 heatmap.3: Visualization in cluster analysis, with evaluation

A.5.1 Examples using mtcars data

`example-heatmap3a.R` (fig. 23) is an example on how to make a heatmap with summary visualization of observations.

```

1  require("GMD")           # load library
2  data(mtcars)             # load data
3  x <- as.matrix(mtcars) # data as a matrix
4
5  dev.new(width=10,height=8)
6  heatmap.3(x)             # default, with reordering and dendrogram
7  heatmap.3(x, Rowv=FALSE, Colv=FALSE) # no reordering and no dendrogram
8  heatmap.3(x, dendrogram="none")      # reordering without dendrogram
9  heatmap.3(x, dendrogram="row")       # row dendrogram with row (and col) reordering
10 heatmap.3(x, dendrogram="row", Colv=FALSE) # row dendrogram with only row reordering
11 heatmap.3(x, dendrogram="col")       # col dendrogram
12 heatmap.3(x, dendrogram="col", Rowv=FALSE) # col dendrogram with only col reordering
13 heatmapOut <-
14   heatmap.3(x, scale="column")        # scaled by column
15 names(heatmapOut)                    # view the list that is returned
16 heatmap.3(x, scale="column", x.center=0) # colors centered around 0
17 heatmap.3(x, scale="column", trace="column") # turn "trace" on
18
19 ## coloring cars (row observations) by brand
20 brands <- sapply(rownames(x), function(e) strsplit(e," ")[[1]][1])
21 names(brands) <- c()
22 brands.index <- as.numeric(as.factor(brands))
23 RowIndividualColors <- rainbow(max(brands.index))[brands.index]
24 heatmap.3(x, scale="column", RowIndividualColors=RowIndividualColors)
25
26 ## coloring attributes (column features) randomly (just for a test :)
27 heatmap.3(x, scale="column", ColIndividualColors=rainbow(ncol(x)))
28
29 ## add a single plot for all row individuals
30 dev.new(width=12,height=8)
31 expr1 <- list(quote(plot(row.data[rowInd,"hp"],1:nrow(row.data),xlab="hp",ylab="",
32   main="Gross horsepower",yaxt="n")),
33   quote(axis(2,1:nrow(row.data),rownames(row.data)[rowInd],las=2)))
34 heatmap.3(x, scale="column", plot.row.individuals=TRUE, row.data=x,
35   plot.row.individuals.list=list(expr1))
36
37 ## add a single plot for all col individuals
38 dev.new(width=12,height=8)
39 expr2 <- list(quote(plot(colMeans(col.data)[colInd], xlab="", ylab="Mean",xaxt="n",
40   main="Mean features",cex=1,pch=19)),
41   quote(axis(1,1:ncol(col.data),colnames(col.data)[colInd],las=2)))
42 heatmap.3(x, scale="column", plot.col.individuals=TRUE, col.data=x,
43   plot.col.individuals.list=list(expr2))
44
45 ## add another single plot for all col individuals
46 dev.new(width=12,height=8)
47 expr3 <- list(quote(op <- par(mar = par("mar")*1.5)),
48   quote(mytmp.data <- apply(col.data,2,function(e) e/sum(e))),
49   quote(boxplot(mytmp.data[,colInd], xlab="", ylab="Value",
50   main="Boxplot of normalized column features")),
51   quote(par(op)))
52 heatmap.3(x, scale="column", plot.col.individuals=TRUE, col.data=x,
53   plot.col.individuals.list=list(expr3))

