# 33-Project-5.R

Suman Paudel

2024-06-09

```r
# Part I: Use NCI60 data of ISLR2 package and page 540 of ISLR2 book to do as follows in R Studio
# to knit PDF output:

# Define nci labels (NCI$labs) as nci.labs and nci data (NCI$data) and nic.data
# Check dimension of nci.data object and interpret it carefully
# Check first four cancer types using nci.labs object
# Fit principal component analysis (PCA) on nci.data with scale = TRUE argument as pr.out object
# Create a plot showing first three PCA components with three different colors
# Get summary of pr.out object and interpret it carefully
# Plot pro.out object and interpret it carefully
# Create custom scatterplots with principal components in x-axis and proportion variance explained
# (PVE) in y-axis for the first plot and cumulative PVE in the y-axis for the second plot
# and interpret them carefully
# Perform PCA with varimax rotation and compare it with the PCA result obtained above
# Write summary of the results and conclusion based on your findings


# a. Define nci labels (NCI$labs) as nci.labs and nci data (NCI$data) as nic.data

library(ISLR2)
nci.labs <- NCI60$labs
nci.data <- NCI60$data

# b. Check dimension of nci.data object and interpret it carefully
dim(nci.data)# The data has 64 rows and 6,830 columns.
```

```
## [1]   64 6830
```

```r
# c. Check first four cancer types using nci.labs object
nci.labs[1:4]
```

```
## [1] "CNS"   "CNS"   "CNS"   "RENAL"
```

```r
table(nci.labs)
```

```
## nci.labs
##      BREAST         CNS       COLON K562A-repro K562B-repro    LEUKEMIA
##           7           5           7           1           1           6
## MCF7A-repro MCF7D-repro    MELANOMA       NSCLC     OVARIAN    PROSTATE
##           1           1           8           9           6           2
##       RENAL     UNKNOWN
##           9           1
```
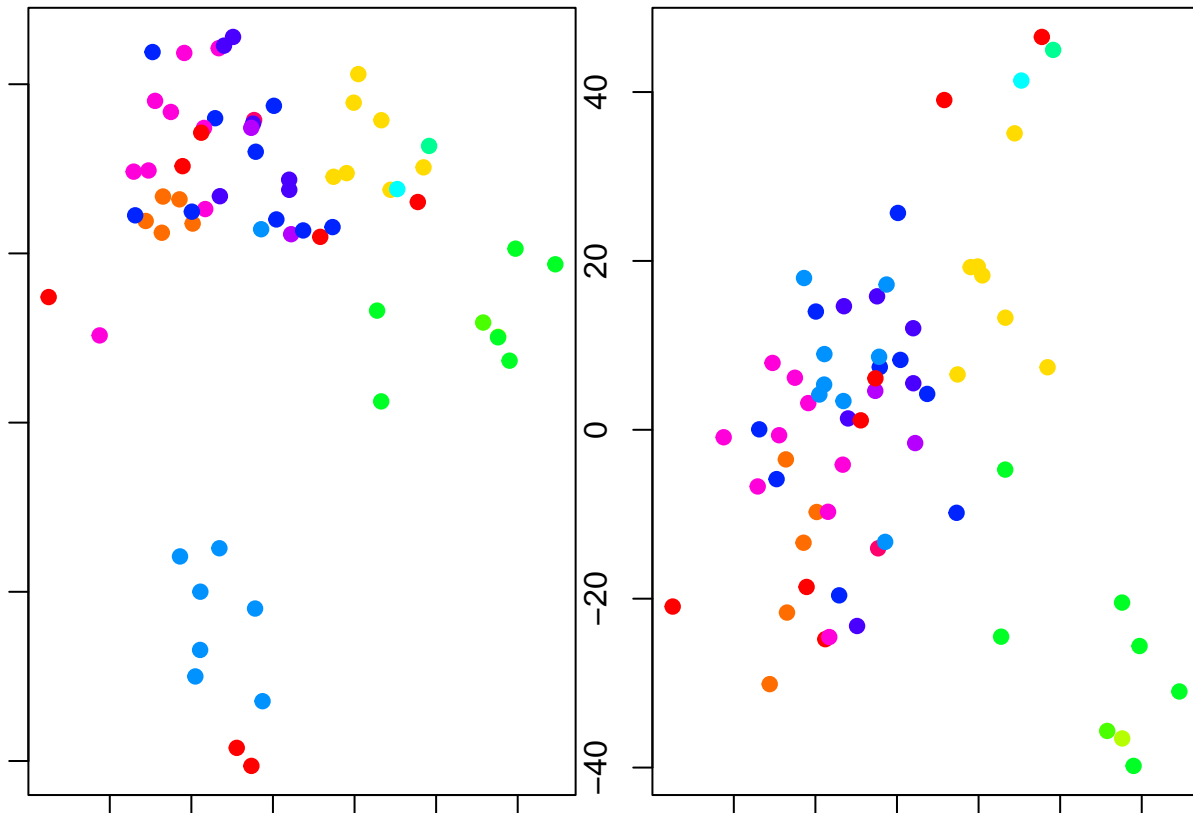
```r
# d. Fit principal component analysis (PCA) on nci.data with scale = TRUE argument as
# pr.out object
pr.out <- prcomp(nci.data, scale = TRUE)


# e. Create a plot showing first three PCA components with three different colors
Cols <- function(vec) {
  cols <- rainbow(length(unique(vec)))
  return(cols[as.numeric(as.factor(vec))])
}

# Adjusting figure margins
par(mar = c(1, 1, 1, 1))
par(mfrow = c(1, 2))
plot(pr.out$x[, 1:2], col = Cols(nci.labs), pch = 19,
     xlab = "Z1", ylab = "Z2")
plot(pr.out$x[, c(1, 3)], col = Cols(nci.labs), pch = 19,
     xlab = "Z1", ylab = "Z3")
```



