# Market Segmentation – A Case Study

By

Hardik Sharma
Mridul Jain
Shyamshree Ghorai
Vansh Sonavane

Date: 14th July 2023

MARKET SEGMENTATION

# STEP 1: CHECKLIST

The provided checklist outlines the initial steps and considerations for deciding whether to proceed with market segmentation. It includes a series of tasks and questions that serve as knock-out criteria. The summary of the checklist is as follows:

The checklist begins by assessing the organization's culture and willingness to change. If the organization is market-oriented and open to new ideas, the process can proceed. Long-term perspective, good communication across units, and the ability to make structural changes are also important factors.

Financial resources play a crucial role, as the organization should have sufficient funds to support a market segmentation strategy. Visible commitment from senior management is necessary, along with their active involvement and required financial support.

Understanding the market segmentation concept and its implications is essential. Training should be conducted until these aspects are fully grasped. Forming a segmentation team with expertise in marketing, data analysis, and data management is important.

An advisory committee representing all affected organizational units should be set up. Clear objectives for the market segmentation analysis should be established, and a structured process should be developed and followed. Responsibilities should be assigned to team members according to the process.

Sufficient time should be allocated for the market segmentation analysis, without time pressure.

In summary, the checklist emphasizes the importance of a market-oriented culture, willingness to change, long-term perspective, open communication, financial resources, and senior management commitment. It also highlights the need for a knowledgeable team, clear objectives, and a structured process to conduct the market segmentation analysis.

# STEP 2: SPECIFYING THE IDEAL TARGET SEGMENT

In Step 2 of market segmentation analysis, the focus is on specifying the ideal target segment by establishing segment evaluation criteria. User input plays a crucial role throughout the process, as their involvement is needed to ensure the usefulness of the analysis. The organization must contribute conceptually to the market segmentation analysis, guiding subsequent steps.

Two sets of segment evaluation criteria are defined: knock-out criteria and attractiveness criteria. Knock-out criteria are essential and non-negotiable features that segments must possess to be considered as potential targets. They include factors such as homogeneity, distinctiveness, size, match with organizational strengths, identifiability, and reachability.

On the other hand, attractiveness criteria are used to assess the relative attractiveness of segments that comply with the knock-out criteria. There is a wide range of proposed criteria in the literature, and the segmentation team selects a subset of no more than six criteria that are most relevant to their specific situation.

A structured process is recommended for evaluating market segments. A popular approach involves using a segment evaluation plot, which shows segment attractiveness and organizational competitiveness. The criteria for both attractiveness and competitiveness need to be negotiated and agreed upon by the segmentation team, preferably with input from representatives of different organizational units in the advisory committee.



During this step, the team determines the segment attractiveness criteria and assigns weights to each criterion, indicating its relative importance. The weights are typically decided through team discussions and negotiations, ensuring agreement among team members. The proposed criteria and weights are then presented to the advisory committee for further discussion and adjustment if necessary.

By the end of Step 2, the segmentation team should have a list of segment attractiveness criteria, each with its assigned weight. This groundwork facilitates data collection in Step 3 and simplifies the selection of a target segment in Step 8.

The checklist for Step 2 includes tasks such as convening a segmentation team meeting, discussing and agreeing on knock-out criteria, presenting them to the advisory committee, studying and selecting attractiveness criteria, distributing weights, and presenting the selected criteria and weights to the advisory committee for further discussion and adjustment if needed.

# STEP 3: DATA COLLECTION

Step 3 of market segmentation involves collecting data to identify and describe market segments. Empirical data is essential for both commonsense and data-driven segmentation. In commonsense segmentation, a single characteristic, such as gender, is used as the segmentation variable to divide the sample into segments. Other personal characteristics serve as descriptor variables to describe the segments in detail.

Data-driven segmentation, on the other hand, utilizes multiple segmentation variables to identify naturally existing or artificially created market segments. These variables can be characteristics or benefits sought by consumers. The quality of empirical data is crucial in developing valid segmentation solutions and accurately describing the segments.

Data for segmentation studies can be obtained from various sources such as surveys, observations (e.g., scanner data), or experimental studies. Survey data is commonly used but may have limitations in reflecting actual behavior, especially for socially desirable actions. Therefore, alternative data sources should be explored to reflect consumer behavior accurately.

Segmentation criteria are the basis for market segmentation and can be geographic, socio-demographic, psychographic, or behavioral. Geographic segmentation uses location as the criterion, while socio-demographic segmentation considers age, gender, income, and education. Psychographic segmentation focuses on psychological criteria, such as beliefs, interests, preferences, and benefits sought. Behavioral segmentation analyzes actual behavior or reported behavior, such as purchase history or information search behavior.

When collecting data through surveys, the response options provided to respondents play a crucial role. Binary or metric response options are preferred for segmentation analysis as they enable distance measures and statistical procedures. Ordinal response options have an ordered scale but lack a clearly defined distance between adjacent options.

Survey data is susceptible to response biases and styles, which can impact the segmentation results. Response biases include extreme or midpoint responses, while response styles are consistent biases shown by respondents over time. Researchers should minimize these biases to ensure accurate segment identification.

Sample size is another important consideration. Larger sample sizes improve the accuracy of segment extraction. The recommended sample size depends on the number of segmentation variables and should be sufficient to enable correct identification of segments.

Internal data from organizations, such as scanner data or online purchase data, can be valuable for segmentation analysis as they reflect actual consumer behavior. However, caution should be exercised to avoid bias towards existing customers and ensure representation of potential future customers.

Experimental studies can provide data for segmentation analysis, particularly through tests on consumer responses to advertisements or choice experiments that assess preferences for specific product attributes.

In summary, collecting high-quality data is crucial for effective market segmentation. Careful consideration should be given to the selection of segmentation variables, data sources, response options, sample size, and potential biases to ensure accurate segment identification and description.

# Step 4: Exploring Data

After data collection, this step involves the following steps:

1) Dataset and Features
2) Data Cleaning
3) Data Analysis
4) Data Pre-Processing and Preparation
5) Principal Component Analysis

This step first cleans the data and then pre-process if necessary and then provides insights about the suitability of different segmentation methods for extracting the market segments.

Data Description:

To illustrate this step, we use Australian travel motives dataset. The dataset contains 20 travel motives reported by 1000 Australian residents. The dataset consists of 1000 rows and 32 columns.

In R, using Col names () command, we obtained features of the dataset:

```
1  install.packages('tidyverse')
2  install.packages(c('quantmod','ff','foreign','R.matlab'),dependency=T)
3
4  suppressPackageStartupMessages(library(tidyverse))
5  data1 <- read.csv('vacation_complete_dataset.csv')
```

```
· colnames(data1)
 [1] "Gender"                      "Age"
 [3] "Education"                    "Income"
 [5] "Income2"                      "Occupation"
 [7] "State"                        "Relationship.Status"
 [9] "Obligation"                   "Obligation2"
[11] "NEP"                          "Vacation.Behaviour"
[13] "rest.and.relax"               "luxury...be.spoilt"
[15] "do.sports"                    "excitement..a.challenge"
[17] "not.exceed.planned.budget"    "realise.creativity"
[19] "fun.and.entertainment"        "good.company"
[21] "health.and.beauty"            "free.and.easy.going"
[23] "entertainment.facilities"     "not.care.about.prices"
[25] "life.style.of.the.local.people" "intense.experience.of.nature"
[27] "cosiness.familiar.atmosphere" "maintain.unspoilt.surroundings"
[29] "everything.organised"         "unspoilt.nature.natural.landscape"
[31] "cultural.offers"              "change.of.surroundings"
```

In R, using summary () command, we obtained summary of the dataset:

```
> summary(data1)
    Gender               Age           Education          Income
 Length:1000        Min.   : 18.00   Min.   :1.000    Length:1000
 Class :character   1st Qu.: 32.00   1st Qu.:3.000    Class :character
 Mode  :character   Median : 42.00   Median :6.000    Mode  :character
                    Mean   : 44.17   Mean   :4.814
                    3rd Qu.: 57.00   3rd Qu.:7.000
                    Max.   :105.00   Max.   :8.000
                                     NA's   :8
    Income2             Occupation           State          Relationship.Status
 Length:1000        Length:1000        Length:1000        Length:1000
 Class :character   Class :character   Class :character   Class :character
 Mode  :character   Mode  :character   Mode  :character   Mode  :character


   Obligation        Obligation2            NEP            Vacation.Behaviour
 Min.   :1.000    Length:1000        Min.   :1.733      Min.   :1.233
 1st Qu.:3.367    Class :character   1st Qu.:3.267      1st Qu.:2.467
 Median :3.800    Mode  :character   Median :3.667      Median :2.944
 Mean   :3.735                       Mean   :3.649      Mean   :2.963
 3rd Qu.:4.200                       3rd Qu.:4.067      3rd Qu.:3.429
 Max.   :5.000                       Max.   :5.000      Max.   :4.900
                                                        NA's   :25
 rest.and.relax     luxury...be.spoilt  do.sports          excitement..a.challenge
 Length:1000        Length:1000        Length:1000        Length:1000
 Class :character   Class :character   Class :character   Class :character
 Mode  :character   Mode  :character   Mode  :character   Mode  :character
```

In R, using unique () command, we obtained unique values for 4 features:

```
> unique(data1$Age)
 [1]  25  31  21  18  61  63  58  41  36  56  30  40  75  50  62  33  55  51  69  23  42
[22]  35  37  57  45  26  27  43  38  39  22  48  70  34  28  32  46  20  64  54  68  60
[43]  29  66  65  59  76  49  24  52  71  44  74  84  47  53  19  67  72  73 105
> unique(data1$Gender)
[1] "Female" "Male"
> unique(data1$Income)
 [1] "$30,001 to $60,000"   "$120,001 to $150,000" "$90,001 to $120,000"
 [4] "Less than $30,000"    "$60,001 to $90,000"   NA
 [7] "$180,001 to $210,000" "more than $240,001"   "$150,001 to $180,000"
[10] "$210,001 to $240,000"
> unique(data1$Income2)
[1] "30-60k"  ">120k"   "90-120k" "<30k"    "60-90k"  NA
```

