# 1. Coding with R

2. The current working directory is displayed by the RStudio IDE within the title region of the Console pane. You can also check your current working directory by running the command getwd() in the console

3.

4. Open the Excel file containing your data: select and copy the data (ctrl + c)

5. Remember to import separately qualitative data Y and qualitative data X

6.

7. IMPORTANT: Rename LABELS of VARIABLES using the code:

8. **V1  V2  V3  V4  V5  V6  V7  V8  V9  V10 V11 V12 V13 V14 V15 V16 V17 V18 V19 V20 V21 V22 V23  V24**

9.

10. **Rename also the CATEGORIES using the CODE a1, a2, a3 ect**

11. **if you have more than one Qualitative Column rename Using a1, a2, a3 next category b1, b2, b3, ect…...**

12. #Type the R code below to import the copied data from the clipboard into R and store the data in a data frame (my_data):

13.

14. X <- read.table(file = "clipboard",  sep = "\t", header=TRUE)

15. X<- read.csv("filename.csv", row.names = 1, header= TRUE)

# Change column name

colnames(XY)[1]  <- "A"

16. #rConvert first column in data.frame to row index [duplicate]

17.

18. > X <- data.frame(X[,-1], row.names = X$Samples)

19. Remove object from Project

20. rm(X)

21. #Select specific row or Column

22. #only columns

23. X1<-data.frame(X[c(1,3:5,7)])

24. #for rows and specific columns

25. X<-data.frame(X[c(1,3:5,7),(3:7)])

26. #Missing data

27. #In R, missing values are represented by the symbol NA (not available). Impossible values (e.g., dividing by zero) are represented by the symbol NaN (not a number).

28. https://cran.r-project.org/web/packages/finalfit/vignettes/missing.html

29.

30. SORT and MATCH and JOIN DATA

31. #Methods to sort a dataframe:

32. #order() function (increasing and decreasing order)

33. #arrange() function from dplyr package

34. #setorder() function from data.table package

35.

36.

37. # sort by colname

38. Xsort <- data.table::setorder(X, namecol)

39. df <-data_frame[order(data_frame$c1),]

40. # Sort by c3 and c4

41. df <-data_frame[order(data_frame$c3, data_frame$c4),]

42. # Sort by c3(descending) and c4(acending)

43. df <-data_frame[order(-data_frame$c3, data_frame$c4),]

44.

45.

46.

47. #match column name

48. x<-colnames(X)

49. y<-colnames(Y)

50. myx<-match(y, x)

51.

52.

53. # Reorder columns based on values in a particular row (df=data frame or your matrix X or XY).

54. x <- structure(list(aa = c(3L, 5L, 7L, 33L), bb = c(4L, 4L, 8L, 63L),

55.   cc = c(5L, 3L, 6L, 55L)), .Names = c("aa", "bb", "cc"),

56.   class = "data.frame", row.names = c("1", "2", "3", "100"))

57. x[,order(-x[nrow(x),])]

58.

59. ### Reorder the columns of the dataframe in Alphabetical order

60. df[,order(colnames(df))]

61. #### Reorder the columns of the dataframe in Alphabetical order

62. df[,order(colnames(df),decreasing = TRUE)]

63.

64. # sort df by row name (first row in example)

65. dfR<-df[ ,order(df[which(rownames(df) == '1'),], decreasing=TRUE)]

66.

67. #Reorder the columns of the dataframe in data type order

68.

69. #list of data type

70. classes = data.frame(sapply(df,class))

71. #transpose the data.frame

72. classes <-t(classes)

73. #bind rows

74. dfVar<-rbind(classes, df)

75. rownames(dfVar) <- NULL

76. # Reorder the columns of the dataframe in data type order

77. dfVarO<-dfVar[, order(dfVar[which(rownames(dfVar) == '1'), ]) ]

78.

79.

80.

81.

82. #Converting character matrix (chr, qualitative data) to numerical matrix for pls

83. XY<-data.matrix(C, rownames.force = NA)

84.

85. #Import data in R  from  excel  going to file->import dataset->excel…….. Follow the instruction in the wizard window

86. #Analyzing Qualitative data Y  on Excel using PIVOT command or with R following the instructions in the following websites:

87. https://uc-r.github.io/descriptives_categorical

88. https://cran.r-project.org/web/packages/vcdExtra/vignettes/vcd-tutorial.pdf

89. https://www.r-bloggers.com/2019/06/exploratory-data-analysis-with-categorical-data

90. # counts for gender categories

91. table(XY$Cat1)

92. # cross classication counts for gender by marital status

93. table(X$Cat1, X$Cat2)

94. #percentages of gender categories

95. table2 <- table(X$Cat1)

96. prop.table(table2)

97.

98.

99.

## 100. #Check the structure

101. str(X)

102.

## 103. #Do summary statistics of the quantitative data matrix X using

104. summary(X)

105.

## 106. #Do summary statistics with more parameters  of the quantitative data matrix X using

107. #(remember always activate the library you need if you dont have library -> install the package)

108. library(psych) #also HMiSC does describe so deactivate it

109. describe(X)

110.

## 111. #Do boxplot of the quantitative data matrix X using

112. https://www.r-graph-gallery.com/boxplot.html

113.

## 114. **#Do normalization using excel or with mdatools**

115. library(mdatools)

116. # autoscaling (mean centering and standardization)

117. XN = prep.autoscale(X, center = TRUE, scale = TRUE)

118.

**119.**

**120. (stripchart applied to boxplot!!) #All lines below together**

121. https://stackoverflow.com/questions/23675735/how-to-add-boxplots-to-scatterplot-with-jitter

122.

123. boxplot(XN)

124. stripchart(XN, method = "jitter", pch = 19, col = 4, vertical = TRUE,  add = TRUE)

125.

126. boxplot(X)

127. stripchart(X, method = "jitter", pch = 19, col = 4, vertical = TRUE,  add = TRUE)

128.

129. BoxData<- boxplot(X)

130.

## 131. #DO analysis of outlierS

132. https://www.r-bloggers.com/2020/01/how-to-remove-outliers-in-r/

133. https://www.r-bloggers.com/2021/09/how-to-remove-outliers-in-r-3/

134. https://statsandr.com/blog/outliers-detection-in-r/

135.

136. boxplot(X)$out

137. library(outliers)

138. outlier(X,  logical = TRUE)

139.

## 140. #Merge matrices

141. Use command merge

142. https://www.rdocumentation.org/packages/Matrix.utils/versions/0.9.8/topics/merge.Matrix

143. X <- data.frame(X1, X2, X3, X4, X5, X6, X7, X8)

144. Xm <- as.matrix(X)

145.

## 146. #Transpose Matrix

147. t() function. As previously, mentioned, you can use t(YourMatrix) to get the transpose of your matrix "YourMatrix".

## 148. #ANOVA using multiple variable (where V1, V2, V3, V4, V5, V6, V7  are the variable and Cat the groups vector, ALL in the same matrix)

149.

150. anX_multi<- aov(formula = cbind(V1, V2, V3, V4, V5, V6, V7) ~ Cat, data = XY)

151.  summary(anX_multi)

152.

153. #ANOVA for matrix X using one variable (where Letters are  the variables and Cat the groups vector, both in the same matrix XY that in this case is a mix of qualitative Y and quantitative X)

154.

155. anV1X <- aov(V1 ~ Cat, data = XY)

156. summary(anV1X)

```r
157.    TukeyHSD(anV1X, which = "Cat")
158.
159.    anV2X <- aov(V2 ~ Cat, data = XY)
160.    summary(anV2X)
161.    TukeyHSD(anV2X, which = "Cat")
162.
163.    anV3X <- aov(V3 ~ Cat, data = XY)
164.    summary(anV3X)
165.    TukeyHSD(anV3X, which = "Cat")
166.
167.    anV4X <- aov(V4 ~ Cat, data = XY)
168.    summary(anV4X)
169.    TukeyHSD(anV4X, which = "Cat")
170.
171.    anV5X <- aov(V5 ~ Cat, data = XY)
172.    summary(anV5X)
173.    TukeyHSD(anV5X, which = "Cat")
174.
175.    anV6X <- aov(V6 ~ Cat, data = XY)
176.    summary(anV6X)
177.    TukeyHSD(anV6X, which = "Cat")
178.
179.    anV7X <- aov(V7 ~ Cat, data = XY)
180.    summary(anV7X)
181.    TukeyHSD(anV7X, which = "Cat")
182.
183.
```

184. For visualizing boxplot with categorical data (this example is for the matrix "X")

```r
185.    library("ggpubr")
186.
187.    x <- which(names(XY) == "Cat") # name of grouping variable
188.    y <- which(names(XY) == "V1" # names of variables to test
189.    | names(XY) == "V2"
190.    | names(XY) == "V3"
191.    | names(XY) == "V4"
192.    | names(XY) == "V5"
193.    | names(XY) == "V6"
194.    | names(XY) == "V7")
195.    method1 <- "anova" # one of "anova" or "kruskal.test"
196.    method2 <- "t.test" # one of "wilcox.test" or "t.test"
197.    my_comparisons <- list(c("a1", "a2"), c("a1", "a3"), c("a1", "a4"), c("a2", "a3"), c("a2", "a4"), c("a3", "a4")) # comparisons for post-hoc tests
```

198.

