

# Small RNA dataset

Luis Ángel Rodríguez García

25-05-2022

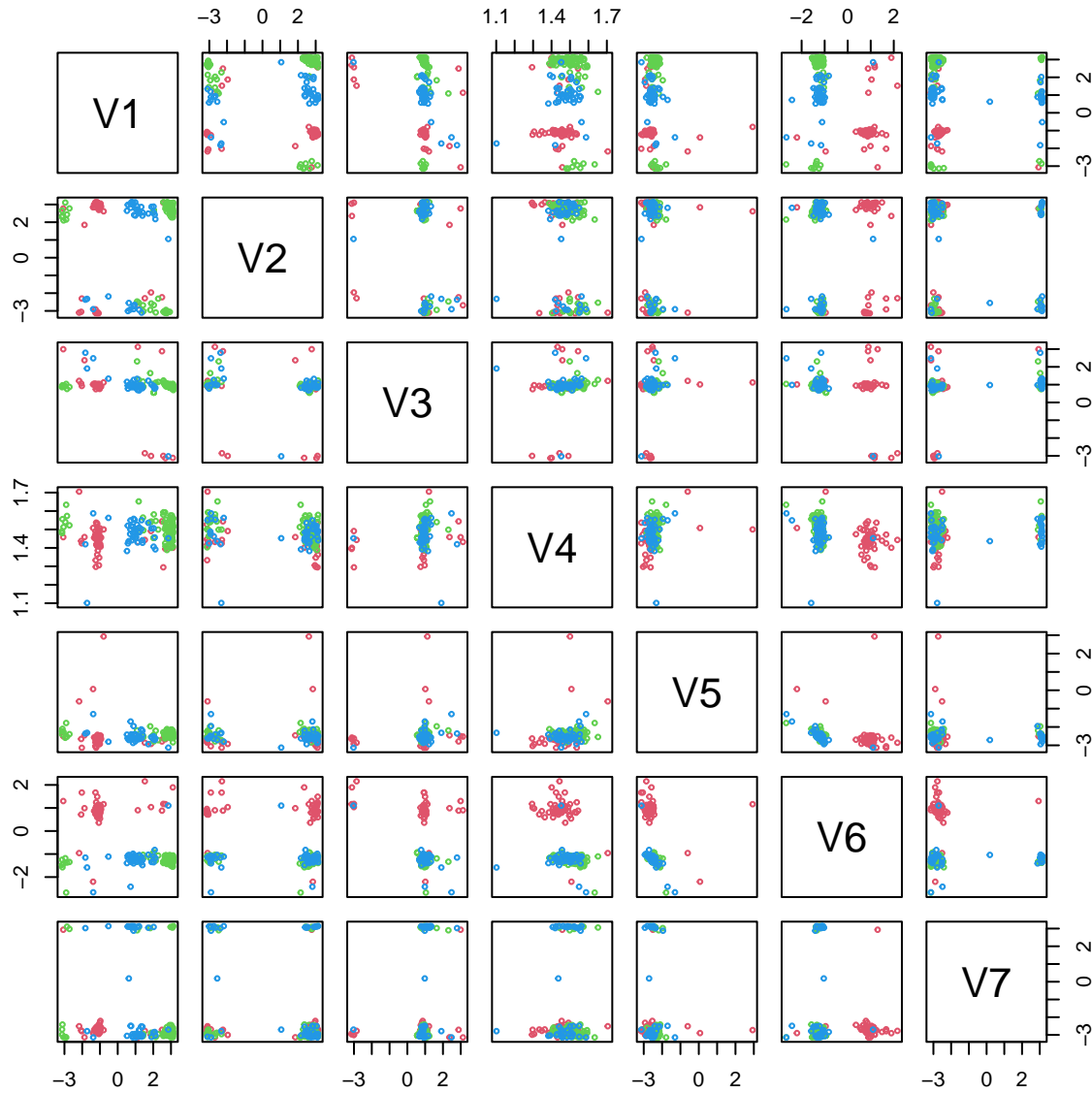
The objective in this case study is to recover clusters identified in `smallrna$clusters` using only the information on `smallrna$angles`, a 7-dimensional matrix of angles (i.e., data on  $(\mathbb{S}^1)^7$ ). The clusters have been constructed using the information in `smallrna$torsion`. If a dimension-reduction technique is able to successfully identify clusters, then it will be doing a good job in terms of identifying the underlying structure of the data. Section 5.2 in Zoubouoglou et al. (2021) describes the history of the “Small RNA” dataset and its construction.