**Data Cleaning**: Before we conduct data analysis, we have to go through the data cleaning step. This step involves to clean unnecessary data and to check the consistent labels for the levels of categorical variable have been used. From the summery of the dataset, we can say that there is no cleaning required for the variable Age and Gender. But categories of Income2 variable are not sorted properly. So, we have sorted them by the following approach.

```
[10]    $210,001 to $240,000
> unique(data1$Income2)
[1] "30-60k"  ">120k"    "90-120k" "<30k"     "60-90k"  NA
> inc2 <- data1$Income2
> unique(inc2)
[1] "30-60k"  ">120k"    "90-120k" "<30k"     "60-90k"  NA
> lev <- unique(inc2)
> lev
[1] "30-60k"  ">120k"    "90-120k" "<30k"     "60-90k"  NA
> lev[c(4,1,5,3,2)]
[1] "<30k"    "30-60k"   "60-90k"  "90-120k" ">120k"
> inc2 <- factor(inc2, labels = lev[c(4,1,5,3,2)])
> table(orig = data1$Income2, new = inc2)
            new
orig      <30k 30-60k 60-90k 90-120k >120k
  <30k     150      0      0       0     0
  >120k      0    140      0       0     0
  30-60k     0      0    265       0     0
  60-90k     0      0      0     233     0
  90-120k    0      0      0       0   146
> data1$Income2 <- inc2
> |
```
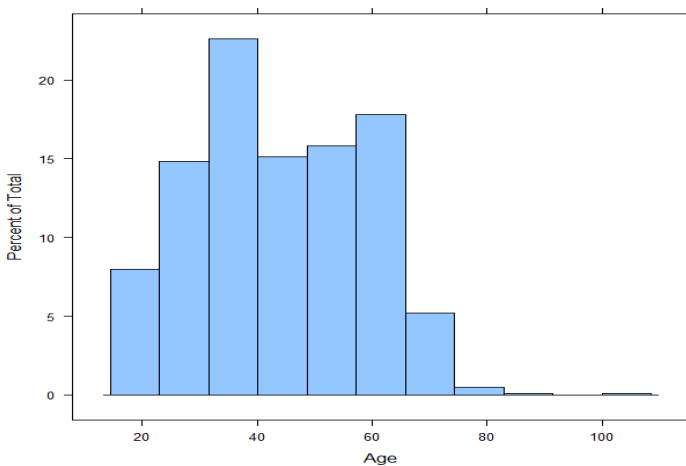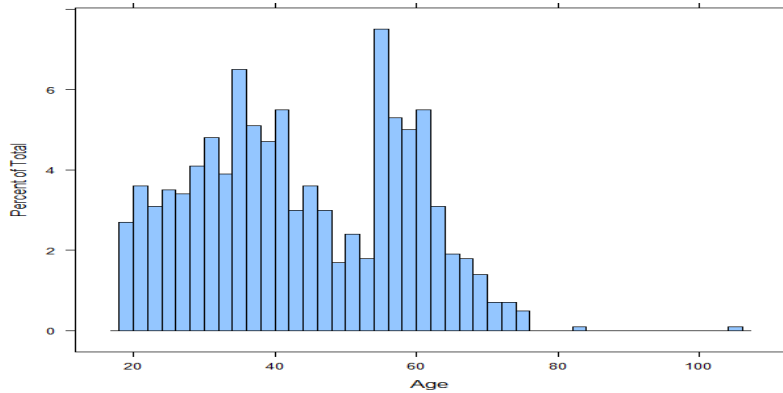
Descriptive Analysis:

Descriptive analysis helps to get meaning insights about the data through analysis. There are various plot which can help us to analyse data like histogram, box-plot, scatter plot etc. Histogram is a graphical representation of the distribution of the data. It helps us to visualize the distribution of the numeric variable. It also helps us to visualize the frequency of the observations within a certain range. We can check weather a distribution of a variable is symmetric or skewed using histogram. To create histogram, first we have to create bins (i.e, categories of values).

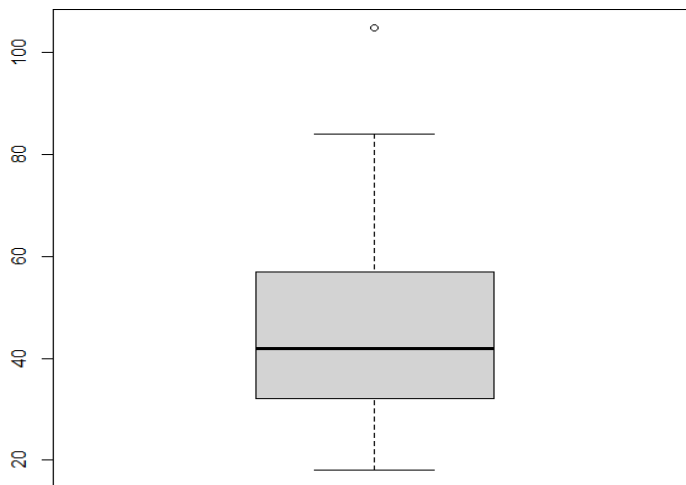In R, We use package lattice. We created histogram for Age variable:

We created 50 bins for the above histogram of the Age variable. From the above figure, we can say that the distribution is bimodal with many respondents aged around 35-40 and around 60 years.
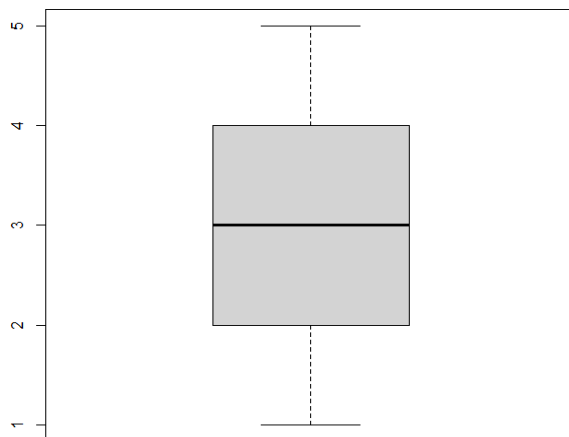
Another useful visualization tool is boxplot. This plot compresses a data into minimum, first quartile, median, third quartile and maximum. These five numbers often called five-point summary.

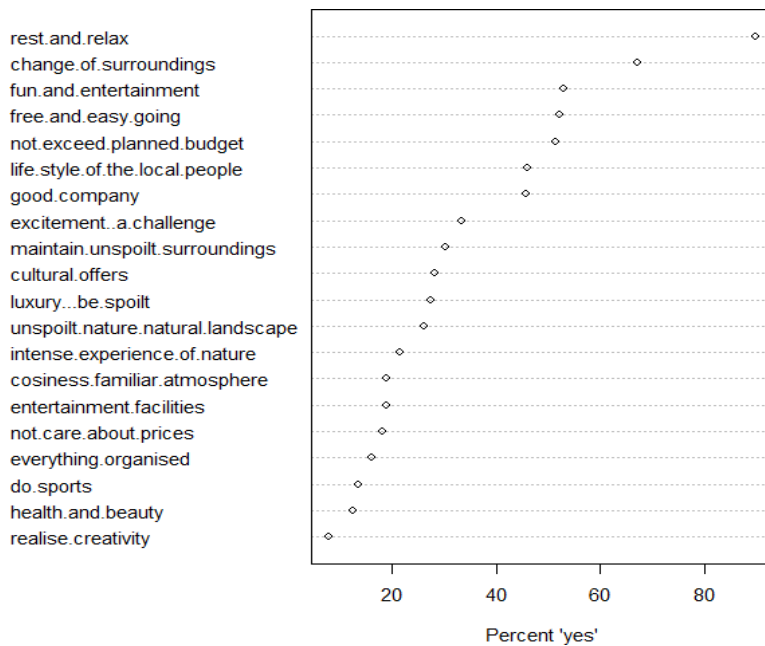In R, using boxplot () command we created boxplot for Age variable.



Interpretation: The 105-year-old respondent is clearly an outlier.

We created boxplot for Income2 variable.

Interpretation: There is no outlier for Income2 variable.

Columns 13 to 32 represents the travel motive of the respondent. We created dot chart of these travel motives with "yes". "Yes" means that the motive does apply.



Pre-Processing:

Categorical variables

Two pre-processing procedures are often used for categorical features. First one is merging the levels of categorical variables and second one is converting the categorical variables to numerical variables. Merging all the categories of Income variable, results the new variable Income2, which is much more balanced frequencies.

```
> # Pre-Processing
> #
> sort(table(data1$Income))

$210,001 to $240,000    more than $240,001 $180,001 to $210,000 $150,001 to $180,000
                  10                     11                   15                   32
$120,001 to $150,000  $90,001 to $120,000     Less than $30,000     $60,001 to $90,000
                  72                    146                  150                  233
  $30,001 to $60,000
                 265
> table(data1$Income2)

  <30k  30-60k  60-90k 90-120k    >120k
   150     140     265     233      146
```

Binary variables can always be converted numeric variables. Most statistical procedures work correctly if there are only two categories. We converted to a numeric matrix with 0 and 1 for No and YES.

```
   150     140     265     233      146
> vacmot <- (data1[, 13:32] == "yes") + 0
```

Numerical Variables:

The range of values of a segmentation variable affects its relative influence in distance-based methods of segment extraction. To balance the influence of segmentation variables on segmentation results, variables can be standardised. Is transform the variable in such a way that puts them on a common scale. The default standardisation method in statistics subtracts the empirical mean and divides by the empirical standard deviation.

In R, using scale () command, we can standardise the data.

```
> vacmot <- (data1[, 13:32] ==  yes ) + 0
> vacmot.scaled <- scale(vacmot)
> vacmot.pca <- prcomp(vacmot)
```

**Principal Components Analysis:**

Principal component analysis is a dimensionality reduction technique that is often used to reduce the dimensionality of the large dataset into the smaller one that still contains the most information about the large dataset. Since principal component analysis reduces the number of variables, it naturally comes at the expense of accuracy. For conducting PCA, first we need to standardize the range of the continuous variable. Mathematically, this can be done by,

(Value-mean)/standard deviation

Next, we must calculate covariance matrix to see if there is any relationship between the variables. Principal components are new variables that are constructed as linear combination of the initial variables. These combinations are done in such a way new variables (principal components) are uncorrelated and most of the information within the initial variables is compressed into the first components. So, the idea is 10-dimensional data gives us 10 principal components, but PCA puts most of the information into the component, then maximum remaining information into the second principal component and so on.

We applied PCA on the Australian travel motives dataset and get the following result.

```
> vacmot.pca <- prcomp(vacmot)
> vacmot.pca
Standard deviations (1, .., p=20):
 [1] 0.8123768 0.5735039 0.5285558 0.5094231 0.4695718 0.4549552 0.4314122 0.4197407
 [9] 0.4054634 0.3754545 0.3638569 0.3596946 0.3484495 0.3320907 0.3283236 0.3196988
[17] 0.3059506 0.2974572 0.2813886 0.2434328

Rotation (n x k) = (20 x 20):
                                       PC1          PC2          PC3          PC4
rest.and.relax                  -0.06277267  0.011968683  0.134476275 -0.07714032
luxury...be.spoilt              -0.10920691  0.393220137 -0.116667935 -0.38147871
do.sports                       -0.09455162  0.145580945 -0.045647512 -0.09360000
excitement..a.challenge         -0.27740140  0.222698527 -0.210311857  0.13855126
not.exceed.planned.budget       -0.28553655 -0.156119478  0.583064849  0.19164919
realise.creativity              -0.10951057 -0.012162975 -0.015277762 -0.02537091
fun.and.entertainment           -0.27924352  0.520509339  0.086515129  0.20396019
good.company                    -0.28435966 -0.009653647  0.129145586  0.30978323
health.and.beauty               -0.13972865  0.050941101  0.003911895 -0.18411602
free.and.easy.going             -0.31720483  0.057544926  0.244457844 -0.18710631
entertainment.facilities        -0.11769062  0.320726761  0.005007925 -0.05438875
not.care.about.prices           -0.04947810  0.239679194 -0.298809744 -0.16315215
life.style.of.the.local.people  -0.35272940 -0.267243059 -0.398233311  0.28716052
intense.experience.of.nature    -0.24129906 -0.213250514 -0.076280334 -0.19484703
cosiness.familiar.atmosphere    -0.13244960 -0.013284092  0.201726959 -0.19766207
maintain.unspoilt.surroundings  -0.30651478 -0.336128427  0.005248829 -0.34938475
everything.organised            -0.09160870  0.164922912  0.078046269 -0.12922382
unspoilt.nature.natural.landscape -0.26852625 -0.183104682 -0.055577530 -0.38654783
cultural.offers                 -0.26006232 -0.115982891 -0.428186679  0.18241933
change.of.surroundings          -0.25917233  0.091930993  0.104343189  0.25725834
```

The column PC1 indicates how the first principal component is composed of the initial variables. This shows that the first principal component separates the two answer tendencies "almost no motives apply" and "all motives apply". For the second principal component, the variables loading highest are FUN and ENTERTAINMENT, LUXURY / BE SPOILT and to MAINTAIN AN UNSPOILT SURROUNDING. For the third principal component not exceeding the planned budget, cultural offers, and the life style of the local people are important variables. For the fourth principal component unspoilt nature/natural landscape, luxury, maintain an unspoilt surrounding and good company are important variables.

