


```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import datetime
import warnings
warnings.filterwarnings('ignore')
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)

In [2]: sun_df = pd.read_csv('sun.csv')

In [34]: sun_df.head()
```

Out[34]:	Unnamed: 0	CustomerID	Gender	Age_group	Ufly_membership_status	Card_holder?	Total_trips	Total_amount_spent	#Discounts	Preferred_tr
0	0	4120414C52484D414E44696420493F7C20676574207468...	M	Young Adults	Not member	No	1	174.0	0	
1	1	414142454C44696420493F7C2067657420746869732072...	M	Youth	Not member	No	1	231.9	0	
2	2	4141424552472042524F4F4B5344696420493F7C206765...	F	Middle Aged	Not member	No	1	294.9	0	
3	3	41414245524744696420493F7C20676574207468697320...	M	Middle Aged	Not member	No	2	0.0	2	
4	4	41414245524744696420493F7C20676574207468697320...	M	Young Adults	Standard	No	2	973.6	0	

Random Sampling

We have used block randomization technique to take sample of data. This helped us wholistically understand every type of Sun Country Airlines customer.

```
In [43]: new1 = sun_df[sun_df['Ufly_membership_status']=='Elite'].sample(n=1300,random_state=2)
new2 = sun_df[sun_df['Ufly_membership_status']=='Standard'].sample(n=1000,random_state=2)
new3 = sun_df[sun_df['Card_holder?']=='Yes'].sample(n=2400,random_state=2)
new4 = sun_df[sun_df['Card_holder?']=='No'].sample(n=1200,random_state=2)
new5 = sun_df[sun_df['Preferred_travel_class']=='Coach'].sample(n=800,random_state=2)
new6 = sun_df[sun_df['Preferred_travel_class']=='First Class'].sample(n=1400,random_state=2)
new7 = sun_df[sun_df['Age_group']=='Young Adults'].sample(n=800,random_state=2)
new8 = sun_df[sun_df['Age_group']=='Children'].sample(n=1100,random_state=2)
new9 = sun_df[sun_df['Age_group']=='Youth'].sample(n=1100,random_state=2)
new10 = sun_df[sun_df['Age_group']=='Middle Aged'].sample(n=800,random_state=2)
new11 = sun_df[sun_df['Age_group']=='Senior'].sample(n=800,random_state=2)
new12 = sun_df[sun_df['Preferred_source-booking']=='SCA Website Booking'].sample(n=400,random_state=2)
new13 = sun_df[sun_df['Preferred_source-booking']=='Outside Booking'].sample(n=400,random_state=2)
new14 = sun_df[sun_df['Preferred_source-booking']=='Airport'].sample(n=800,random_state=2)
new15 = sun_df[sun_df['Preferred_source-booking']=='Tour Operator Portal'].sample(n=800,random_state=2)
new16 = sun_df[sun_df['Preferred_source-booking']=='Reservations Booking'].sample(n=800,random_state=2)
new17 = sun_df[sun_df['Preferred_source-booking']=='No Preference'].sample(n=800,random_state=2)
new18 = sun_df[sun_df['Preferred_source-booking']=='SY Vacation'].sample(n=800,random_state=2)

In [44]: new_df = pd.concat([new1, new2, new3, new4, new5, new6, new7, new8, new9, new10, new11, new12, new13, new14, new15, \
                             new16, new17, new18])

In [45]: new_df.drop_duplicates(inplace=True)

In [46]: new_df.count()

Out[46]: Unnamed: 0      17345
CustomerID      17345
Gender          17345
Age_group       17345
Ufly_membership_status 17345
Card_holder?    17345
Total_trips     17345
Total_amount_spent 17345
#Discounts      17345
Preferred_travel_class 17345
#Upgrades       17345
Preferred_source-booking 17345
dtype: int64

In [47]: # We did not add the feature - "total amount spent" in X because it is correlated with "total trips".
# This helps in better analysis and faster execution.
X = new_df[['Gender', 'Age_group', 'Ufly_membership_status', 'Card_holder?', 'Total_trips', \
            '#Discounts', '#Upgrades', 'Preferred_travel_class', 'Preferred_source-booking']]

In [48]: X.head()

Out[48]:
```

	Gender	Age_group	Ufly_membership_status	Card_holder?	Total_trips	#Discounts	#Upgrades	Preferred_travel_class	Preferred_source-booking
748819	F	Middle Aged	Elite	No	2	2	0.0	First Class	Reservations Booking
1453483	M	Senior	Elite	No	6	0	0.0	Coach	Outside Booking
403441	M	Young Adults	Elite	No	5	5	0.0	First Class	Airport
780044	M	Middle Aged	Elite	No	2	2	0.0	First Class	SCA Website Booking
22251	F	Young Adults	Elite	No	1	1	0.0	Coach	SCA Website Booking

```
In [49]: import gower
from sklearn_extra.cluster import KMedoids

In [50]: # We have used gower distance as distance metric
gower_dist = gower.gower_matrix(X)
```

