```
import pandas as pd
```

```
path = "/content/mcdonalds.csv"
import pandas as pd
from sklearn.decomposition import PCA
import numpy as np
from sklearn import preprocessing
from bioinfokit.visuz import cluster
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.utils import resample
from sklearn.metrics import adjusted_rand_score
```

```
pip install bioinfokit
```

```
Requirement already satisfied: bioinfokit in /usr/local/lib/python3.10/dist-packages (2.1.4)
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (2.0.3)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (1.25.2)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (3.7.1)
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (1.11.4)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (1.2.2)
Requirement already satisfied: seaborn in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (0.13.1)
Requirement already satisfied: matplotlib-venn in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (0.11.10)
Requirement already satisfied: tabulate in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (0.9.0)
Requirement already satisfied: statsmodels in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (0.14.2)
Requirement already satisfied: textwrap3 in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (0.9.2)
Requirement already satisfied: adjustText in /usr/local/lib/python3.10/dist-packages (from bioinfokit) (1.2.0)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (1.2.1)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (4.53.1)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (1.4.5)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (24.1)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (9.4.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (3.1.2)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib->bioinfokit) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas->bioinfokit) (2023.4)
Requirement already satisfied: tzdata>=2022.1 in /usr/local/lib/python3.10/dist-packages (from pandas->bioinfokit) (2024.1)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn->bioinfokit) (1.4.2)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn->bioinfokit) (3.5.0)
Requirement already satisfied: patsy>=0.5.6 in /usr/local/lib/python3.10/dist-packages (from statsmodels->bioinfokit) (0.5.6)
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (from patsy>=0.5.6->statsmodels->bioinfokit) (1.16.0)
```

```
data=pd.read_csv("mcdonalds.csv")
data1=pd.read_csv("mcdonalds.csv")
data.columns.values.tolist()
```

```
['yummy',
 'convenient',
 'spicy',
 'fattening',
 'greasy',
 'fast',
 'cheap',
 'tasty',
 'expensive',
 'healthy',
 'disgusting',
 'Like',
 'Age',
 'VisitFrequency',
 'Gender']
```

```
data.shape
```

```
(1453, 15)
```

```
data.head(6)
```

| | yummy | convenient | spicy | fattening | greasy | fast | cheap | tasty | expensive | healthy | disgusting | Like | Age | VisitFrequency | Gender |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | No | Yes | No | Yes | No | Yes | Yes | No | Yes | No | No | -3 | 61 | Every three months | Female |
| 1 | Yes | Yes | No | Yes | Yes | Yes | Yes | Yes | Yes | No | No | +2 | 51 | Every three months | Female |
| 2 | No | Yes | Yes | Yes | Yes | Yes | No | Yes | Yes | Yes | No | +1 | 62 | Every three months | Female |
| 3 | Yes | Yes | No | Yes | Yes | Yes | Yes | Yes | No | No | Yes | +4 | 69 | Once a week | Female |

```python
MD=data.iloc[:,0:11].replace("Yes",1).replace("No",0)
mean=round(MD.mean(),2)
mean
```

```
yummy         0.55
convenient    0.91
spicy         0.09
fattening     0.87
greasy        0.53
fast          0.90
cheap         0.60
tasty         0.64
expensive     0.36
healthy       0.20
disgusting    0.24
dtype: float64
```

```python
pca = PCA()
MD_pca=pca.fit_transform(MD)
MD_p=pca.fit(MD)

SD=np.sqrt(pca.explained_variance_)
PV=pca.explained_variance_ratio_
index=[]
for i in range(len(SD)):
    i=i+1
    index.append("PC{}".format(i))

sum=pd.DataFrame({
    "Standard deviation":SD,"Proportion of Variance":PV,"Cumulative Proportion":PV.cumsum()
},index=index)
sum
```