We can conduct further analysis on the fitted object with the summary function.

```
> print(summary(vacmot.pca), digits = 2)
Importance of components:
                        PC1  PC2   PC3   PC4  PC5   PC6   PC7   PC8   PC9  PC10  PC11
Standard deviation     0.81 0.57 0.529 0.509 0.47 0.455 0.431 0.420 0.405 0.375 0.364
Proportion of Variance 0.18 0.09 0.077 0.071 0.06 0.057 0.051 0.048 0.045 0.039 0.036
Cumulative Proportion  0.18 0.27 0.348 0.419 0.48 0.536 0.587 0.635 0.681 0.719 0.756
                        PC12  PC13 PC14 PC15  PC16  PC17  PC18  PC19  PC20
Standard deviation     0.360 0.348 0.33 0.33 0.320 0.306 0.297 0.281 0.243
Proportion of Variance 0.035 0.033 0.03 0.03 0.028 0.026 0.024 0.022 0.016
Cumulative Proportion  0.791 0.824 0.85 0.88 0.912 0.938 0.962 0.984 1.000
> library("flexclust")
```
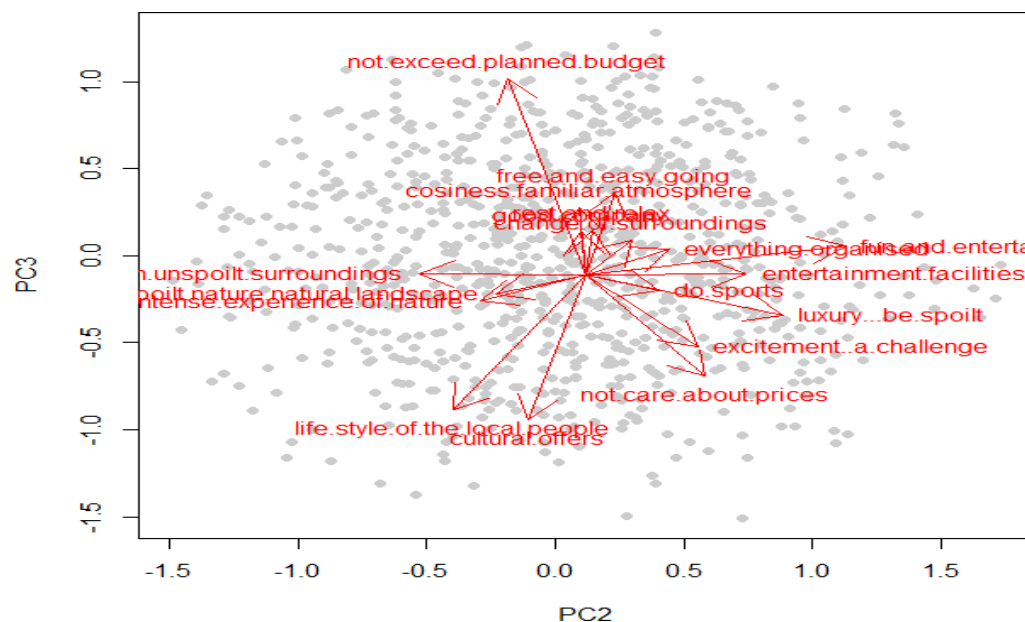
Interpretation: Principal component 1 explains 18% of the variation in the original data. Principal component 2 explains 9% of the variation in the original data.

Now, we want to plot the data into two-dimensional space. Inspecting the rotation matrix reveals that the first principal component does not differentiate well between motives because all motives load on it negatively. So, we consider principal component 2 and principal component 3 to create perceptual map. Function projAxes plots how the principal components are composed of the original variables and visualises the rotation matrix.

```
> library("flexclust")
> plot(predict(vacmot.pca)[, 2:3], pch = 16,
+       col = "grey80")
> projAxes(vacmot.pca, which = 2:3)
```
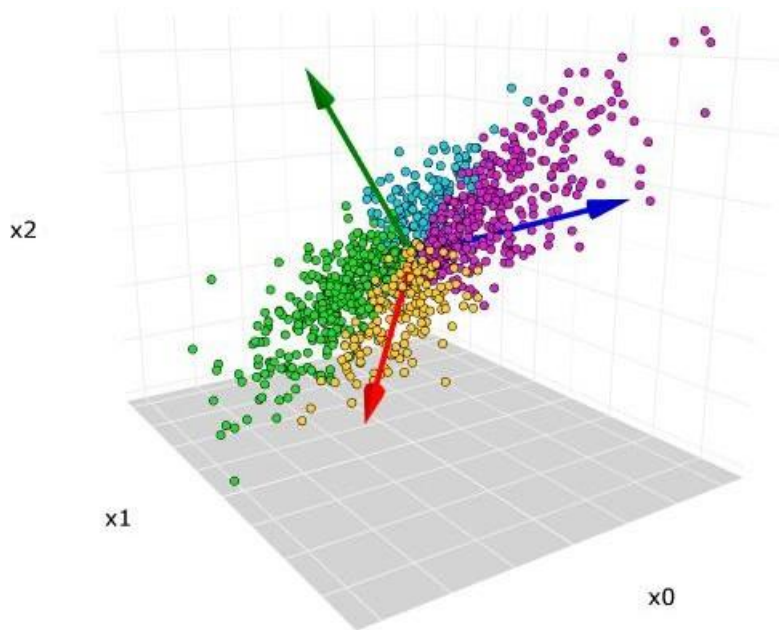


Interpretation: NOT EXCEEDING THE PLANNED BUDGET (represented by the arrow pointing in the top slightly left direction) is a travel motive that is quite unique. On the other hand,

LIFESTYLE OF LOCAL PEOPLE, and interest in CULTURAL OFFERS available at destinations often occur simultaneously (as indicated by the two arrows both pointing to the left bottom).

We can use principal component analysis to reduce the number of segmentation variable before extracting the market segment from the consumer data. It is helpful because after reducing the number of variables it provides meaningful insights which is easy handle.

# Step 5: Extracting Segments

Step 5 is where we extract segments. To illustrate a range of extraction techniques, we subdivide this step into sections. In the first section, we will use standard k-means analysis.

## 5.1 Grouping Consumers:

For extracting market segments from data mostly clustering analysis is used with a variety of approaches. Selecting a suitable clustering method requires matching the data analytic featuresof the resulting clustering with the context-dependent requirements that are desired by the researcher of the market.

Some of Data set and segment characteristics informing extraction algorithm selection given:

Data set characteristics: – Size (number of consumers, number of segmentation variables) – Scale level of segmentation variables (nominal, ordinal, metric, mixed) – Special structure, additional information.

Segment characteristics: – Similarities of consumers in the same segment – Differences between consumers from different segments – Number and size of segments
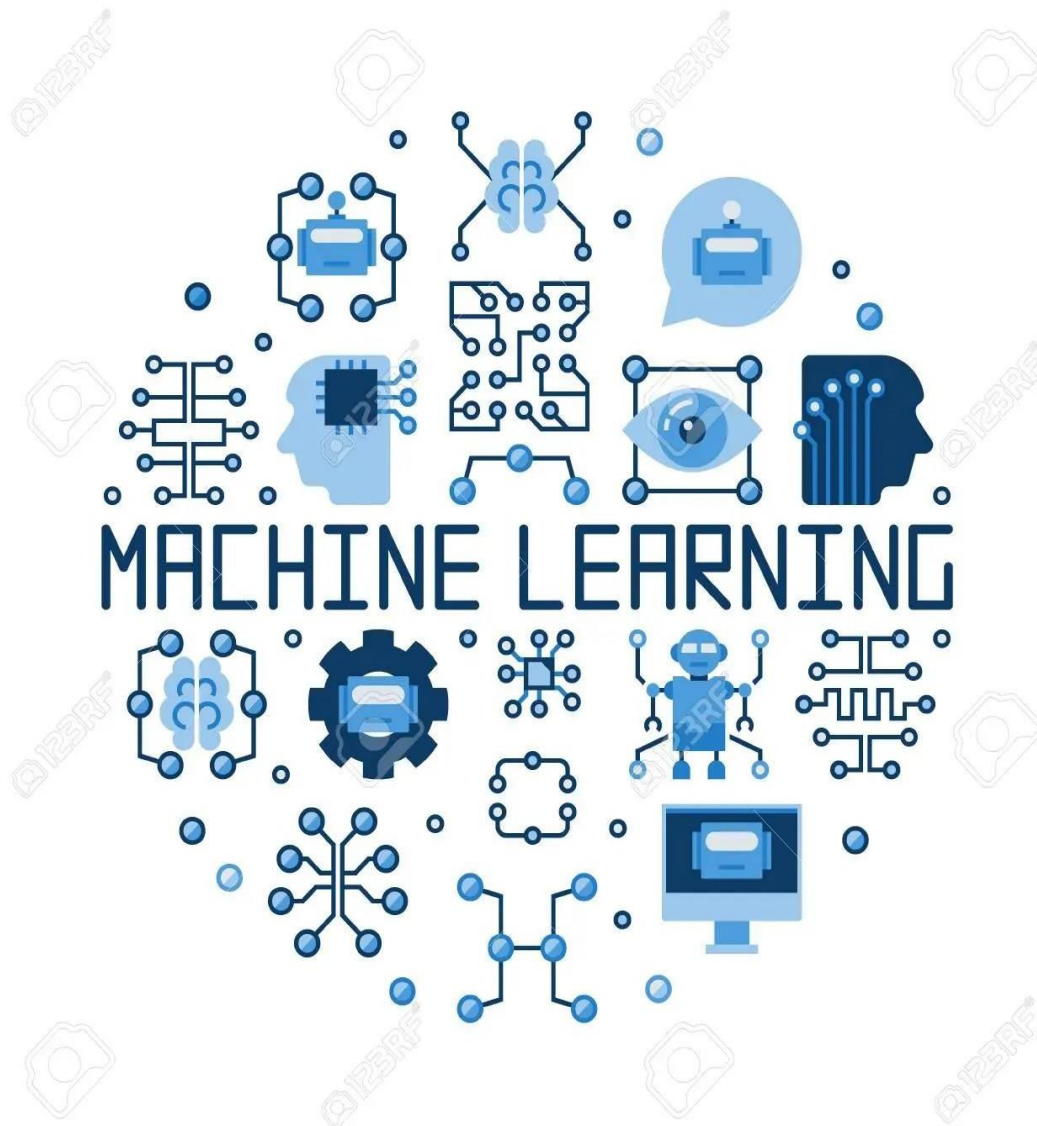
None of these methods outperform other methods in all situations. Rather, each method has pros and cons.

