

ADS505 -- Assignment 5.1

Segmenting Consumers of Bath Soap

Stephen Kuc

10/10/22

```
In [89]: # Importing necessary libraries

import numpy as np
import pandas as pd

from sklearn import preprocessing
from sklearn.preprocessing import Normalizer, StandardScaler
from sklearn.model_selection import import train_test_split, cross_val_score, GridSearchCV
from sklearn.metrics import pairwise, accuracy_score
from scipy.cluster.hierarchy import dendrogram, linkage, fcluster
from sklearn.cluster import KMeans
from pandas.plotting import parallel_coordinates
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier

import dmba
import classificationSummary, plotDecisionTree

import matplotlib.pyplot as plt
import seaborn as sns

In [18]: # Import required packages for this chapter

import pandas as pd
import numpy as np

from scipy.spatial.distance import cosine
from sklearn.metrics.pairwise import cosine_similarity

from surprise import Dataset
from surprise import Reader
from surprise import SVDBasic

In [19]: # If you're warnings
import warnings
warnings.filterwarnings('ignore')
```

Overall problem -- CRISA wants to segment the market based on (1) purchase behavior and (2) basis of purchase; and then use this to deploy promotional budgets more effectively, among other benefits and uses. Measuring brand loyalty will be an important part of this as well.

Before diving in, let's load the data and do some basic EDA

```
In [20]: #loading data

df = pd.read_csv("C:/Users/steph/OneDrive/Documents/USB/ADS505/Data/BathSoapHousehold.csv")
df.head()
```

Member Id	SEC	FEH	MT	SEX	AGE	EDU	HS	CHILD	CS	...	PrpCat 6	PrpCat 7	PrpCat 8	PrpCat 9	PrpCat 10	PrpCat 11	PrpCat 12	PrpCat 13	
0	1010010	3	2	1	4	2	4	2	4	1	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.028037	0.0	
1	1010020	3	2	1	0	2	4	4	2	1	0.347048	0.026834	0.016100	0.014311	0.0	0.059034	0.000000	0.0	
2	1014020	2	3	1	0	2	4	5	6	4	1	0.121212	0.033550	0.010823	0.006658	0.0	0.000000	0.016234	0.0
3	1014030	4	0	0	0	4	0	0	5	0	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.0	
4	1014190	4	1	1	0	2	4	4	3	1	0.000000	0.000000	0.048193	0.000000	0.0	0.000000	0.000000	0.0	

5 rows x 46 columns

```
In [21]: print(df.dtypes)

Member Id      int64
SEC            int64
FEH            int64
MT             int64
SEX            int64
AGE            int64
EDU            int64
HS             int64
CHILD          int64
CS             int64
Affluence Index int64
No. of Brands  int64
Brand Runs    int64
Total Volume  int64
No. of Trans  int64
Value         float64
Trans / Brand Runs float64
Vol/Tran      float64
Avg. Price    float64
Pur Vol Mo Promo - % float64
Pur Vol Promo 6 % float64
Pur Vol Other Promo % float64
Br. Cd. 57, 144 float64
Br. Cd. 55    float64
Br. Cd. 272   float64
Br. Cd. 481   float64
Br. Cd. 352   float64
Br. Cd. 5    float64
Others 999    float64
Pr Cat 1      float64
Pr Cat 2      float64
Pr Cat 3      float64
Pr Cat 4      float64
PrpCat 5      float64
PrpCat 6      float64
PrpCat 7      float64
PrpCat 8      float64
PrpCat 9      float64
PrpCat 10     float64
PrpCat 11     float64
PrpCat 12     float64
PrpCat 13     float64
PrpCat 14     float64
PrpCat 15     float64
dtype: object

In [22]: df.isnull().sum()

Member Id      0
SEC            0
FEH            0
MT             0
SEX            0
AGE            0
EDU            0
HS             0
CHILD          0
CS             0
Affluence Index 0
No. of Brands  0
Brand Runs     0
Total Volume   0
No. of Trans   0
Value          0
Trans / Brand Runs 0
Avg. Price     0
Pur Vol Mo Promo - % 0
Pur Vol Promo 6 % 0
Pur Vol Other Promo % 0
Br. Cd. 57, 144 0
Br. Cd. 55      0
Br. Cd. 272     0
Br. Cd. 481     0
Br. Cd. 352     0
Br. Cd. 5       0
Others 999      0
Pr Cat 1        0
Pr Cat 2        0
Pr Cat 3        0
Pr Cat 4        0
PrpCat 5        0
PrpCat 6        0
PrpCat 7        0
PrpCat 8        0
PrpCat 9        0
PrpCat 10       0
PrpCat 11       0
PrpCat 12       0
PrpCat 13       0
PrpCat 14       0
PrpCat 15       0
dtype: int64
```

Since we'll be using K-means on different sets of variables, let's examine each set, and then recombine them when the time comes for that.