| | Standard deviation | Proportion of Variance | Cumulative Proportion |
|---|---|---|---|
| PC1 | 0.757050 | 0.299447 | 0.299447 |
| PC2 | 0.607456 | 0.192797 | 0.492244 |
| PC3 | 0.504619 | 0.133045 | 0.625290 |
| PC4 | 0.398799 | 0.083096 | 0.708386 |
| PC5 | 0.337405 | 0.059481 | 0.767866 |
| PC6 | 0.310275 | 0.050300 | 0.818166 |
| PC7 | 0.289697 | 0.043849 | 0.862015 |
| PC8 | 0.275122 | 0.039548 | 0.901563 |
| PC9 | 0.265251 | 0.036761 | 0.938323 |
| PC10 | 0.248842 | 0.032353 | 0.970677 |
| PC11 | 0.236903 | 0.029323 | 1.000000 |

```python
print("Standard Deviation:\n",SD.round(1))

load = (pca.components_)
i=0
rot_matrix = MD_p.components_.T

rot_df = pd.DataFrame(rot_matrix, index=MD.columns.values, columns=index)
rot_df=round(-rot_df,2)
rot_df
```

Standard Deviation:
[0.8 0.6 0.5 0.4 0.3 0.3 0.3 0.3 0.3 0.2 0.2]

|            | PC1   | PC2   | PC3   | PC4   | PC5   | PC6   | PC7   | PC8   | PC9   | PC10  | PC11  |
|------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| yummy      | 0.48  | -0.36 | 0.30  | -0.06 | 0.31  | -0.17 | 0.28  | -0.01 | -0.57 | 0.11  | -0.05 |
| convenient | 0.16  | -0.02 | 0.06  | 0.14  | -0.28 | 0.35  | 0.06  | 0.11  | 0.02  | 0.67  | 0.54  |
| spicy      | 0.01  | -0.02 | 0.04  | -0.20 | -0.07 | 0.36  | -0.71 | -0.38 | -0.40 | 0.08  | -0.14 |
| fattening  | -0.12 | 0.03  | 0.32  | 0.35  | 0.07  | 0.41  | 0.39  | -0.59 | 0.16  | 0.01  | -0.25 |
| greasy     | -0.30 | 0.06  | 0.80  | -0.25 | -0.36 | -0.21 | -0.04 | 0.14  | 0.00  | -0.01 | -0.00 |
| fast       | 0.11  | 0.09  | 0.06  | 0.10  | -0.11 | 0.59  | 0.09  | 0.63  | -0.17 | -0.24 | -0.34 |
| cheap      | 0.34  | 0.61  | 0.15  | -0.12 | 0.13  | 0.10  | 0.04  | -0.14 | -0.08 | -0.43 | 0.49  |
| tasty      | 0.47  | -0.31 | 0.29  | 0.00  | 0.21  | 0.08  | -0.36 | 0.07  | 0.64  | -0.08 | -0.02 |
| expensive  | -0.33 | -0.60 | -0.02 | -0.07 | 0.00  | 0.26  | 0.07  | -0.03 | -0.07 | -0.45 | 0.49  |
| healthy    | 0.21  | -0.08 | -0.19 | -0.76 | -0.29 | 0.18  | 0.35  | -0.18 | 0.19  | 0.04  | -0.16 |
| disgusting | -0.37 | 0.14  | 0.09  | -0.37 | 0.73  | 0.21  | 0.03  | 0.17  | 0.07  | 0.29  | 0.04  |