Let’s begin by importing the data.

```
# Load package
library(pscsne)
stopifnot(packageVersion("pscsne") >= "0.0.1.900002")

# Load dataset
data("smallrna")

# Visualize original data
pairs(smallrna$angles, col = smallrna$clusters + 1, cex = 0.5)
```



We can now run psc-SNE. First, we transform the data and obtain the  $\rho$ 's giving the prescribed perplexity.

```
# Data to Cartesian coordinates
smallrna_X <- sphunif::Theta_to_X(Theta = smallrna$angles)

# Obtain rhos for given perplexity
rho_psc_list <- rho_optim_bst(x = smallrna_X, perp_fixed = 30)
```

```
## Time difference of 10.1769 secs
```

We run psc-SNE for  $d = 1$  with its default  $\eta$ .

```
# Default
fit_1 <- psc_sne(X = smallrna_X, d = 1, rho_psc_list = rho_psc_list,
  eta = 200, maxit = 1e3, tol = 1e-6, show_prog = 10,
  colors = smallrna$clusters)
```

```
## It: 1; obj: 1.196e+01; abs: 0.000e+00; rel: 0.000e+00; norm: 3.156e-01; mom: 0.000e+00;
## best it: 1; best obj: 1.196e+01
```

```
## It: 100; obj: 1.048e+01; abs: 1.892e+00; rel: 1.529e-01; norm: 2.554e-01; mom: 7.510e+00;
```

```

## best it: 14; best obj: 1.027e+01
## It: 200; obj: 1.150e+00; abs: 1.523e-02; rel: 1.342e-02; norm: 8.878e-02; mom: 4.454e+00;
## best it: 101; best obj: 1.078e+00

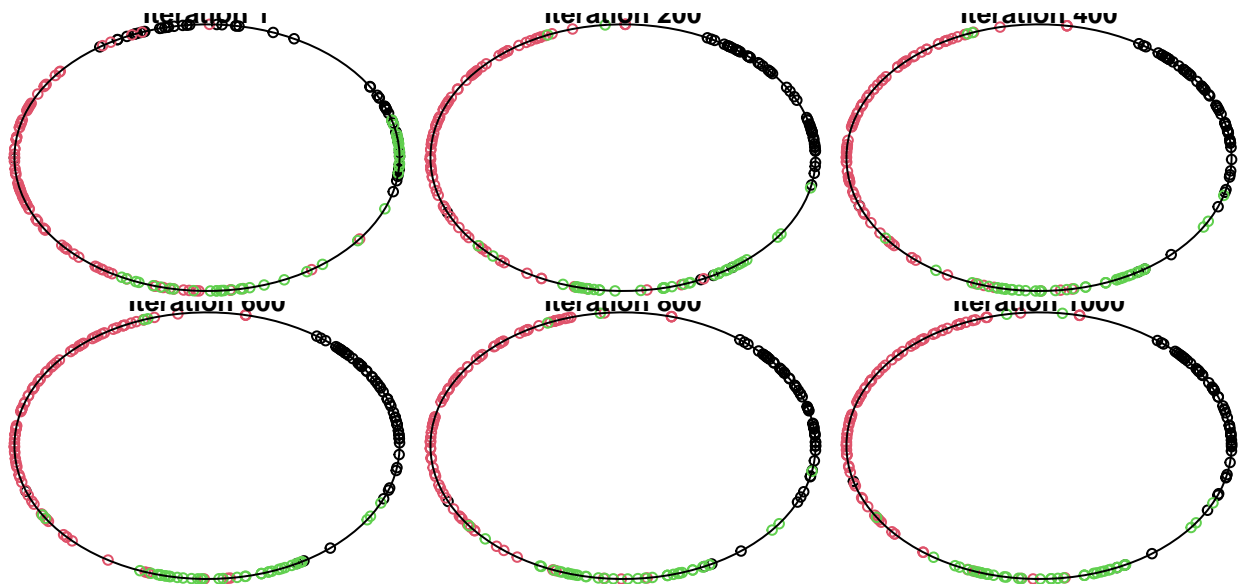
