

Creating Customer Segments

In this project you, will analyze a dataset containing annual spending amounts for internal structure, to understand the variation in the different types of customers that a wholesale distributor interacts with.

Instructions:

- Run each code block below by pressing **Shift+Enter**, making sure to implement any steps marked with a TODO.
- Answer each question in the space provided by editing the blocks labeled "Answer:".
- When you are done, submit the completed notebook (.ipynb) with all code blocks executed, as well as a .pdf version (File > Download as).

```
In [1]: # Import libraries: NumPy, pandas, matplotlib
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# Tell iPython to include plots inline in the notebook
%matplotlib inline

# Read dataset
data = pd.read_csv("wholesale-customers.csv")
print "Dataset has {} rows, {} columns".format(*data.shape)
print data.head() # print the first 5 rows

#Exploratory Data Analysis
print "\nCorrelation Matrix: \n"
pd.set_option('precision', 2)
print data.corr()

print "\nVariance for each product: \n"
```

Dataset has 440 rows, 6 columns

	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicatessen
0	12669	9656	7561	214	2674	1338
1	7057	9810	9568	1762	3293	1776
2	6353	8808	7684	2405	3516	7844
3	13265	1196	4221	6404	507	1788
4	22615	5410	7198	3915	1777	5185

Correlation Matrix:

	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicatessen
Fresh	1.00	0.10	-0.01	0.35	-0.10	0.24
Milk	0.10	1.00	0.73	0.12	0.66	0.41
Grocery	-0.01	0.73	1.00	-0.04	0.92	0.21
Frozen	0.35	0.12	-0.04	1.00	-0.13	0.39
Detergents_Paper	-0.10	0.66	0.92	-0.13	1.00	0.07
Delicatessen	0.24	0.41	0.21	0.39	0.07	1.00

Variance for each product:

##Feature Transformation

1) In this section you will be using PCA and ICA to start to understand the structure of the data. Before doing any computations, what do you think will show up in your computations? List one or two ideas for what might show up as the first PCA dimensions, or what type of vectors will show up as ICA dimensions.

Answer: Based on some exploratory data analysis, my guess is that 'Fresh' and 'Grocery' would be a potential candidates for the first PCA dimensions since they presumably have the highest variance and hence can potentially differentiate between different customer segments better.

It is also interesting to note that 'Grocery' and 'Detergent_Papers' are highly correlated, meaning it's likely that they would be forming an independent component together under ICA analysis.

###PCA

```
In [2]: # TODO: Apply PCA with the same number of dimensions as variables in the dataset
from sklearn.decomposition import PCA
pca = PCA(n_components = 6)
pca.fit(data)

# Print the components and the amount of variance in the data contained in each dimension
print "PCA Components:"
print pca.components_

print "\nPCA Explained Variance Ratio:"
print pca.explained_variance_ratio_

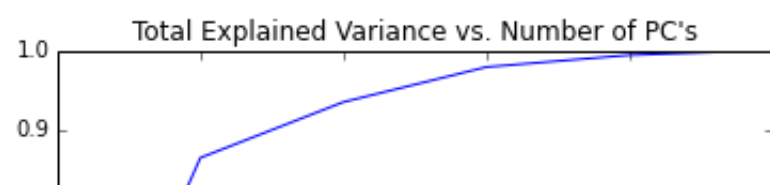
print "\nThe figure below shows total explained variance of data over the number of PC's"
x = np.arange(6)
plt.title('Total Explained Variance vs. Number of PC\'s')
plt.plot(x+1, np.cumsum(pca.explained_variance_ratio_), '-')
```

```
PCA Components:
[[-0.97653685 -0.12118407 -0.06154039 -0.15236462  0.00705417 -0.06810471]
 [-0.11061386  0.51580216  0.76460638 -0.01872345  0.36535076  0.05707921]
 [-0.17855726  0.50988675 -0.27578088  0.71420037 -0.20440987  0.28321747]
 [-0.04187648 -0.64564047  0.37546049  0.64629232  0.14938013 -0.02039579]
 [ 0.015986    0.20323566 -0.1602915   0.22018612  0.20793016 -0.91707659]
 [-0.01576316  0.03349187  0.41093894 -0.01328898 -0.87128428 -0.26541687]]
```

```
PCA Explained Variance Ratio:
[ 0.45961362  0.40517227  0.07003008  0.04402344  0.01502212  0.00613848]
```

The figure below shows total explained variance of data over the number of PC's

```
Out[2]: [<matplotlib.lines.Line2D at 0x1123e1fd0>]
```



2) How quickly does the variance drop off by dimension? If you were to use PCA on this dataset, how many dimensions would you choose for your analysis? Why?

Answer: As observed in `pca.explained_variance_ratio_`, about 46% of the data variance could be explained by the first PC, with an additional 40% variance explained by PC2. Meaning that the first two PC's can represent ~86% of the variance in data, therefore, I would pick 2 PCs for the analysis so that it's also easier to visualize the data points and clusters.

3) What do the dimensions seem to represent? How can you use this information?

Answer: Looking at the first and second rows of `pca.components_` and their most significant coefficients, it seems that the first PC mainly focuses on 'Fresh' with some smaller correlation with 'Frozen' and 'Milk', while the second PC focuses on 'Grocery', 'Milk', and 'Detergents_Paper' in that order.

