**King Saud University**

**College of Computer and Information Sciences**

**Information Technology department**

# IT 326: Data Mining
# Course Project

# McDonald menu

**Project Report**

Group #: 8                                              LAB Day-Time: Monday (8-10)

Group members:

| Name | ID | Section |
|---|---|---|
| Nada Almutairi | | |
| | | |

[Pick the date]

# 1 Problem

The availability of fast foods, fatty food, and unhealthy snacks in the living environment of children is assumed to give a very big reason of obesity among children.

In particular, it is very clear that food retailers are focusing on kid's using different strategies to dazzle them. We are using McDonald's menu to investigate :

(1) the clustering of food facts in the whole menu .

 (2) the effects of particular food item in mac's big menu on children.

McDonald menu contains breakfasts and fast food like burgers and sandwiches and we would like to analyze it with respect to the content of food and quality. Also, we want to know all the good and bad sides of McDonald's most famous menu items such as Mc Nuggets, Big Mac and McChicken.
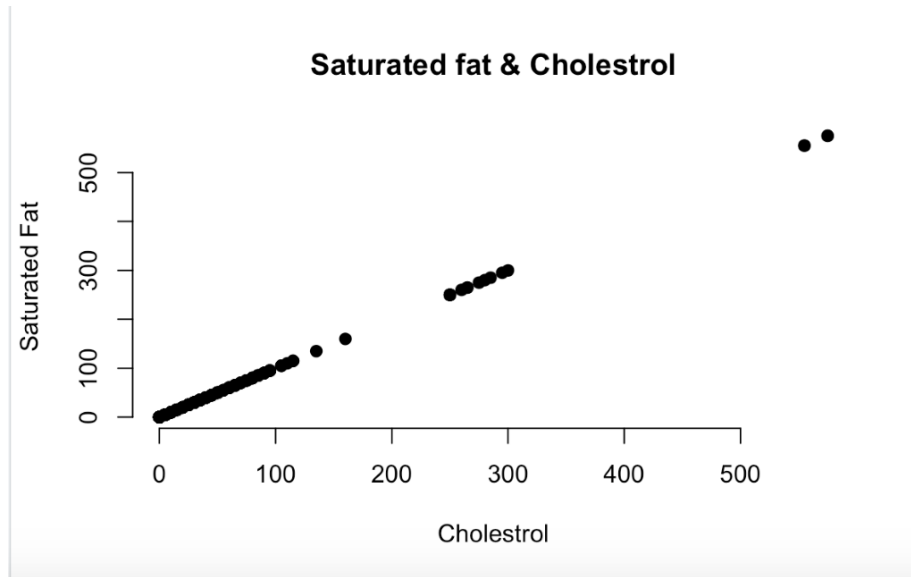
# 2 Data Mining Task

For this problem we will use clustering on food facts in the overall menu of US McDonald's to predict reasons of obesity among children.
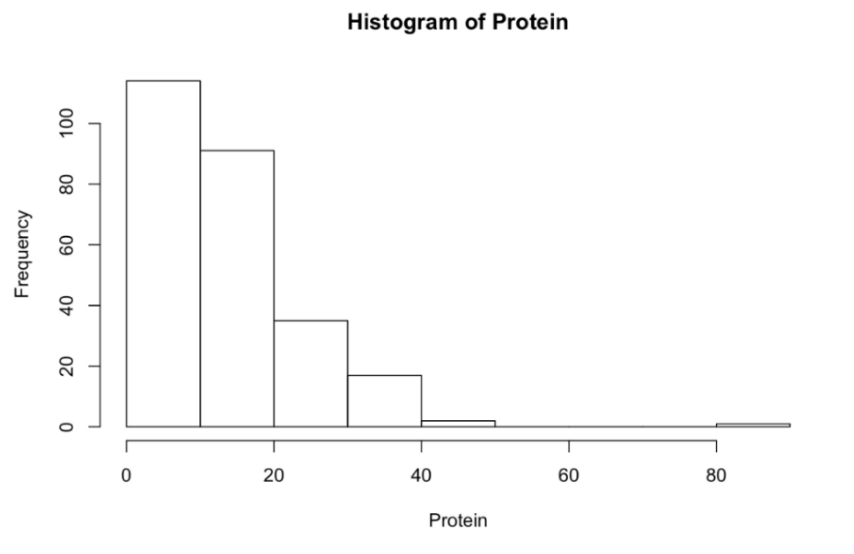
# 3 Data