199. To produce graph with p-value

200. for (i in y) {

201.     for (j in x) {

202.         p <- ggboxplot(XY,

203.             x = colnames(XY[j]), y = colnames(XY[i]),

204.             color = colnames(XY[j]),

205.             legend = "none",

206.             palette = "npg",

207.             add = "jitter"

208.         )

209.         print(

210.             p + stat_compare_means(aes(label = paste0(..method.., ", p-value = ", ..p.format..)),

211.                 method = method1, label.y = max(XY[, i], na.rm = TRUE)

212.             )

213.             + stat_compare_means(comparisons = my_comparisons, method = method2, label = "p.format") # remove if p-value of ANOVA or Kruskal-Wallis test >= alpha

214.         )

215.     }

216. }

217.

218.

219.

220.

221.

222.   ANOVA follow the example on http://www.sthda.com/english/wiki/two-way-anova-test-in-r

223.   ANOVA follow the example https://statsandr.com/blog/how-to-do-a-t-test-or-anova-for-many-variables-at-once-in-r-and-communicate-the-results-in-a-better-way/

224.

# 225.     If you have problems with plots digit that

226.   par(mar=c(1,1,1,1))

227.   **Or work with the zoom command from menu**

# 228.    BIVARIATE ANALYSIS -----Do correlation plot of the quantitative data matrix X using

**229.**

## 230. #Do normalization using excel or with mdatools

```
231.    library(mdatools)
232.    # autoscaling (mean centering and standardization)
233.    XN = prep.autoscale(X, center = TRUE, scale = TRUE)
234.
235.    plot(X)
236.
237.    plot(XN)
238.
```

## 239. Do correlation plot with graph and r values of the quantitative data matrix X using

```
240.    (remember always activate the library you need if you dont have library
        install the package)
241.    library("PerformanceAnalytics")
242.    chart.Correlation(X, histogram=TRUE, pch=19)
243.
```

## 244. Store correlation and p value matrix.. (Hmisc package and read the output as matrix.)

```
245.     library(Hmisc)
246.    corX<- rcorr(as.matrix(X),type="pearson")
247.
```

## 248. Export correlation matrix to  manage by excel using

```
249.     write.csv(corX$P, file = "corX.csv")
```

## 250. MULTIVARIATE ANALYSIS with mdatools library

## 251. DO PCA

```
252.    library(mdatools)
253.    XPCA = pca(X, 7, scale = TRUE, info = "X PCA model")
254.    summary(XPCA)
255.    plot(XPCA, show.labels = TRUE)
256.    plotBiplot(XPCA, show.labels = TRUE)
257.
```

258.

## 259. for making Scatterplot Matrices

260.   pairs(XPCA$calres$scores, bg = c("red", "green3", "blue", "black")[unclass(XY$Cat1)])

261.   #with group colours

262.   l <- length(unique(XY$A))

263.   pairs(XPCA$calres$scores, col = hcl.colors(l, "Temps")[XY$A])

264.

265.

## 266. CHECKING DISTANCES

267.   **#orthogonal distance Q** (The distance shows how well the object is fitted by the PCA model and allows to detect objects that do not follow a common trend, captured by the PCs) and **Hotelling T2 distance or a score distance**. (The T2 distance allows to see extreme objects — which are far from the origin)

268.   c = categorize(XPCA)

269.   plotResiduals(XPCA, show.labels = TRUE, cgroup = c)

270.

## 271. # change both levels to 1% (default are 5%)

272.   XPCA1 = setDistanceLimits(XPCA, alpha = 0.01, gamma = 0.01)

273.   plotResiduals(XPCA1, show.labels = TRUE, cgroup = c)

274.

275.

276.

277.   #Follow this example to do PCA 2D (3D doesn't work) with cluster analysis in very quirky way

278.   https://planspace.org/2013/02/03/pca-3d-visualization-and-clustering-in-r

279.

280.   The package pca3d quickly generates 2D and 3D graphics of PCA.

281.   https://cran.r-project.org/web/packages/pca3d/vignettes/pca3d.pdf

282.

## 283. MAKING 3D GRAPH

284.   https://r-graph-gallery.com/index.html

285.   https://r-graph-gallery.com/3d.html

286.

287.

288.   #3D graph

289.

290.   ```{r}

291.

292.   #library

293.   library(rgl)

```r
# This is to output a rgl plot in a rmarkdown document.
# setupKnitr()

# Data: the iris data is provided by R
data <- data.frame(XPCA$res$cal$scores[,1:3])


# Add a new column with color
mycolors <- c('royalblue1', 'darkcyan', 'oldlace')
data$color <- mycolors[ as.numeric(XY$A) ]

# Plot
plot3d(
  x=data$Comp.1, y=data$Comp.2, z=data$Comp.3,
  col = data$color,
  type = 's',
  radius = .1,
  xlab="PC1", ylab="PC2", zlab="PC3", row)

# To display in an R Markdown document:
# rglwidget()

# To save to a file:
htmlwidgets::saveWidget(rglwidget(width = 520, height = 520),
                file = "HtmlWidget/3dscatter.html",
                libdir = "libs",
                selfcontained = FALSE
                )


```

#3d surface plot with R and plotly
```{r}

# Library
library(plotly)

# Data: volcano is provided by plotly data have to be as.matrix
volc<-volcano
data2<-as.matrix(X [,5:9])
```

```r
338.    # Plot
339.    p <- plot_ly(z =data2, type = "surface")
340.    p
341.
342.    # save the widget
343.    # library(htmlwidgets)
344.    # saveWidget(p, file=paste0( getwd(), "/HtmlWidget/3dSurface.html"))
345.
346.    ```
347.
348.
349.
350.    #animated
351.
352.    ```{r}
353.    library(rgl)
354.    library(magick)
355.
356.    # Let's use the iris dataset
357.    # iris
358.
359.    # This is ugly
360.    colors <- c("royalblue1", "darkcyan", "oldlace")
361.    iris$color <- colors[ as.numeric( as.factor(iris$Species) ) ]
362.
363.    # Static chart
364.    plot3d( iris[,1], iris[,2], iris[,3], col = iris$color, type = "s", radius = .2 )
365.
366.    # We can indicate the axis and the rotation velocity
367.    play3d( spin3d( axis = c(0, 0, 1), rpm = 20), duration = 10 )
368.
369.    # Save like gif
370.    movie3d(
371.      movie="3dAnimatedScatterplot",
372.      spin3d( axis = c(0, 0, 1), rpm = 7),
373.      duration = 10,
374.      dir = "~/Desktop",
375.      type = "gif",
376.      clean = TRUE
377.    )
378.
379.
380.
381.    ```
```

382.

383.

# 384. Install package pca3d then

385.    library(pca3d)

386.    pcaX <- prcomp(X,, scale.=TRUE)

387.

388.   Coloring samples based on categories remember to use Y matrix for categories

389.    XCat <- factor(XY[,1])

390.

391.    pca3d(XPCA$loadings, show.labels = TRUE)

392.     pca3d(XPCA$calres$scores, show.labels = TRUE)

393.

394.

395.    #Go to folder where project is working (where you saved the project and see the picture)

396.    The next coding is to see ellipses

397.    pca3d(pcaX, group=XCat, show.labels = TRUE,  show.ellipses=TRUE, ellipse.ci=0.75, show.plane=FALSE)

398.

399.    snapshotPCA3d(file="ellipses.png")

400.

# 401. CLUSTER ANALYSIS

402.    https://www.r-bloggers.com/2021/04/cluster-analysis-in-r

403.    or

404.    http://www.sthda.com/english/articles/25-clusteranalysis-in-r-practical-guide/

405.

406.

407.    Normalization is very important in cluster analysis, sometimes we have variables in different scales, need to normalized based on scale function before clustering the data sets.

408.    Normalization is mandatory for cluster analysis.

409.

410.   #Hierarchical agglomerative clustering (XN was normalized using mdatools or excel)

411.

412.    disXN = dist(XN)

413.    hcXN = hclust(disXN)

414.

415.    plot(hcXN)

416.