Distance-based methods use a particular notion of similarity or distance between observations (consumers), and try to find groups of similar observations (market segments). For distance- based methods, the choice of the distance measure depends on the scale level of the data.



Model-based methods formulate a concise stochastic model for the market segments where mainly distributions and probabilistic approach is used. Model based approaches are more probability based considering parameters of segment size and characteristics consumer's probability of getting fit into segment is derived and best solution is provided in end. If the data set contains repeated measurements

of consumers over time, for example, an algorithm that takes this longitudinal natureof the data into account is needed. Such data generally requires a model-

based approach. If the data contains purchase histories and price information, and market segments are based on similar price sensitivity levels, regression models are needed. This, in turn, calls for theuse of a model-based segment extraction algorithm.

In the case of binary segmentation variables, another aspect needs to be considered. We may want consumers in the same segments to have both the presence and absence of segmentation variables in common, here these variables would be symmetrical (with 0s and 1s treated equally).



Alternatively, we may be concerned about segmentation variables consumers have in common, herethese variables would be symmetrical (with only common 1s being of interest). Biclustering uses binary information asymmetrically. Distance-based methods can use distance measures that accountfor this asymmetry, and extract segments characterized by common 1s.

Data-driven market segmentation analysis is exploratory by nature. Consumer data sets are typically not well structured. Consumers come in all shapes and forms a two-dimensional plot of consumers' product preferences typically does not contain clear groups of consumers. Rather, consumer preferences are spread across the entire plot. The combination of exploratory methods andunstructured consumer data that results from any method used to extract market segments from suchdata will strongly depend on the assumptions made on the structure of the segments implied by the method. The result of a market segmentation analysis, therefore, is determined as much by the underlying data as it is by the extraction algorithm chosen.

There are few types of **Distance-based Methods:**

## 5.2 Distance Methods:

### 1) Euclidean distance

Euclidean Distance represents the shortest distance between two points. Most machine learning algorithms including K-Means use this distance metric to measure the similarity between observations

$$d = ((p_1 - q_1)^2 + (p_2 - q_2)^2)^{1/2}$$
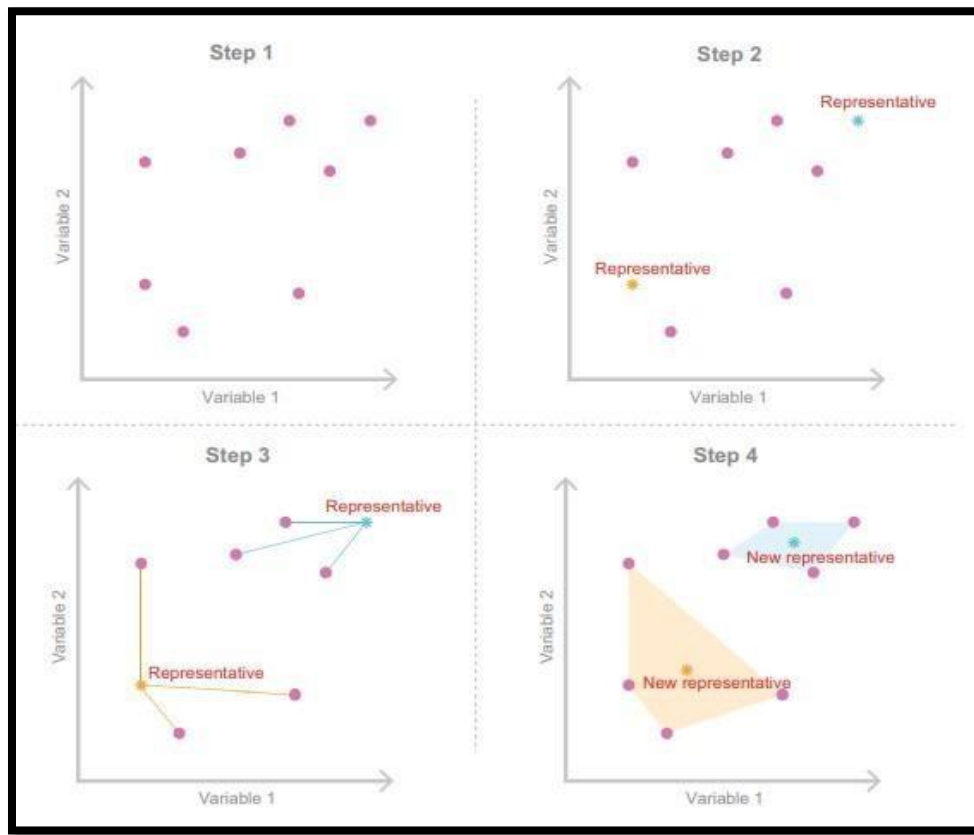
### 2) Hierarchical measures

Hierarchical clustering methods are the most intuitive way of grouping data because they mimic how a human would approach the task of dividing a set of n observations(consumers) into k groups (segments). If the aim is to have one large market segment (k = 1), the only possible solution is one big market segment containing all consumers in data X.

$$l(X, Y) = \frac{1}{|X||Y|} \sum_{x \in X} \sum_{y \in Y} d(x, y),$$

### 3) Partitioning method

Hierarchical clustering methods are particularly well suited for the analysis of small data sets with up

to a few hundred observations. For larger data sets, dendrograms are hard to read, and the matrix of pairwise distances usually does not fit into computer memory. For data sets containing more than 1000 observations (consumers), clustering methods creating a single partition are more suitable thana nested sequence of partitions. This means that – instead of computing all distances between all pairsof observations in the data set.
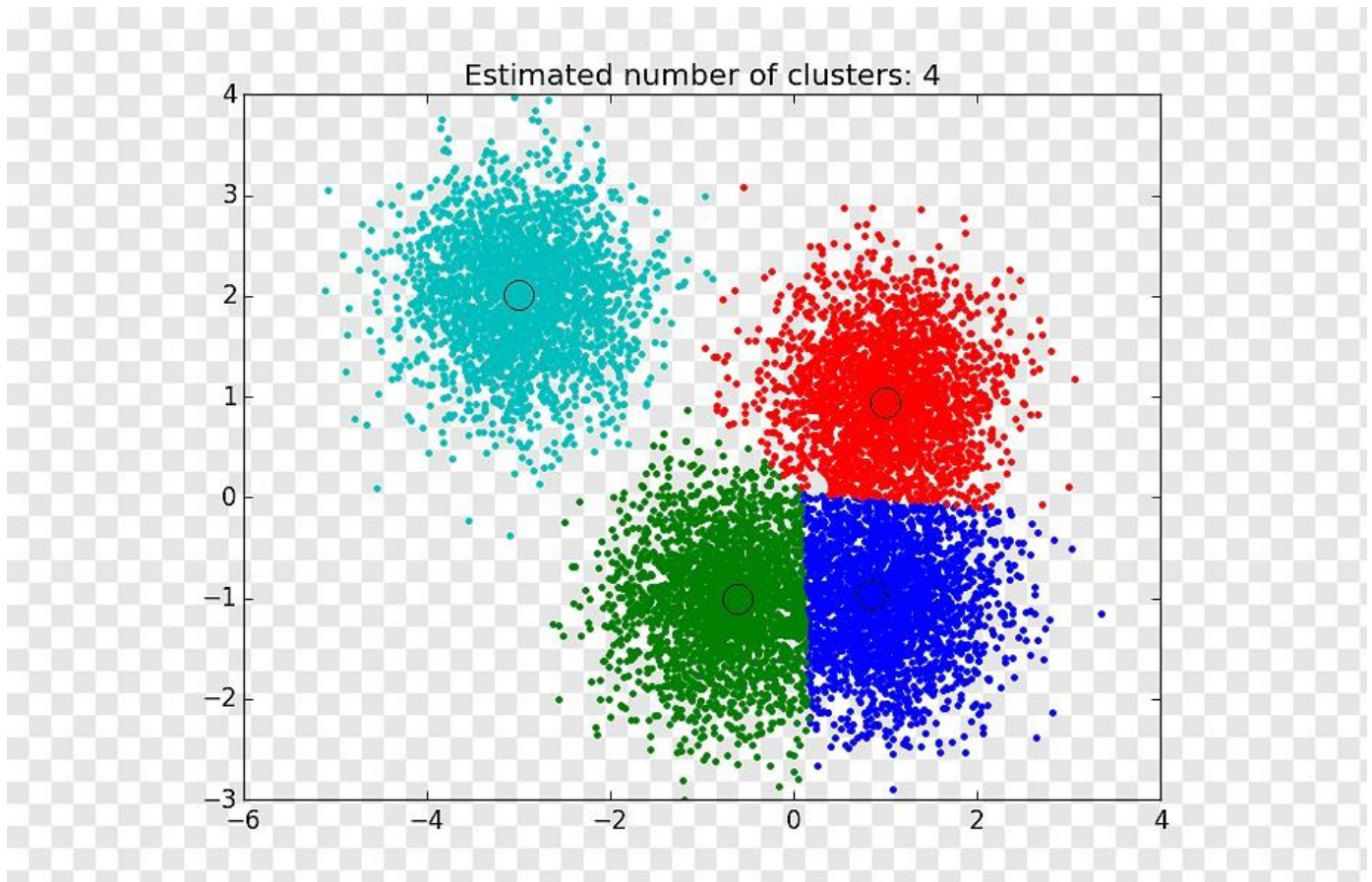