## It: 300; obj: 1.084e+00; abs: 3.854e-03; rel: 3.542e-03; norm: 8.592e-02; mom: 6.150e+00;
## best it: 248; best obj: 1.022e+00
## It: 400; obj: 1.088e+00; abs: 1.494e-03; rel: 1.371e-03; norm: 8.640e-02; mom: 6.108e+00;
## best it: 248; best obj: 1.022e+00

## It: 500; obj: 1.082e+00; abs: 1.050e-02; rel: 9.610e-03; norm: 8.634e-02; mom: 6.133e+00;
## best it: 248; best obj: 1.022e+00
## It: 600; obj: 1.083e+00; abs: 1.018e-02; rel: 9.309e-03; norm: 8.623e-02; mom: 6.173e+00;
## best it: 248; best obj: 1.022e+00

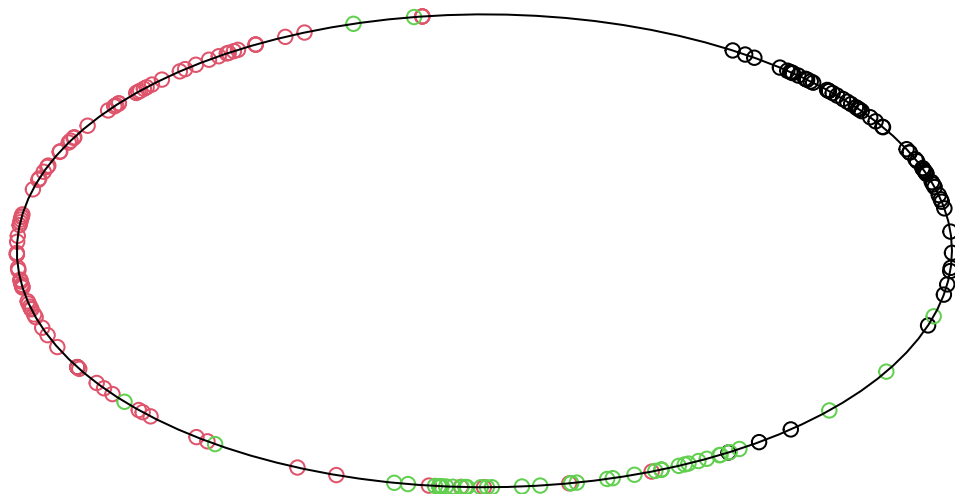
## It: 700; obj: 1.088e+00; abs: 5.319e-04; rel: 4.885e-04; norm: 8.626e-02; mom: 6.081e+00;
## best it: 248; best obj: 1.022e+00
## It: 800; obj: 1.085e+00; abs: 8.191e-03; rel: 7.495e-03; norm: 8.668e-02; mom: 6.122e+00;
## best it: 248; best obj: 1.022e+00

## It: 900; obj: 1.082e+00; abs: 7.312e-03; rel: 6.712e-03; norm: 8.591e-02; mom: 6.155e+00;
## best it: 248; best obj: 1.022e+00
## It: 1000; obj: 1.089e+00; abs: 2.103e-04; rel: 1.932e-04; norm: 8.629e-02; mom: 6.071e+00;
## best it: 248; best obj: 1.022e+00