417. Look at all elements in the object X.hclust using
418. hcXN$order
419.
420. #Cluster membership
421. member = cutree(hcXN,4)
422. table(member)
423.
424. #HeatMAP
425. hm<- data.matrix(XN)
426. heatmap(hm)
427. ht<- heatmap(hm)
428.
429.


# 430. PLS-DA

431.
432.   https://mdatools.com/docs/index.html
433.
434. calibrate the models:
435. library(mdatools)
436. daX = plsda(X, XY$Cat, ncomp.selcrit = "min", scale = TRUE, cv = 1)
437. daX = plsda(X, XY$Cat, ncomp.selcrit = "min", ncomp = 10, scale = TRUE, cv = list('ven', nseg = 5))
438.
439. Check Number of Components:
440. daX$ncomp
441. plot(daX, ncomp = 7)
442.
443.
444. Check Number of Components that gives the minimus misclassification:
445. summary(daX, ncomp = 7)
446.
447. summary(daX$calres, ncomp = 7)
448.
449. plot(daX$calres, ncomp = 7)
450.
451. plotMisclassified(daX$calres, ncomp = 7)
452.
453. summary(daX$cvres, ncomp = 7)

454.

455.    plot(daX$cvres, ncomp = 7)

456.

457.    plotMisclassified(daX$cvres, ncomp = 7)

458.

459.

**460. In all Multivariate test you calibrate and validate …. using PLSDA with have to show two set** looking at the summary for each of the mode (Cal and Val).

461.    As you can see, summary for multi class PLS-DA simply shows one set of results for each class. The performance statistics include explained X and Y variance (cumulative), values for confusion matrix (true positives, false positives, true negatives, false negatives) as well as specificity, sensitivity and accuracy values.

462.    summary(daX, ncomp = 4)

463.

464.    #you can also get a confusion matrix for particular result.

465.     getConfusionMatrix(daX$res$cal, ncomp = 4)

466.     getConfusionMatrix(daX$res$cv, ncomp = 4)

467.

468.    #you can show how sensitivity, specificity and total amount of misclassified samples depending on number of components by using corresponding plots.

469.

470.    par(mfrow = c(3, 2))

471.    plotMisclassified(daX, ncomp = 4)

472.    plotSensitivity(daX, ncomp = 4)

473.    plotSpecificity(daX, ncomp = 4)

474.

475.    #HAving details of the samples classified and misclassified in CAL or CV

476.

477.    print(daX$res$cal$misclassified)

478.    print(daX$res$cal$y.pred)

479.

480.    print(daX$res$cv$misclassified)

481.    print(daX$res$cv$y.pred)

482.

# 483.    #SVM

484.

485.    library(e1071)

```
486.  library(rpart)
487.  data(Glass, package="mlbench")
488.   ## split data into a train and test set
489.  index <- 1:nrow(Glass)
490.  testindex <- sample(index, trunc(length(index)/3))
491.  testset <- Glass[testindex,]
492.  trainset <- Glass[-testindex,]
493.
494.  # svm
495.  svm.model <- svm(Type ~ ., data = trainset, cost = 100, gamma = 1)
496.  svm.pred <- predict(svm.model, testset[,-10])
497.
498.  # rpart
499.  rpart.model <- rpart(Type ~ ., data = trainset)
500.  rpart.pred <- predict(rpart.model, testset[,-10], type = "class")
501.
502.  ## compute svm confusion matrix
503.  table(pred = svm.pred, true = testset[,10])
504.
505.  ## compute rpart confusion matrix
506.  table(pred = rpart.pred, true = testset[,10])


507.  #K-fold cross-validation nested PLS-DA
508.
509.  #libraries
510.
511.  library(mdatools)
512.  library(psych)
513.  library(rms)
514.  library(ggpubr)
515.  library(PerformanceAnalytics)
516.  library(Hmisc)
517.  library(pca3d)
518.  library(rpart)
519.  library(e1071)
520.
521.  #split dataset in training and test
522.
523.  index <- 1:nrow(h144N)
524.  testindex <- sample(index, trunc(length(index)/3))
```

```
525.   testset1<- h144N[testindex,]
526.   trainset1<- h144N[-testindex,]
527.   index <- 1:nrow(h144N)
528.   testindex <- sample(index, trunc(length(index)/3))
529.   testset2<- h144N[testindex,]
530.   trainset2<- h144N[-testindex,]
531.   index <- 1:nrow(h144N)
532.   testindex <- sample(index, trunc(length(index)/3))
533.   testset3<- h144N[testindex,]
534.   trainset3<- h144N[-testindex,]
535.   index <- 1:nrow(h144N)
536.   testindex <- sample(index, trunc(length(index)/3))
537.   testset4<- h144N[testindex,]
538.   trainset4<- h144N[-testindex,]
539.   index <- 1:nrow(h144N)
540.   testindex <- sample(index, trunc(length(index)/3))
541.   testset5<- h144N[testindex,]
542.   trainset5<- h144N[-testindex,]
543.   index <- 1:nrow(h144N)
544.   testindex <- sample(index, trunc(length(index)/3))
545.   testset6<- h144N[testindex,]
546.   trainset6<- h144N[-testindex,]
547.   index <- 1:nrow(h144N)
548.   testindex <- sample(index, trunc(length(index)/3))
549.   testset7<- h144N[testindex,]
550.   trainset7<- h144N[-testindex,]
551.   index <- 1:nrow(h144N)
552.   testindex <- sample(index, trunc(length(index)/3))
553.   testset8<- h144N[testindex,]
554.   trainset8<- h144N[-testindex,]
555.   index <- 1:nrow(h144N)
556.   testindex <- sample(index, trunc(length(index)/3))
557.   testset9<- h144N[testindex,]
558.   trainset9<- h144N[-testindex,]
559.   index <- 1:nrow(h144N)
560.   testindex <- sample(index, trunc(length(index)/3))
561.   testset10<- h144N[testindex,]
562.   trainset10<- h144N[-testindex,]
563.
564.  #PLSDA model on training
565.
566.   daT1 = plsda(trainset1[,-1], trainset1[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
```

```r
567.  daT2 = plsda(trainset2[,-1], trainset2[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
568.  daT3 = plsda(trainset3[,-1], trainset3[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
569.  daT4 = plsda(trainset4[,-1], trainset4[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
570.  daT5 = plsda(trainset5[,-1], trainset5[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
571.  daT6 = plsda(trainset6[,-1], trainset6[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
572.  daT7 = plsda(trainset7[,-1], trainset7[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
573.  daT8 = plsda(trainset8[,-1], trainset8[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
574.  daT9 = plsda(trainset9[,-1], trainset9[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
575.  daT10 = plsda(trainset10[,-1], trainset10[,1],  scale = TRUE, cv =
      list("rand",nseg=4,nrep=10))
576.
577.  #prediction on test
578.
579.  T1= predict(daT1, testset1[,-1], testset1[,1])
580.  T2= predict(daT2, testset2[,-1], testset1[,1])
581.  T3= predict(daT3, testset3[,-1], testset1[,1])
582.  T4= predict(daT4, testset4[,-1], testset1[,1])
583.  T5= predict(daT5, testset5[,-1], testset1[,1])
584.  T6= predict(daT6, testset6[,-1], testset1[,1])
585.  T7= predict(daT7, testset7[,-1], testset1[,1])
586.  T8= predict(daT8, testset8[,-1], testset1[,1])
587.  T9= predict(daT9, testset9[,-1], testset1[,1])
588.  T10= predict(daT10, testset10[,-1], testset1[,1])
589.



590.
591.
```

# 592. DOE ------ experimental design

597.

598. install.packages("daewr")

599. install.packages("DoE.base")

600. install.packages("FrF2")

601. install.packages("rsm")

602.

603. library(daewr)

604. library(DoE.base)

605. library(FrF2)

606. library(rsm)

607.

608.

609.

# 610. FULL vs Factorial Factorial design

611.

612. ff23 <- FrF2( 8, 3, randomize = FALSE)

613. ff23

614.

615.

616.

617.

618. ff23f <- FrF2( 4, 3, randomize = FALSE)

619. ff23f

620. y <- runif(4, 0, 1)

621. aliases( lm( y~ (.)^4, data = ff23f))

622.

623. design <- FrF2( 16, 5, generators = "ABCD", randomize = FALSE)

624. design

625. y <- runif(16, 0, 1)

626. aliases( lm( y~ (.)^4, data = design))

627.

## 628. Plackett-Burman designs

629. it can be created easily using the FrF2 package. The

630. example below illustrates the use of the pb function in that package to create

631. the design with 11 factors using 12 runs.

632. library(FrF2)

633. pb( nruns = 12, randomize=FALSE)

634.

635.

636. Exercise full factorial

637. volt

638. modv <- lm( y ~ A*B*C, data=volt, contrast=list(A=contr.FrF2, B=contr.FrF2, C=contr.FrF2))

639. summary(modv)

640. par( mfrow = c(2,2) )

641. IAPlot(modv)

642. IAPlot(modv, select = c(1,3))

643.