```r
# On the whole, cell lines corresponding to a single cancer type do tend to have similar
# values on the first few principal component score vectors. This indicates that cell
# lines from the same cancer type tend to have pretty similar gene expression levels.
# f. Get summary of pr.out object and interpret it carefully
summary(pr.out)


## Importance of components:
```

```
##                                PC1      PC2      PC3      PC4      PC5      PC6
## Standard deviation         27.8535 21.48136 19.82046 17.03256 15.97181 15.72108
## Proportion of Variance      0.1136  0.06756  0.05752  0.04248  0.03735  0.03619
## Cumulative Proportion       0.1136  0.18115  0.23867  0.28115  0.31850  0.35468
##                                PC7      PC8      PC9     PC10     PC11     PC12
## Standard deviation         14.47145 13.54427 13.14400 12.73860 12.68672 12.15769
## Proportion of Variance      0.03066  0.02686  0.02529  0.02376  0.02357  0.02164
## Cumulative Proportion       0.38534  0.41220  0.43750  0.46126  0.48482  0.50646
##                               PC13     PC14     PC15     PC16     PC17     PC18
## Standard deviation         11.83019 11.62554 11.43779 11.00051 10.65666 10.48880
## Proportion of Variance      0.02049  0.01979  0.01915  0.01772  0.01663  0.01611
## Cumulative Proportion       0.52695  0.54674  0.56590  0.58361  0.60024  0.61635
##                               PC19     PC20     PC21     PC22    PC23    PC24
## Standard deviation         10.43518 10.3219 10.14608 10.0544 9.90265 9.64766
## Proportion of Variance      0.01594  0.0156  0.01507  0.0148 0.01436 0.01363
## Cumulative Proportion       0.63229  0.6479  0.66296  0.6778 0.69212 0.70575
##                               PC25    PC26    PC27   PC28    PC29    PC30    PC31
## Standard deviation         9.50764 9.33253 9.27320 9.0900 8.98117 8.75003 8.59962
## Proportion of Variance     0.01324 0.01275 0.01259 0.0121 0.01181 0.01121 0.01083
## Cumulative Proportion      0.71899 0.73174 0.74433 0.7564 0.76824 0.77945 0.79027
##                               PC32    PC33    PC34    PC35    PC36    PC37    PC38
## Standard deviation         8.44738 8.37305 8.21579 8.15731 7.97465 7.90446 7.82127
## Proportion of Variance     0.01045 0.01026 0.00988 0.00974 0.00931 0.00915 0.00896
## Cumulative Proportion      0.80072 0.81099 0.82087 0.83061 0.83992 0.84907 0.85803
##                               PC39    PC40    PC41   PC42    PC43   PC44    PC45
## Standard deviation         7.72156 7.58603 7.45619 7.3444 7.10449 7.0131 6.95839
## Proportion of Variance     0.00873 0.00843 0.00814 0.0079 0.00739 0.0072 0.00709
## Cumulative Proportion      0.86676 0.87518 0.88332 0.8912 0.89861 0.9058 0.91290
##                               PC46    PC47    PC48    PC49    PC50    PC51    PC52
## Standard deviation         6.8663 6.80744 6.64763 6.61607 6.40793 6.21984 6.20326
## Proportion of Variance     0.0069 0.00678 0.00647 0.00641 0.00601 0.00566 0.00563
## Cumulative Proportion      0.9198 0.92659 0.93306 0.93947 0.94548 0.95114 0.95678
##                               PC53    PC54    PC55    PC56    PC57   PC58    PC59
## Standard deviation         6.06706 5.91805 