```



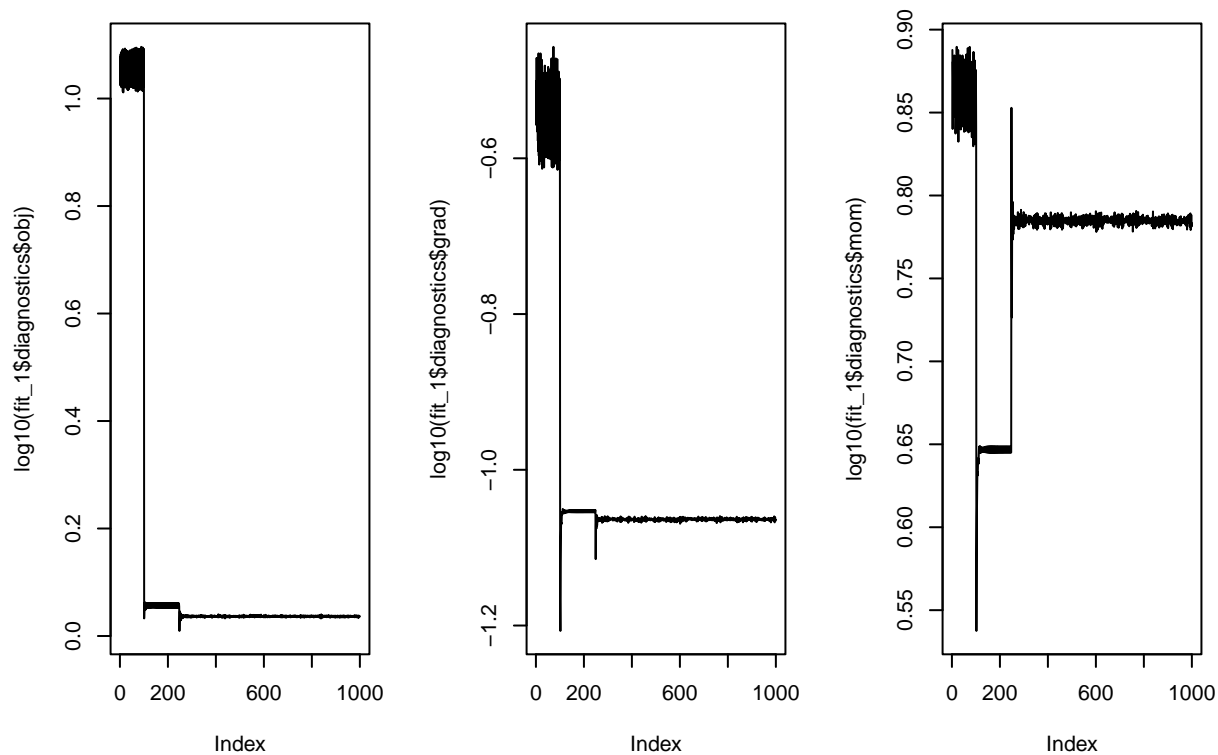
## Iteration 248



```
# Does not converge
fit_1$convergence
```

```
## [1] FALSE
```

```
par(mfrow = c(1, 3))
plot(log10(fit_1$diagnostics$obj), type = "l")
plot(log10(fit_1$diagnostics$grad), type = "l")
plot(log10(fit_1$diagnostics$mom), type = "l")
```



The employed  $\eta$  seems to be too large. Let's reduce it.

```
# Lower eta
```

```
fit_2 <- psc_sne(X = smallrna_X, d = 1, rho_psc_list = rho_psc_list,
  eta = 10, maxit = 1e3, tol = 1e-6, show_prog = 10,
  colors = smallrna$clusters)
```

```
## It: 1; obj: 1.026e+01; abs: 0.000e+00; rel: 0.000e+00; norm: 3.156e-01; mom: 0.000e+00;
```

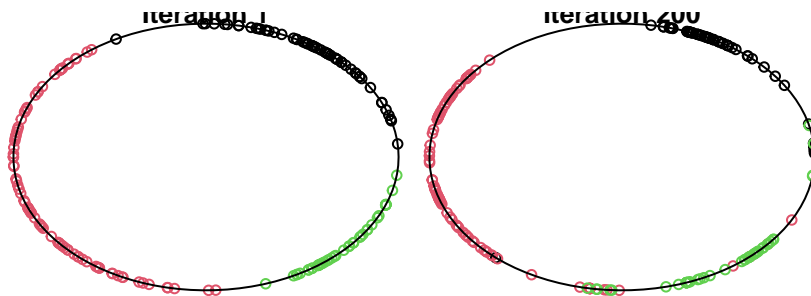
```
## best it: 1; best obj: 1.026e+01
```

```
## It: 100; obj: 9.646e+00; abs: 3.490e-04; rel: 3.618e-05; norm: 9.664e-04; mom: 9.981e-03;
```

```
## best it: 13; best obj: 9.567e+00
```

```
## It: 200; obj: 8.597e-01; abs: 4.874e-10; rel: 5.669e-10; norm: 4.820e-06; mom: 5.699e-05;
```

```
## best it: 200; best obj: 8.597e-01
```