Cluster quality analysis

```
In [56]: from sklearn.metrics import silhouette_samples, silhouette_score

In [68]: print(silhouette_score(gower_dist, X['cluster']))

0.22189361

Cluster quality was measured using silhouette coefficient which was highest for 5 cluster solution.

Silhouette score = 0.22
```

Cluster creation

```
In [52]: clusterer = KMedoids(n_clusters = 5, random_state = 10, method = 'pam')
X['cluster'] = clusterer.fit_predict(gower_dist)
```

Visualizing Clusters



Understanding CLusters

```
In [54]: new_df['cluster'] = X['cluster']

In [55]: # Summary statistics by cluster

print('gender')
print(new_df.groupby('cluster')['Gender'].describe())
print('Age_Group')
print(new_df.groupby('cluster')['Age_group'].describe())
print('UflyMemberStatus')
print(new_df.groupby('cluster')['Ufly_membership_status'].describe())
print('CardHolder')
print(new_df.groupby('cluster')['Card_holder?'].describe())
print('NumTrips')
print(new_df.groupby('cluster')['Total_trips'].describe())
print('TotalDocAmt')
print(new_df.groupby('cluster')['Total_amount_spent'].describe())
print('# Discounts')
print(new_df.groupby('cluster')['#Discounts'].describe())
print('Preferred_travel_class')
print(new_df.groupby('cluster')['Preferred_travel_class'].describe())
print('#Upgrades')
print(new_df.groupby('cluster')['#Upgrades'].describe())
print('Preferred_source-booking')
print(new_df.groupby('cluster')['Preferred_source-booking'].describe())

gender
cluster
count unique top freq
0 2391 2 F 2214
1 3214 1 F 3214
2 3429 1 M 3429
3 3278 1 F 3278
4 5033 2 M 5032
Age_Group
count unique top freq
cluster
0 2391 5 Senior 1366
1 3214 5 Young Adults 973
2 3429 5 Senior 1473
3 3278 5 Children 927
4 5033 5 Young Adults 1306
UflyMemberStatus
count unique top freq
cluster
0 2391 3 Standard 2144
1 3214 3 Not member 2937
2 3429 3 Standard 2577
3 3278 3 Not member 2375
4 5033 3 Not member 4838
CardHolder
count unique top freq
cluster
0 2391 2 Yes 1471
1 3214 2 No 3211
2 3429 2 No 2210
3 3278 1 No 3278
4 5033 2 No 5028
NumTrips
count mean std min 25% 50% 75% max
cluster
0 2391.0 4.790882 6.897803 1.0 2.0 2.0 5.0 94.0
1 3214.0 2.277225 2.159958 1.0 2.0 2.0 2.0 45.0
2 3429.0 4.907262 7.468770 1.0 2.0 2.0 4.0 103.0
3 3278.0 2.298658 1.705742 1.0 2.0 2.0 2.0 49.0
4 5033.0 2.362911 3.087374 1.0 2.0 2.0 2.0 62.0
TotalDocAmt
count mean std min 25% 50% 75% \
cluster
0 2391.0 1948.815337 3317.402165 0.0 496.4000 943.20 2040.56
1 3214.0 642.542866 1072.959548 0.0 138.8000 496.00 836.00
2 3429.0 2002.656113 3373.671886 0.0 479.6000 922.00 1871.60
3 3278.0 767.290912 790.607066 0.0 337.8475 626.65 975.41
4 5033.0 766.905786 1391.741479 0.0 178.9000 536.00 927.60
max
cluster
0 43721.47
1 23756.18
2 46556.78
3 17613.70
4 26748.80
# Discounts
count mean std min 25% 50% 75% max
cluster
0 2391.0 2.583438 4.299806 0.0 0.0 2.0 3.0 53.0
1 3214.0 0.911325 1.259107 0.0 0.0 0.0 2.0 17.0
2 3429.0 2.362788 4.235852 0.0 0.0 2.0 2.0 52.0
3 3278.0 1.024100 1.171982 0.0 0.0 1.0 2.0 14.0
4 5033.0 0.973376 1.539260 0.0 0.0 0.0 2.0 43.0
Preferred_travel_class
count unique top freq
cluster
0 2391 2 Coach 1593
1 3214 6 Outside Booking 1911
2 3429 7 SCA Website Booking 2291
3 3278 7 SCA Website Booking 2053
4 5033 7 Outside Booking 1810
#Upgrades
count mean std min 25% 50% 75% max
cluster
0 2391.0 0.997909 3.849849 -3.0 0.0 0.0 0.0 69.0
1 3214.0 0.130989 0.961418 -2.0 0.0 0.0 0.0 26.0
2 3429.0 1.132983 3.999977 -4.0 0.0 0.0 1.0 70.0
3 3278.0 0.135143 0.747926 -4.0 0.0 0.0 0.0 25.0
4 5033.0 0.175641 1.306606 -2.0 0.0 0.0 0.0 46.0
Preferred_source-booking
count unique top freq
cluster
0 2391 7 SCA Website Booking 1683
1 3214 6 Outside Booking 1911
2 3429 7 SCA Website Booking 2291
3 3278 7 SCA Website Booking 2053
4 5033 7 Outside Booking 1810
```