5.91233 5.73539 5.47261 5.2921 5.02117
## Proportion of Variance     0.00539 0.00513 0.00512 0.00482 0.00438 0.0041 0.00369
## Cumulative Proportion      0.96216 0.96729 0.97241 0.97723 0.98161 0.9857 0.98940
##                               PC60    PC61    PC62    PC63     PC64
## Standard deviation         4.68398 4.17567 4.08212 4.04124 1.951e-14
## Proportion of Variance     0.00321 0.00255 0.00244 0.00239 0.000e+00
## Cumulative Proportion      0.99262 0.99517 0.99761 1.00000 1.000e+00
```

```r
# g. Plot pro.out object and interpret it carefully
plot(pr.out)
# the height of each bar in the bar plot is given by squaring the corresponding
# element of pr.out$sdev.



# h. Create custom scatter plots with principal components in x-axis and proportion
# variance explained (PVE) in y-axis for the first plot and cumulative PVE in the y-axis
# for the second plot and interpret them carefully
pve <- 100 * pr.out$sdev^2 / sum(pr.out$sdev^2)
par(mfrow = c(1, 2))
```
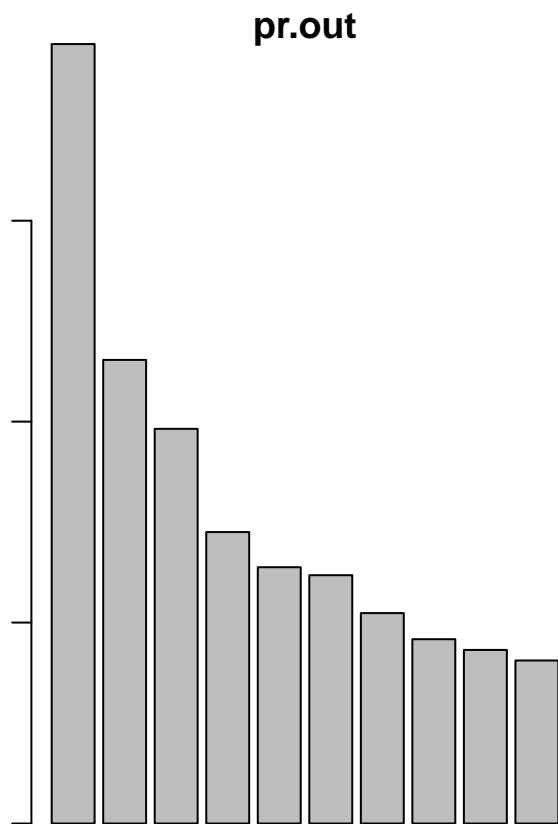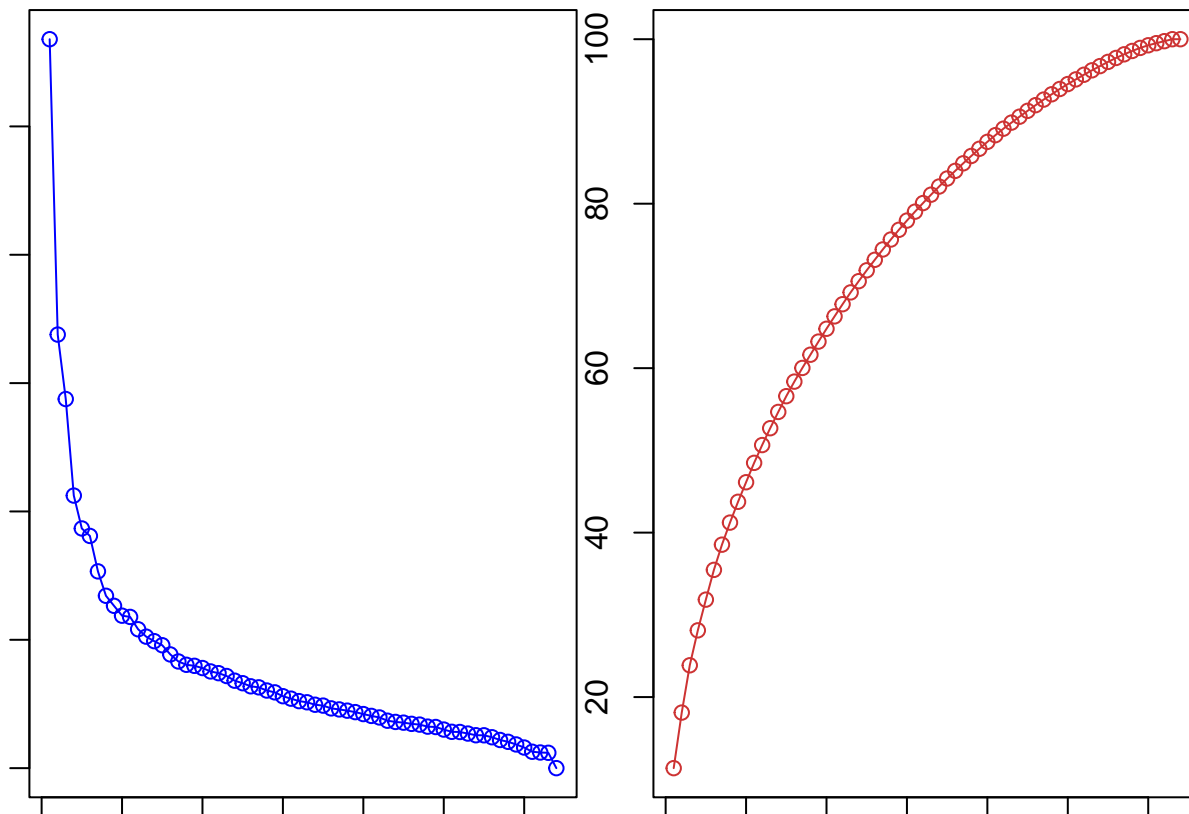
## pr.out



```r
plot(pve, type = "o", ylab = "PVE",
     xlab = "Principal Component", col = "blue")
plot(cumsum(pve), type = "o", ylab = "Cumulative PVE",
     xlab = "Principal Component", col = "brown3")
```

```
# Interpretation: We see that together, the first seven principal components explain
# around 40 % of the variance in the data. This is not a huge amount of the variance.
# However, looking at the scree plot, we see that while each of the first seven principal
# components explain a substantial amount of variance, there is a marked decrease in
# the variance explained by further principal components. That is, there is an elbow in
# the plot after approximately the seventh principal component. This suggests that there
# may be little benefit to examining more than seven or so principal components (though
# even examining seven principal components may be difficult).

# i. Perform PCA with varimax rotation and compare it with the PCA result obtained above
# library(psych)
#
# pca.varimax <- principal(nci.data, nfactors = 7, rotate = "varimax", scores = TRUE)
# pca.varimax

# Comparing PCA with varimax rotation to the original PCA results:
# Varimax rotation maximizes the sum of the variance of the squared loadings, making
# the interpretation easier.
# However, it does not change the amount of variance explained by the components.


# j. Write summary of the results and conclusion based on your findings
# Summary:
# 1. The NCI60 data set contains gene expression levels for 64 cancer cell lines and
# 6,830 genes.
# 2. PCA was performed on the data, revealing that the first few principal components
# explain a substantial amount of variance.
```

```
# 3. A scree plot and cumulative PVE plot indicate that the first seven principal
# components account for around 40% of the variance in the data.
# 4. The first three principal components were visualized, showing that cell lines fr
# om the same cancer type tend to cluster together.
# 5. PCA with varimax rotation was also performed for comparison, confirming the
# interpretability of rotated components without changing the variance explained.