644. Exercise full factorial

645. chem

646. modf <- lm( y ~ A*B*C*D, data = chem)

647.  summary(modf)

648. fullnormal(coef(modf)[-1],alpha=.025)

649.

650. Exercise full factorial

651. data(BoxM)

652. Gaptest(BoxM)

653.

654. Exercise fractional factorial

655.

656. soup <- FrF2(16, 5, generators = "ABCD", factor.names =list(Ports=c(1,3), Temp=c("Cool","Ambient"), MixTime=c(60,80),BatchWt=c(1500,2000), delay=c(7,1)), randomize = FALSE)

657. y <- c(1.13, 1.25, .97, 1.70, 1.47, 1.28, 1.18, .98, .78, 1.36, 1.85, .62, 1.09, 1.10, .76, 2.10)

658. soup <- add.response( soup , y )

659. mod1 <- lm( y ~ (.)^2, data = soup)

660. mod1

661. summary(mod1)

662.

663. soupc<-FrF2(16,5,generators="ABCD",randomize=FALSE)

664. soupc<-add.response(soupc, y)

665.

```
666.   soupc<-FrF2(16,5,generators="ABCD",randomize=FALSE)
667.   soupc<-add.response(soupc, y)
668.
669.   modc<-lm(y~(.)^2, data=soupc)
670.   LGB(coef(modc)[-1], rpt = FALSE)
671.
```

## 672. #response surface experiment

```
673.
```

## 674. #Exercise central composite designs 2 factors

## 675. #https://www.youtube.com/watch?v=5Zb-3gZIL1E

```
676.   # Lez74 http://www.lithoguru.com/scientist/statistics/course.html
677.
678.   #------------------------------------#
679.   #--- Response Surface Modeling in R ---#
680.   #------------------------------------#
681.
682.   #First, install and load the "rsm" package
683.
684.   # install.packages("rsm")
685.   library(rsm)
686.
687.   # Example generating a Box-Behnken design with three factors and two
       center points (no)
688.   bbd(3, n0 = 2, coding = list(x1 ~ (Force - 20)/3, x2 ~ (Rate - 50)/10, x3 ~
       Polish - 4))
689.
690.
691.   # Example data set
692.   data = ChemReact
693.   plot(data)
694.
695.
696.   # The data set was collected in two blocks.
697.   # Block1 is a 2-level, two-factor factorial design with three repeated center
       points.
698.   # Block 2 is the Central Composite Design (circomscribed) with 3 center
       points.
699.   # The variables are Time = 85 +/- 5 and Temp = 175 +/- 5,
700.   # Thus, the coded variables are  x1 = (Time-85)/5 and x2 = (Temp-175)/5
701.   CR <- coded.data(ChemReact, x1 ~ (Time - 85)/5, x2 ~ (Temp - 175)/5)
702.   CR[1:7,]
703.
```

```
704.   # Note: If the data are already coded, use as.coded.data() to convert to the
       proper coded data object
705.
706.   # Let's work as though the first block (full factorial) has been finished,
707.   # and we'll fit a linear model, first order (FO), to it (Yield is the response)
708.   CR.rsm1 <- rsm(Yield ~ FO(x1, x2), data = CR, subset = (Block == "B1"))
709.   summary(CR.rsm1)
710.
711.   #The fit is not very good.  Let's include the interaction term (TWI) and
       update the model, or start over with a new model (these two lines do the
       same thing)
712.   CR.rsm1.5 <- update(CR.rsm1, . ~ . + TWI(x1, x2))
713.   CR.rsm1.5 <- rsm(Yield ~ FO(x1, x2)+TWI(x1, x2), data = CR, subset =
       (Block == "B1"))
714.   summary(CR.rsm1.5)
715.   #This is no better!  The reason is the strong quadratic response, with the
       peak near the center.
716.
717.   # Now let's assume the second block has been collected.  We use the SO
       (second order) function, which includes FO and TWI
718.   CR.rsm2 <- rsm(Yield ~ Block + SO(x1, x2), data = CR)
719.   summary(CR.rsm2)
720.
721.   # The secondary point is a maximum (both eigenvalues are negative) and
       within the experimental design range (no extrapolation)
722.
723.   # Also note that the block is significant, meaning that the processes shifted
       between the first set of data and the second.  This is not good.  The
       coefficient is -4.5, meaning the yield shifted down by 4.5% between the two
       blocks - a more significant effect than either temperatue or time!  This is most
       easily seen by looking at the repeat center points.
724.
725.   # We can plot the fitted response as a contour plot.
726.   contour(CR.rsm2, ~ x1 + x2, at = summary(CR.rsm2)$canonical$xs)
727.
728.
729.   #OTHER example
730.   library(rsm)
731.   cube (2, n0 = 4)
732.   ccd.pick(k=2)
733.
734.   #Copy this matrix for example
735.   x1 x2     Time  Temp  y
736.   -1 -1     80     179   76.5
```

```
737.  -1  1       80      180     77
738.  1   0.1     90      179     78
739.  1   1       90      180     79.5
740.  0   0       85      175     79.9
741.  0   0       85      175     80.3
742.  0   0       85      175     80
743.  0   0       85      175     79.7
744.
745.  ccd1 <- read.table(file = "clipboard",  sep = "\t", header=TRUE)
746.  ccd1<- as.coded.data(cdd1, x1 ~ (Time-85)/5, x2 ~ (Temp-175)/5)
747.
748.  Model_Y1<- rsm(y ~ SO(x1, x2), data = ccd1)
749.  Model_Y3<- rsm(y ~ FO(x1, x2) + PQ(x1, x2), data = ccd1)
750.
751.  par(mfrow = c(1, 2))
752.  contour(Model_Y3, ~ x1+x2, image = TRUE, yagp=c(168,182, 2),
      xlabs=c("Time", "Temp"))
753.  points(ccd1$Time, ccd1$Temp)
754.   persp(Model_Y3,  x1~x2, col = terrain.colors(50), contours = "colors")
755.
756.  max <- data.frame(x1 = 0.361, x2 = 0.257)
757.
758.
759.
760.
```

## 761. central composite designs 3 factors

```
762.  library(rsm)
763.  ccd.pick(k=3)
764.
765.  data(Treb)
766.  treb.quad <- rsm(y ~ SO(x1, x2, x3), data = Treb)
767.  summary(treb.quad)
768.
769.  par (mfrow=c(2,2))
770.  contour(treb.quad, ~ x1+x2+x3 )
771.
```

## 772. 3D response surface experiment

```
773.
774.  par (mfrow=c(1,3))
775.  persp(treb.quad, ~ x1+x2+x3, zlab="Distance", contours=list(z="bottom"))
776.
```

## 777. 3D response surface experiment one factor

```
778.
```

```
779.    par (mfrow=c(1,3))
780.    contour(treb.quad, x1~x3, at=list(x2=1))
781.    persp(treb.quad, x1~x3, at=list(x2=1),zlab="Distance",
        contours=list(z="bottom"))
782.
```

783. DETERMINING OPTIMUM OPERATING CONDITIONS

```
784.    ridge<-steepest(treb.quad, dist=seq(0, 1.412, by=.1), descent=FALSE)
785.
786.    ridge
787.
788.
```

# 789. MIXTURE EXPERIMENTS

```
790.
791.    library(mixexp)
792.    SLD(3,2)
793.    des<-SLD(3,3)
794.    DesignPoints(des)

795.    library(daewr)
796.    data(pest)
797.    DesignPoints(pest)
798.     spc <- lm(y ~ x1 + x2 + x3 + x1:x2 + x1:x3 + x2:x3 + x1:x2:x3 -1, data =
        pest)
799.    summary(spc)
800.
801.    MixturePlot(des = pest,mod = 2)
802.
803.    EffPlot(des=pest,mod=2,dir=1)
804.
```

805. Other example of ANALYSIS OF MIXTURE EXPERIMENTS

```
806.
807.    data(polvdat)
808.    sqm <- lm(y ~ x1 + x2 + x3 + x1:x2 + x1:x3 + x2:x3 + x1:x2:x3 -1, data =
        polvdat)
809.    summary(sqm)
810.
811.    MixturePlot(des = polvdat, mod = 4, lims=c(0,.8,.1,.95, .05, .50),
        constrts=TRUE, pseudo=TRUE)
812.
813.
```