| Dats set information | |
|---|---|
| **This data set contains nutrition facts for mcDonalds' menu items** | |
| **Source** | Kaggle<br>https://www.kaggle.com/mcdonalds/nutrition-facts/data |
| **Number of observations** | 260 |
| **Number of attributes** | 24 |
| **Missing values** | No missing values |

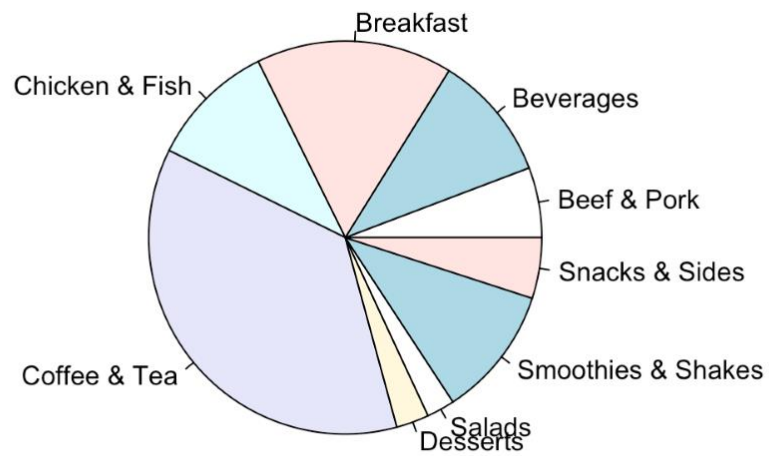| Attribute | Attribute information | Data type |
|---|---|---|
| Category | The types of meals | String |
| Item | Items of food sold from the McDonalds' US menu | String |
| Serving size | Recommended serving size for each item | String |
| Calories | The amount of calories in each item | Integer |
| Calories from fat | The amount of calories in each item that comes from fat | Integer |
| Total fat | The total fat for each item | Integer |
| Total fat (% daily value) | How much nutrition value in % from fat in an item | Integer |
| Saturated fat | How much saturated fat in an item | Integer |
| Saturated fat(% daily value) | How much nutrition value in % from saturated fat in an item | Integer |
| Trans fat | How much trans fat in an item | Integer |
| Cholesterol | How much choleterol in an item | Integer |
| Cholesterol (% daily value) | How much nutrition value in % from cholesterol in an item | Integer |
| Sodium | How much sodium in an item | Integer |
| Sodium (% daily value) | How much nutrition value in % from sodium in an item | Integer |
| Carbohydrates | How much carbohydrates in an item | Integer |
| Carbohydrates (% daily value) | How much nutrition value in % from carbohydrates in an item | Integer |
| Dietary fiber | How much dietary fiber in an item | Integer |
| Dietary fiber (% daily value) | How much nutrition value in % from dietary fiber in an item | Integer |
| Sugars | How much sugars in an item | Integer |
| Protein | How much protein in an item | Integer |
| Vitamin A (% daily value) | How much nutrition value in % from vitamin A in an item | Integer |
| Vitamin C (% daily value) | How much nutrition value in % from vitamin C in an item | Integer |
| Calcium (% daily value) | How much nutrition value in % from calcium in an item | Integer |
| Iron (% daily value) | How much nutrition value in % from iron in an item | Integer |

# 4 Data preprocessing



**Saturated fat & Cholestrol**

The graph shows that the more number you have of saturated fat in a meal the more the cholesterol in the product as increases , which also might  increase the risk of getting heart disease and stroke.