rot_df

|            | PC1   | PC2   | PC3   | PC4   | PC5   | PC6   | PC7   | PC8   | PC9   | PC10  | PC11  |
|------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| yummy      | 0.48  | -0.36 | 0.30  | -0.06 | 0.31  | -0.17 | 0.28  | -0.01 | -0.57 | 0.11  | -0.05 |
| convenient | 0.16  | -0.02 | 0.06  | 0.14  | -0.28 | 0.35  | 0.06  | 0.11  | 0.02  | 0.67  | 0.54  |
| spicy      | 0.01  | -0.02 | 0.04  | -0.20 | -0.07 | 0.36  | -0.71 | -0.38 | -0.40 | 0.08  | -0.14 |
| fattening  | -0.12 | 0.03  | 0.32  | 0.35  | 0.07  | 0.41  | 0.39  | -0.59 | 0.16  | 0.01  | -0.25 |
| greasy     | -0.30 | 0.06  | 0.80  | -0.25 | -0.36 | -0.21 | -0.04 | 0.14  | 0.00  | -0.01 | -0.00 |
| fast       | 0.11  | 0.09  | 0.06  | 0.10  | -0.11 | 0.59  | 0.09  | 0.63  | -0.17 | -0.24 | -0.34 |
| cheap      | 0.34  | 0.61  | 0.15  | -0.12 | 0.13  | 0.10  | 0.04  | -0.14 | -0.08 | -0.43 | 0.49  |
| tasty      | 0.47  | -0.31 | 0.29  | 0.00  | 0.21  | 0.08  | -0.36 | 0.07  | 0.64  | -0.08 | -0.02 |
| expensive  | -0.33 | -0.60 | -0.02 | -0.07 | 0.00  | 0.26  | 0.07  | -0.03 | -0.07 | -0.45 | 0.49  |
| healthy    | 0.21  | -0.08 | -0.19 | -0.76 | -0.29 | 0.18  | 0.35  | -0.18 | 0.19  | 0.04  | -0.16 |
| disgusting | -0.37 | 0.14  | 0.09  | -0.37 | 0.73  | 0.21  | 0.03  | 0.17  | 0.07  | 0.29  | 0.04  |

```python
from bioinfokit.visuz import cluster
```

```python
cluster.biplot(cscore=MD_pca, loadings=-load, labels=data.columns.values,var1=0,var2=0, show=True, dim=(10, 10))
```

```
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
WARNING:matplotlib.font_manager:findfont: Font family 'Arial' not found.
```
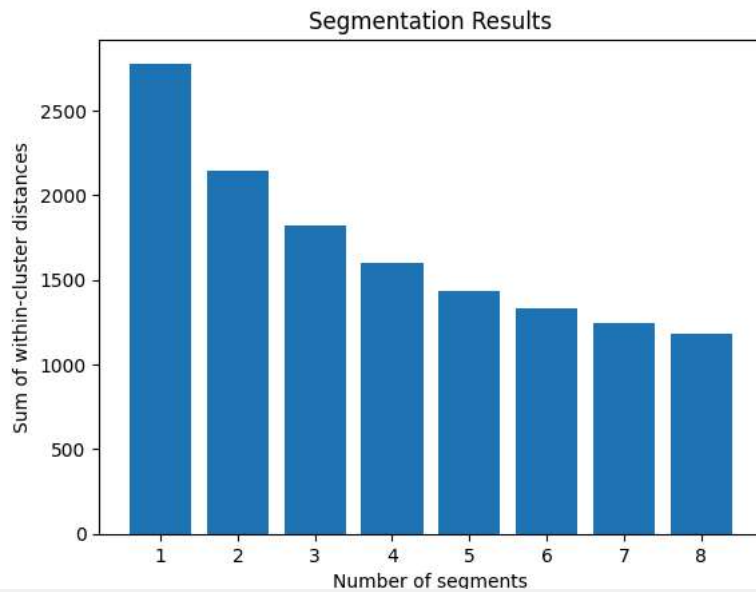


```python
np.random.seed(1234)

nrep = 10

num_segments = range(1, 9)
within_cluster_distances = []
MD_km28 = {}

for k in num_segments:
    kmeans = KMeans(n_clusters=k, n_init=nrep, random_state=1234)
    kmeans.fit(MD)
    within_cluster_distances.append((kmeans.inertia_))
    MD_km28[str(k)] = kmeans

plt.bar(num_segments, within_cluster_distances)
plt.xlabel("Number of segments")
plt.ylabel("Sum of within-cluster distances")
plt.title("Segmentation Results")
plt.show()
```

```
num_segments = range(2, 9)

segment_stability = []
for segment in range(2, 9):
    labels_segment = MD_km28[str(segment)].predict(MD)
    segment_stability.append(labels_segment)

plt.figure(figsize=(8, 6))
for i, segment in enumerate(range(2, 9)):
    plt.plot(num_segments, [np.mean(segment_stability[i] == labels) for labels in segment_stability], marker='o', label=f'Segment {segment}'