# Conclusion:
# PCA is a valuable tool for reducing dimensionality and identifying patterns in high
# -dimensional data such as gene expression levels. The results suggest that a few
# principal components capture significant patterns in the data, and varimax rotation aids
# in interpretability. Further analysis can be conducted to understand the biological
# significance of these patterns.


# Part II: Use the distance between 10 US cities provided below in R Studio to knot
# PDF output:

# a. Get dissimilarity distance as city.dissimilarity object
# Step 1: Define the distance matrix
# Distance matrix for 10 US cities

city_distances <- matrix(c(
  0, 587, 1212, 701, 1936, 604, 748, 2139, 2182, 543,
  587, 0, 920, 940, 1745, 1188, 713, 2182, 2234, 597,
  1212, 920, 0, 879, 1949, 1726, 1631, 949, 1021, 1494,
  701, 940, 879, 0, 2394, 968, 1420, 2420, 2442, 597,
  1936, 1745, 1949, 2394, 0, 2300, 1645, 347, 403, 2339,
  604, 1188, 1726, 968, 2300, 0, 781, 2372, 2420, 1121,
  748, 713, 1631, 1420, 1645, 781, 0, 1923, 1960, 688,
  2139, 2182, 949, 2420, 347, 2372, 1923, 0, 214, 2571,
  2182, 2234, 1021, 2442, 403, 2420, 1960, 214, 0, 2534,
  543, 597, 1494, 597, 2339, 1121, 688, 2571, 2534, 0),
  nrow = 10, byrow = TRUE)

# City names
city_names <- c("Atlanta", "Chicago", "Denver", "Houston", "Los Angeles", "Miami",
                "New York", "San Francisco", "Seattle", "Washington")
# Assign row and column names to the distance matrix
rownames(city_distances) <- city_names
colnames(city_distances) <- city_names
# Convert to a distance object
(city.dissimilarity <- as.dist(city_distances))
```

```
##               Atlanta Chicago Denver Houston Los Angeles Miami New York
## Chicago           587
## Denver           1212     920
## Houston           701     940    879
## Los Angeles      1936    1745   1949    2394
## Miami             604    1188   1726     968        2300
## New York          748     713   1631    1420        1645   781
## San Francisco    2139    2182    949    2420         347  2372     1923
## Seattle          2182    2234   1021    2442         403  2420     1960
```

```
## Washington              543    597   1494    597           2339  1121        688
##                  San Francisco Seattle
## Chicago
## Denver
## Houston
## Los Angeles
## Miami
## New York
## San Francisco
## Seattle                         214
## Washington                     2571   2534
```

```r
# b. Fit a classical multidimensional model using the city.dissimilarity object
city_mds <- cmdscale(city.dissimilarity, eig = TRUE, k = 2)
# c. Get the summary of the model and interpret it carefully
# Get the MDS coordinates
mds_coordinates <- city_mds$points
# Print the summary of the model
summary(city_mds)
```

```
##        Length Class  Mode
## points 20     -none- numeric
## eig    10     -none- numeric
## x       0     -none- NULL
## ac      1     -none- numeric
## GOF     2     -none- numeric
```

```r