**Histogram of Protein**

This histogram made us find out  that the data is skewed on the right-hand side having highest frequency of items with low protein while protein is an essential element to build muscles especially for children

From this pie chart we can tell that the most category that its items varies the most is coffee&tea, which is not the children's first choice anyway, and the salads has the least amount of menu items which is less options to choose from. and it is one of the healthiest categories.

To analysis our dataset, we had to do some of preprocessing.
so, we start of deleting the first column which is the category. Because it is just a meal information that we will not be needing it in our study.
Also, we converted the first two columns to numeric , Since we still have two more character attributes Items and serving size.
finally, we want to make sure if there were any missing value or not.

- Important note :

   we can not transform the column Item to numeric because we need the items' value but not

   the category itself ,so we transformed to row names to perform data mining task.

```r
1   #read the dataset
2   macData2 = read.csv('Mac.csv')
3   #Cleaning the data
4   install.packages("dplyr")
5   library(dplyr)
6   #Delete the first column
7   macData2$Category <-NULL
8
9   #Now we will convert column 2 to numeric
10  macData2[, 2] <-sapply(macData2[,2], as.numeric)
11
12  #Remove duplicates
13  CleanMenuData <- macData2 %>% distinct(Item , .keep_all = TRUE)
14
15  #Now we cheak if there is any missing value
16  sum(is.na(CleanMenuData))
17
18  #Set item as row names so that we can delete that column
19  rownames(CleanMenuData) <- CleanMenuData[,1]
20  #Since we do Clusting method ,we deleted the column Item after retaining it as Row names
21  CleanMenuData$Item <- NULL
22
23
24  ScaledMenuData <- scale(CleanMenuData, center = FALSE)
25
26
```
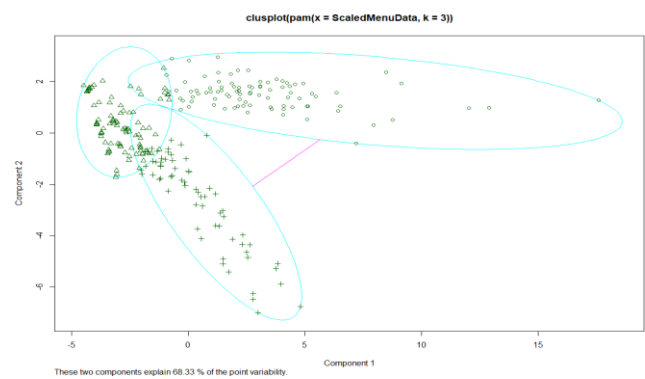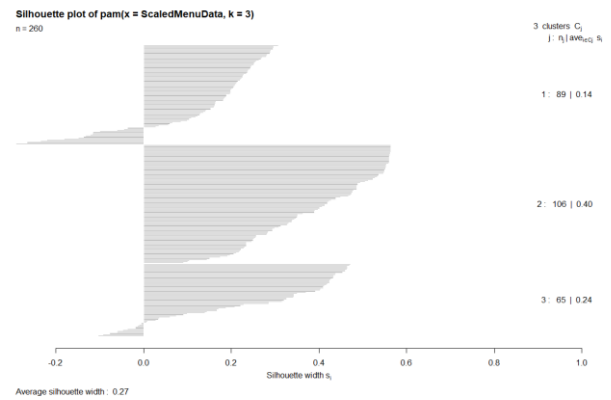
# 5 Data Mining Technique