plt.xlabel('Number of Segments')
plt.ylabel('Segment Level Stability')
plt.title('Segment Level Stability Across Solutions (SLSA) Plot')
plt.xticks(num_segments)
plt.legend()
plt.grid(True)

plt.show()
```

```python
segment_solutions = ["2", "3", "4", "5"]
segment_labels = {}
segment_similarities = {}

for segment in segment_solutions:
    segment_labels[segment] = MD_km28[segment].predict(MD)
    segment_similarities[segment] = MD_km28[segment].transform(MD).min(axis=1)


segment_stability_values = []
for segment in segment_solutions:
    similarities = segment_similarities[segment]
    normalized_similarities = similarities / np.max(similarities)
    segment_stability_values.append(normalized_similarities)

plt.boxplot(segment_stability_values, whis=1.5)
plt.xlabel("Segment Number")
plt.ylabel("Segment Stability")
plt.xticks(range(1, len(segment_solutions) + 1), segment_solutions)
plt.ylim(0, 1)
plt.title("Segment Level Stability within Solutions")

plt.show()
```



```python
from scipy.stats import entropy

np.random.seed(1234)
k_values = range(2, 9)
MD_m28 = []

for k in k_values:
    model = KMeans(n_clusters=k, random_state=1234)
    model.fit(MD.values)
    iter_val = model.n_iter_
    converged = True
    k_val = k
    k0_val = k
    log_likelihood = -model.inertia_
    n_samples, _ = MD.shape
    aic = -2 * log_likelihood + 2 * k
    bic = -2 * log_likelihood + np.log(n_samples) * k
    labels = model.labels_
    counts = np.bincount(labels)
    probs = counts / float(counts.sum())
    class_entropy = entropy(probs)
    icl = bic - class_entropy

    MD_m28.append((iter_val, converged, k_val, k0_val, log_likelihood, aic, bic, icl))
MD_m28 = pd.DataFrame(MD_m28, columns=['iter', 'converged', 'k', 'k0', 'logLik', 'AIC', 'BIC', 'ICL'])

print(MD_m28)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
  warnings.warn(
   iter  converged  k  k0      logLik          AIC          BIC          ICL
0     6       True  2   2 -2145.503727  4295.007454  4305.570225  4304.904547
1     6       True  3   3 -1818.717659  3643.435318  3659.279475  3658.209672
2     9       True  4   4 -1604.107292  3216.214583  3237.340126  3235.989403
3    11       True  5   5 -1434.610417  2879.220835  2905.627763  2904.035975
4     8       True  6   6 -1331.652440  2675.304880  2706.993194  2705.228430
5     6       True  7   7 -1248.417887  2510.835774  2547.805474  2545.884829