```

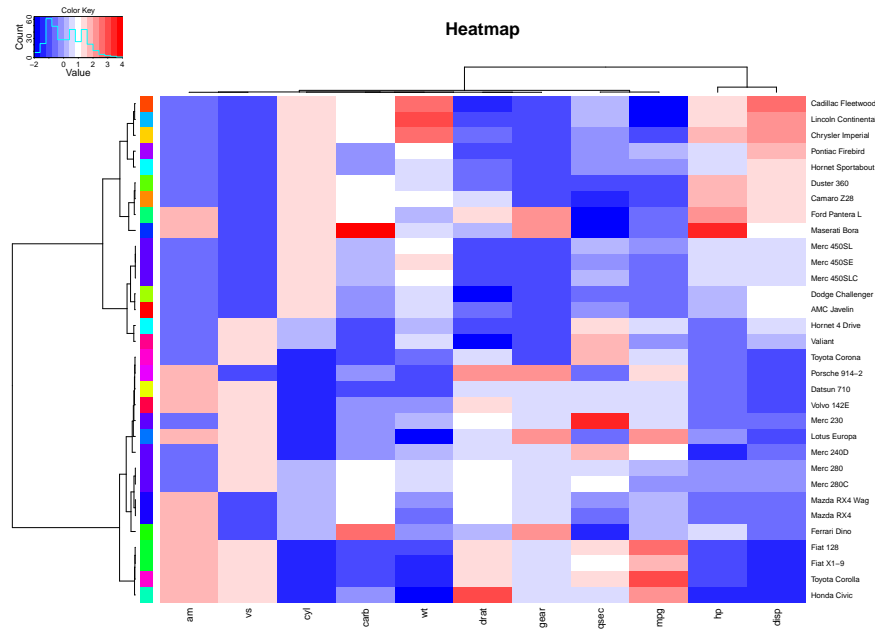


Figure 24: Graphical output 1 of source code “example-heatmap3a.R”.

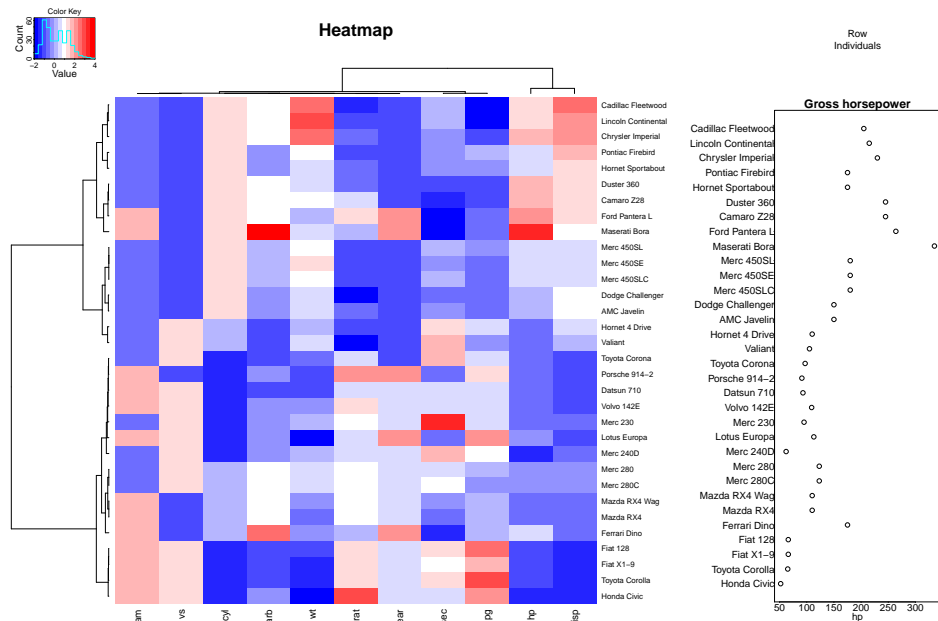


Figure 25: Graphical output 2 of source code “example-heatmap3a.R”.