# 814. #Analyzing literatures blibliometrix

```
815.
816.    library(bibliometrix)
```

```
817.
818.    M <- convert2df(file = "HABSscopus.bib" , dbsource = "scopus", format =
        "bibtex")
819.    results <- biblioAnalysis(M, sep = ";")
820.    S <- summary(object = results, k = 10, pause = FALSE)
821.    plot(x = results, k = 10, pause = FALSE)
822.    ID <- cocMatrix(M, Field = "ID", sep = ";")
823.    sort(Matrix::colSums(ID), decreasing = TRUE)[1:20]
824.    rm(A)
825.    AB <- cocMatrix(M, Field = "AB", sep = " ")
826.    sort(Matrix::colSums(AB), decreasing = TRUE)[1:100]
827.
828.    # Create keyword co-occurrences network
829.
830.    NetMatrix <- biblioNetwork(M, analysis = "co-occurrences", network =
        "keywords", sep = ";")
831.
832.    # Plot the network
833.    net=networkPlot(NetMatrix, normalize="association", weighted=T, n = 30,
        Title = "Keyword Co-occurrences", type = "fruchterman", size=T,edgesize =
        5,labelsize=0.7)
834.
835.    # Conceptual Structure using keywords (method="CA")
836.
837.    CS <- conceptualStructure(M,field="ID", method="CA", minDegree=4,
        clust=5, stemming=FALSE, labelsize=10, documents=10)
838.
```