First, a set of variables that describe purchase behavior, including brand loyalty.

Purchase behavior also included -- volume, frequency, susceptibility to discounts, and brand loyalty.

Viewing the features:

No. of Brands, Brand Runs, Total Volume, Sum of Value, Trans/Brand Runs, Vol/Trans, Pur Vol, No Promo %, Pur Vol Promo %, Pur Vol Other Promo %, and all brandwise purchase volume (which may be further analyzed on its own)

```
In [23]: # demographics dataframe, that may not be used
demographics = ['SEC', 'FEH', 'MT', 'SEX', 'AGE', 'EDU', 'HS', 'SHILD', 'CS', 'Affluence_Index']

# purchase behavior dataframe, initially without brand
pb_df = df[['No. of Brands', 'Brand Runs', 'Total Volume',
            'Value', 'Trans / Brand Runs', 'Vol/Tran',
            'Pur Vol Mo Promo - %', 'Pur Vol Promo 6 %',
            'Pur Vol Other Promo %']]

# Brand dataframe - will likely transform
brand_pct = df[['Br. Cd. 57, 144', 'Br. Cd. 55', 'Br. Cd. 272', 'Br. Cd. 286',
               'Br. Cd. 24', 'Br. Cd. 481', 'Br. Cd. 352', 'Br. Cd. 5', 'Others 999']]
brand_pct.head()
```

	Br. Cd. 57, 144	Br. Cd. 55	Br. Cd. 272	Br. Cd. 286	Br. Cd. 24	Br. Cd. 481	Br. Cd. 352	Br. Cd. 5	Others 999
0	0.376947	0.130841	0.0	0.000000	0.0	0.000000	0.0	0.000000	0.492212
1	0.021467	0.075134	0.0	0.000000	0.0	0.059034	0.0	0.014490	0.699463
2	0.025974	0.545455	0.0	0.030303	0.0	0.000000	0.0	0.019481	0.378788
3	0.400000	0.600000	0.0	0.000000	0.0	0.000000	0.0	0.000000	0.000000
4	0.048193	0.144578	0.0	0.000000	0.0	0.000000	0.0	0.000000	0.722229

Since we're interested in measuring brand loyalty, as mentioned in a note in question, a good way to see if one customer does really prefer one brand over another, we can just check for the max in the brand_pct dataframe, without using the Others 999 column, as that encompasses all other brands, not just one singular.

Before that, let's at least try to find the distances between brands

```
In [24]: # before we simply use this, let's see if a distance metric can be found from the brand percentages

brand_pct_sp = brand_pct.drop('Others 999', axis = 1)
brand_pct_sp = brand_pct_sp.apply(lambda xx: x.astype(float))

brand_distance = pairwise.pairwise_distances(brand_pct_sp, metric = 'euclidean')
brand_distance = pd.DataFrame(brand_distance)
brand_distance
```

	0	1	2	3	4	5	6	7	8	9	...	590	591	592	5
0	0.000000	0.392366	0.544412	0.469725	0.329041	0.306842	0.114276	0.739209	0.131092	0.508000	...	0.163751	0.359799	1.043335	0.4408
1	0.392366	0.000000	0.094812	0.657772	0.173258	0.166288	0.459365	0.732421	0.402890	0.271948	...	0.485995	0.191845	0.978114	0.2564
2	0.544412	0.094812	0.000000	0.279695	0.403104	0.483485	0.658258	0.248089	0.654030	0.574144	...	0.707680	0.556860	1.108527	0.5784
3	0.469725	0.657772	0.379695	0.000000	0.575480	0.624051	0.557281	0.404554	0.600186	0.708862	...	0.604710	0.698034	1.203885	0.7456
4	0.329041	0.173258	0.403104	0.575480	0.000000	0.083679	0.416172	0.645765	0.366574	0.403563	...	0.541018	0.190259	0.975994	0.2416

Since we're interested in measuring brand loyalty, as mentioned in a note in question, a good way to see if one customer does really prefer one brand over another, we can just check for the max in the brand_pct dataframe, without using the Others 999 column, as that encompasses all other brands, not just one singular.