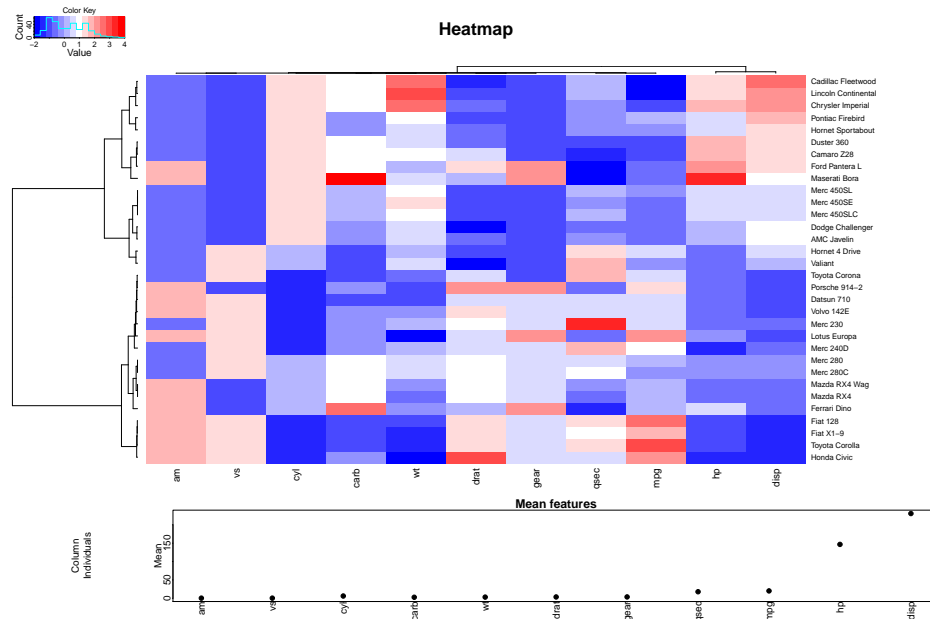


Figure 26: Graphical output 3 of source code “example-heatmap3a.R”.

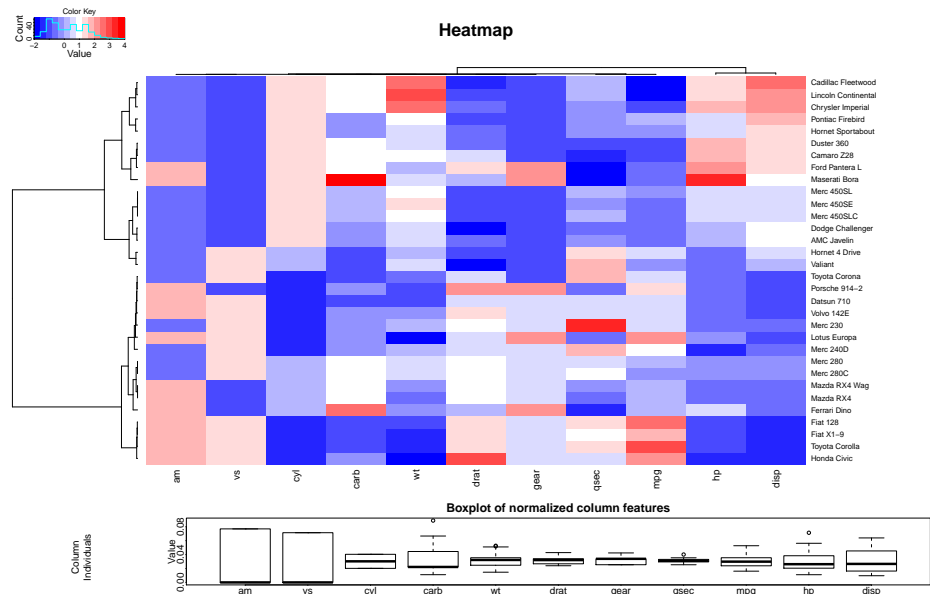


Figure 27: Graphical output 4 of source code “example-heatmap3a.R”.