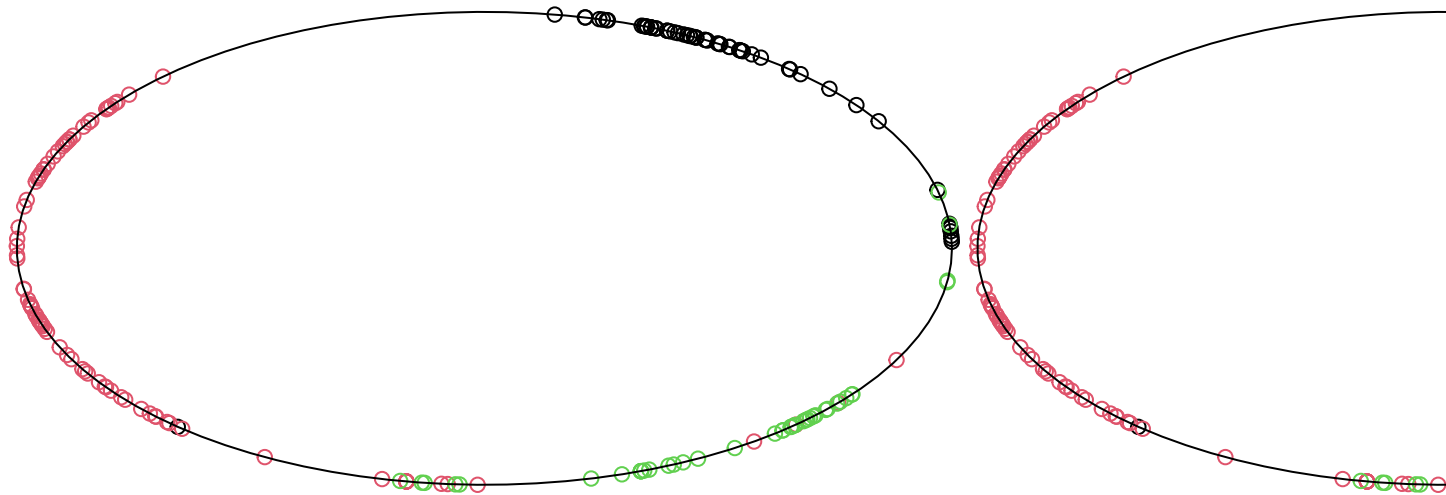


```
## It: 227; obj: 8.597e-01; abs: 1.978e-11; rel: 2.301e-11; norm: 9.828e-07; mom: 1.072e-05;
```

```
## best it: 227; best obj: 8.597e-01
```

**Iteration 227**

**Iteration**



```
# Converges
```

```
fit_2$convergence
```

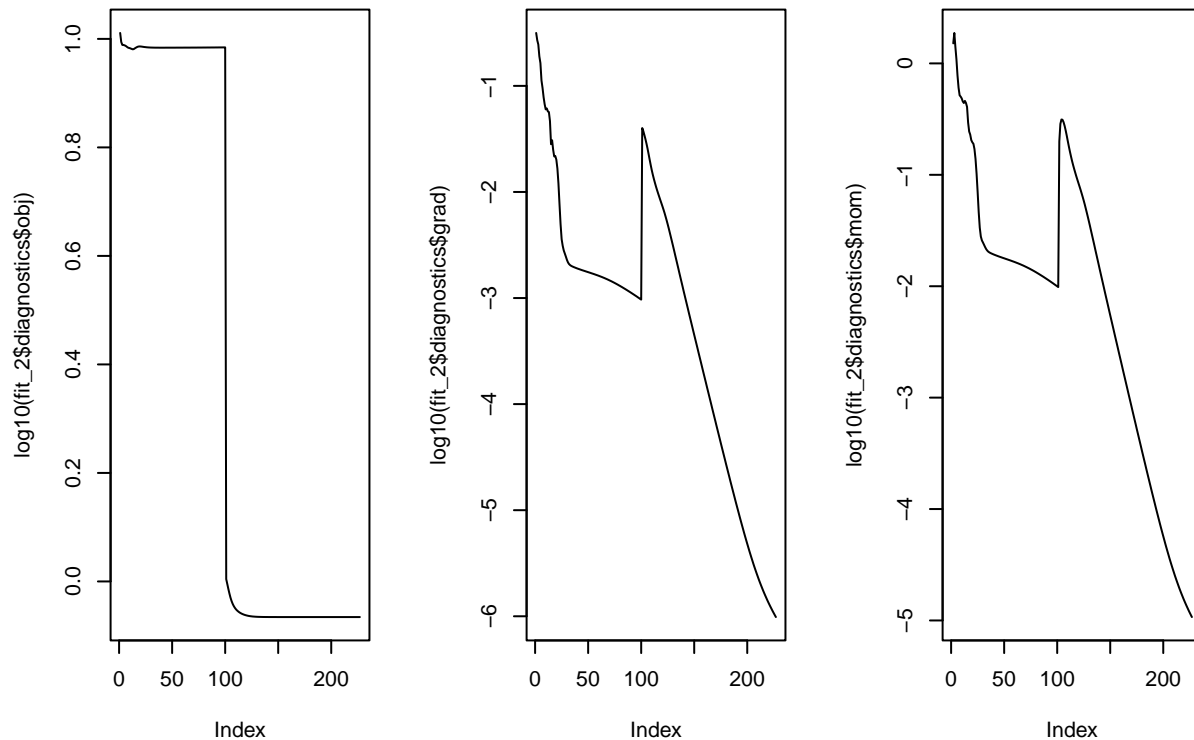
```
## [1] TRUE
```

```
par(mfrow = c(1, 3))
```

```
plot(log10(fit_2$diagnostics$obj), type = "l")
```

```
plot(log10(fit_2$diagnostics$grad), type = "l")
```

```
plot(log10(fit_2$diagnostics$mom), type = "l")
```