# Interpretation: The summary includes the eigenvalues, which indicate the amount of
# variance captured by each dimension.
# eigenvalues <- city_mds$eig
# variance_explained <- eigenvalues / sum(eigenvalues) * 100
# variance_explained


# d. Get the bi-plot of the model and interpret it carefully
# Plot the MDS result
plot(mds_coordinates, type = "n", xlab = "Coordinate 1", ylab = "Coordinate 2",
     main = "Classical MDS of US Cities")
text(mds_coordinates, labels = rownames(mds_coordinates), col = "blue")



# Part III: Part I: Use NCI60 data of ISLR2 package and page 543 of ISLR2 book to
# do as follows in R Studio to knit PDF output:

# a. Scale the nci.data as sd.data object
sd.data <- scale(nci.data)

# b. Fit hierarchical cluster analysis on the sd.data using complete, average and single
# linkage methods, show the results with dendrogram and interpret them carefully
par(mfrow = c(1, 3))
```
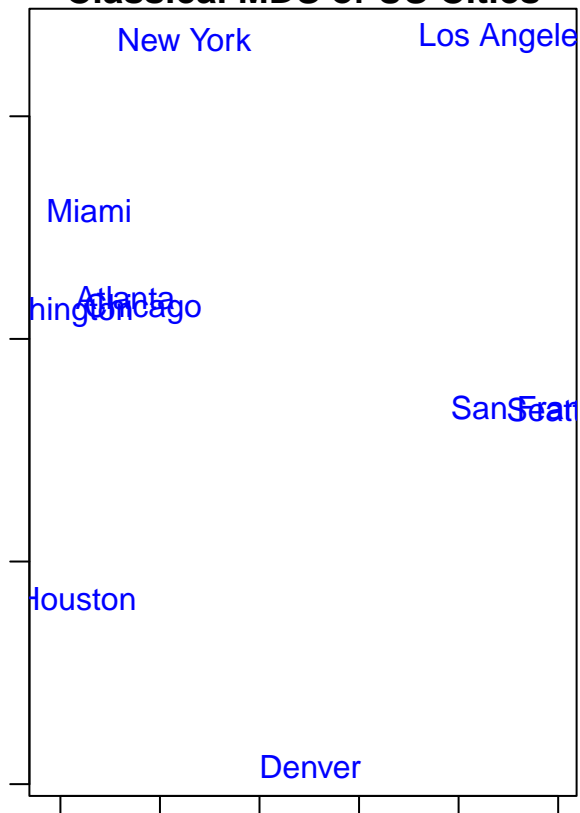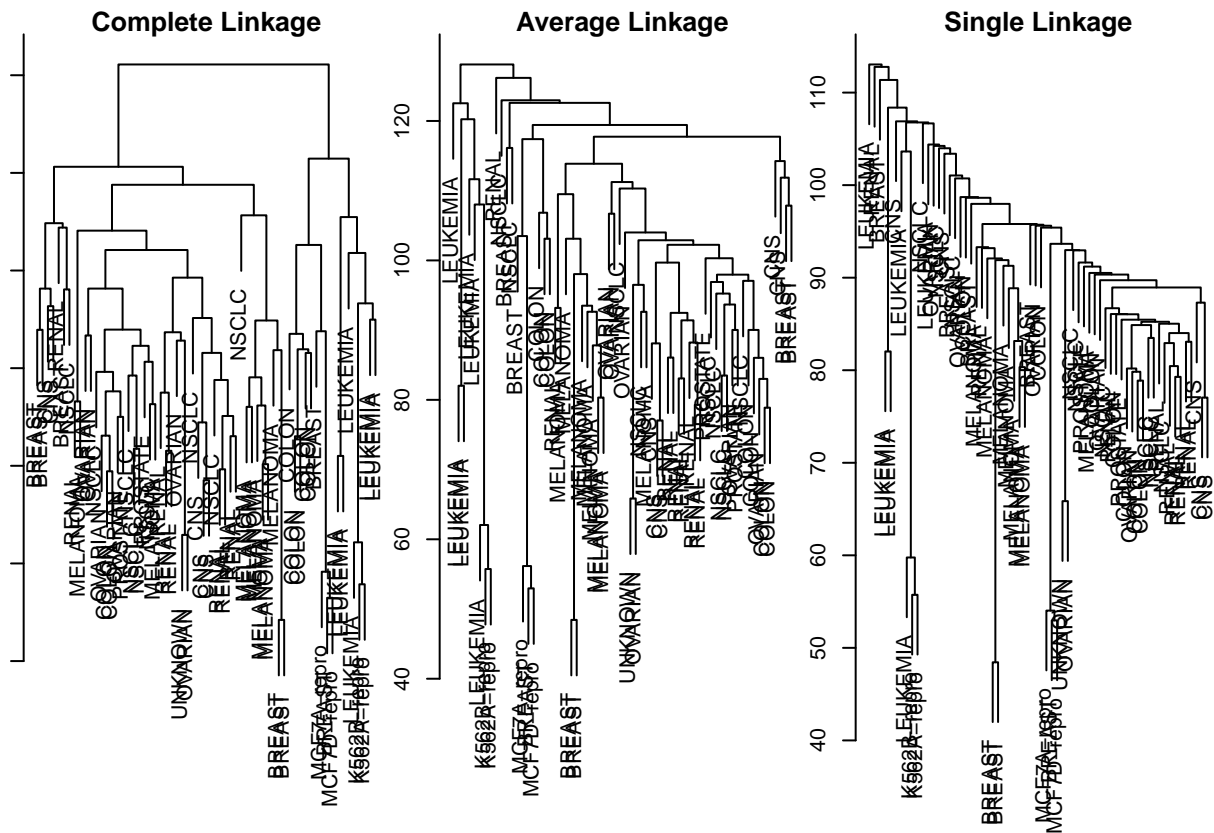
## Classical MDS of US Cities

New York                Los Angele

Miami

hington Atlanta Chicago

San Sean Fran

Houston

Denver

```
data.dist <- dist(sd.data)
plot(hclust(data.dist), xlab = "", sub = "", ylab = "",
     labels = nci.labs, main = "Complete Linkage")
plot(hclust(data.dist, method = "average"),
     labels = nci.labs, main = "Average Linkage",
     xlab = "", sub = "", ylab = "")
plot(hclust(data.dist, method = "single"),
     labels = nci.labs, main = "Single Linkage",
     xlab = "", sub = "", ylab = "")
```

```
# Interpretation:The choice of linkage certainly does affect the results obtained.
# Typically, single linkage will tend to yield trailing clusters: very large clusters
# onto which individual observations attach one-by-one. On the other hand, complete and
# average linkage tend to yield more balanced, attractive clusters. For this reason,
# complete and average linkage are generally preferred to single linkage. Clearly cell lines
# within a single cancer type do tend to cluster together, although the clustering is
# not perfect. We will use complete linkage hierarchical clustering for the analysis that follows
# We can cut the dendrogram at the height that will yield a particular number of clusters, say four.


# c. Find the best number for clusters using "cutree" function with best distance value
# Let's use the complete linkage method for this example

hc.out <- hclust(dist(sd.data))
hc.clusters <- cutree(hc.out, 4)
table(hc.clusters, nci.labs)
```
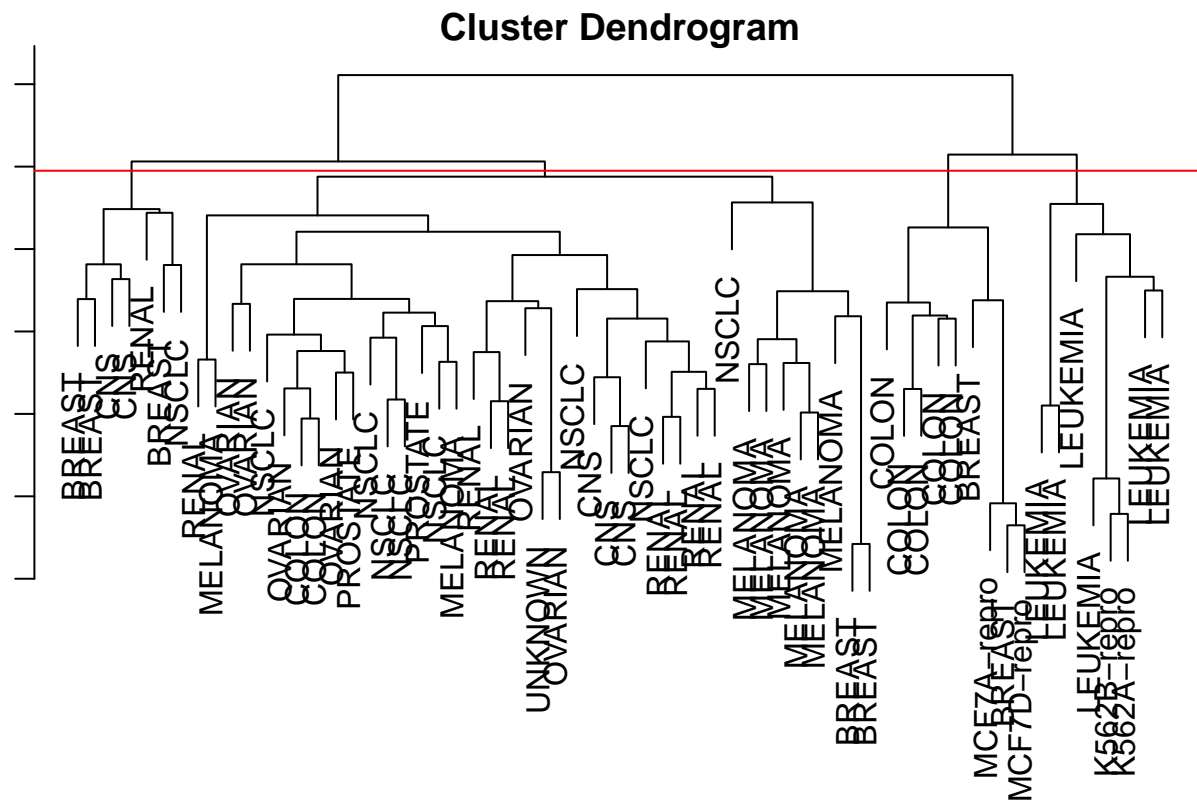
```
##           nci.labs
## hc.clusters BREAST CNS COLON K562A-repro K562B-repro LEUKEMIA MCF7A-repro
##          1      2   3     2           0           0        0           0
##          2      3   2     0           0           0        0           0
##          3      0   0     0           1           1        6           0
##          4      2   0     5           0           0        0           1
##           nci.labs
## hc.clusters MCF7D-repro MELANOMA NSCLC OVARIAN PROSTATE RENAL UNKNOWN
##          1           0        8     8       6        2     8       1
##          2           0        0     1       0        0     1       0
```

```