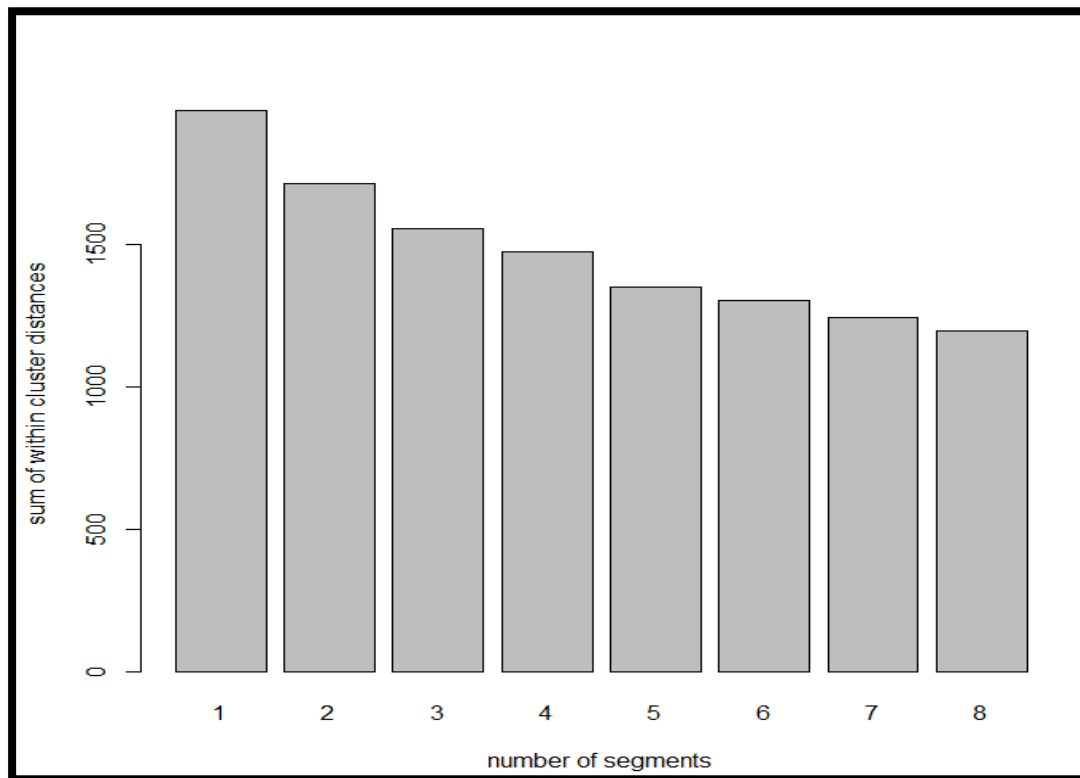**Simplified visualisation of the k-means clustering algorithm**

## 5.3 Using k-Means

We calculate solutions for two to eight market segments using standard k-means analysis with ten

random restarts. We then relabel segment numbers such that they are consistent across segmentations.

We extract between two and eight segments because we do not know in advance what the best numberof market segments is. If we calculate a range of solutions, we can compare them and choose the onewhich extracts segments containing similar consumers which are distinctly different from members of other segments. We compare different solutions using a scree plot:
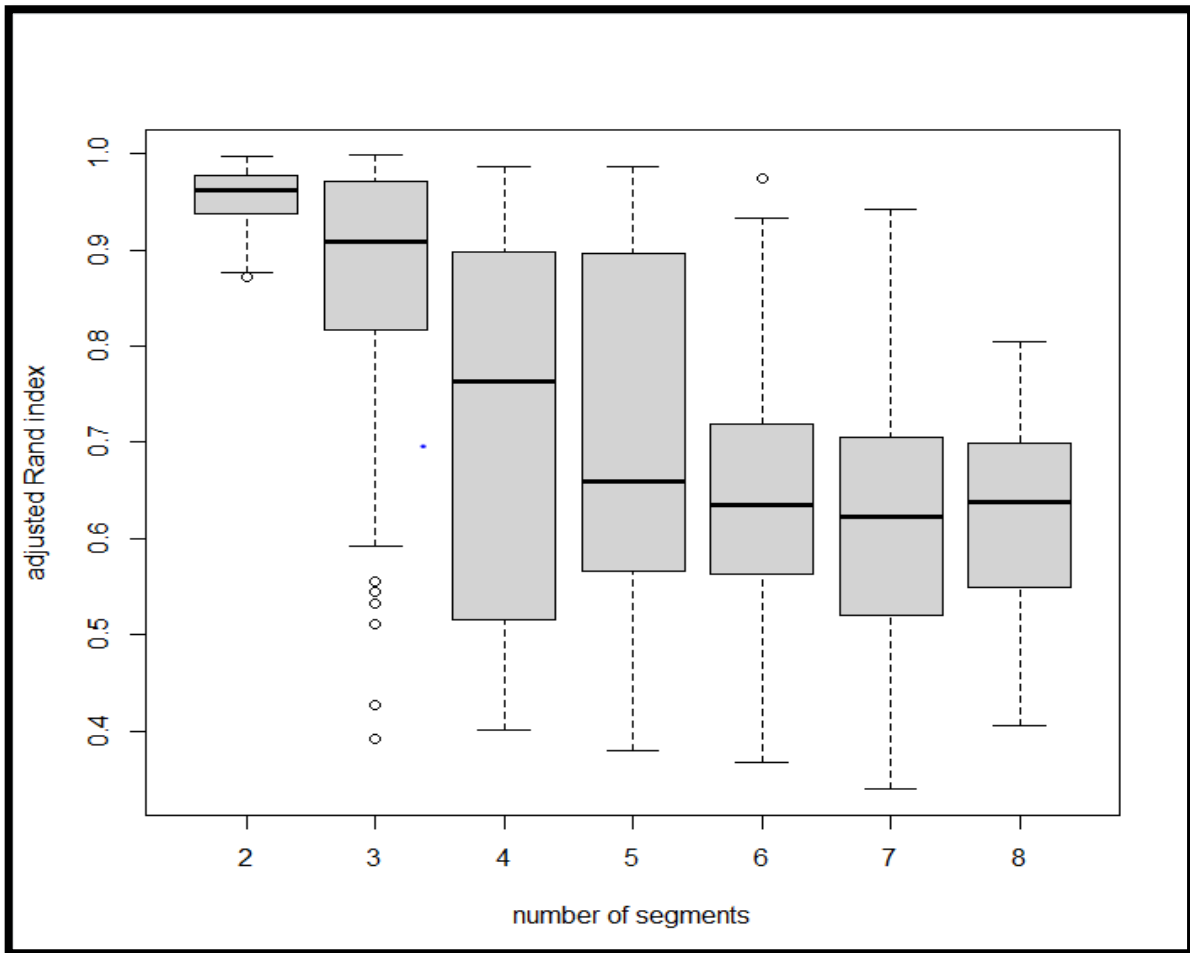
**Scree plot for the McDonald's data set**

The scree plot has no distinct elbow: the sum of distances within market segments drops slowly as the number of market segments increases. We expect the values to decrease because more market segments automatically mean that the segments are smaller and, as a consequence, that segment members are more similar to one another. But the much-anticipated point where the sum of distancesdrops dramatically is not visible. This scree plot does not provide useful guidance on the number of market segments to extract.

A second approach to determining a good number of segments is to use stability-based data structureanalysis. Stability-based data structure analysis also indicates whether market segments occurnaturally in the data, or if they have to be artificially constructed. Stability-based data structure analysis uses stability across replications as criterion to offer this guidance. Imagine using a market segmentation solution which cannot be reproduced. Such a solution would give McDonald's management little confidence in terms of investing substantial resources into a market segmentationstrategy.

**Global stability of k-means segmentation solutions for the McDonald's data set**

The vertical boxplots show the distribution of stability for each number of segments. The median is indicated by the fat black horizontal line in the middle of the box. Higher stability is better. Inspecting points to the two-, three- and four-segment solutions as being quite stable. However, the two- and three-segment solutions do not offer a very differentiated view of the market.

We gain further insights into the structure of the four-segment solution with a gorge plot:

Spending Score (1-100) vs Annual Income (k$)

```mermaid
graph TD
    ML[Machine Learning]
    ML --> SL[Supervised Learning]
    ML --> RL[Reinforcement Learning]
    ML --> UL[Unsupervised Learning]
    SL --> Classification
    SL --> Regression
    UL --> Association
    UL --> Clustering
    Clustering --> Hierarchical
    Clustering --> Overlapping
    Clustering --> Exclusive
    Exclusive --> KMeans[K-Means]
```

- Machine Learning
  - Supervised Learning
    - Classification
    - Regression
  - Reinforcement Learning
  - Unsupervised Learning
    - Association
    - Clustering
      - Hierarchical
      - Overlapping
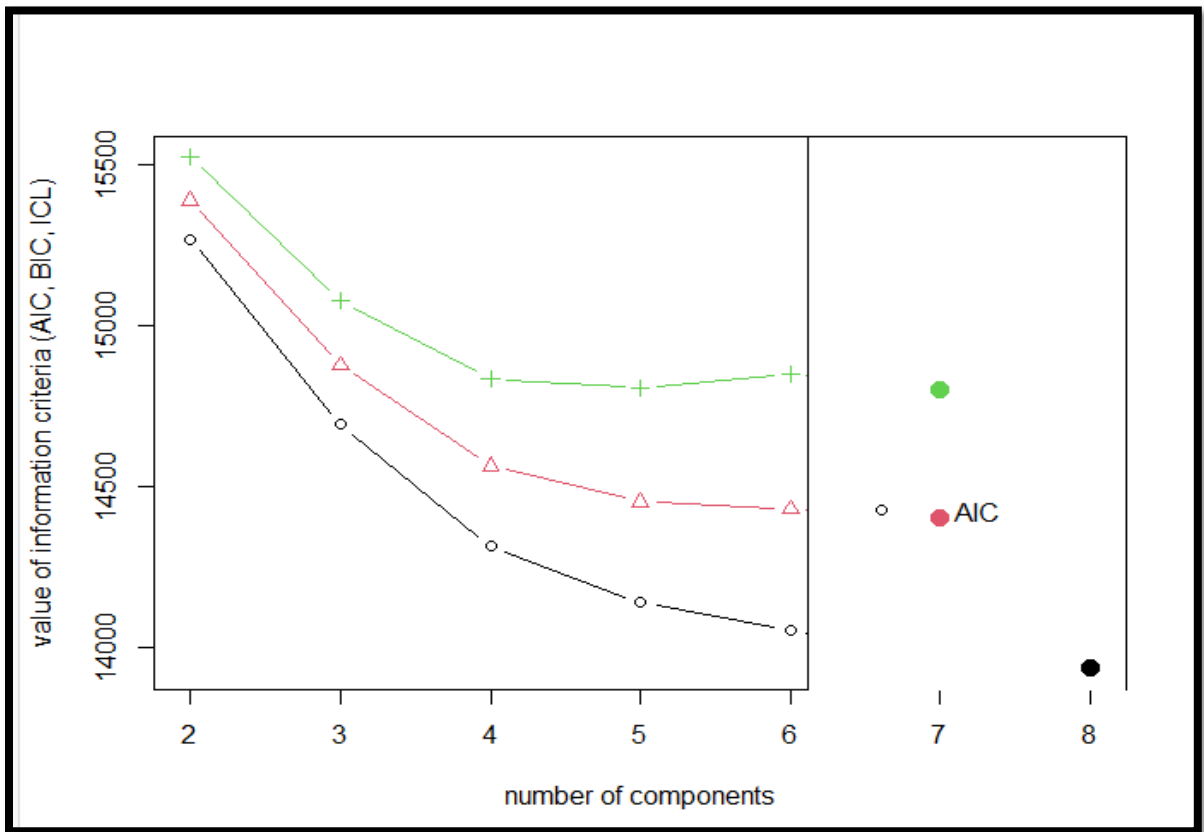      - Exclusive
        - K-Means

## 5.4 Using Mixtures of Distributions:

We calculate latent class analysis using a finite mixture of binary distributions. The mixture model maximizes the likelihood to extract segments (as opposed to minimizing squared Euclidean distance, as is the case for k-means).



