

Statistisches Data Mining (StDM)

Woche 3

Aufgabe 1 Non-Metrical MDS

- a) Load the data set `voting.rda`. This data set has been taken from the HSAUR2 package and contains the number of times two congressmen voted differently on 19 environmental bills in New Jersey. Use `isoMDS` and plot the results. The party of congressman can be obtained by

```
party = as.factor(c('R','R','D','D','R','R','R','D','D','D','D','D','R','R','D','D'))
```

- b) To illustrate that only the order of the distances matter add 1 to each distance and then take the logarithm. Does the isoMDS significantly change?
- c) Repeat a) b) with Sammon Mapping
- d) Calculate the stress in the case of sammon mapping and compare against the result from c)
- e) To show the iterative nature of the `isoMDS` procedure, we start with a random initial configuration `Y = matrix(rnorm(2 * 15), ncol = 2)` and then perform only three iterations `isoMDS(voting, y = Y, maxit = 3)`. We plot the result and use it as a new starting point. We repeat this 10 times and look at the resulting plots. If you like, you can create an animation by storing the pngs and stick them together.

Note: Normally `isoMDS` uses `cmdscale` as a starting point and therefore usually converges much faster.

Aufgabe 2 Visualizing Images

The file `training_48x48_aligned.gz` contains images and labels of the faces of several people. Use the following code to load the images, replace filename appropriately. If you like you can of course can create your pictures.

```
filename = file.path(baseDir, 'training_48x48_aligned.gz')
dumm <- as.matrix(read.table(filename, sep=" ", stringsAsFactors = FALSE))
X = dumm[, -1] #226 examples, 48^2 pixels
y = as.factor(dumm[, 1]) #The label of the person from 0 to 5
N = sqrt(dim(X)[2])
par(mfrow=c(1,4))
par(mai=c(0.1,0.1,0.1,0.1))
```

```
for (i in c(1,50,100,200)) {
  m <- matrix(rev(X[i,]), nrow = N, ncol = N)
  image(m, useRaster = TRUE, axes = FALSE, col=gray((0:255)/255))
}
```



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

- Sidetrack Eigenfaces (optional). Perform a PCA on the transposed matrix **X** and plot the first 16 **scores** as images. That is take the first, second, ... 2304 dimensional score vector and create an image 48x48 image using e.g. the **matrix** command as above.
- Now use the first 2 **loadings** of the PCA and plot them in a scatter plot together with we color by the labels **y**.
- Now perform a non-Metrical MDS using the **isoMDS** and the **sannon** scaling, using euclidian distances, between the image.
- Now perform a **tSNE** analysis using euclidian distances, between the image.