A.5.2 Examples using ruspini data

example-heatmap3b.R (fig. 28) is an example on how to make a heatmap with summary visualization of clusters.

```
example-heatmap3b.R
1  ## load library
2  require("GMD")
3  require(cluster)
4
5  ## load data
6  data(ruspini)
7
8  ## heatmap on a `dist' object
9  x <- gdist(ruspini)
10 main <- "Heatmap of Ruspini data"
11 dev.new(width=10,height=10)
12 heatmap.3(x, main=main, mapratio=1) # default with a title and a map in square!
13 heatmap.3(x, main=main, revC=TRUE) # reverse column for a symmetric look
14 heatmap.3(x, main=main, kr=2, kc=2) # show partition by predefined number of clusters
15
16 ## show partition by elbow
17 css.multi.obj <- css.hclust(x,hclust(x))
18 elbow.obj <- elbow.batch(css.multi.obj,ev.thres=0.90,inc.thres=0.05)
19 heatmap.3(x, main=main, revC=TRUE, kr=elbow.obj$k, kc=elbow.obj$k)
20 heatmap.3(x, main=main, sub=sub("\n"," ",attr(elbow.obj,"description")), cex.sub=1.25,
21           revC=TRUE,kr=elbow.obj$k, kc=elbow.obj$k) # show elbow info as subtitle
22
23 ## side plot for every row clusters
24 dev.new(width=10,height=10)
25 expr1 <- list(quote(plot(do.call(rbind,i.x),xlab="x",ylab="y",
26                               xlim=range(ruspini$x),ylim=range(ruspini$y),)))
27 heatmap.3(x, main=main, revC=TRUE, kr=elbow.obj$k, kc=elbow.obj$k, trace="none",
28           row.data=as.list(data.frame(t(ruspini))),
29           plot.row.clusters=TRUE,plot.row.clusters.list=list(expr1))
30
31 ## side plot for every col clusters
32 dev.new(width=10,height=10)
33 expr2 <- list(quote(plot(do.call(rbind,i.x),xlab="x",ylab="y",
34                               xlim=range(ruspini$x),ylim=range(ruspini$y),)))
35 heatmap.3(x, main=main, revC=TRUE, kr=elbow.obj$k, kc=elbow.obj$k, trace="none",
36           col.data=as.list(data.frame(t(ruspini))),
37           plot.col.clusters=TRUE,plot.col.clusters.list=list(expr2))
38
```

Figure 28: Source code of “example-heatmap3b.R”.

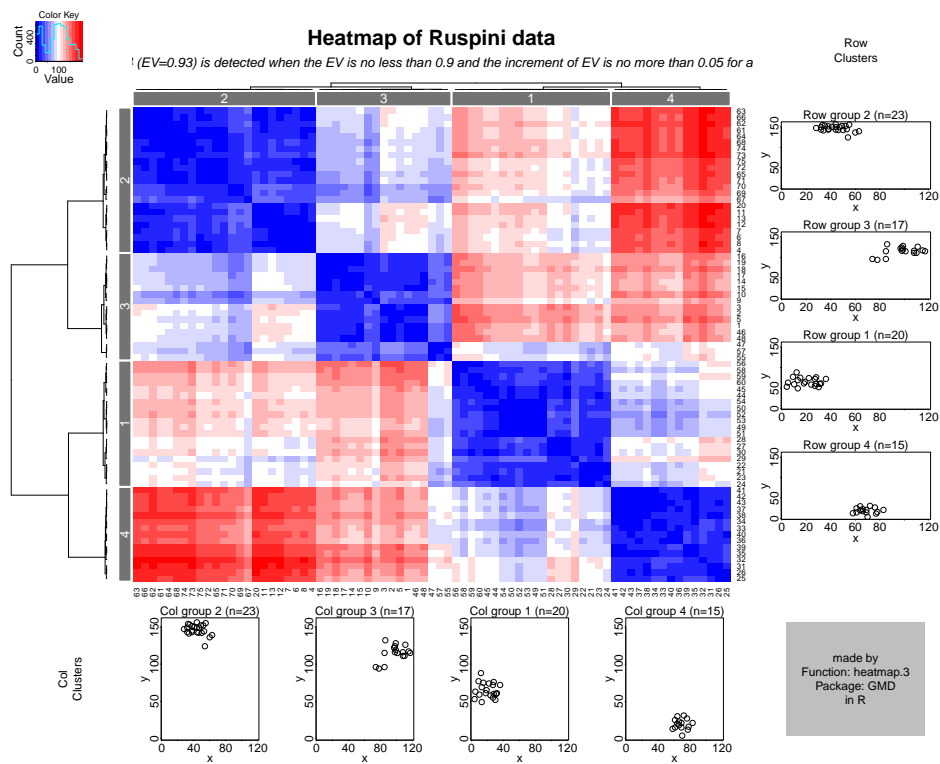


Figure 29: Graphical output 1 of source code “example-heatmap3b.R”.