##          3         0        0        0        0        0        0        0
##          4         1        0        0        0        0        0        0
```

```r
# There are some clear patterns. All the leukemia cell lines fall in cluster 3, while
# the breast cancer cell lines are spread out over three diferent clusters. We can plot
# the cut on the dendrogram that produces these four clusters:
par(mfrow = c(1, 1))
plot(hc.out, labels = nci.labs)
abline(h = 139, col = "red")
```

## Cluster Dendrogram



```r
# Printing the output of hclust
hc.out
```

```
##
## Call:
## hclust(d = dist(sd.data))
##
## Cluster method   : complete
## Distance         : euclidean
## Number of objects: 64
```

```r
# d. Use your roll number as set.seed and perform k-means clustering on sd.data with
# the best number of clusters/distance value with nstart=20
set.seed(30)
km.out <- kmeans(sd.data, 4, nstart = 20)
km.clusters <- km.out$cluster
table(km.clusters, hc.clusters)
```

10

```
##          hc.clusters
## km.clusters  1  2  3  4
##           1  9  0  0  0
##           2  0  0  8  0
##           3 20  7  0  0
##           4 11  0  0  9
```

```r
# Interpretation: The four clusters obtained using hierarchical clustering and Kmeans
# clustering are somewhat different. Cluster 2 in K-means clustering is identical to
# cluster 1 in hierarchical clustering. However, the other clusters differ: for instance,
# cluster 3 in K-means clustering contains a portion of the observations assigned to
# cluster 4 by hierarchical clustering, as well as all of the observations assigned to
# cluster 3 by hierarchical clustering. Rather than performing hierarchical clustering on
# the entire data matrix, we can simply perform hierarchical clustering on the frst few
# principal component score vectors, as follows:


# e. Get summary of the k-means clustering and interpret them carefully
summary(km.out)
```

```
##              Length Class  Mode
## cluster         64  -none- numeric
## centers      27320  -none- numeric
## totss            1  -none- numeric
## withinss         4  -none- numeric
## tot.withinss     1  -none- numeric
## betweenss        1  -none- numeric
## size             4  -none- numeric
## iter             1  -none- numeric
## ifault           1  -none- numeric
```

```r
# f. Plot this k-means results using base r plot and cluster package and interpret them carefully
par(mfrow = c(1, 2))
plot(sd.data, col = km.out$cluster, main = "K-means Clustering Results", pch = 20)
points(km.out$centers, col = 1:4, pch = 8, cex = 2)
par(mfrow=c(1,1))



# Part IV: Use "Groceries" data available in the "datasets" package to do as follows
# in R Studio to knit PDF output

# a. Load "arules" and "arulesViz" libraries
library(arules)
```
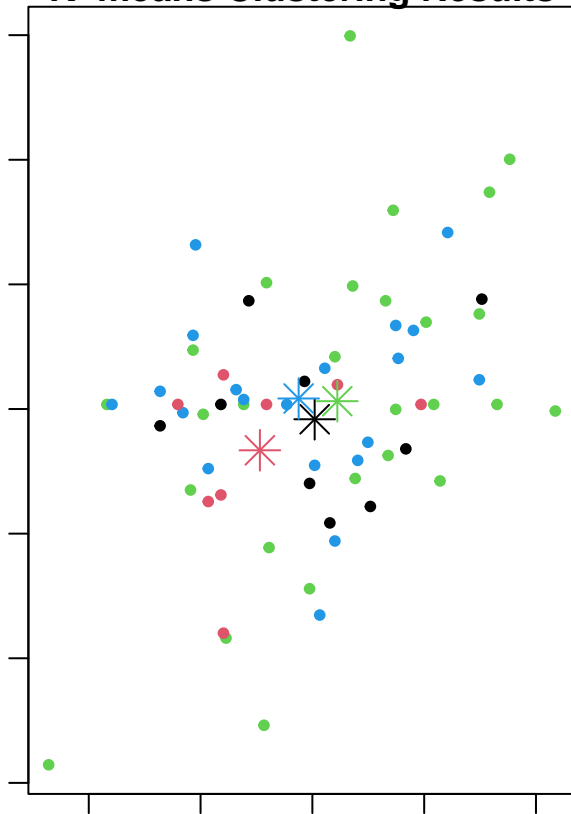
```
## Loading required package: Matrix


##
## Attaching package: 'arules'


## The following objects are masked from 'package:base':
##
##     abbreviate, write
```

## K−means Clustering Results



```r
library(arulesViz)
# b. Load "Groceries" data, check its structure and interpret it carefully
data("Groceries")
str(Groceries)
```

