

Market Segmentation for Airlines

Market segmentation is a strategy that divides a broad target market of customers into smaller, more similar groups, and then designs a marketing strategy specifically for each group. Clustering is a common technique for market segmentation since it automatically finds similar groups given a data set.

In this problem, we'll see how clustering can be used to find similar groups of customers who belong to an airline's frequent flyer program. The airline is trying to learn more about its customers so that it can target different customer segments with different types of mileage offers.

The file [AirlinesCluster.csv](#) contains information on 3,999 members of the frequent flyer program. This data comes from the textbook "Data Mining for Business Intelligence," by Galit Shmueli, Nitin R. Patel, and Peter C. Bruce.

There are seven different variables in the dataset, described below:

- **Balance** = number of miles eligible for award travel
 - **QualMiles** = number of miles qualifying for TopFlight status
 - **BonusMiles** = number of miles earned from non-flight bonus transactions in the past 12 months
 - **BonusTrans** = number of non-flight bonus transactions in the past 12 months
 - **FlightMiles** = number of flight miles in the past 12 months
 - **FlightTrans** = number of flight transactions in the past 12 months
 - **DaysSinceEnroll** = number of days since enrolled in the frequent flyer program
-

Q1. Which TWO variables have (on average) the smallest values and largest values?

Q2. In this problem, we will normalize our data before we run the clustering algorithms. In the normalized data, which variable has the largest maximum and smallest minimum value? (*Hint: Use the `preProcess` and `predict` function from CARET package to normalize the data*).

Q3. **Hierarchical clustering:** Compute the distances between data points (using euclidean distance) and then run the Hierarchical clustering algorithm (using `method="ward.D"`) on the normalized data.

Then, plot the dendrogram of the hierarchical clustering process. Suppose the airline is looking for somewhere between 2 and 10 clusters. According to the dendrogram, which of the following is NOT a good choice for the number of clusters?

Q4. Suppose that after looking at the dendrogram and discussing with the marketing department, the airline decides to proceed with 5 clusters. Divide the data points into 5 clusters by using the `cutree` function. How many data points are in Cluster 1?

Q5. Compute the average values in each of the variables for the 5 clusters (the centroids of the clusters). You may want to compute the average values of the unnormalized data so that it is easier to interpret. Provide a Business Interpretation of the all the variables, as accordance to the clusters.

Q6. K-Means Clustering: Now run the k-means clustering algorithm on the normalized data, again creating 5 clusters. Set the seed to 88 right before running the clustering algorithm, and set the argument `iter.max` to 1000. How many clusters have more than 1000 observations?

Q7. Compute the average values in each of the variables for the 5 clusters (the centroids of the clusters) for the output from k-means clustering. You may want to compute the average values of the unnormalized data so that it is easier to interpret. Provide a Business Interpretation of the all the variables, as accordance to the clusters.