

Math 583      HW 4 Fall 2020                      Due Friday, Sept. 25.  
Quiz 4 on Wednesday will have problems on DD plot and maybe MVN.  
Note that Exam 2 is now Friday Oct. 23.

Problem numbers are from Olive (2020).

**3.4.** Suppose that

$$\mathbf{X} \sim (1 - \gamma)EC_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g_1) + \gamma EC_p(\boldsymbol{\mu}, c\boldsymbol{\Sigma}, g_2)$$

where  $c > 0$  and  $0 < \gamma < 1$ . Following Example 3.2, show that  $\mathbf{X}$  has an elliptically contoured distribution assuming that all relevant expectations exist.

*R* Problems. Every time you log into *R* for the Math 583 homework, copy and paste the two source commands

```
source("http://parker.ad.siu.edu/Olive/rpack.txt")  
source("http://parker.ad.siu.edu/Olive/robdata.txt")
```

into *R*. These source commands are near the top of the robRhw file (<http://parker.ad.siu.edu/Olive/robRhw.txt>). This file has *R* commands for the homework.

**3.29a)** a) Assuming that you have done the two source commands (and the *R* command *library(MASS)*), type the command *ddcomp(buxx)* (just copy and paste the commands for this problem into *R*). This will make 4 DD plots based on the DGK, FCH, FMCD, and median ball estimators. The DGK and median ball estimators are the two attractors used by the FCH estimator. With the leftmost mouse button, move the cursor to an outlier and click. This data is the Buxton (1920) data and cases with numbers 61, 62, 63, 64, and 65 were the outliers with head lengths near 5 feet. After identifying at least three outliers in each plot, hold the rightmost mouse button down (and in *R* click on *Stop*) to advance to the next plot. When done, hold down the *Ctrl* and *c* keys to make a copy of the plot. Then paste the plot in *Word*.

**3.30.** The *concmv* function illustrates concentration with  $p = 2$  and a scatterplot of  $X_1$  versus  $X_2$ . The outliers are such that the MBA and FCH estimators can not always detect them. Type the command *concmv()*. Hold the rightmost mouse button down (and in *R* click on *Stop*) to see the DD plot after one concentration step. The start uses the coordinatewise median and  $diag([MAD(X_i)]^2)$ . Repeat 4 more times to see the DD plot based on the attractor. The outliers have large values of  $X_2$  and the highlighted cases have the smallest distances. Repeat the command *concmv()* several times. Sometimes the start will contain outliers but the attractor will be clean (none of the highlighted cases will be outliers), but sometimes concentration causes more and more of the highlighted cases to be outliers, so that the attractor is worse than the start. Copy one of the DD plots where none of the outliers are highlighted into *Word*.

**3.31.** The *ddmv* function illustrates concentration with the DD plot. The outliers are highlighted. The first graph is the DD plot after one concentration step. Hold the rightmost mouse button down (and in *R* click on *Stop*) to see the DD plot after two concentration steps. Repeat 4 more times to see the DD plot based on the attractor. In this problem, try to determine the proportion of outliers *gam* that the DGK estimator can detect for  $p = 2, 4, 10$  and  $20$ . Make a table of  $p$  and *gam*. For example the command *ddmv(p=2,gam=.4)* suggests that the DGK estimator can tolerate nearly 40% outliers with  $p = 2$ , but the command *ddmv(p=4,gam=.4)* suggest that *gam* needs to be lowered (perhaps by 0.1 or 0.05). Try to make  $0 < gam < 0.5$  as large as possible.

**3.38.** a) Copy and paste the commands for this problem into *R*.)

b) The commands are used to create a bivariate data set with outliers and to obtain a classical and robust RMVN covering ellipsoid. Include the two plots in *Word*.

**3.35.** a) Copy and paste the commands for this problem into *R*.)

b) Using the function *ddsims* for  $p = 2, 3, 4$ , determine how large the sample size  $n$  should be in order for the RFCH DD plot of  $n N_p(\mathbf{0}, \mathbf{I}_p)$  cases to cluster tightly about the identity line with high probability. Table your results. (Hint: type the command *ddsims(n=20,p=2)* and increase  $n$  by 10 until most of the 20 plots look linear. Then repeat for  $p = 3$  with the  $n$  that worked for  $p = 2$ . Then repeat for  $p = 4$  with the  $n$  that worked for  $p = 3$ .)

**3.36.** a) Copy and paste the commands for this problem into *R*.)

b) A numerical quantity of interest is the correlation between the  $MD_i$  and  $RD_i$  in a RFCH DD plot that uses  $n N_p(\mathbf{0}, \mathbf{I}_p)$  cases. Using the function *corrsims* for  $p = 2, 3, 4$ , determine how large the sample size  $n$  should be in order for 9 out of 10 correlations to be greater than 0.9. (Try to make  $n$  small.) Table your results. (Hint: type the command *corrsims(n=20,p=2,nruns=10)* and increase  $n$  by 10 until 9 or 10 of the correlations are greater than 0.9. Then repeat for  $p = 3$  with the  $n$  that worked for  $p = 2$ . Then repeat for  $p = 4$  with the  $n$  that worked for  $p = 3$ .)

**3.37\*.** a) Copy and paste the commands for this problem into *R*.)

b), c) The following commands to make generate data from the EC distribution  $(1 - \epsilon)N_p(\mathbf{0}, \mathbf{I}_p) + \epsilon N_p(\mathbf{0}, 25 \mathbf{I}_p)$  where  $p = 3$  and  $\epsilon = 0.4$ . Use the command *ddplot(x)* to make a DD plot and include the plot in *Word*. What is the slope of the line followed by the plotted points?