Convergence is attained in the second run, yet it is weird that the objective function takes exactly the zero value.

Let's see the recovery of the clusters.

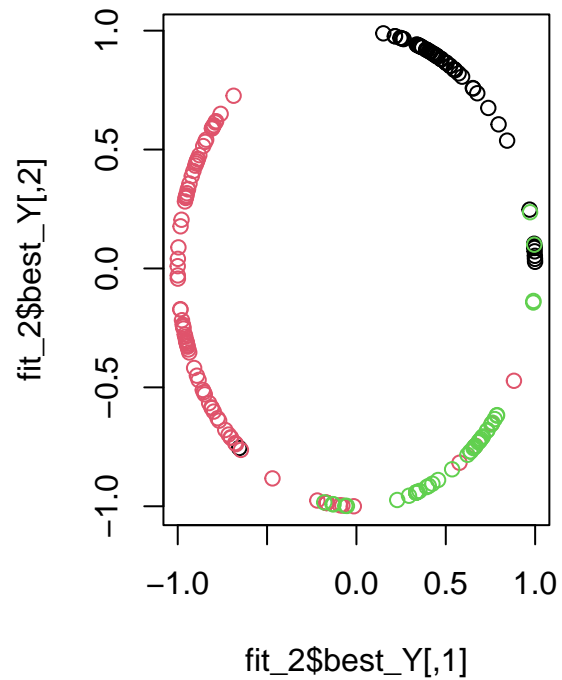
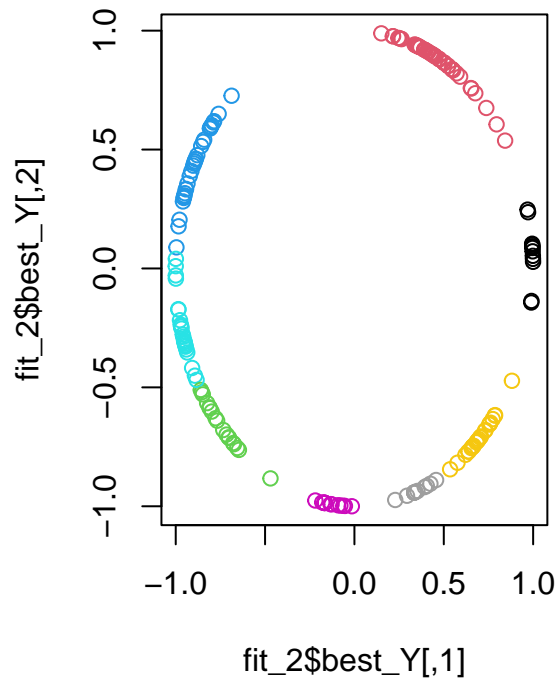
```
# Kernel mean shift clustering
n <- nrow(fit_2$best_Y)
d <- ncol(fit_2$best_Y) - 1
fit_mix <- DirStats::bic_vmf_mix(fit_2$best_Y, kappa_max = 1e3)
h <- DirStats::bw_dir_emi(data = fit_2$best_Y, fit_mix = fit_mix)$h_opt *
  n^(1 / (d + 4)) * n^(-1 / (d + 6))
kms <- kms_dir(x = fit_2$best_Y, data = fit_2$best_Y, h = h)
```

```
## |
```

```
# Detects 8 clusters by splitting the 3 original clusters
length(unique(kms$cluster))
```