## 839. #bio3d working with PDB files

```
840.    # List of bio3d functions with brief description
841.    help(package=bio3d)
842.
843.    # Detailed help on a particular function, e.g. 'pca.xyz'
844.    help(pca.xyz)
845.
846.    pdb <- read.pdb("4q21")
847.
848.    attributes(pdb)
849.
850.    head(pdb$atom)
851.
852.    # Print $atom data for the first two atoms
853.    pdb$atom[1:2, ]
854.
```

```
855.   # Print a subset of $atom data for the first two atoms
856.   pdb$atom[1:2, c("eleno", "elety", "x","y","z")]
857.
858.   # Note that individual $atom records can also be accessed like this
859.   pdb$atom$elety[1:2]
860.
861.   # Which allows us to do the following (see Figure 1.)
862.   plot.bio3d(pdb$atom$b[pdb$calpha], sse=pdb, typ="l", ylab="B-factor")
863.
864.   # Examine the row and column dimensions
865.   dim(pdb$xyz)
866.
867.   pdb$xyz[ 1, atom2xyz(1:2) ]
868.
869.   # Select all C-alpha atoms (return their indices)
870.   ca.inds <- atom.select(pdb, "calpha")
871.   ca.inds
872.
873.   # Print details of the first few selected atoms
874.   head( pdb$atom[ca.inds$atom, ] )
875.
876.   # And selected xyz coordinates
877.   head( pdb$xyz[, ca.inds$xyz] )
878.
879.   # Select chain A
880.   a.inds <- atom.select(pdb, chain="A")
881.
882.   # Select C-alphas of chain A
883.   ca.inds <- atom.select(pdb, "calpha", chain="A")
884.
885.   # We can combine multiple selection criteria to return their intersection
886.   cab.inds <- atom.select(pdb, elety=c("CA","CB"), chain="A", resno=10:20)
887.
888.   # Select all atoms except waters
889.   nowat.inds <- atom.select(pdb, "water", inverse=TRUE)
890.
891.   # Select protein + GDP
892.   sele <- atom.select(pdb, "protein", resid="GDP", operator="OR")
893.
894.   sele <- atom.select(pdb, "protein", elety=c("N", "CA", "C"), resno=50:60,
       verbose=T)
895.
896.   a.inds <- atom.select(pdb, "protein")
897.   b.inds <- atom.select(pdb, resid="GDP")
```

```r
898.    sele <- combine.select(a.inds, b.inds, operator="OR")
899.
900.    aa321(aa3)
901.
902.    # Output a backbone only PDB file to disc
903.    b.inds <- atom.select(pdb, "back")
904.    backpdb <- trim.pdb(pdb, b.inds)
905.    write.pdb(backpdb, file="4q21_back.pdb")
906.
907.    # Selection statements can be passed directly to trim.pdb()
908.    backpdb <- trim.pdb(pdb, "backbone")
909.
910.    # The 'value=TRUE' option of atom.select() will result in a PDB object
        being returned
911.    backpdb <- atom.select(pdb, "backbone", value=TRUE)
912.
913.    # Renumber all residues
914.    write.pdb(backpdb, resno=backpdb$atom$resno+10)
915.
916.    # Assign chain B to all residues
917.    write.pdb(backpdb, chain="B")
918.
919.    pdb <- read.pdb("4lhy")
920.
921.    # select chains A, E and F
922.    inds <- atom.select(pdb, chain=c("A", "E", "F"))
923.
924.    # trim PDB to selection
925.    pdb2 <- trim.pdb(pdb, inds)
926.
927.    # assign new chain identifiers
928.    pdb2$atom$chain[ pdb2$atom$chain=="E" ] <- "B"
929.    pdb2$atom$chain[ pdb2$atom$chain=="F" ] <- "C"
930.
931.    # re-number chain B and C
932.    pdb2$atom$resno[ pdb2$atom$chain=="B" ] <- pdb2$atom$resno[
        pdb2$atom$chain=="B" ] - 156
933.    pdb2$atom$resno[ pdb2$atom$chain=="C" ] <- pdb2$atom$resno[
        pdb2$atom$chain=="C" ] - 156
934.
935.    # assign the GDP residue a residue number of 500
936.    pdb2$atom$resno[ pdb2$atom$resid=="GDP" ] <- 500
937.
938.    # use chain D for the GDP residue
```

```r
939.    pdb2$atom$chain[ pdb2$atom$resid=="GDP" ] <- "D"
940.
941.    # Center, to the coordinate origin, and orient, by principal axes,
942.    # the coordinates of a given PDB structure or xyz vector.
943.    xyz <- orient.pdb(pdb2)
944.
945.    # write the new pdb object to file
946.    write.pdb(pdb2, xyz=xyz, file="4LHY_AEF-oriented.pdb")
947.
948.    # read two G-protein structures
949.    a <- read.pdb("4q21")
950.    b <- read.pdb("4lhy")
951.
952.    a1 <- trim.pdb(a, chain="A")
953.
954.    b1 <- trim.pdb(b, chain="A")
955.    b2 <- trim.pdb(b, chain="E")
956.    b3 <- trim.pdb(b, chain="F")
957.
958.
959.  # Redundant testing excluded FASTA
960.    # Read sequence alignment
961.    file <- system.file("examples/kif1a.fa",package="bio3d")
962.    aln <- read.fasta(file)
963.    # Read aligned PDBs
964.    pdbs <- read.fasta.pdb(aln)
965.    # Structure/sequence names/ids
966.    basename( pdbs$id )
967.    # Alignment positions 335 to 339
968.    pdbs$ali[,335:339]
969.    pdbs$resid[,335:339]
970.    pdbs$resno[,335:339]
971.    pdbs$b[,335:339]
972.    # Alignment C-alpha coordinates for these positions
973.    pdbs$xyz[, atom2xyz(335:339)]
974.    # See 'fit.xyz()' function for actual coordinate superposition
975.    # e.g. fit to first structure
976.    # xyz <- fit.xyz(pdbs$xyz[1,], pdbs)
977.    # xyz[, atom2xyz(335:339)]
978.
979.  # concatenate PDBs
980.    new <- cat.pdb(a1, b1, b2, b3, rechain=TRUE)
981.    unique(new$atom$chain)
```

```r
982.
983.   # write new PDB object to file
984.   write.pdb(new, file="4Q21-4LHY.pdb")
985.
986.   # Align and superpose two or more structures
987.   pdbs <- pdbaln(c("4q21", "521p"), fit=TRUE)
988.
989.   # read two G-protein structures
990.   a <- read.pdb("4q21")
991.   b <- read.pdb("4lhy")
992.
993.   # perform iterative alignment
994.   aln <- struct.aln(a, b)
995.
996.   # store new coordinates of protein B
997.   b$xyz <- aln$xyz
998.
999.   # indices at which the superposition should be based
1000.  a.ind <- atom.select(a, chain="A", resno=87:103, elety="CA")
1001.  b.ind <- atom.select(b, chain="A", resno=93:109, elety="CA")
1002.
1003.  # perform superposition
1004.  xyz <- fit.xyz(fixed=a$xyz, mobile=b$xyz,
1005.                 fixed.inds=a.ind$xyz,
1006.                 mobile.inds=b.ind$xyz)
1007.
1008.  # write coordinates to file
1009.  write.pdb(b, xyz=xyz, file="4LHY-at-4Q21.pdb")
1010.
1011.  # read G-protein structure
1012.  pdb <- read.pdb("4q21")
1013.  bs <- binding.site(pdb)
1014.
1015.  # residue names of identified binding site
1016.  print(bs$resnames)
1017.
1018.  b <- read.pdb("4lhy")
1019.
1020.  # atom selection
1021.  a.inds <- atom.select(b, chain="A")
1022.  b.inds <- atom.select(b, chain=c("E", "F"))
1023.
1024.  # identify interface residues
```

```
1025.  bs <- binding.site(b, a.inds=a.inds, b.inds=b.inds)
1026.
1027.  # use b-factor column to store interface in PDB file
1028.  b$atom$b[ bs$inds$atom ] <- 1
1029.  b$atom$b[ -bs$inds$atom ] <- 0
1030.
1031.  # write to file
1032.  write.pdb(b, file="4LHY-interface.pdb")
1033.
1034.  b <- read.pdb("4lhy")
1035.
1036.  # atom selection
1037.  a.inds <- atom.select(b, chain="A")
1038.  b.inds <- atom.select(b, chain=c("E", "F"))
1039.
1040.  # identify interface residues
1041.  bs <- binding.site(b, a.inds=a.inds, b.inds=b.inds)
1042.
1043.  # use b-factor column to store interface in PDB file
1044.  b$atom$b[ bs$inds$atom ] <- 1
1045.  b$atom$b[ -bs$inds$atom ] <- 0
1046.
1047.  # write to file
1048.  write.pdb(b, file="4LHY-interface.pdb")
1049.
1050.  # The xyz component contains 20 frames
1051.  pdb$xyz
1052.
1053.  # Select a subset of the protein
1054.  ca.inds <- atom.select(pdb, "calpha")
1055.
1056.  # Access C-alpha coordinates of the first 5 models
1057.  #pdb$xyz[1:5, ca.inds$xyz]
1058.
1059.  library(bio3d.geostas)
1060.
1061.  # Domain analysis
1062.  gs  <- geostas(pdb)
1063.
1064.  # Fit all frames to the 'first' domain
1065.  domain.inds <- gs$inds[[1]]
1066.
1067.  xyz <- pdbfit(pdb, inds=domain.inds)
1068.
```

```r
1069.   # write fitted coordinates