```
## Formal class 'transactions' [package "arules"] with 3 slots
##   ..@ data       :Formal class 'ngCMatrix' [package "Matrix"] with 5 slots
##   .. .. ..@ i       : int [1:43367] 13 60 69 78 14 29 98 24 15 29 ...
##   .. .. ..@ p       : int [1:9836] 0 4 7 8 12 16 21 22 27 28 ...
##   .. .. ..@ Dim     : int [1:2] 169 9835
##   .. .. ..@ Dimnames:List of 2
##   .. .. .. ..$ : NULL
##   .. .. .. ..$ : NULL
##   .. .. ..@ factors : list()
##   ..@ itemInfo   :'data.frame':  169 obs. of  3 variables:
##   .. ..$ labels: chr [1:169] "frankfurter" "sausage" "liver loaf" "ham" ...
##   .. ..$ level2: Factor w/ 55 levels "baby food","bags",..: 44 44 44 44 44 44 44 42 42 41 ...
##   .. ..$ level1: Factor w/ 10 levels "canned food",..: 6 6 6 6 6 6 6 6 6 6 ...
##   ..@ itemsetInfo:'data.frame':  0 obs. of  0 variables
```
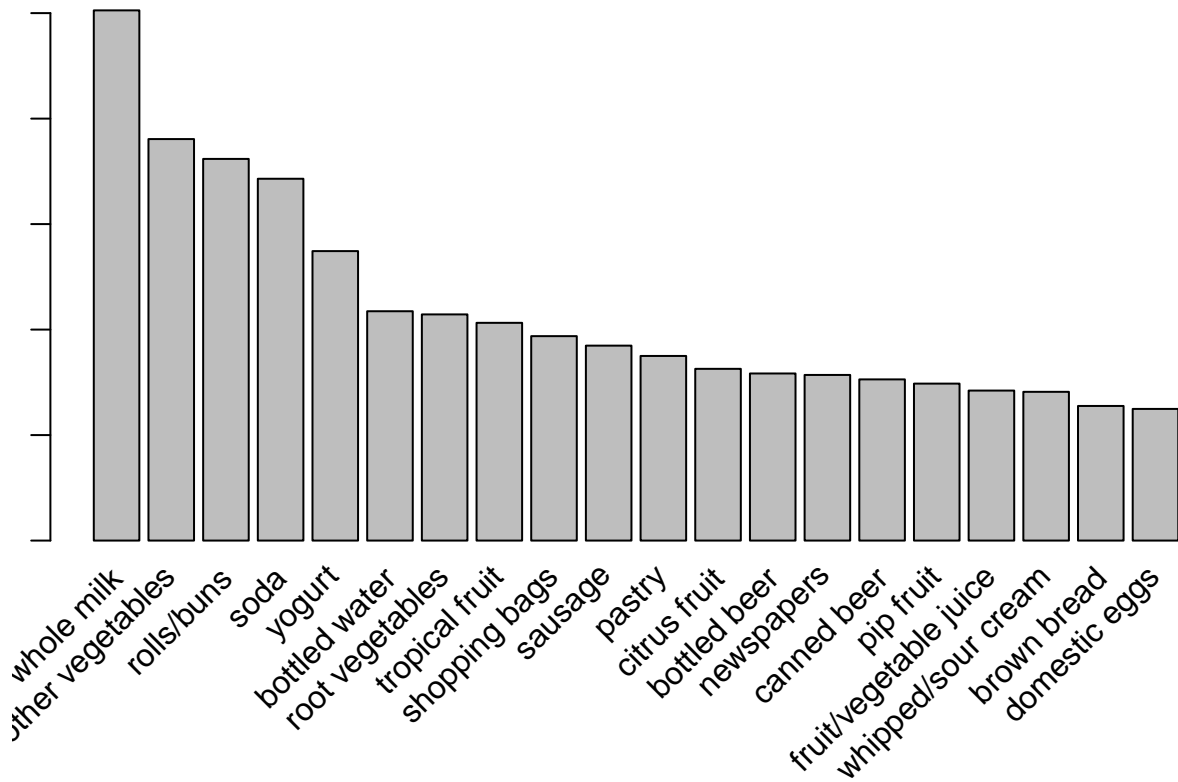
```r
summary(Groceries)
```

```
## transactions as itemMatrix in sparse format with
##  9835 rows (elements/itemsets/transactions) and
##  169 columns (items) and a density of 0.02609146
```

```
## 
## most frequent items:
##       whole milk other vegetables       rolls/buns            soda
##             2513             1903             1809            1715
##           yogurt          (Other)
##             1372            34055
## 
## element (itemset/transaction) length distribution:
## sizes
##     1     2     3     4     5     6     7     8     9    10    11    12    13    14    15    16
## 2159  1643  1299  1005   855   645   545   438   350   246   182   117    78    77    55    46
##    17    18    19    20    21    22    23    24    26    27    28    29    32
##    29    14    14     9    11     4     6     1     1     1     1     3     1
## 
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   1.000   2.000   3.000   4.409   6.000  32.000
## 
## includes extended item information - examples:
##        labels  level2            level1
## 1 frankfurter sausage meat and sausage
## 2     sausage sausage meat and sausage
## 3  liver loaf sausage meat and sausage
```

```
# Interpretation: The Groceries dataset is a sparse matrix of transactions, where each
# row represents a transaction and each column represents an item. The summary gives
# an overview of the number of transactions, the number of items, and some basic statistics.
# c. Get Frequent Item frequencies using itemFrequencyPlot function and interpret it carefully.
```

```r
itemFrequencyPlot(Groceries, topN = 20, type = "absolute")
```

```
# Interpretation: This plot shows the absolute frequencies of the top 20 items. Items
# like "whole milk" and "other vegetables" are among the most frequently purchased items

# d. Set a priori rule with support = 0.001 and confidence = 0.8 and interpret the output carefully
rules <- apriori(Groceries, parameter = list(support = 0.001, confidence = 0.8))
```