6     9       True  8   8 -1182.100019  2380.200037  2422.451123  2420.433939
```

```python
num_segments = MD_m28["k"]
AIC_values = MD_m28["AIC"]
BIC_values = MD_m28["BIC"]
ICL_values = MD_m28["ICL"]

plt.plot(num_segments, AIC_values, marker='o', label='AIC')
plt.plot(num_segments, BIC_values, marker='o', label='BIC')
plt.plot(num_segments, ICL_values, marker='o', label='ICL')

plt.xlabel('Number of Segments')
plt.ylabel('Value of Information Criteria')
plt.title('Information Criteria (AIC, BIC, ICL)')
plt.legend()
plt.grid(True)

plt.show()
```
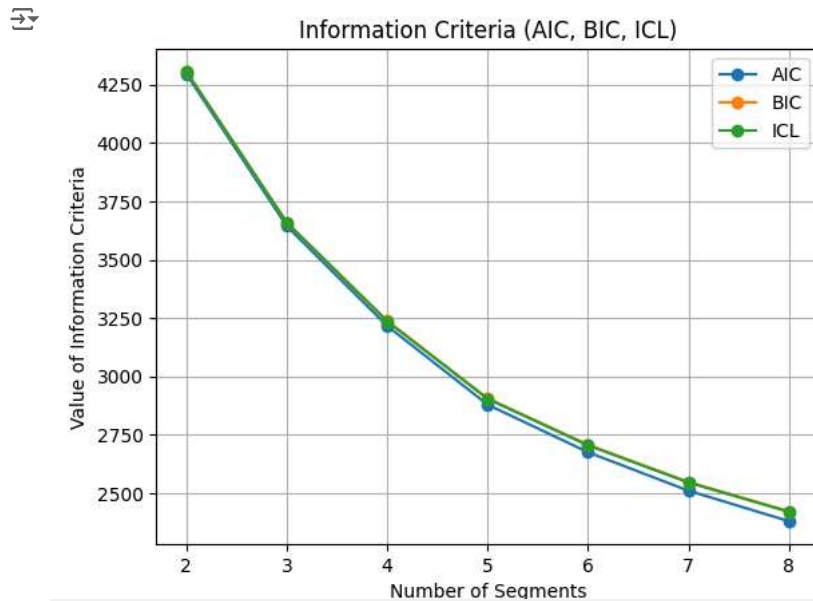
```python
from sklearn.mixture import GaussianMixture
k = 4
kmeans = KMeans(n_clusters=k, random_state=1234)
kmeans.fit(MD)
kmeans_clusters = kmeans.predict(MD)

gmm = GaussianMixture(n_components=k, random_state=1234)
gmm.fit(MD)
gmm_clusters = gmm.predict(MD)

results = pd.DataFrame({'kmeans': kmeans_clusters, 'mixture': gmm_clusters})

MD_m4 = MD[results['mixture'] == 3]

k4_m4 = KMeans(n_clusters=k, random_state=1234)
k4_m4.fit(MD_m4)
k4_m4_clusters = k4_m4.predict(MD_m4)

results_m4 = pd.DataFrame({'kmeans': k4_m4_clusters, 'mixture': 3})

print(pd.crosstab(results['kmeans'], results['mixture']))
print(pd.crosstab(results['kmeans'], results_m4['kmeans']))
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 1€
  warnings.warn(
mixture   0    1    2    3
kmeans
0        20  302    0  191
1         2  116   59  189
2        90   30   14  108
3        17  150    1  164
kmeans    0    1    2    3
kmeans
0        76   52   73   32
1        40   39   51   27
2        27   30   31   11
3        43   35   51   34
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 1€
  warnings.warn(
```

```python
like_mapping = {
    'I HATE IT!-5': -5,
    '-4': -4,
    '-3': -3,
    '-2': -2,
    '-1': -1,
    '0': 0,
    '1': 1,
    '2': 2,
    '3': 3,
    '4': 4,
    'I LOVE IT!+5': 5
}

data['Like.n'] = data['Like'].map(like_mapping)


like_n_counts = data['Like.n'].value_counts()


print(like_n_counts)
```

```
Like.n
 0.0    169
-3.0     73
-4.0     71
-2.0     59
-1.0     58
Name: count, dtype: int64
```

```python
from patsy import dmatrices

independent_vars = data.columns[0:11]

formula_str = ' + '.join(independent_vars)

formula_str = 'Like ~ ' + formula_str


f = dmatrices(formula_str, data=data)[1]

print(f)
```

```
[[1. 0. 1. ... 1. 0. 0.]
 [1. 1. 1. ... 1. 0. 0.]
 [1. 0. 1. ... 1. 1. 0.]
 ...
 [1. 1. 1. ... 1. 0. 0.]
 [1. 1. 1. ... 0. 1. 0.]
 [1. 0. 1. ... 1. 0. 1.]]
```

