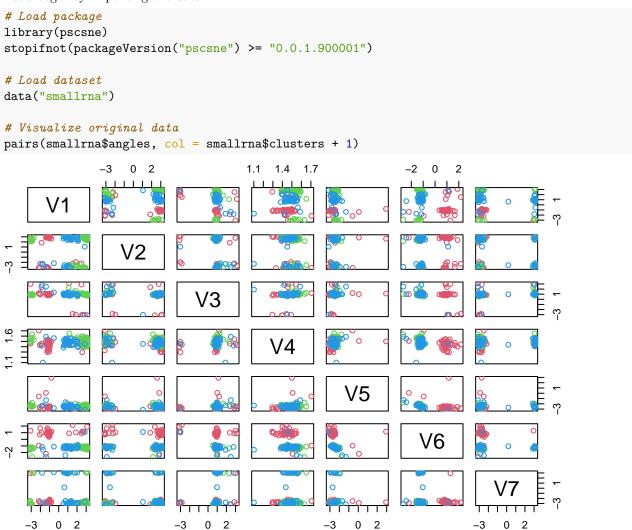
Small RNA dataset

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 Zoubouloglou et al. (2021) describes the history of the "Small RNA" dataset and its construction.

Let's begin by importing the data.



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

```
# Data to Cartesian coordinates
smallrna_X <- sphunif::Theta_to_X(Theta = smallrna$angles)</pre>
# Obtain rhos for given perplexity
rho_psc_list <- rho_optim_bst(x = smallrna_X, perp_fixed = 30)</pre>
## Time difference of 10.45816 secs
We run psc-SNE for d=1 with its default \eta.
# Default
void_1 <- capture.output(</pre>
  fit_1 <- psc_sne(X = smallrna_X, d = 1, rho_psc_list = rho_psc_list,</pre>
                       eta = 200, num_iteration = 1e3, tol = 1e-6,
                       show_prog = FALSE, colors = smallrna$clusters))
# 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$grad), type = "1")
plot(log10(fit_1$diagnostics$mom), type = "1")
    1.0
                                                                              0.85
                                                                              0.80
    0.8
                                                                          og10(fit_1$diagnostics$mom)
                                    og10(fit_1$diagnostics$grad)
og10(fit_1$diagnostics$obj)
                                                                              0.75
                                         9.0
    9.0
                                                                              0.70
    0.4
                                                                              0.65
                                         -1.0
                                                                              0.60
    0.2
                                                                              0.55
                                         -1.2
            200
                    600
                             1000
                                                 200
                                                                  1000
                                                                                      200
                                                                                              600
                                                                                                       1000
         0
                                              0
                                                          600
                                                                                   0
                  Index
                                                       Index
                                                                                            Index
```

The employed η seems to be too large. Let's reduce it.

```
show_prog = FALSE, colors = smallrna$clusters))
# Converges
fit_2$convergence
## [1] TRUE
par(mfrow = c(1, 3))
plot(log10(fit_2$diagnostics$obj), type = "1")
plot(log10(fit_2$diagnostics$grad), type = "1")
plot(log10(fit_2$diagnostics$mom), type = "1")
     1.0
                                              ī
                                                                                       7
                                              7
                                                                                 log10(fit_2$diagnostics$mom)
                                         log10(fit_2$diagnostics$grad)
log10(fit_2$diagnostics$obj)
     9.0
                                                                                       7
                                              _{\rm L}
     0.4
                                                                                       က
                                              4
     0.2
                                              2
                                                                                       4
     0.0
                                              ၯ
                                                                                       2
          0
              50
                   100
                              200
                                                   0
                                                       50
                                                            100
                                                                       200
                                                                                            0
                                                                                                50
                                                                                                     100
                                                                                                               200
                                                             Index
                                                                                                      Index
```

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

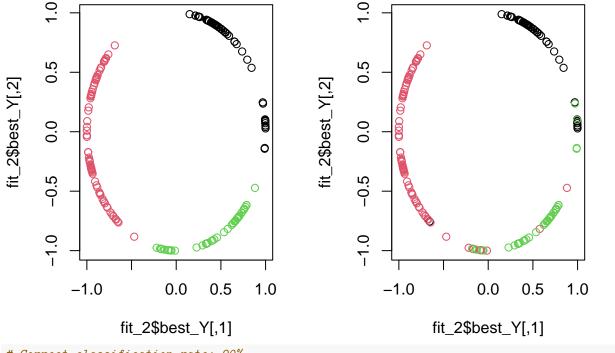
Let's see the recovery of the clusters.

[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)
      0.5
                                                           0.5
fit_2$best_Y[,2]
                                                    fit_2$best_Y[,2]
                                           0
      0.0
                                                           0.0
                                           0
      -0.5
                                                           -0.5
      -1.0
                                                           -1.0
                                                                              -1.0
                                                                -1.0
                           0.0
                                  0.5
                                          1.0
                                                                                0.0
                                                                                       0.5
                                                                                               1.0
                   fit_2$best_Y[,1]
                                                                        fit_2$best_Y[,1]
```

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)</pre>
```



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.