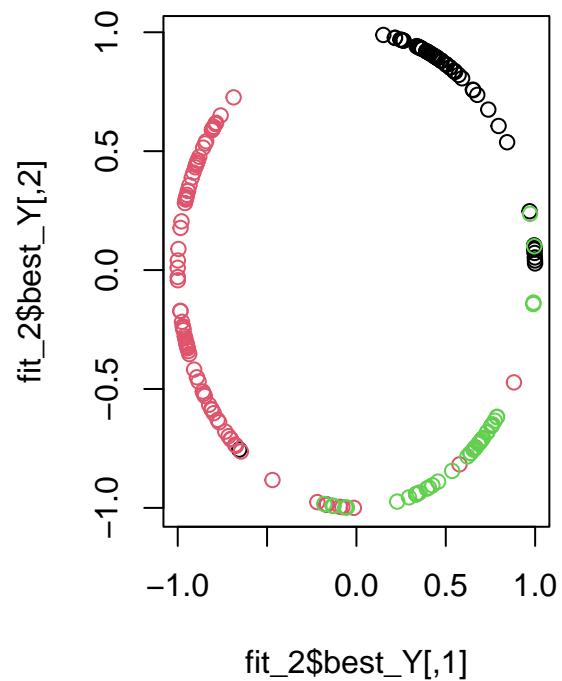
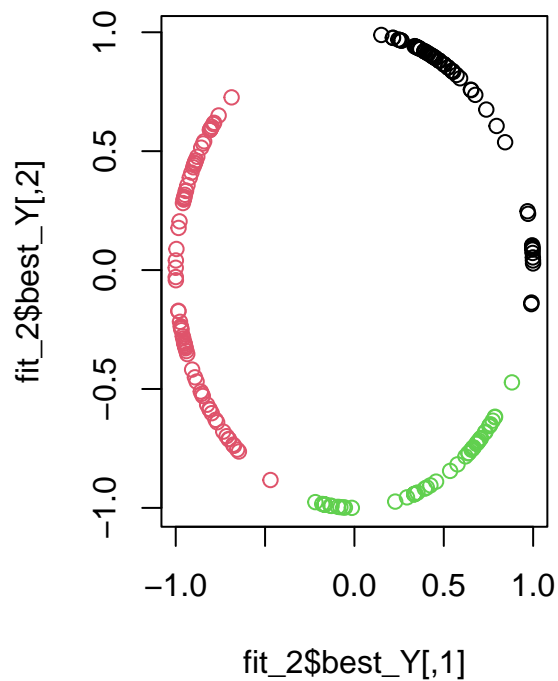
```
## [1] 8
```

```
# Fully-automatically recovered clusters vs. real clusters
par(mfrow = c(1, 2))
plot(fit_2$best_Y, col = kms$cluster)
plot(fit_2$best_Y, col = smallrna$clusters)
```



The original clusters are not fully recovered, in the sense that more clusters are obtained. However, the three-cluster structure is present, as the new clusters appear dividing the three main ones. This can be checked by cutting the hierarchical clustering tree behind kernel mean shift clustering exactly at three groups. Or, in other words, by merging the 8 groups into 3.

```
# Recovered clusters with three clusters vs. real clusters
par(mfrow = c(1, 2))
labels <- cutree(kms$tree, k = 3)
plot(fit_2$best_Y, col = labels)
plot(fit_2$best_Y, col = smallrna$clusters)
```



```
# Correct classification rate: 90%  
mean(labels == smallrna$clusters)
```

```
## [1] 0.9
```

```
# 19 incorrectly classified observations  
sum(labels != smallrna$clusters)
```

```
## [1] 19
```

The classification accuracy is on-par with Zoubouloglou et al. (2021), which misclassifies 16 points and has a classification rate of 0.916.