We used clustering technique which is used to find the structure of unlabeled data. Since that the type of data we have , and in order to discover groups of objects where the average distances between the members of each cluster shall be closer than to members in other clusters. All of that is for better understanding of the data by categorizing them.
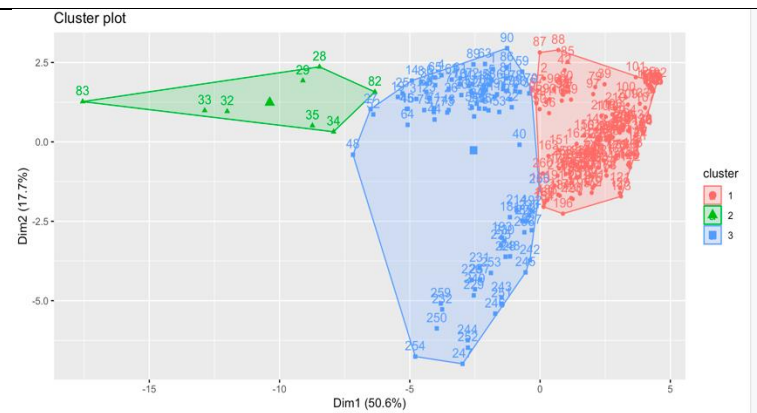
we used dplyr package for preprocessing , and factoextra , cluster , NbClust and ggplot2 for clustering .

| | 3 clusters |
|---|---|
| **Silhouette width for 3 cluster** |  |
| **Visualization** | <br><br>clusters are way close to each other which is not what we want, only cluster 2,3 have a good distance between them |

| | 4 clusters |
|---|---|
| **Silhouette width for 4 cluster** |  |
| **Visualization** | \n\nwe now can tell that the more the number of clusters increases the better the distance between clusters get |

| | 5 clusters |
|---|---|
| **Silhouette width for 5 cluster** |  |
| **Visualization** | 

clusters are getting more spaced out |

| | 6 clusters |
|---|---|
| **Silhouette width for 6 cluster** | **Silhouette plot of pam(x = ScaledMenuData, k = 6)**<br>n = 260<br><br>6 clusters $C_j$<br>j : $n_j$ \| $ave_{i \in C_j}$ $s_i$<br><br>1 : 58 \| 0.27<br>2 : 25 \| 0.005<br>3 : 23 \| -0.07<br>4 : 65 \| 0.38<br>5 : 65 \| 0.28<br>6 : 24 \| 0.47<br><br>-0.2  0.0  0.2  0.4  0.6  0.8  1.0<br>Silhouette width $s_i$<br>Average silhouette width : 0.26<br><br>**clusplot(pam(x = ScaledMenuData, k = 6))**<br>Component 2<br>Component 1<br>These two components explain 68.33 % of the point variability. |
| **Visualization** | Cluster plot |

- **Silhouette width for all cluster**



Optimal number of clusters
Silhouette method

# 6  Findings

As we see the optimal number of cluster is 4. Now we will performs k-means clustering. Since the data has more than two attributes , we have to choose two attributes and studying them based on our problem.