Before that, let's at least try to find the distances between brands

```
In [24]: # before we simply use this, let's see if a distance metric can be found from the brand percentages

brand_pct_sp = brand_pct.drop('Others 999', axis = 1)
brand_pct_sp = brand_pct_sp.apply(lambda xx: x.astype(float))

brand_distance = pairwise.pairwise_distances(brand_pct_sp, metric = 'euclidean')
brand_distance = pd.DataFrame(brand_distance)
brand_distance
```

	0	1	2	3	4	5	6	7	8	9	...	590	591	592	5
0	0.000000	0.392366	0.544412	0.469725	0.329041	0.306842	0.114276	0.739209	0.131092	0.508000	...	0.163751	0.359799	1.043335	0.4408
1	0.392366	0.000000	0.094812	0.657772	0.173258	0.166288	0.459365	0.732421	0.402890	0.271948	...	0.485995	0.191845	0.978114	0.2564
2	0.544412	0.094812	0.000000	0.279695	0.403104	0.483485	0.658258	0.248089	0.654030	0.574144	...	0.707680	0.556860	1.108527	0.5784
3	0.469725	0.657772	0.379695	0.000000	0.575480	0.624051	0.557281	0.404554	0.600186	0.708862	...	0.604710	0.698034	1.203885	0.7456
4	0.329041	0.173258	0.403104	0.575480	0.000000	0.083679	0.416172	0.645765	0.366574	0.403563	...	0.541018	0.190259	0.975994	0.2416

Since we're interested in measuring brand loyalty, as mentioned in a note in question, a good way to see if one customer does really prefer one brand over another, we can just check for the max in the brand_pct dataframe, without using the Others 999 column, as that encompasses all other brands, not just one singular.

Before that, let's at least try to find the distances between brands

```
In [24]: # before we simply use this, let's see if a distance metric can be found from the brand percentages

brand_pct_sp = brand_pct.drop('Others 999', axis = 1)
brand_pct_sp = brand_pct_sp.apply(lambda xx: x.astype(float))

brand_distance = pairwise.pairwise_distances(brand_pct_sp, metric = 'euclidean')
brand_distance = pd.DataFrame(brand_distance)
brand_distance
```

	0	1	2	3	4	5	6	7	8	9	...	590	591	592	5
0	0.000000	0.392366	0.544412	0.469725	0.329041	0.306842	0.114276	0.739209	0.131092	0.508000	...	0.163751	0.359799	1.043335	0.4408
1	0.392366	0.000000	0.094812	0.657772	0.173258	0.166288	0.459365	0.732421	0.402890	0.271948	...	0.485995	0.191845	0.978114	0.2564
2	0.544412	0.094812	0.000000	0.279695	0.403104	0.483485	0.658258	0.248089	0.654030	0.574144	...	0.707680	0.556860	1.108527	0.5784
3	0.469725	0.657772	0.379695	0.000000	0.575480	0.624051	0.557281	0.404554	0.600186	0.708862	...	0.604710	0.698034	1.203885	0.7456
4	0.329041	0.173258	0.403104	0.575480	0.000000	0.083679	0.416172	0.645765	0.366574	0.403563	...	0.541018	0.190259	0.975994	0.2416

Since we're interested in measuring brand loyalty, as mentioned in a note in question, a good way to see if one customer does really prefer one brand over another, we can just check for the max in the brand_pct dataframe, without using the Others 999 column, as that encompasses all other brands, not just one singular.

Before that, let's at least try to find the distances between brands

```
In [24]: # before we simply use this, let's see if a distance metric can be found from the brand percentages

brand_pct_sp = brand_pct.drop('Others 999', axis = 1)
brand_pct_sp = brand_pct_sp.apply(lambda xx: x.astype(float))

brand_distance = pairwise.pairwise_distances(brand_pct_sp, metric = 'euclidean')
brand_distance = pd.DataFrame(brand_distance)
brand_distance
```

	0	1	2	3	4	5	6	7	8	9	...	590	591	592	5
0	0.000000	0.392366	0.544412	0.469725	0.329041	0.306842	0.114276	0.739209	0.131092	0.508000	...	0.163751	0.359799	1.043335	0.4408
1	0.392366	0.000000	0.094812	0.657772	0.173258	0.166288	0.459365	0.732421	0.402890	0.271948	...	0.485995	0.191845	0.978114	0.2564
2	0.544412	0.094812	0.000000	0.279695	0.403104	0.483485	0.658258	0.248089	0.654030	0.574144	...	0.707680	0.556860	1.108527	0.5784
3	0.469725	0.657772	0.379695	0.000000	0.575480	0.624051	0.557281	0.404554	0.600186	0.708862	...	0.604710	0.698034	1.203885	0.7456
4	0.329041	0.173258	0.403104	0.575480	0.000000	0.083679	0.416172	0.645765	0.366574	0.403563	...	0.541018	0.190259	0.975994	0.2416

Since we're interested in measuring brand loyalty, as mentioned in a note in question, a good way to see if one customer does really prefer one brand over another, we can just check for the max in the brand_pct dataframe, without using the Others 999 column, as that encompasses all other brands, not just one singular.