**Information criteria for the mixture models of binary distributions with 2 to 8 components (segments) for the McDonald's data set**

$$
\begin{aligned}
K &= \text{number of mixture components} \\
N &= \text{number of observations} \\
\theta_{i=1\ldots K} &= \text{parameter of distribution of observation associated with component } i \\
\phi_{i=1\ldots K} &= \text{mixture weight, i.e., prior probability of a particular component } i \\
\phi &= K\text{-dimensional vector composed of all the individual } \phi_{1\ldots K}; \text{ must sum to 1} \\
z_{i=1\ldots N} &= \text{component of observation } i \\
x_{i=1\ldots N} &= \text{observation } i \\
F(x|\theta) &= \text{probability distribution of an observation, parametrized on } \theta \\
z_{i=1\ldots N} &\sim \text{Categorical}(\phi) \\
x_{i=1\ldots N} | z_{i=1\ldots N} &\sim F(\theta_{z_i})
\end{aligned}
$$

# Step 6: Profiling Segments

## What is profiling segments?

The term 'profiling' refers to analysis of market segments. This step involves regarding data-driven market segmentation. In case of data-driven market segmentation, after extraction step, profiling is required. On the other hand, this step is not necessary for common-sense segmentation because in this case profiles are predefined. When data-driven market segmentation is conducted, profiling helps to interpret the solution of the market segmentation. Consumer may want to extract the solution of the data-driven market segmentation based on their requirements. In this situation, this step is required.

## Different approaches to profiling market segments

### 1. Traditional Approach

In this approach, data-driven market segmentation solution is represented by a large table. The table contains percentage of segmentation variables for each segment. As a result, it is very difficult to extract the key insights from the table quickly.

**Example**: We use the Australian vacation motives data set. Segments were extracted from this dataset using the neural gas clustering algorithm with number of segments varied from 3 to 8 and with 20 random restarts. We reload the segmentation solution.

```
88  data("vacmot", package = "flexclust")
89  set.seed(1234)
90  vacmot.k38 <- stepcclust(vacmot, k = 3:8,
91                           method = "neuralgas", nrep = 20, save.data = TRUE,
92                           verbose = FALSE)
93  vacmot.k38 <- relabel(vacmot.k38)
94  vacmot.k6 <- vacmot.k38[["6"]]
95  save(vacmot.k38, vacmot.k6,
96       + file = "vacmot-clusters.RData")
97  vacmot.r6 <- slswFlexclust(vacmot, vacmot.k6,
98                           method = "neuralgas", FUN = "cclust")
99
00  data("vacmot", package = "flexclust")
01  load("vacmot-clusters.RData")
```

The below table shows the mean values of the segmentation variables by segment (extracted from the return object using parameters (vacmot.k6)), together with the overall mean values.

Because the travel motives are binary, the segment means are equal to the percentage of segment members engaging in each activity.

|  | Seg.1 | Seg.2 | Seg.3 | Seg.4 | Seg.5 | Seg.6 | Total |
|---|---|---|---|---|---|---|---|
| Rest and relax | 83 | 96 | 89 | 82 | 98 | 96 | 90 |
| Change of surroundings | 27 | 82 | 73 | 82 | 87 | 77 | 67 |
| Fun and entertainment | 7 | 71 | 81 | 60 | 95 | 37 | 53 |
| Free-and-easy-going | 12 | 65 | 58 | 45 | 87 | 75 | 52 |
| Not exceed planned budget | 23 | 100 | 2 | 49 | 84 | 73 | 51 |
| Life style of the local people | 9 | 29 | 30 | 90 | 75 | 80 | 46 |
| Good company | 14 | 59 | 40 | 58 | 77 | 55 | 46 |
| Excitement, a challenge | 9 | 17 | 39 | 57 | 76 | 36 | 33 |
| Maintain unspoilt surroundings | 9 | 10 | 16 | 7 | 67 | 95 | 30 |
| Cultural offers | 4 | 2 | 5 | 96 | 62 | 38 | 28 |
| Luxury / be spoilt | 19 | 24 | 39 | 13 | 89 | 6 | 28 |
| Unspoilt nature/natural landscape | 10 | 10 | 13 | 15 | 69 | 64 | 26 |
| Intense experience of nature | 6 | 8 | 9 | 21 | 50 | 58 | 22 |
| Cosiness/familiar atmosphere | 11 | 24 | 12 | 7 | 49 | 25 | 19 |
| Entertainment facilities | 5 | 25 | 30 | 14 | 53 | 6 | 19 |
| Not care about prices | 8 | 7 | 43 | 19 | 29 | 10 | 18 |
| Everything organised | 7 | 21 | 15 | 12 | 46 | 9 | 16 |
| Do sports | 8 | 12 | 13 | 10 | 46 | 7 | 14 |
| Health and beauty | 5 | 8 | 10 | 8 | 49 | 16 | 12 |
| Realise creativity | 2 | 2 | 3 | 8 | 29 | 14 | 8 |

Interpretation: Basic interpretation from the above table can be done like basic characteristic of segment 2 are being motivated by rest and relax, not wanting to exceed planned budget and care about changing of surrounding.

The above table contains the mean value of segmentation variable by segment together with the overall mean values. It is very difficult to find the defining characteristics of the market segments. If we want to compare between the segments by segmentation variable, for each row we need to compare the pair of numbers 15 times. Therefore in total we have to compare the pair of numbers 300 times. So it is a tedious task to understand the defining characteristic of the market segment.

## 2. Graphical statistics approaches

Graphics is an important part in exploratory data analysis because they provide complex insights between the variables. To inspect one or more segmentation solution, visualization tools can help. The process of segmenting data always leads to many alternative solutions. Selecting one of the possible decisions is a critical decision. In this situation, statistical graphs can help us. To order segmentation variables by similarity, we can use clustering method.

A good way to visualize and understanding a segment is to produce a segment profile plot. Segment profile plot is also called panel plot. Each of the six panel (Fig. 1) represents one segment. For each segment, segment profile plot shows the cluster centres. The dots in Fig. 1 are identical in each of the six panels and represents the total mean values for the segmentation variables across all observations in the data set

```
99
100  data("vacmot", package = "flexclust")
101  load("vacmot-clusters.RData")
102  vacmot.vdist <- dist(t(vacmot))
103  vacmot.vclust <- hclust(vacmot.vdist, "ward.D2")
104  barchart(vacmot.k6, shade = TRUE,
105          which = rev(vacmot.vclust$order))
106
```
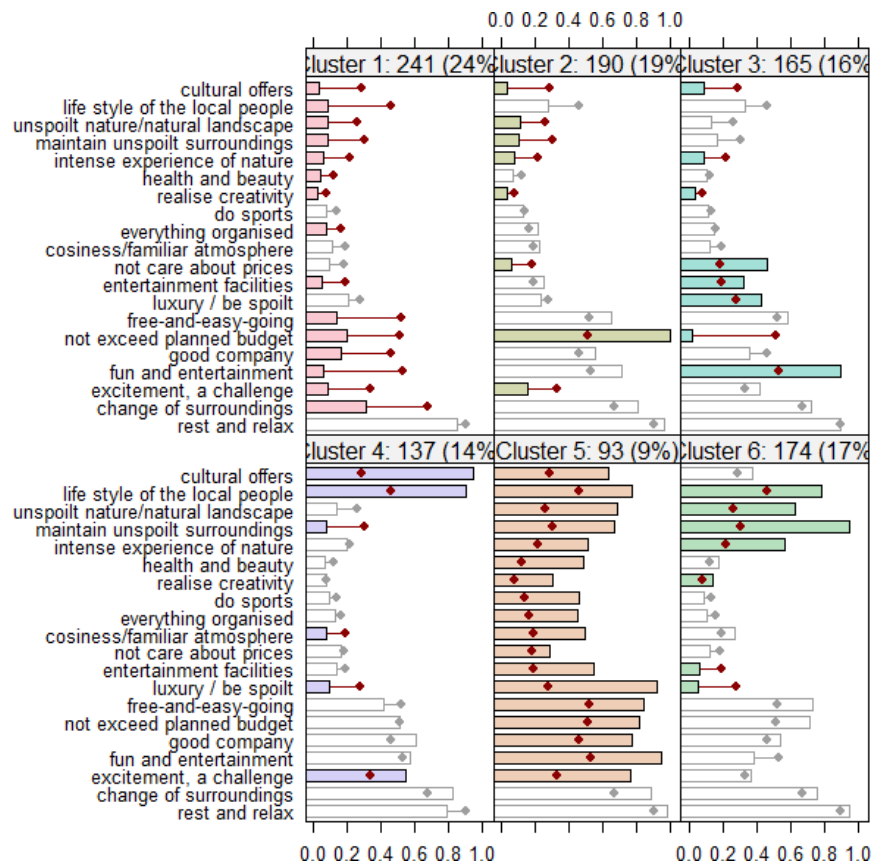
Fig.1 Segment profile plot for the six-segment solution of the Australian travel motives dataset

Interpretation: To interpret the chart, marker variables appear in colour. Marker variables are defined as variables which deviate by more than 0.25 from the overall mean. For example, a variable with a total sample mean of 0.20, and a segment mean of 0.60 qualifies as marker variable (0.20 + 0.25 = 0.45 < 0.60). Such large absolute difference is hard to obtain. A relative difference of 50% from the total mean, therefore, also makes the variable a marker variable. For travel motive Health and Beauty, this segmentation variable has a sample mean of 0.12. This means that only 12% of all the people who participated in the survey indicated that HEALTH AND BEAUTY were a travel motive for them. For segments with HEALTH AND BEAUTY outside of the interval
(0.12 - 0.06, 0.12+0.06) this vacation activity will be considered a marker variable, because 0.06 is 50% of 0.12.

Segment separation can be visualized by segment separation plot. Segment separation plot are very useful if the number of variables is low. But in case of high-dimensional dataset we need to reduce the dimension of the dataset before using segment separation plot.
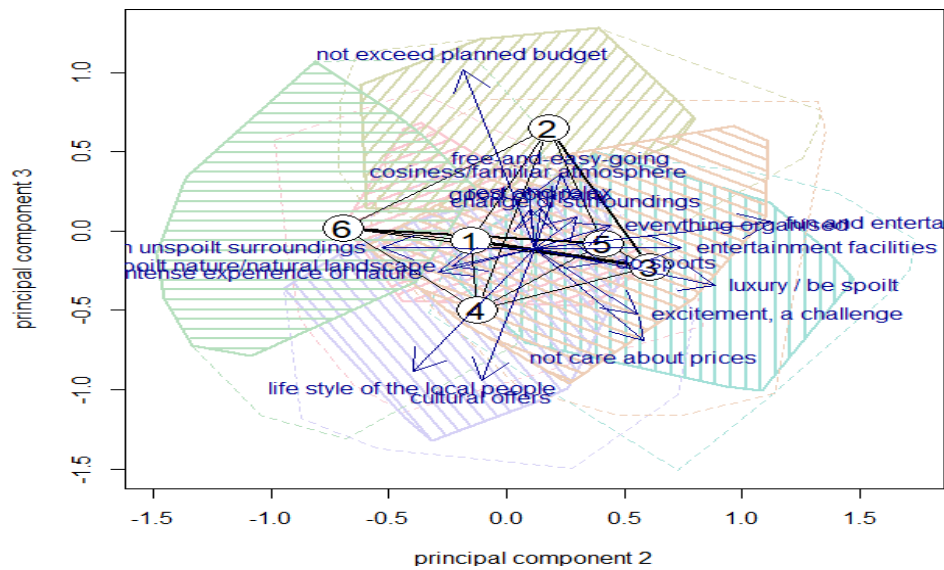
Fig.3 Segment Separation plot for Australian travel motives dataset

Interpretation: We can be interpreting the segment separation plot as follows the segment 6 that cares about maintaining unspoilt surroundings, spoilt nature and want to intense experience to nature.