B Data

B.1 GMD dataset overview

```
> data(package="GMD")
```

```
Data sets in package 'GMD':
```

```
cage          CAGE Data
cagel         CAGE Data
chipseq_hCD4T ChIP-seq Data
chipseq_mES   ChIP-seq Data
```

B.2 CAGE data: cage and cagel

```
> help(cage)
```

```
> data(cage)
> class(cage)
```

```
[1] "list"
```

```
> length(cage)
```

```
[1] 20
```

```
> names(cage)
```

```
[1] "NA (T01F029805F8)"      "Glul (T01F092C2995)"
[3] "Cyp4a14 (T04R06C91673)"  "Stxbp4 (T11R05607FD4)"
[5] "D230039L06Rik (T01FOAA465EB)" "Gas5 (T01F09995479)"
[7] "Rab5c (T11R05FBC6C4)"    "BC003940 (T11R072A6CE0)"
[9] "Tpt1 (T14F04079189)"    "Pcna (T02R07DE319B)"
[11] "DOH4S114 (T18R020553F0)" "Gstol (T19F02D03566)"
[13] "Hsd11b1 (T01R0B8305BD)"  "Csf1 (T03R0672174D)"
[15] "B2m (T02F0743FF05)"     "Alox5ap (T05F08BCF2C4)"
[17] "Pfkfb3 (T02R00AEC2D8)"   "Hig1 (T09R0743763C)"
[19] "Cd72 (T04R028B8BC9)"    "Egln1 (T08R0769239F)"
```

```
> data(cagel)
> names(cagel)
```

```
[1] "NA (T17F05912B83)"      "Tpt1 (T14F04079189)"
[3] "B2m (T02F0743FF05)"     "Grn (T11F0615F289)"
[5] "2600001A11Rik (T12R043A2595)" "Rbbp7 (TOXF091A7ACA)"
[7] "Rpl141 (T10R07AB7138)"   "H2afy (T13R034ACF47)"
[9] "Dscr111 (T17F02802885)"  "Ckap4 (T10R0504CE97)"
[11] "Rab5c (T11R05FBC6C4)"    "Pfkfb3 (T02R00AEC2D8)"
[13] "Csf1 (T03R0672174D)"     "Ctsb (T14F0348EDBA)"
[15] "Crim1 (T17F04928998)"     "3930401E15Rik (T18R02CDD141)"
[17] "Rai17 (T14F014BF473)"    "Hig1 (T09R0743763C)"
[19] "Apbb2 (T05R03E329C8)"    "Ptn (T06R0230806E)"
[21] "Tneff2 (T01F030EC757)"   "Mrps6 (T16F0583C906)"
[23] "Hsd11b1 (T01R0B8305BD)"  "DOH4S114 (T18R020553F0)"
[25] "4930430J02Rik (T09F036E80C6)" "Phtf2 (T05R0125E896)"
[27] "Trpv2 (T11F03B4EBD8)"    "Slco3a1 (T07R03AC06B9)"
[29] "Scd1 (T19R029B5186)"     "Ctxn (T08R0040864A)"
[31] "5730596K20Rik (T19F006DFC1A)" "Arbp (T05F06BBE13B)"
[33] "Klh15 (T05F03CCA673)"    "Gstol (T19F02D03566)"
[35] "NA (T02R07EF5EDA)"       "Srpk1 (T17R019F4A41)"
[37] "Nudt7 (T08F06C3B651)"    "Tnfrsf10b (T14F03AB1306)"
[39] "Egln1 (T08R0769239F)"    "9630050M13Rik (T02F002EC972)"
[41] "BC003940 (T11R072A6CE0)" "Ppia (T11F00604AFF)"
[43] "Alox5ap (T05F08BCF2C4)"  "Pcna (T02R07DE319B)"
[45] "Gch1 (T14R02602138)"     "Yap1 (T09R0079F3FF)"
[47] "Vrk1 (T12F06010C9B)"     "Cd72 (T04R028B8BC9)"
[49] "Wdtdc1 (T04R07DAFEDC)"    "Centg2 (T01F055392D1)"
```