Before that, let's at least try to find the distances between brands

```
In [24]: # before we simply use this, let's see if a distance metric can be found from the brand percentages

brand_pct_sp = brand_pct.drop('Others 999', axis = 1)
brand_pct_sp = brand_pct_sp.apply(lambda xx: x.astype(float))

brand_distance = pairwise.pairwise_distances(brand_pct_sp, metric = 'euclidean')
brand_distance = pd.DataFrame(brand_distance)
brand_distance
```

	0	1	2	3	4	5	6	7	8	9	...	590	591	592	5
0	0.000000	0.392366	0.544412	0.469725	0.329041	0.306842	0.114276	0.739209	0.131092	0.508000	...	0.163751	0.359799	1.043335	0.4408
1	0.392366	0.000000	0.094812	0.657772	0.173258	0.166288	0.459365	0.732421	0.402890	0.271948	...	0.485995	0.191845	0.978114	0.2564
2	0.544412	0.094812	0.000000	0.279695	0.403104	0.483485	0.658258	0.248089	0.654030	0.574144	...	0.707680	0.556860	1.108527	0.5784
3	0.469725	0.657772	0.379695	0.000000	0.575480	0.624051	0.557281	0.404554	0.600186	0.708862	...	0.604710	0.698034	1.203885	0.7456
4	0.329041	0.173258	0.403104	0.575480	0.000000	0.083679	0.416172	0.645765	0.366574	0.403563	...	0.541018	0.190259	0.975994	0.2416

Since we're interested in measuring brand loyalty, as mentioned in a note in question, a good way to see if one customer does really prefer one brand over another, we can just check for the max in the brand_pct dataframe, without using the Others 999 column, as that encompasses all other brands, not just one singular.

Before that, let's at least try to find the distances between brands

```
In [24]: # before we simply use this, let's see if a distance metric can be found from the brand percentages

brand_pct_sp = brand_pct.drop('Others 999', axis = 1)
brand_pct_sp = brand_pct_sp.apply(lambda xx: x.astype(float))

brand_distance = pairwise.pairwise_distances(brand_pct_sp, metric = 'euclidean')
brand_distance = pd.DataFrame(brand_distance)
brand_distance
```

	0	1	2	3	4	5	6	7	8	9	...	590	591	592	5
0	0.000000	0.392366	0.544412	0.469725	0.329041	0.306842	0.114276	0.739209	0.131092	0.508000	...	0.163751	0.359799	1.043335	0.4408
1	0.392366	0.000000	0.094812	0.657772	0.173258	0.166288	0.459365	0.732421	0.402890	0.271948	...	0.485995	0.191845	0.978114	0.2564
2	0.544412	0.094812	0.000000	0.279695	0.403104	0.483485	0.658258	0.248089	0.654030	0.574144	...	0.707680	0.556860	1.108527	0.5784
3	0.469725	0.657772	0.379695	0.000000	0.575480	0.624051	0.557281	0.404554	0.600186	0.708862	...	0.604710	0.698034	1.203885	0.7456
4	0.329041	0.173258	0.403104	0.575480	0.000000	0.083679	0.416172	0.645765	0.366574	0.403563	...	0.541018	0.190259	0.975994	0.2416

Since we're interested in measuring brand loyalty, as mentioned in a note in question, a good way to see if one customer does really prefer one brand over another, we can just check for the max in the brand_pct dataframe, without using the Others 999 column, as that encompasses all other brands, not just one singular.

Before that, let's at least try to find the distances between brands

```
In [24]: # before we simply use this, let's see if a distance metric can be found from the brand percentages

brand_pct_sp = brand_pct.drop('Others 999', axis = 1)
brand_pct_sp = brand_pct_sp.apply(lambda xx: x.astype(float))

brand_distance = pairwise.pairwise_distances(brand_pct_sp, metric = 'euclidean')
brand_distance = pd.DataFrame(brand_distance)
brand_distance
```

	0	1	2	3	4	5	6	7	8	9	...	590	591	592	5
0	0.000000	0.392366	0.544412	0.469725	0.329041										