1070.   write.pdb(pdb, xyz=xyz, chain=gs$atomgrps, file="1d1d_fit-domain1.pdb")
1071.
1072.   # plot geostas results
1073.   plot(gs, contour=FALSE)
1074.
1075.   # Invariant core
1076.   core <- core.find(pdb)
1077.
1078.   # fit to core region
1079.   xyz <- pdbfit(pdb, inds=core)
1080.
1081.   # write fitted coordinates
1082.   write.pdb(pdb, xyz=xyz, file="1d1d_fit-core.pdb")
1083.
1084. # Read PDB from online database
1085.   pdb <- read.pdb("2dn1")
1086.
1087.   # Examine biological unit matrices
1088.   pdb$remark$biomat
1089.
1090.   biounit(pdb)
1091.
1092.   bio <- biounit(pdb)
1093.   names(bio)
1094.
1095.   # Download some example PDB files
1096.   ids <- c("1TND_B","1AGR_A","1FQJ_A","1TAG_A","1GG2_A","1KJY_A")
1097.   raw.files <- get.pdb(ids)
1098.
1099. # Extract and align the chains we are interested in
1100.   files <- pdbsplit(raw.files, ids)
1101.   pdbs <- pdbaln(files,
          exefile="C:/Users/marce/Documents/R/bioinformatics/muscle3.exe")
1102.   # Calculate sequence identity
1103.   pdbs$id <- basename.pdb(pdbs$id)
1104.   seqidentity(pdbs)
1105.
1106.   ## Calculate RMSD
1107.   rmsd(pdbs, fit=TRUE)
1108.
1109. ## Quick PCA (see Figure 9)
1110.   pc <- pca(pdbfit(pdbs), rm.gaps=TRUE)
```

```
1111.  plot(pc)
1112.
1113.  ## Quick NMA of all structures (see Figure 10)
1114.  library(bio3d.nma)
1115.  modes <- nma(pdbs)
1116.  plot(modes, pdbs, spread=TRUE)
1117.
1118.  library(rmarkdown)
1119.  render("Bio3D_pdb.Rmd", "all")
1120.
1121.  # Information about the current Bio3D session
1122.  sessionInfo()
1123.
1124.  attach(transducin)
1125.
1126.  pdb <- read.pdb("1tag")
1127.  seq <- pdbseq(pdb)
1128.  blast <- blast.pdb(seq)
1129.
1130.  # See Figure 2.
1131.  hits <- plot.blast(blast, cutoff=240)
1132.
1133.  head(hits$hits)
1134.
1135.  # Download PDBs and split by chain ID
1136.  files <- get.pdb(hits, path="raw_pdbs", split = TRUE)
1137.
1138.  # Extract and align sequences
1139.  pdbs <- pdbaln(files,
         exefile="C:/Users/marce/Documents/R/bioinformatics/muscle3.exe")
1140.  core <- core.find(pdbs)
1141.
1142.  # See Figure 3.
1143.  col=rep("black", length(core$volume))
1144.  col[core$volume<2]="pink"; col[core$volume<1]="red"
1145.  plot(core, col=col)
1146.
1147.  xyz <- pdbfit( pdbs, core.inds )
1148.
1149.  rd <- rmsd(xyz)
1150.  hist(rd, breaks=40, xlab="RMSD (Å)", main="Histogram of RMSD")
1151.
```

```r
1152.  # RMSD clustering
1153.  hc.rd <- hclust(as.dist(rd))
1154.
1155.  pdbs$id <- substr(basename(pdbs$id), 1, 6)
1156.  hclustplot(hc.rd, colors=annotation[, "color"], labels=pdbs$id, cex=0.5,
1157.          ylab="RMSD (Å)", main="RMSD Cluster Dendrogram",
       fillbox=FALSE)
1158.
1159.  # Ignore gap containing positions
1160.  gaps.res <- gap.inspect(pdbs$ali)
1161.  gaps.pos <- gap.inspect(pdbs$xyz)
1162.
1163.  # Tailor the PDB structure to exclude gap positions for SSE annotation
1164.  id <- grep("1TAG", pdbs$id)
1165.  inds <- atom.select(pdb, resno = pdbs$resno[id, gaps.res$f.inds])
1166.  ref.pdb <- trim.pdb(pdb, inds = inds)
1167.
1168.  # Plot RMSF with SSE annotation and labeled with residue numbers
       (Figure 8.)
1169.  rf <- rmsf(xyz[, gaps.pos$f.inds])
1170.  plot.bio3d(rf, resno=ref.pdb, sse=ref.pdb, ylab="RMSF (Å)",
1171.          xlab="Residue No.", typ="l")
1172.
1173.  tor <- torsion.pdb(pdb)
1174.
1175.  # Basic Ramachandran plot (Figure 9)
1176.  plot(tor$phi, tor$psi, xlab="phi", ylab="psi")
1177.
1178.  # Locate the two structures in pdbs
1179.  ind.a <- grep("1TAG_A", pdbs$id)
1180.  ind.b <- grep("1TND_B", pdbs$id)
1181.
1182.  # Exclude gaps in the two structures to make them comparable
1183.  gaps.xyz2 <- gap.inspect(pdbs$xyz[c(ind.a, ind.b), ])
1184.  a.xyz <- pdbs$xyz[ind.a, gaps.xyz2$f.inds]
1185.  b.xyz <- pdbs$xyz[ind.b, gaps.xyz2$f.inds]
1186.
1187.  # Compare CA based pseudo-torsion angles between the two structures
1188.  a <- torsion.xyz(a.xyz, atm.inc=1)
1189.  b <- torsion.xyz(b.xyz, atm.inc=1)
1190.  d.ab <- wrap.tor(a-b)
1191.  d.ab[is.na(d.ab)] <- 0
1192.
```

```
1193.  # Plot results with SSE annotation
1194.  plot.bio3d(abs(d.ab), resno=pdb, sse=pdb, typ="h", xlab="Residue No.",
1195.          ylab = "Difference Angle")
1196.
1197.  a <- dm.xyz(a.xyz)
1198.  b <- dm.xyz(b.xyz)
1199.
1200.  plot.dmat( (a - b), nlevels=10, grid.col="gray", xlab="1tag", ylab="1tnd")
1201.
1202.  # Do PCA
1203.  pc.xray <- pca.xyz(xyz[, gaps.pos$f.inds])
1204.  pc.xray
1205.
1206.  plot(pc.xray, col=annotation[, "color"])
1207.
1208.  plot(pc.xray, pc.axes=1:2, col=annotation[, "color"])
1209.
1210.  # Left-click on a point to label and right-click to end
1211.  identify(pc.xray$z[,1:2], labels=basename.pdb(pdbs$id))
1212.
1213.  par(mfrow = c(3, 1), cex = 0.75, mar = c(3, 4, 1, 1))
1214.  plot.bio3d(pc.xray$au[,1], resno=ref.pdb, sse=ref.pdb, ylab="PC1")
1215.  plot.bio3d(pc.xray$au[,2], resno=ref.pdb, sse=ref.pdb, ylab="PC2")
1216.  plot.bio3d(pc.xray$au[,3], resno=ref.pdb, sse=ref.pdb, ylab="PC3")
1217.
1218.  hc <- hclust(dist(pc.xray$z[,1:2]))
1219.  grps <- cutree(hc, h=30)
1220.  cols <- c("red", "green", "blue")[grps]
1221.  plot(pc.xray, pc.axes=1:2, col=cols)
1222.
1223.
1224. #Enhanced Normal Modes Analysis with Bio3D
1225.
1226.
1227.  modes <- nma(pdb)
1228.
1229.  print(modes)
1230.
1231.  plot(modes, sse=pdb)
1232.
1233.  # Calculate modes with various force fields
1234.  modes.a <- nma(pdb, ff="calpha")
1235.  modes.b <- nma(pdb, ff="anm")# Make a PDB trajectory
```

```
1236.  mktrj(modes, mode=7)
1237.
1238.  # Vector field representation (see Figure 3.)
1239.  pymol(modes, mode=7)
1240.
1241.
1242.  modes.c <- nma(pdb, ff="pfanm")
1243.  modes.d <- nma(pdb, ff="reach")
1244.  modes.e <- nma(pdb, ff="sdenm")
1245.
1246.  # Root mean square inner product (RMSIP)
1247.  r <- rmsip(modes.a, modes.b)
1248.
1249.  plot(r, xlab="ANM", ylab="C-alpha FF")
1250.
1251.  # Download PDB, calcualte normal modes of the open subunit
1252.  pdb.full   <- read.pdb("1sx4")
1253.  pdb.open   <- trim.pdb(pdb.full, atom.select(pdb.full, chain="A"))
1254.  modes      <- nma(pdb.open)
1255.
1256.
1257.
1258.  # Calculate the cross-correlation matrix
1259.  cm <- dccm(modes)
1260.
1261.  # Plot a correlation map with plot.dccm(cm)
1262.  plot(cm, sse=pdb.open, contour=F, col.regions=bwr.colors(20),
    at=seq(-1,1,0.1) )
1263.
1264.  # View the correlations in the structure (see Figure 5.)
1265.  pymol(cm, pdb.open, type="launch")
1266.
1267.  # Deformation energies
1268.  defe <- deformation.nma(modes)
1269.  defsums <- rowSums(defe$ei[,1:3])
1270.
1271.  # Fluctuations
1272.  flucts <- fluct.nma(modes, mode.inds=seq(7,9))
1273.
1274.  # Write to PDB files (see Figure 6.)
1275.  write.pdb(pdb=NULL, xyz=modes$xyz, file="R-defor.pdb", b=defsums)
1276.  write.pdb(pdb=NULL, xyz=modes$xyz, file="R-fluct.pdb", b=flucts)
1277.
1278.  # Closed state of the subunit
```

```
1279.  pdb.closed <- trim.pdb(pdb.full, atom.select(pdb.full, chain="H"))
1280.
1281.  # Align closed and open PDBs
1282.  aln <- struct.aln(pdb.open, pdb.closed, max.cycles=0)
1283.  pdb.closed$xyz <- aln$xyz
1284.
1285.  # Caclulate a difference vector
1286.  xyz <- rbind(pdb.open$xyz[aln$a.inds$xyz],
       pdb.closed$xyz[aln$a.inds$xyz])
1287.  diff <- difference.vector(xyz)
1288.
1289.  # Calculate overlap
1290.  oa <- overlap(modes, diff)
1291.
1292.  plot(oa$overlap, type='h', xlab="Mode index", ylab="Squared overlap",
       ylim=c(0,1))
1293.  points(oa$overlap, col=1)
1294.  lines(oa$overlap.cum, type='b', col=2, cex=0.5)
1295.  text(c(1,5)+.5, oa$overlap[c(1,5)], c("Mode 1", "Mode 5"), adj=0)
1296.
1297.
```