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

kmeans = MD_km28['4']

labels = kmeans.labels_

MD_mean = MD.groupby(labels).mean()

fig, axs = plt.subplots(2, 2, figsize=(10, 6))
axs[0, 0].barh(range(MD_mean.shape[1]), MD_mean.iloc[0])
axs[0, 0].set_title('Component 1')
axs[0, 1].barh(range(MD_mean.shape[1]), MD_mean.iloc[1])
axs[0, 1].set_title('Component 2')
axs[1, 0].barh(range(MD_mean.shape[1]), MD_mean.iloc[2])
axs[1, 0].set_title('Component 3')
axs[1, 1].barh(range(MD_mean.shape[1]), MD_mean.iloc[3])
axs[1, 1].set_title('Component 4')

for ax in axs.flat:
    ax.set(ylabel='Variable', xlabel='Proportion')
    ax.set_yticks(range(MD_mean.shape[1]))
    ax.set_yticklabels(MD.columns)

for ax in axs.flat:
    ax.label_outer()

fig.suptitle('Segment Profiles')

fig.tight_layout()

plt.show()
```

## Segment Profiles



```
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt

kmeans = KMeans(n_clusters=4)
kmeans.fit(MD)

pca = PCA(n_components=2)
MD_pca = pca.fit_transform(MD)

fig, ax = plt.subplots()

ax.scatter(MD_pca[:, 0], MD_pca[:, 1])
ax.set_xlabel('principal component 1')
ax.set_ylabel('principal component 2')
plt.show()
```
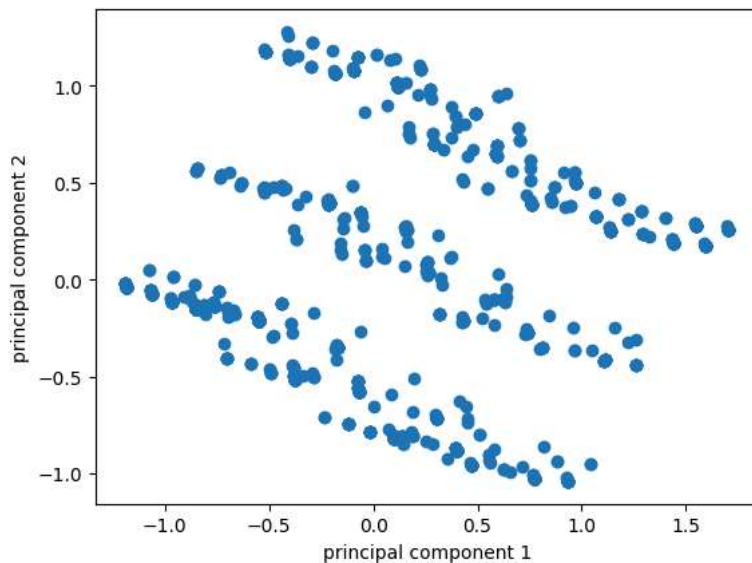
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10
    warnings.warn(

```
from statsmodels.graphics.mosaicplot import mosaic
from itertools import product
#Label encoding for categorical - Converting 11 cols with yes/no

from sklearn.preprocessing import LabelEncoder
def labelling(x):
    data1[x] = LabelEncoder().fit_transform(data1[x])
    return data1


cat = ['yummy', 'convenient', 'spicy', 'fattening', 'greasy', 'fast', 'cheap',
       'tasty', 'expensive', 'healthy', 'disgusting']


for i in cat:
    labelling(i)
data1
df_eleven = data1.loc[:,cat]
df_eleven
kmeans = KMeans(n_clusters=4, init='k-means++', random_state=0).fit(df_eleven)
data1['cluster_num'] = kmeans.labels_
crosstab =pd.crosstab(data1['cluster_num'],data1['Like'])
#Reordering cols
data1
# crosstab = crosstab[['I hate it!-5','-4','-3','-2','-1','0','1','2','3','4','I love it!+5']]
crosstab
plt.rcParams['figure.figsize'] = (7,5)
mosaic(crosstab.stack())
plt.show()
```
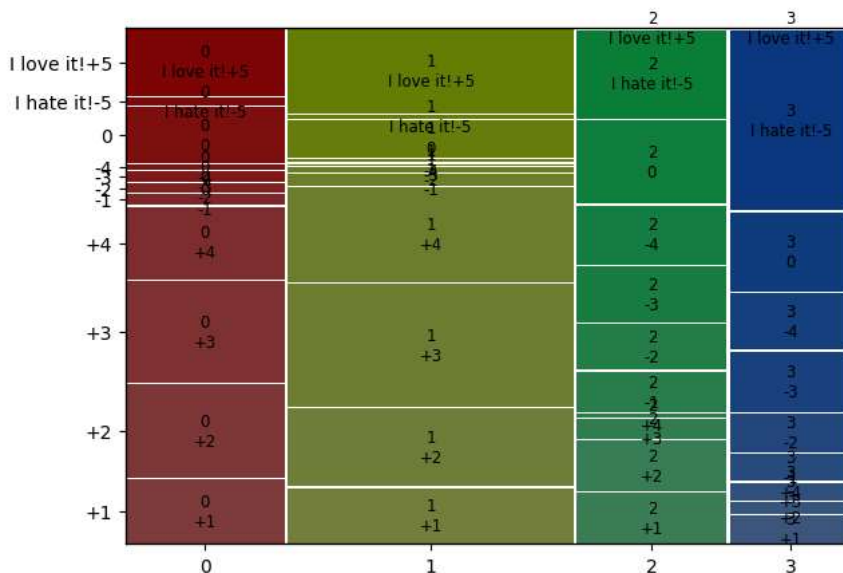
```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 18
    warnings.warn(
```
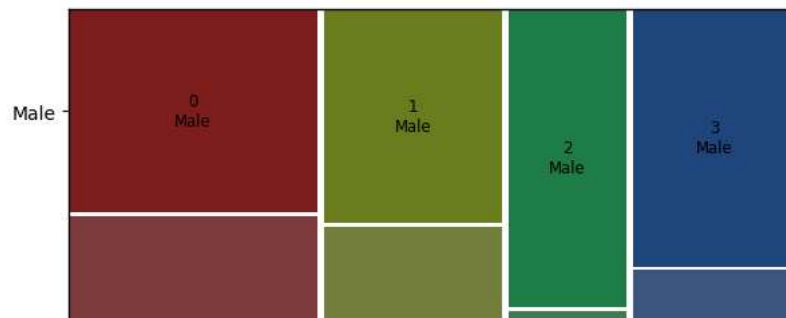