So, our focus is mainly on the calories and total fat, which is the worst things that we can give to the children. And we will try avoid them.

```
> menuKCluster <- kmeans(CleanMenuData[, 3:5], 4)
> print(menuKCluster)
K-means clustering with 4 clusters of sizes 75, 8, 80, 97

Cluster means:
  Calories.from.Fat Total.Fat Total.Fat....Daily.Value.
1          239.8667 26.666667                  41.080000
2          550.0000 61.125000                  94.125000
3          115.0000 12.837500                  19.725000
4           15.0000  1.721649                   2.680412

Clustering vector:
                                              Egg McMuffin
                                                         3
                                         Egg White Delight
                                                         3
                                          Sausage McMuffin
                                                         1
                                 Sausage McMuffin with Egg
                                                         1
                          Sausage McMuffin with Egg Whites
                                                         1
                                       Steak & Egg McMuffin
                                                         1
               Bacon, Egg & Cheese Biscuit (Regular Biscuit)
                                                         1
                 Bacon, Egg & Cheese Biscuit (Large Biscuit)
                                                         1
  Bacon, Egg & Cheese Biscuit with Egg Whites (Regular Biscuit)
                                                         1
    Bacon, Egg & Cheese Biscuit with Egg Whites (Large Biscuit)
                                                         1
                            Sausage Biscuit (Regular Biscuit)
                                                         1
                              Sausage Biscuit (Large Biscuit)
                                                         1
                   Sausage Biscuit with Egg (Regular Biscuit)
                                                         1
                     Sausage Biscuit with Egg (Large Biscuit)
                                                         1
            Sausage Biscuit with Egg Whites (Regular Biscuit)
                                                         1
              Sausage Biscuit with Egg Whites (Large Biscuit)
                                                         1
               Southern Style Chicken Biscuit (Regular Biscuit)
```

```
                                                                    1
                                Strawberry Shake (Large)
                                                                    1
                                Chocolate Shake (Small)
                                                                    3
                                Chocolate Shake (Medium)
                                                                    1
                                Chocolate Shake (Large)
                                                                    1
                                Shamrock Shake (Medium)
                                                                    3
                                Shamrock Shake (Large)
                                                                    1
                  McFlurry with M&Mâ€™s Candies (Small)
                                                                    1
                  McFlurry with M&Mâ€™s Candies (Medium)
                                                                    1
                  McFlurry with M&Mâ€™s Candies (Snack)
                                                                    3
                       McFlurry with Oreo Cookies (Small)
                                                                    3
                       McFlurry with Oreo Cookies (Medium)
                                                                    1
                       McFlurry with Oreo Cookies (Snack)
                                                                    3
        McFlurry with Reese's Peanut Butter Cups (Medium)
                                                                    1
        McFlurry with Reese's Peanut Butter Cups (Snack)
                                                                    3
Within cluster sum of squares by cluster:
[1] 189890.85 325779.75  88138.84  42555.58
 (between_SS / total_SS =  85.3 %)

Available components:

[1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss" "betweenss"     "size"
[8] "iter"         "ifault"
> menuKCluster$centers
  Calories.from.Fat Total.Fat Total.Fat....Daily.Value.
1          239.8667 26.666667                  41.080000
2          550.0000 61.125000                  94.125000
3          115.0000 12.837500                  19.725000
4           15.0000  1.721649                   2.680412
>
```
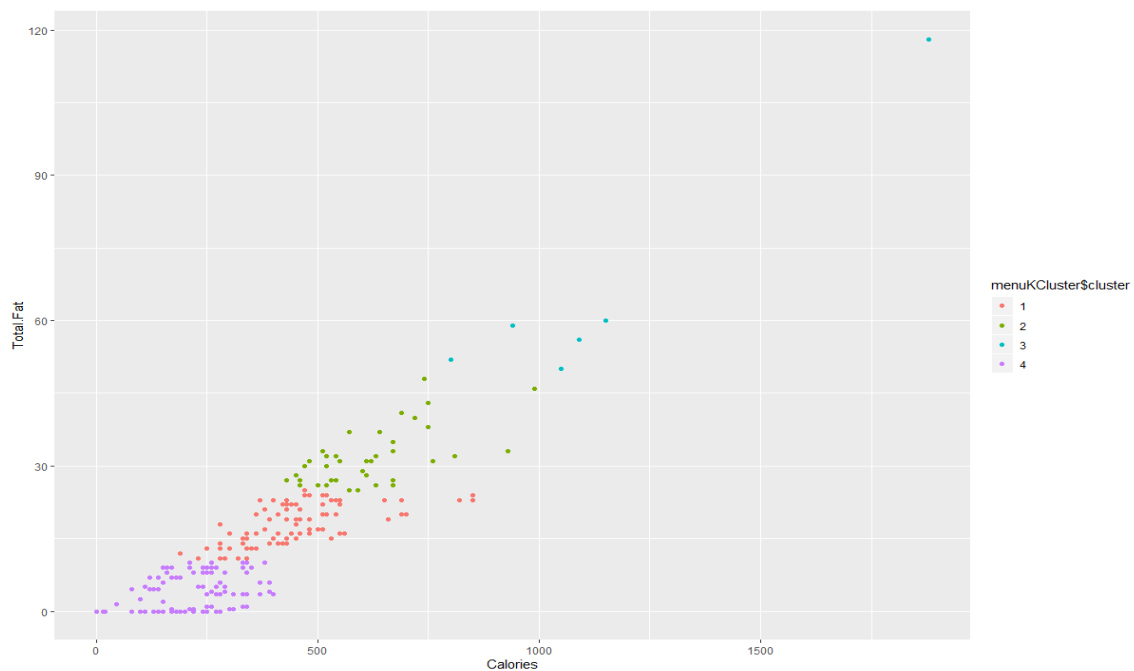
Then by plotting the points, each color represents the cluster.
As we can see there is a positive correlation between calories and total fat. As the calories and total fat increase, another cluster is formatted.