# Step 7: Describing Segments



## 9.1 Developing a Complete Picture Market Segments:

In this section, the goal is to gain a comprehensive understanding of market segments. This involves gathering data and conducting analysis to identify key characteristics and attributes that define each segment.

## 9.2 Using Visualizations to Describe Market Segments:

Visualizations are powerful tools for conveying information about market segments. This subsection explores the use of visual representations, such as charts and graphs, to effectively communicate the characteristics and differences between various segments.

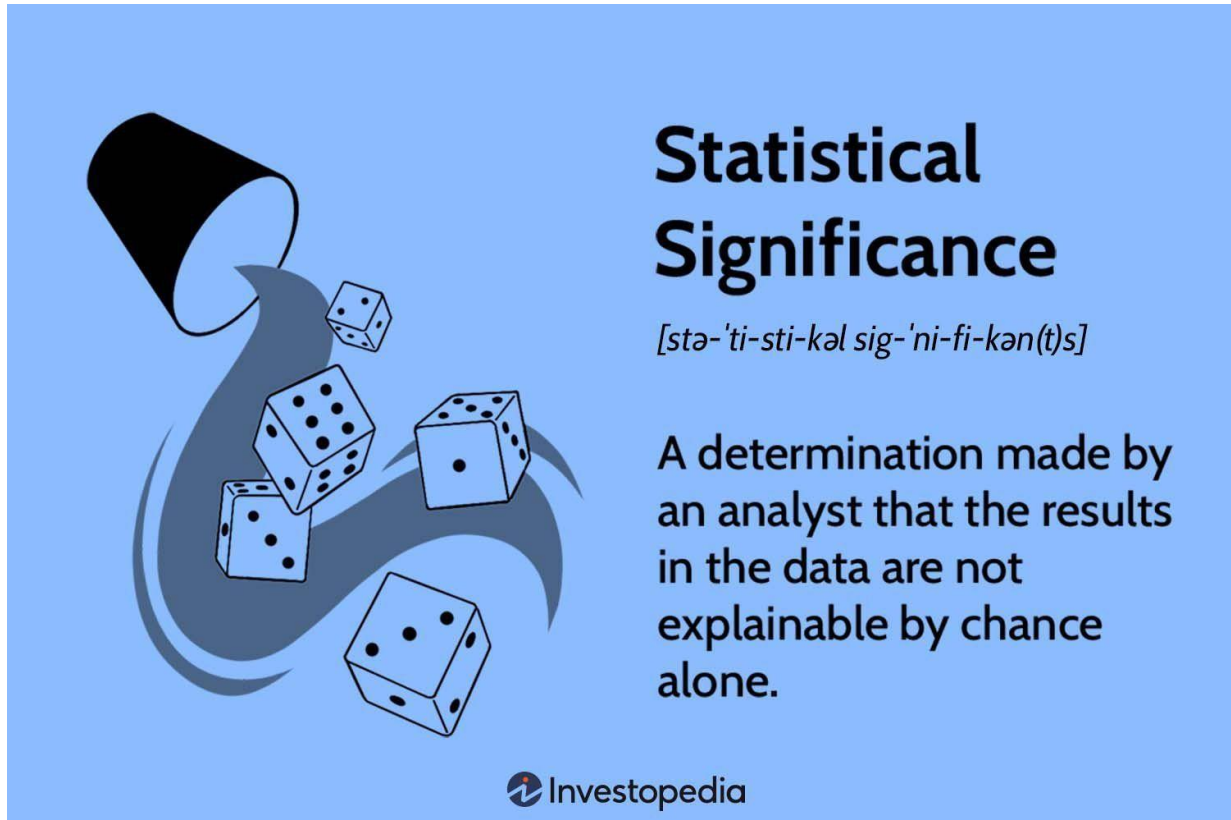## 9.2.1 Nominal and Ordinal Descriptor Variables:

This subsection discusses the use of nominal and ordinal descriptor variables in visualizations. Nominal variables represent categories without a specific order, while ordinal variables have a predetermined order. Visualizations help in presenting these variables in a clear and understandable manner.

## 9.2.2 Metric Descriptor Variables:

Metric descriptor variables are discussed in this section. These variables are quantitative and can be measured on a continuous scale. Visualizations are utilized to showcase relationships, trends, and differences among market segments based on these metric variables.

## 9.3 Testing for Segment Differences in Descriptor Variables:

To determine if there are significant differences between market segments, statistical tests are performed on descriptor variables. This subsection explains various statistical methods that can be employed to test for differences and establish the uniqueness of each segment.

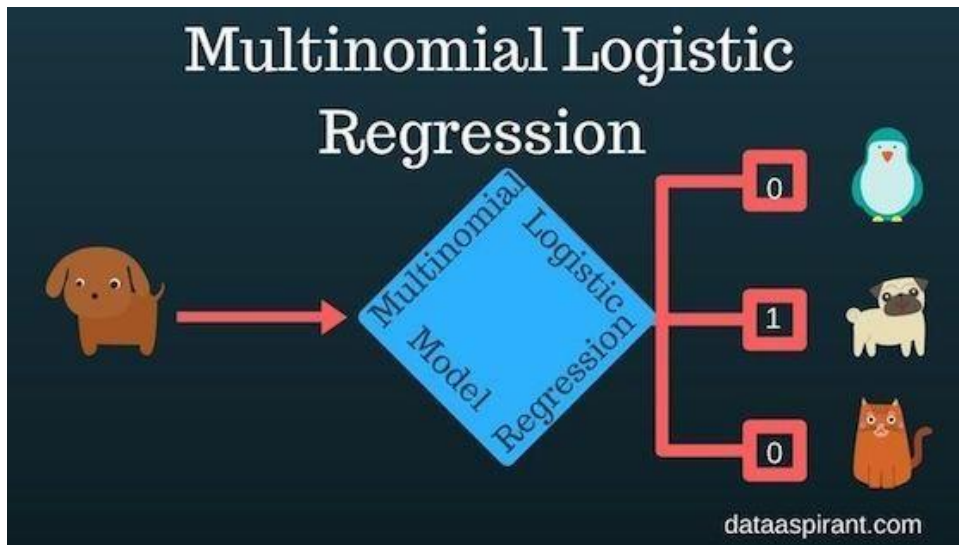## 9.4 Predicting Segments from Descriptor Variables:

This section focuses on predicting market segments using descriptor variables. It introduces binary logistic regression, multinomial logistic regression, and tree-based methods as techniques to predict the likelihood of an individual belonging to a particular segment based on the descriptor variables.

## 9.4.1 Binary Logistic Regression:

Binary logistic regression is a statistical technique used to predict the probability of an event occurring. In the context of market segmentation, it can be used to predict which segment an individual is likely to belong to based on specific descriptor variables.

# 9.4.2 Multinomial Logistic Regression:

Multinomial logistic regression extends binary logistic regression to predict outcomes with more than two categories. This subsection explains how multinomial logistic regression can be used to predict market segments using multiple descriptor variables.

# 9.4.3 Tree-Based Methods:

Tree-based methods, such as decision trees and random forests, are introduced as alternative approaches to predict market segments. These methods create decision rules based on the descriptor variables to classify individuals into different segments



Decision Tree



Random Forest

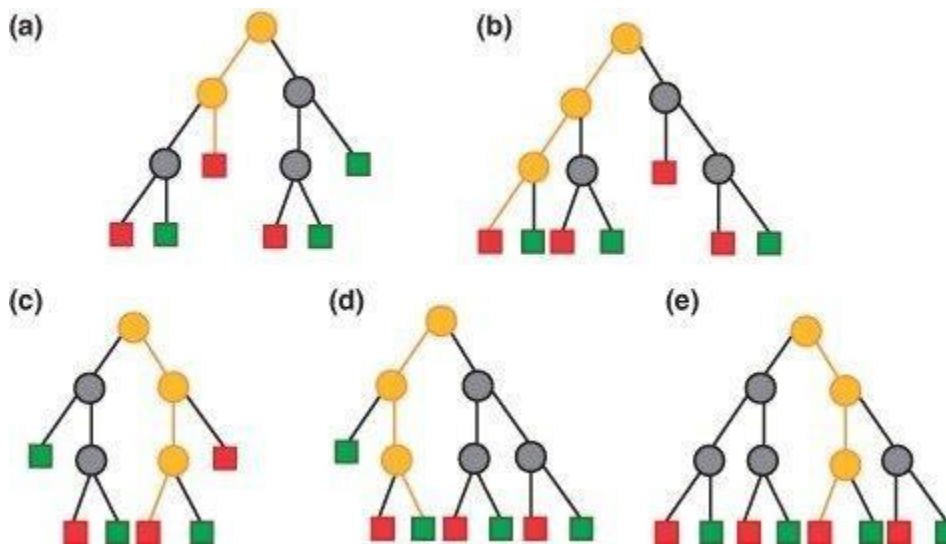# Step 8: Selecting (the) Target Segment(s)

Step 8 of the topic is centered around selecting the target segment(s) for a marketing strategy. This step involves making informed decisions about which market segment(s) to prioritize and focus on. The summary of each subsection is as follows:

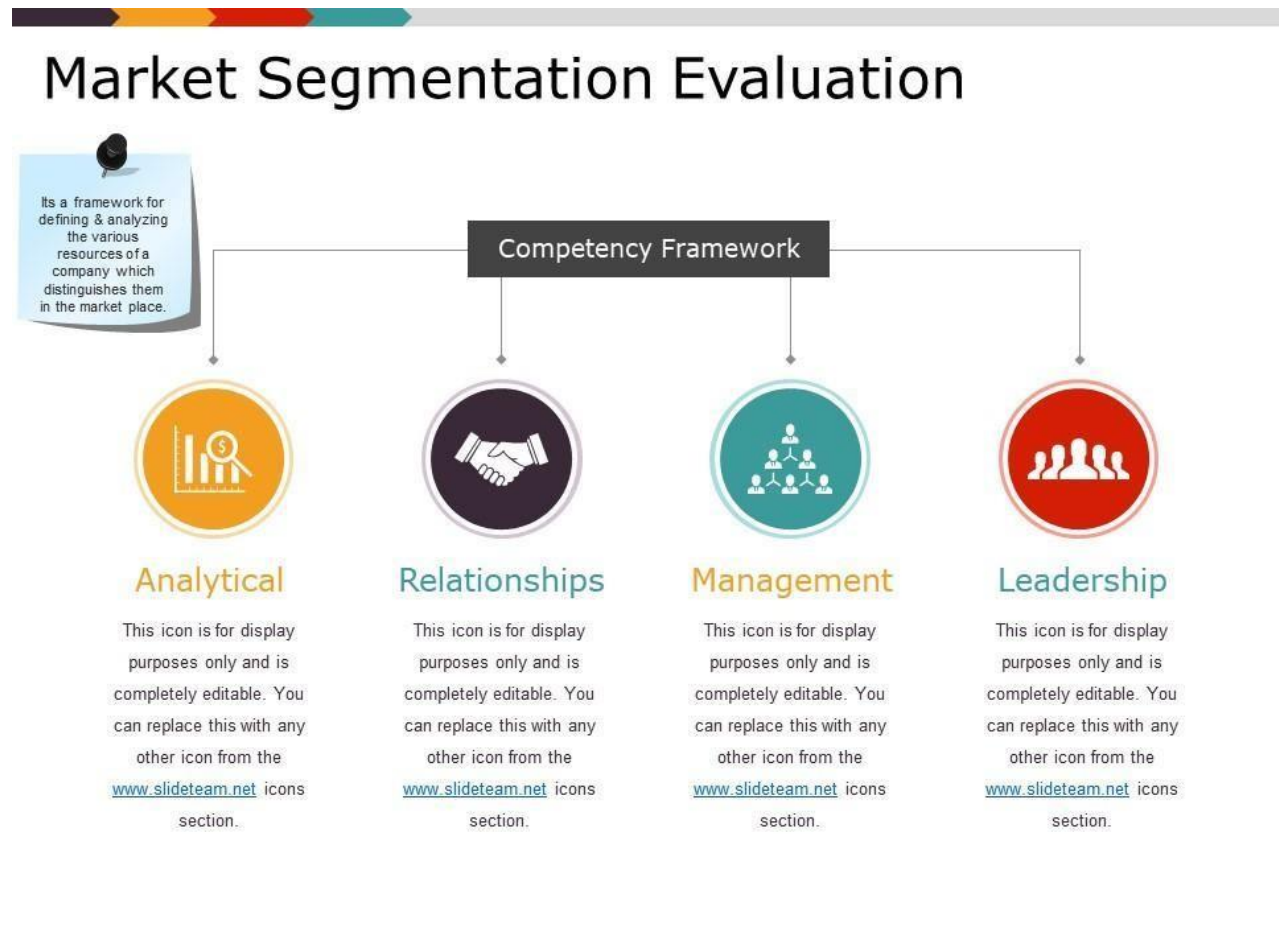## 10.1 The Targeting Decision:

This section emphasizes the importance of the targeting decision in marketing. It highlights the need to choose the most viable and profitable segment(s) to allocate resources and efforts towards. Factors such as segment attractiveness, fit with the company's objectives and capabilities, and potential for long-term growth are considered in the decision-making process.

# 10.2 Market Segment Evaluation:

The subsection delves into the evaluation of market segments to assess their suitability as target segments. It outlines various criteria and methods to evaluate segments, including market size, growth potential, competitive intensity, segment accessibility, and compatibility with the company's offerings. By thoroughly evaluating each segment, marketers can identify the most promising opportunities.



Market Segmentation Evaluation

Its a framework for defining & analyzing the various resources of a company which distinguishes them in the market place.

Competency Framework

**Analytical**

This icon is for display purposes only and is completely editable. You can replace this with any other icon from the www.slideteam.net icons section.

**Relationships**

This icon is for display purposes only and is completely editable. You can replace this with any other icon from the www.slideteam.net icons section.

**Management**

This icon is for display purposes only and is completely editable. You can replace this with any other icon from the www.slideteam.net icons section.

**Leadership**

This icon is for display purposes only and is completely editable. You can replace this with any other icon from the www.slideteam.net icons section.

# STEP 9 – CUSTOMIZING THE MARKETING MIX

In the past, marketing was seen as a toolbox of different strategies toachieve sales results. This toolbox included things like product planning pricing advertising distribution
One common model for the marketing mix is the 4Ps: Product, Price,Promotion, and Place.

Market segmentation is not a standalone strategy but works together withother strategic areas like competition and positioning. The segmentation- targeting-positioning (STP) approach is often used, where segmentation is
the first step, followed by targeting a specific segment, and then positioningthe product in a distinct way.

## Know Your Competition

Evaluate their marketing mix

1

2 Analyze their target market

Examine their messaging

3

4 Study their channels

Keep an eye on their new initiatives

5

When selecting a target segment, it is important to customize the marketingmix accordingly. This means adjusting the product, price, promotion, and
place to meet the needs and preferences of the chosen segment. For example, if a segment is interested in cultural activities, a company maydesign a product specifically tailored to their interests, offer relevant promotions, and choose appropriate distribution channels.

Each element of the marketing mix can be influenced by the target segment.
For instance, product design may be modified to better meet customer
needs, pricing decisions can be adjusted based on segment preferences, andpromotional messages can be tailored to resonate with the target segment.

Overall, the content emphasizes the importance of aligning the marketingmix with the chosen target segment to effectively meet their needs and

increase the chances of success in the market.

## Bi Clustering the Price

Bi Clustering is the process of dividing the rows and columns in the form ofclusters, here's the code which was originally written in R but have been converted to Python:

```python
import numpy as np
from sklearn.cluster import SpectralBiclustering
from sklearn.datasets import load_iris

# Load the data
data = load_iris()
X = data.data

# Perform biclustering
bicluster = SpectralBiclustering(n_clusters=12, random_state=0)
bicluster.fit(X)

# Assign bicluster labels to rows
row_labels = np.zeros(X.shape[0])
for i, rows in enumerate(bicluster.rows_):
    row_labels[rows] = i + 1

# Count the number of rows in each bicluster
row_counts = np.bincount(row_labels.astype(int))

print(row_counts)
```

### Importance of Bi Clustering the Price aspect of segmentation:

Identify Price Sensitive Segments: Identify those segments which respondto the market in a similar way

1. Customize Pricing Strategies: Make changes in the pricing strategies by knowing/getting the unique patterns of buying of the customers
2. Optimize Pricing Structures
3. Bi Clustering can reveal the customer buying patterns and price preferences
4. Price Positioning

# PLACE:

The segment being referred to as "segment 3," which consists of customerswho are interested in a destination with a rich cultural heritage. To better understand the booking preferences of these customers, a survey was conducted during their last domestic holiday. The survey allowed respondents to select multiple options for how they booked their accommodation. This information is valuable for the destination, as it helpsthem ensure that their "MUSEUMS, MONUMENTS & MUCH, MUCH MORE"product is available through the preferred distribution channels of the segment.

```python
Code Used:
import pandas as pd

import matplotlib.pyplot as plt

# Load the data into a DataFrame (assuming 'ausActivDesc' is the data)
ausActivDesc = pd.read_csv('your_data_file.csv')

# Extract the columns starting with "book" using regular expressions
book_columns = ausActivDesc.filter(regex=r'^book')

# Specify the segment membership (assuming 'cl12.3' is the segment membership)
segment_membership = cl12_3

# Create a bar chart showing the proportions of booking behavior
```

```
plt.bar(book_columns.columns,
ausActivDesc.groupby(segment_membership)[book_columns.columns].mean().values[2],
width=0.5)
plt.xlabel("Booking Behavior")
plt.ylabel("Percent")
plt.xlim(-2, 102)
plt.show()
```

## Benefits:

1) Distribution Channel Insight

2) Decision Making-For Direct Sales or Intermediaries


## PROMOTION:

In simple terms, the content explains the importance of promotion decisionsin the marketing mix. It discusses the need to develop an advertising
message that resonates with the target market and identifies effective waysof communicating this message. Other promotion tools, such as public

relations, personal selling, and sponsorship, are also mentioned.


The goal is to determine the best information sources to reach these customers and inform them about the "MUSEUMS, MONUMENTS & MUCH, MUCH MORE" product. This is done by comparing the information sources they used for their last domestic holiday and investigating their preferred TVstations.


To visualize the use of different information sources, a plot is generated using the same command as before, but this time using variables starting with "info". The resulting plot (Fig. 11.4) shows that members of segment 3rely more on information provided by tourist centers when deciding whereto spend their vacation compared to other tourists. This insight can be used

to design the promotion component of the marketing mix, such as creatingspecific information packs for the product available in hard copy at local tourist information centers and online on the tourist information center's website.

Additionally, a mosaic plot (Fig. 11.5) is used to display TV channel preferences. This plot helps understand the preferred TV channels ofcustomers in segment 3.

## Benefits of Promotion in Marketing Segmentation with respect to the mentioned code:

1. Targeted Advertising: By understanding the preferred information sources and TV channel preferences of customers in segment 3, businesses can tailor their advertising messages to resonate with this specific market segment.
This increases the effectiveness of promotional efforts and enhances the chances of reaching the target audience with the right message.

2. Customized Information Packs: The insight gained from the informationsources analysis allows businesses to create customized information packsfor the "MUSEUMS, MONUMENTS & MUCH, MUCH MORE" product.
Providing these packs both in hard copy at local tourist information centersand online on the tourist information center's website caters to the
preferences of segment 3 customers, ensuring they have access to relevantinformation through their preferred channels.

3. Improved Communication Channels: Knowing the preferred TV channelsof customers in segment 3 helps in selecting the most effective communication channels for promotional activities. Businesses can allocateresources towards advertising on these preferred channels, maximizing thereach and impact of their promotional messages.

By leveraging promotion strategies based on market segmentation insights, businesses can effectively communicate their message, engage with the target audience, and increase the likelihood of customer engagement,

ultimately leading to higher conversions and sales.

## Code Used:

```python
import pandas as pd
import matplotlib.pyplot as plt

# Set the rotation of the x-axis labels to 2 (las = 2)
plt.xticks(rotation=2)

# Create a contingency table of segment membership and TV channel preference
table_data = pd.crosstab(cl12_3, ausActivDesc['TV.channel'])

# Create a mosaic plot of the contingency table
plt.title("")
plt.xlabel("")
plt.mosaic(table_data, labelizer=lambda k: "")
plt.show()
```

# Market Segmentation Case Study on McDonalds Dataset

Kindly refer to any of the following GitHub links for the complete code implementation.

| Name | GitHub Link |
| --- | --- |
| **Hardik Sharma** | https://github.com/hardiksharmmaaaa/Feynn-Labs-Report |
| **Shyamshree Ghorai** | https://github.com/ShyamashreeGhorai1/McdonaldsCaseStudy |
| **Mridul Jain** | https://github.com/Spinachboul/Feynn-Labs-report |
| **Vansh Sonvane** | https://github.com/vansh1903/McDonalds-Market-Segmentation |