B.3 ChIP-seq data: chipseq_mES and chipseq_hCD4T

```
> help(chipseq)

> data(chipseq_mES)
> class(chipseq_mES)

[1] "list"

> length(chipseq_mES)

[1] 6

> names(chipseq_mES)

[1] "H3K27me3" "H3K36me3" "H3K4me1" "H3K4me2" "H3K4me3" "H3K9me3"

> data(chipseq_hCD4T)
> names(chipseq_hCD4T)

[1] "CTCF"      "H2AK5ac"    "H2AK9ac"    "H2AZ"      "H2BK120ac"  "H2BK12ac"
[7] "H2BK20ac"  "H2BK5ac"    "H2BK5me1"   "H3K14ac"   "H3K18ac"    "H3K23ac"
[13] "H3K27ac"   "H3K27me1"   "H3K27me2"   "H3K27me3"  "H3K36ac"    "H3K36me1"
[19] "H3K36me3"  "H3K4ac"     "H3K4me1"    "H3K4me2"   "H3K4me3"    "H3K79me1"
[25] "H3K79me2"  "H3K79me3"   "H3K9ac"     "H3K9me1"   "H3K9me2"    "H3K9me3"
[31] "H3R2me1"   "H3R2me2"    "H4K12ac"    "H4K16ac"   "H4K20me1"   "H4K20me3"
[37] "H4K5ac"    "H4K8ac"     "H4K91ac"    "H4R3me2"
```

References

- [1] Artem Barski, Suresh Cuddapah, Kairong Cui, Tae-Young Roh, Dustin E Schones, Zhibin Wang, Gang Wei, Iouri Chepelev, and Keji Zhao. High-resolution profiling of histone methylations in the human genome. *Cell*, 129(4):823–837, May 2007.
- [2] Piero Carninci, Albin Sandelin, Boris Lenhard, Shintaro Katayama, Kazuro Shimokawa, Jasmina Ponjavic, Colin A M Semple, Martin S Taylor, Pr G Engstrm, Martin C Frith, Alistair R R Forrest, Wynand B Alkema, Sin Lam Tan, Charles Plessy, Rimantas Kodzius, Timothy Ravasi, Takeya Kasukawa, Shiro Fukuda, Mutsumi Kanamori-Katayama, Yayoi Kitazume, Hideya Kawaji, Chikatoshi Kai, Mari Nakamura, Hideaki Konno, Kenji Nakano, Salim Mottagui-Tabar, Peter Arner, Alessandra Chesi, Stefano Gustincich, Francesca Persichetti, Harukazu Suzuki, Sean M Grimmond, Christine A Wells, Valerio Orlando, Claes Wahlestedt, Edison T Liu, Matthias Harbers, Jun Kawai, Vladimir B Bajic, David A Hume, and Yoshihide Hayashizaki. Genome-wide analysis of mammalian promoter architecture and evolution. *Nat Genet*, 38(6):626–635, Jun 2006.
- [3] Sung-Hyuk Cha and Sargur N. Srihari. On measuring the distance between histograms. *Pattern Recognition*, 35(6):1355–1370, 2002.
- [4] Tarjei S Mikkelsen, Manching Ku, David B Jaffe, Biju Issac, Erez Lieberman, Georgia Giannoukos, Pablo Alvarez, William Brockman, Tae-Kyung Kim, Richard P Koche, William Lee, Eric Mendenhall, Aisling O’Donovan, Aviva Presser, Carsten Russ, Xiaohui Xie, Alexander Meissner, Marius Wernig, Rudolf Jaenisch, Chad Nusbaum, Eric S Lander, and Bradley E Bernstein. Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature*, 448(7153):553–560, Aug 2007.
- [5] R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2011. ISBN 3-900051-07-0.
- [6] Albin Sandelin, Piero Carninci, Boris Lenhard, Jasmina Ponjavic, Yoshihide Hayashizaki, and David A Hume. Mammalian rna polymerase ii core promoters: insights from genome-wide studies. *Nat Rev Genet*, 8(6):424–436, Jun 2007.
- [7] Xiaobei Zhao, Eivind Valen, Brian J Parker, and Albin Sandelin. Systematic clustering of transcription start site landscapes. *PLoS ONE*, 6(8):e23409, August 2011.
- [8] Vicky W Zhou, Alon Goren, and Bradley E Bernstein. Charting histone modifications and the functional organization of mammalian genomes. *Nat Rev Genet*, 12(1):7–18, Jan 2011.