```
## Apriori
##
## Parameter specification:
##  confidence minval smax arem  aval originalSupport maxtime support minlen
##         0.8    0.1    1 none FALSE            TRUE       5   0.001      1
##  maxlen target  ext
##      10  rules TRUE
##
## Algorithmic control:
##  filter tree heap memopt load sort verbose
##     0.1 TRUE TRUE  FALSE TRUE    2    TRUE
##
## Absolute minimum support count: 9
##
## set item appearances ...[0 item(s)] done [0.00s].
## set transactions ...[169 item(s), 9835 transaction(s)] done [0.00s].
## sorting and recoding items ... [157 item(s)] done [0.00s].
## creating transaction tree ... done [0.00s].
## checking subsets of size 1 2 3 4 5 6 done [0.01s].
## writing ... [410 rule(s)] done [0.00s].
## creating S4 object  ... done [0.00s].
```

```
summary(rules)
```

```
## set of 410 rules
##
## rule length distribution (lhs + rhs):sizes
##   3   4   5   6
##  29 229 140  12
##
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   3.000   4.000   4.000   4.329   5.000   6.000
##
## summary of quality measures:
##     support          confidence         coverage            lift
##  Min.   :0.001017   Min.   :0.8000   Min.   :0.001017   Min.   : 3.131
##  1st Qu.:0.001017   1st Qu.:0.8333   1st Qu.:0.001220   1st Qu.: 3.312
##  Median :0.001220   Median :0.8462   Median :0.001322   Median : 3.588
##  Mean   :0.001247   Mean   :0.8663   Mean   :0.001449   Mean   : 3.951
##  3rd Qu.:0.001322   3rd Qu.:0.9091   3rd Qu.:0.001627   3rd Qu.: 4.341
##  Max.   :0.003152   Max.   :1.0000   Max.   :0.003559   Max.   :11.235
##      count
##  Min.   :10.00
##  1st Qu.:10.00
##  Median :12.00
##  Mean   :12.27
```

```
##  3rd Qu.:13.00
##  Max.   :31.00
##
## mining info:
##       data ntransactions support confidence
##  Groceries         9835   0.001        0.8
##                                                            call
##  apriori(data = Groceries, parameter = list(support = 0.001, confidence = 0.8))
```

```r
# Interpretation: The apriori algorithm generates association rules based on the specified
# minimum support and confidence. The summary provides the number of rules generated
# and basic statistics like minimum, maximum, and average support and confidence.

# e. Show the top five rules using inspect and round the results to two digits
inspect(head(rules, 5))
```

```
##     lhs                     rhs            support     confidence
## [1] {liquor, red/blush wine} => {bottled beer} 0.001931876 0.9047619
## [2] {curd, cereals}          => {whole milk}   0.001016777 0.9090909
## [3] {yogurt, cereals}        => {whole milk}   0.001728521 0.8095238
## [4] {butter, jam}            => {whole milk}   0.001016777 0.8333333
## [5] {soups, bottled beer}    => {whole milk}   0.001118454 0.9166667
##     coverage    lift      count
## [1] 0.002135231 11.235269 19
## [2] 0.001118454  3.557863 10
## [3] 0.002135231  3.168192 17
## [4] 0.001220132  3.261374 10
## [5] 0.001220132  3.587512 11
```

```r
# f. Sort the rule by confidence in decreasing order
rules <- sort(rules, by = "confidence", decreasing = TRUE)
inspect(head(rules, 5))
```

```
##     lhs                rhs            support confidence   coverage     lift count
## [1] {rice,
##      sugar}            => {whole milk} 0.001220132          1 0.001220132 3.913649    12
## [2] {canned fish,
##      hygiene articles} => {whole milk} 0.001118454          1 0.001118454 3.913649    11
## [3] {root vegetables,
##      butter,
##      rice}             => {whole milk} 0.001016777          1 0.001016777 3.913649    10
## [4] {root vegetables,
##      whipped/sour cream,
##      flour}            => {whole milk} 0.001728521          1 0.001728521 3.913649    17
## [5] {butter,
##      soft cheese,
##      domestic eggs}    => {whole milk} 0.001016777          1 0.001016777 3.913649    10
```

```r
# g. Use "whole milk" as target item and show the items in "lhs" with decreasing order
# of confidence and show the top five rules

rules_lhs <- subset(rules, lhs %pin% "whole milk")
```

```r
rules_lhs <- sort(rules_lhs, by = "confidence", decreasing = TRUE)
inspect(head(rules_lhs, 5))
```

```
##      lhs                    rhs                  support confidence    coverage     lift count
## [1] {tropical fruit,
##      grapes,
##      whole milk,
##      yogurt}            => {other vegetables} 0.001016777  1.0000000 0.001016777 5.168156    10
## [2] {ham,
##      tropical fruit,
##      pip fruit,
##      whole milk}        => {other vegetables} 0.001118454  1.0000000 0.001118454 5.168156    11
## [3] {whole milk,
##      rolls/buns,
##      soda,
##      newspapers}        => {other vegetables} 0.001016777  1.0000000 0.001016777 5.168156    10
## [4] {root vegetables,
##      whole milk,
##      yogurt,
##      oil}               => {other vegetables} 0.001423488  0.9333333 0.001525165 4.823612    14
## [5] {citrus fruit,
##      tropical fruit,
##      root vegetables,
##      whole milk,
##      yogurt}            => {other vegetables} 0.001423488  0.9333333 0.001525165 4.823612    14
```

```r
# h. Use "whole milk" as target item and show the items in "rhs" with decreasing order
# of confidence and show the top five rules
rules_rhs <- subset(rules, rhs %pin% "whole milk")
rules_rhs <- sort(rules_lhs, by = "confidence", decreasing = TRUE)
inspect(head(rules_rhs, 5))
```

```
##      lhs                    rhs                  support confidence    coverage     lift count
## [1] {tropical fruit,
##      grapes,
##      whole milk,
##      yogurt}            => {other vegetables} 0.001016777  1.0000000 0.001016777 5.168156    10
## [2] {ham,
##      tropical fruit,
##      pip fruit,
##      whole milk}        => {other vegetables} 0.001118454  1.0000000 0.001118454 5.168156    11
## [3] {whole milk,
##      rolls/buns,
##      soda,
##      newspapers}        => {other vegetables} 0.001016777  1.0000000 0.001016777 5.168156    10
## [4] {root vegetables,
##      whole milk,
##      yogurt,
##      oil}               => {other vegetables} 0.001423488  0.9333333 0.001525165 4.823612    14
## [5] {citrus fruit,
##      tropical fruit,
##      root vegetables,
##      whole milk,
```

```
##      yogurt}          => {other vegetables} 0.001423488   0.9333333 0.001525165 4.823612      14
```

```
# i. Write summary and conclusion based on your findings above
# Based on the analysis, the "Groceries" data reveals several interesting associations:
# 1.Frequent Items: Items like "whole milk" and "other vegetables" are among the most
# frequently purchased.
# 2. Association Rules: Setting a minimum support of 0.001 and confidence of 0.8, we
# generated a significant number of rules, indicating strong associations among items.
# 3. Top Rules: By inspecting the top rules, we can see high-confidence associations,
# such as items frequently bought together with "whole milk."
# 4. Targeted Rules: When "whole milk" is the target item in the LHS or RHS, the rule
# s show high confidence and lift, highlighting items that are commonly bought with "whole milk."
# Overall, the analysis provides valuable insights into customer purchasing behavior,
# which can be leveraged for marketing strategies and inventory management.
```