# 7 Code

Pre processing Code :

```r
1   #read the dataset
2   macData2 = read.csv('Mac.csv')
3   #Cleaning the data
4   install.packages("dplyr")
5   library(dplyr)
6   #Delete the first column
7   macData2$Category <-NULL
8
9   #Now we will convert column 2 to numeric
10  macData2[, 2] <-sapply(macData2[,2], as.numeric)
11
12  #Remove duplicates
13  CleanMenuData <- macData2 %>% distinct(Item , .keep_all = TRUE)
14
15  #Now we cheak if there is any missing value
16  sum(is.na(CleanMenuData))
17
18  #Set item as row names so that we can delete that column
19  rownames(CleanMenuData) <- CleanMenuData[,1]
20  #Since we do Clusting method ,we deleted the column Item after retaining it as Row names
21  CleanMenuData$Item <- NULL
22
23
24  ScaledMenuData <- scale(CleanMenuData, center = FALSE)
25
26
```

# Clustering Code :

```r
27  #Now we will do kmeans clustering to find appropatie No. of clusters
28  install.packages('factoextra')
29  library(factoextra)
30
31  set.seed(150)
32  kmeans3.result <-kmeans(ScaledMenuData,3)
33  kmeans3.result
34  fviz_cluster(kmeans3.result, data = ScaledMenuData)
35
36  set.seed(200)
37  kmeans4.result <-kmeans(ScaledMenuData,4)
38  kmeans4.result
39  fviz_cluster(kmeans4.result, data = ScaledMenuData)
40
41  set.seed(122)
42  kmeans5.result <-kmeans(ScaledMenuData,5)
43  kmeans5.result
44  fviz_cluster(kmeans5.result, data = ScaledMenuData)
45
46  set.seed(346)
47  kmeans6.result <-kmeans(ScaledMenuData,6)
48  kmeans6.result
49  fviz_cluster(kmeans6.result, data = ScaledMenuData)
50  #clustering with PAM
51  install.packages('cluster')
52  library(cluster)
53  # Trying with all four clusters
54  pam.result <- pam(ScaledMenuData, 3)
55  plot(pam.result)
56
57  pam.result <- pam(ScaledMenuData, 4)
58  plot(pam.result)
59
60  pam.result <- pam(ScaledMenuData, 5)
61  plot(pam.result)
62
63  pam.result <- pam(ScaledMenuData, 6)
64  plot(pam.result)
65  #Fainlly we will see the optimal number of clustering
66  install.packages('NbClust')
67  library(NbClust)
68  fviz_nbclust(ScaledMenuData, kmeans, method = "silhouette")+ labs(subtitle = "Silhouette method")
```

```
70
71  #View(CleanMenuData)
72
73  #Now we will performs k-means clustering with the best number of clusters which is 4
74  #Also with the two attributes calories and total fat (3 and 5)
75  set.seed(30)
76  menuKCluster <- kmeans(CleanMenuData[, 3:5], 4)
77  print(menuKCluster)
78  menuKCluster$centers
79  menuKCluster$cluster <- as.factor(menuKCluster$cluster)
80  #Now we want to visualize the result
81  library(ggplot2)
82  ggplot(CleanMenuData, aes(Calories, Total.Fat, color = menuKCluster$cluster), environment=environment())
83  + geom_point()
84
85
86
87
88
89
90  |
91
92
90:1    (Top Level) ↕                                                                                          R Script ↕
```

# 8   References

**E-books**

[1]      Practical guide to cluster analysis in R (Edition 1)

**General Internet Site**

[2]      http://www.rdocumentation.org/

[3]      https://www.datanovia.com/