## 1298.    #Mining text Analyzing text tidy approach for articles details

## 1299.    IMPORTANT FOR LINUX UBUNTU

```
1300.  # for working with library(bibliometrix):
1301.  # install dependencies
1302.  # sudo apt-get update -y
1303.  #sudo apt-get install -y r-cran-factominer
1304.
1305.# some text
1306.  text <- c("Because I could not stop for Death -",
1307.        "He kindly stopped for me -",
1308.        "The Carriage held but just Ourselves -",
1309.        "and Immortality")
```

## 1310.    #IMPORT documents from pdf

```
1311.
1312.  library(tidyverse)
1313.  library(pdftools)
1314.#document split in pages
1315.  HA1_2015 <-
       pdf_text("F:/Work/ARTtoDO/ClampH1N1/2015RBS_HA1.pdf")
```

```
1316.#document split in pages
1317.  HA1_2015str <-
       pdf_text("F:/Work/ARTtoDO/ClampH1N1/2015RBS_HA1.pdf") %>%
       str_split("\n")
1318.
1319.  #All line together page by page
1320.  HA1<-str_squish("HA1_2015")
1321.
1322.#final data frame
1323.  final_df <- as_tibble(HA1_2015)
1324.
1325.
1326.#Working with multiple .pdf files
1327.  #sew wd to where files are located
1328.  setwd("F:/Work/ARTtoDO/ClampH1N1")
1329.
1330.  #create character list of all files ending with .pdf in folder
1331.  files <- list.files("F:/Work/ARTtoDO/ClampH1N1/", pattern = "pdf$")
1332.
1333.  #use lapply() to apply pdf_text or other pdftools function interactively cross
       each of the files
1334.  all_pdf<- lapply(files, pdf_text)#
1335.  #export as excel table each pdf file
1336.  library(writexl)
1337.  df<-data.frame(all_pdf[[1]])
1338.  write_xlsx(df, "C:/Users/marce/Documents/R/scientometrics/name.xlsx")
1339.
1340.#working with Text Mining(TM) package
1341.  library(tm)
1342.  files <- list.files("F:/Work/ARTtoDO/ClampH1N1/", pattern = "pdf$")
1343.  corp <- Corpus(URISource(files),
1344.            readerControl = list(reader = readPDF))
1345.
1346.  # create a term-document matrix, or TDM for short. A TDM stores counts
       of terms for each document.
1347.
1348.  opinions.tdm <- TermDocumentMatrix(corp,
1349.                        control =
1350.                          list(removePunctuation = TRUE,
1351.                               stopwords = TRUE,
1352.                               tolower = TRUE,
1353.                               stemming = TRUE,
1354.                               removeNumbers = TRUE,
```

```
1355.                              bounds = list(global = c(3, Inf))))
1356.
1357.
1358.
1359. #the removePunctuation function
1360.  corp <- tm_map(corp, removePunctuation, ucp = TRUE)
1361.  opinions.tdm <- TermDocumentMatrix(corp,
1362.                        control =
1363.                          list(stopwords = TRUE,
1364.                            tolower = TRUE,
1365.                            stemming = TRUE,
1366.                            removeNumbers = TRUE,
1367.                            bounds = list(global = c(3, Inf))))
1368.  inspect(opinions.tdm[1:10,])
1369.
1370.  findFreqTerms(opinions.tdm, lowfreq = 100, highfreq = Inf)
1371.
1372.  ft <- findFreqTerms(opinions.tdm, lowfreq = 100, highfreq = Inf)
1373.  as.matrix(opinions.tdm[ft,])
1374.
1375.  ft.tdm <- as.matrix(opinions.tdm[ft,])
1376.  sort(apply(ft.tdm, 1, sum), decreasing = TRUE)
1377.
1378.      # more examples Working with multiple pdf files
1379.  library(pdftools)
1380.  library(tm)
1381.  library(wordcloud)
1382.  library(RColorBrewer)
1383.  files <- list.files("~/R/scientometrics", pattern = "pdf$") # Vector of pdf file
    names
1384.  pdfs <- lapply(files, pdf_text)# loads all three files
1385.  length(pdfs)# verify how many files are in the pdfs
1386.  lapply(pdfs, length) # check the length of each pdf file
1387.  #Create a corpus with the vector that has the three files
1388.  pdfdatabase <- Corpus(URISource(files), readerControl = list(reader =
    readPDF))# creating a PDF database
1389.
1390. #Lets clean up the corpus
1391.  #create your own list of stopwords, it has to be performed on the Corpus
1392.  pdfdatabase<- tm_map(pdfdatabase, removeWords, c("abuse", "access",
    "affect"))
1393.  #remove english stopwords
```

```r
1394. pdfdatabase <- tm_map(pdfdatabase, removeWords,
      stopwords("english"))
1395. # Remove numbers
1396. pdfdatabase <- tm_map(pdfdatabase, removeNumbers)
1397.#convert Vcorpus to TIBBLE and data.frame
1398.
1399.
1400.
1401.
1402.#Only words that appear at least two times are counted in this
      example.
1403. pdfs.tdm <- TermDocumentMatrix(pdfdatabase,control =
      list(removePunctuation = TRUE, stopwords = TRUE, tolower = TRUE,
      stemming = FALSE, removeNumbers = TRUE, bounds = list(global =
      c(2,Inf))))
1404. #Examine  10 words at a time in across documents. The range below
      specifies the first 10.
1405.
1406. inspect(pdfs.tdm[1:10,])
1407.
1408. #Frequent terms that appear at least 20 times across all documents
1409.
1410. findFreqTerms(pdfs.tdm, lowfreq = 20, highfreq = Inf)
1411.
1412.#Compare the distribution of frequently appearing words
      across the three documents.
1413.
1414. ft <- findFreqTerms(pdfs.tdm, lowfreq = 20, highfreq = Inf)
1415. as.matrix(pdfs.tdm[ft,])
1416.
1417. #Sum the count of all frequently occurring words
1418.
1419. ft.tdm <- as.matrix(pdfs.tdm[ft,])
1420. sort(apply(ft.tdm, 1, sum), decreasing = TRUE)
1421.
1422. #Lets complete the second part of this exercise by conducting correlation
      analysis and creating graphs and charts. Lets find frequent terms that appear
      at least 10 times.
1423.
1424. findFreqTerms(pdfs.tdm, lowfreq = 10)
1425.
```

```
1426.  #Examine frequent Terms and their association. In this example we are
       looking at the frequent terms related to bullying. The correlation limit that is
       being examined is a correlation of 75% or greater.
1427.
1428.  findAssocs(pdfs.tdm, terms = "bullying", corlimit = 0.75)
1429.
1430.  #To create a word cloud or bar chart you must convert the term
       document matrix to a data frame.
1431.
1432.  m <- as.matrix(pdfs.tdm)
1433.  v <- sort(rowSums(m),decreasing=TRUE)
1434.  d <- data.frame(word = names(v),freq=v)
1435.
1436.  #Create word cloud
1437.
1438.  set.seed(1234)
1439.  wordcloud(words = d$word, freq = d$freq, min.freq = 1,
1440.        max.words=200, random.order=FALSE, rot.per=0.35,
1441.        colors=brewer.pal(8, "Dark2"))
1442.
1443.  #Create Bar chart
1444.
1445.  barplot(d[1:11,]$freq, las = 2, names.arg = d[1:11,]$word,
1446.        col ="lightblue", main ="Most frequent words",
1447.        ylab = "Word frequencies")
1448.
1449.  #Converting PDFs to data.frame
1450.  library(dplyr)
1451.  text_df <- tibble(line = 1:2, text = dfvc)
1452.
1453.
1454.
1455.
1456.  #Tidy text dataset, we first need to put it into a data frame
1457.  library(dplyr)
1458.  text_df <- tibble(line = 1:4, text = text)
1459.
1460.  #Split in single word
1461.  library(tidytext)
1462.  library(janeaustenr)
1463.  library(dplyr)
1464.  library(stringr)
1465.
```

```
1466.
1467. allPDF <- read.table(file = "clipboard",  sep = "\t", header=TRUE)
1468.
1469. #Art is the name of the column
1470. pdf_T<- allPDF %>%  unnest_tokens(word, Art)
1471.
1472. #Remove stop words such as "the", "of", "to", and so forth in English
1473. data(stop_words)
1474. tidyT <- pdf_T %>% anti_join(stop_words)
1475. #find the most common words in all the books as a whole.
1476. tidyT %>% count(word, sort = TRUE)
1477.
1478.
1479. #Plot words change filter number according to the previous step
1480. library(ggplot2)
1481. tidyT %>%
1482.   count(word, sort = TRUE) %>%
1483.   filter(n > 10) %>%
1484.   mutate(word = reorder(word, n)) %>%
1485.   ggplot(aes(n, word)) +
1486.   geom_col() +
1487.   labs(y = NULL)
1488.
1489.#Wordclouds
1490. library(wordcloud)
1491. tidyT %>%
1492.   anti_join(stop_words) %>%
1493.   count(word) %>%
1494.   with(wordcloud(word, n, max.words = 100))
1495.
1496. #work with physics books
1497. library(gutenbergr)
1498. physics <- gutenberg_download(c(37729, 14725, 13476, 30155),
1499.                   meta_fields = "author")
1500.
1501.#Split using two words changing the n = we can work with 3
      words (trigram) ect
1502. #Art is the name of the column
1503. pdf_2T<- allPDF %>%  unnest_tokens(bigram, Art, token = "ngrams", n =
      2)
1504. pdf_2T %>%
1505.   count(bigram, sort = TRUE)
1506.
```

```
1507.#working with 2 words by removing "stop-words"
1508. # https://www.tidytextmining.com/ngrams.html
1509. bigrams_separated <- pdf_2T %>%
1510.   separate(bigram, c("word1", "word2"), sep = " ")
1511.
1512. bigrams_filtered <- bigrams_separated %>%
1513.   filter(!word1 %in% stop_words$word) %>%
1514.   filter(!word2 %in% stop_words$word)
1515.
1516. # new bigram counts:
1517. bigram_counts <- bigrams_filtered %>%
1518.   count(word1, word2, sort = TRUE)
1519.
1520. bigram_counts
1521.
1522. #reunite  words after cleaning
1523. bigrams_united <- bigrams_filtered %>%
1524.   unite(bigram, word1, word2, sep = " ")
1525.
1526.# filter for only relatively common combinations
1527. library(igraph)
1528. bigram_graph <- bigram_counts %>%
1529.   filter(n > 2) %>%
1530.   graph_from_data_frame()
1531. library(ggraph)
1532. set.seed(2017)
1533. bigram_graph
1534.
1535.#Plotting 1 interaction between words
1536. ggraph(bigram_graph, layout = "fr") +
1537.   geom_edge_link() +
1538.   geom_node_point() +
1539.   geom_node_text(aes(label = name), vjust = 1, hjust = 1)
1540.
1541. #Plotting 2 interaction between words
1542. set.seed(2020)
1543. a <- grid::arrow(type = "closed", length = unit(.15, "inches"))
1544.
1545. ggraph(bigram_graph, layout = "fr") +
1546.   geom_edge_link(aes(edge_alpha = n), show.legend = FALSE,
1547.              arrow = a, end_cap = circle(.07, 'inches')) +
1548.   geom_node_point(color = "lightblue", size = 5) +
1549.   geom_node_text(aes(label = name), vjust = 1, hjust = 1) +
```

```
1550.    theme_void()
1551.
1552.
1553.#collect it into a function so that we easily perform it on other
     text datasets.
1554.  #To make it easy to use the count_bigrams() and visualize_bigrams()
1555.
1556.  library(dplyr)
1557.  library(tidyr)
1558.  library(tidytext)
1559.  library(ggplot2)
1560.  library(igraph)
1561.  library(ggraph)
1562.
1563.  count_bigrams <- function(dataset) {
1564.    dataset %>%
1565.      unnest_tokens(bigram, text, token = "ngrams", n = 2) %>%
1566.      separate(bigram, c("word1", "word2"), sep = " ") %>%
1567.      filter(!word1 %in% stop_words$word,
1568.             !word2 %in% stop_words$word) %>%
1569.      count(word1, word2, sort = TRUE)
1570.  }
1571.
1572.  visualize_bigrams <- function(bigrams) {
1573.    set.seed(2016)
1574.    a <- grid::arrow(type = "closed", length = unit(.15, "inches"))
1575.
1576.    bigrams %>%
1577.      graph_from_data_frame() %>%
1578.      ggraph(layout = "fr") +
1579.      geom_edge_link(aes(edge_alpha = n), show.legend = FALSE, arrow =
    a) +
1580.      geom_node_point(color = "lightblue", size = 5) +
1581.      geom_node_text(aes(label = name), vjust = 1, hjust = 1) +
1582.      theme_void()
1583.  }
1584.  #then applying functions
1585.  library(stringr)
1586.  tri_bigrams <- pdf_3T  %>%
1587.    count_bigrams()
1588.
1589.  tri_bigrams %>%
1590.    filter(n > 40,
```

```
1591.         !str_detect(word1, "\\d"),
1592.         !str_detect(word2, "\\d")),
1593.         !str_detect(word2, "\\d")) %>%
1594.          visualize_bigrams()
1595.
1596.
```