```
from statsmodels.graphics.mosaicplot import mosaic

MD_k4=MD_km28['4']
k4 = MD_k4.labels_

ct = pd.crosstab(k4, data['Gender'])
ct
mosaic(ct.stack(),gap=0.01)

plt.show()
```
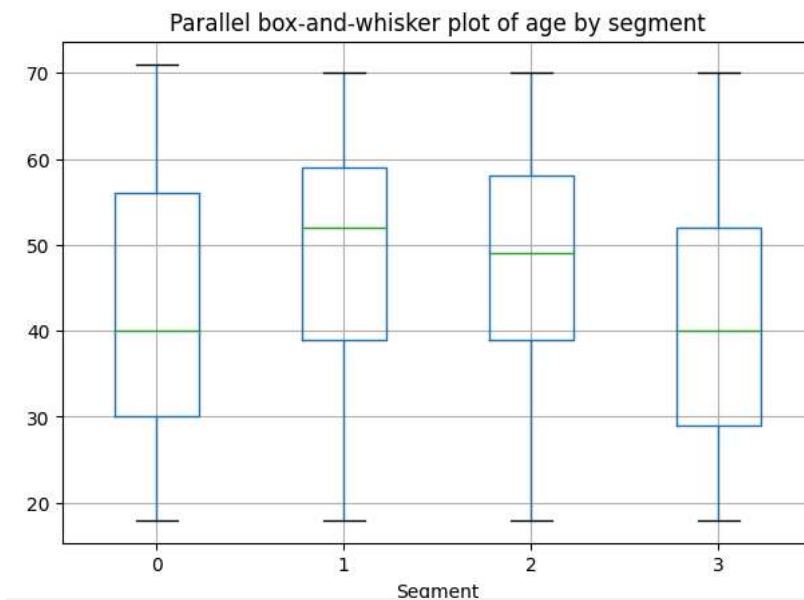
```
df = pd.DataFrame({'Segment': k4, 'Age': data['Age']})

df.boxplot(by='Segment', column='Age')
plt.title('Parallel box-and-whisker plot of age by segment')
plt.suptitle('')
plt.show()
```



```
data1['VisitFrequency'] = LabelEncoder().fit_transform(data1['VisitFrequency'])
visit = data1.groupby('cluster_num')['VisitFrequency'].mean()
visit = visit.to_frame().reset_index()
visit
```

|   | cluster_num | VisitFrequency |
|---|---|---|
| **0** | 0 | 2.547988 |
| **1** | 1 | 2.584483 |
| **2** | 2 | 2.822368 |
| **3** | 3 | 2.654472 |

```
#Like
data1['Like'] = LabelEncoder().fit_transform(data1['Like'])
Like = data1.groupby('cluster_num')['Like'].mean()
```