It is also worth noting that:

- the 'Delicatessen' turns out to be a non-important feature if we are only using two dimensions, in that it would not provide helpful enough information for customer segmentation.
- 'Grocery' and 'Detergents_Paper' are grouped together in PC2 given that they have high correlation as shown in EDA.

PCA1 and PCA2 seem to be almost equally significant, given that their explained variance ratio is 46% and 41% respectively, while 'Delicatessen' feature seems non-influential in customer segmentation. Hence, I'd recommend the business managers to track more detailed features and measurements along PCA1 and PCA2 to capture even more variance along those PC's.

###ICA

```
In [3]: # TODO: Fit an ICA model to the data
# Note: Adjust the data to have center at the origin first!
from sklearn.decomposition import FastICA
ica = FastICA(n_components = 6, random_state=1)
origin_centered_data = data.copy()
origin_centered_data -= origin_centered_data.mean(axis=0)
origin_centered_data /= origin_centered_data.std(axis=0)
ica.fit(origin_centered_data)

# Print the independent components
ica_sources_scaled = np.apply_along_axis(lambda x: x/np.max(abs(x)), 1, ica.components_)
np.set_printoptions(precision=4, suppress=True)
print ica_sources_scaled

# Use line below to put switch back to default printoptions
#np.set_printoptions(edgeitems=3, infstr='inf', linewidth=75, nanstr='nan', precision=8, suppress=False, th
```

```
[[-0.0281  0.1269  0.8514 -0.0528 -1.         -0.1202]
 [ 1.         -0.1262 -0.1187 -0.0654  0.1962 -0.0584]
 [ 0.2022  0.0192 -0.1359 -1.         0.0489  0.3102]
 [-0.0268 -1.         0.7601  0.0243 -0.2179  0.2351]
 [-0.0958 -0.0318 -0.1118 -0.0497  0.0476  1.         ]
 [-0.0434  0.2269 -1.         -0.0324  0.0663  0.0674]]
```

4) For each vector in the ICA decomposition, write a sentence or two explaining what sort of object or property it corresponds to. What could these components be used for?

Answer:

Vector 1: mainly corresponds to features 1 and 5 ('Fresh' and 'Detergents_Paper'). Customer type: fresh and healthy buyer

Vector 2: mainly corresponds to feature 3 ('Grocery'). Customer type: regular grocery shopper

Vector 3: mainly corresponds to features 3, 2, and 6 ('Milk', 'Grocery', and 'Delicatessen'). Customer type: family feeders

Vector 4: mainly corresponds to feature 4 ('Delicatessen'). Customer type: on-the-go

Vector 5: mainly corresponds to features 5 and 3 ('Detergents_Paper' and 'Grocery'). Customer type: here we go again

Vector 6: mainly corresponds to features 4 and 6 ('Frozen' and 'Delicatessen') Customer type: Bulk food shopper

##Clustering

In this section you will choose either K Means clustering or Gaussian Mixed Models clustering, which implements expectation-maximization. Then you will sample elements from the clusters to understand their significance.

###Choose a Cluster Type

5) What are the advantages of using K Means clustering or Gaussian Mixture Models?

Answer: Advantages of K Means clustering:

- hard assignment
- quickly converges to local optima
- does not make any assumptions on the Gaussian-ness of underlying data

Advantages of GMM:

- soft assignment. It doesn't have to 100% choose one or another cluster. It works on boundaries, and assigns intermediate probabilities around 0.5 to points in the middle of decision boundary, which is nice because we get probabilities and not yes / no answers because there is always some underlying chance the the point could be from the other gaussian (cluster).
- it will not bias the clusters to have specific structures
- it gives you significantly more structural information (e.g. being able to measure how 'wide' each cluster is)

For this analysis, I would pick K Means clustering due to the hard assignment for easier visualization and better interpretation of data.

6) Below is some starter code to help you visualize some cluster data. The visualization is based on [this demo \(http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_digits.html\)](http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_digits.html) from the sklearn documentation.

```
In [4]: # Import clustering modules
from sklearn.cluster import KMeans
from sklearn.mixture import GMM
```

```
In [5]: # TODO: First we reduce the data to two dimensions using PCA to capture variation
pca_2 = PCA(n_components = 2)
pca_2.fit(data)
reduced_data = pca_2.fit_transform(data)
# print pca.components_[0]
# print data.loc[0]
# print type(ndarray(data.loc[0]))
# print type(pca.components_[0])
print reduced_data[:10] # print upto 10 elements
```

```
[[ -650.0221  1585.5191]
 [ 4426.805  4042.4515]
 [ 4841.9987  2578.7622]
 [ -990.3464 -6279.806 ]
 [-10657.9987 -2159.7258]
 [ 2765.9616 -959.8707]
 [ 715.5509 -2013.0023]
 [ 4474.5837  1429.497 ]
 [ 6712.0954 -2205.9092]
 [ 4823.6344 13480.5592]]
```

```
In [6]: # TODO: Implement your clustering algorithm here, and fit it to the reduced data for visualization
# The visualizer below assumes your clustering object is named 'clusters'
```

```
# Creating different number of clusters (2,3, and 4) for comparison on what would be a good number
# of clusters
clusters_2 = KMeans(n_clusters=2, random_state=1) #2 clusters
clusters_3 = KMeans(n_clusters=3, random_state=1) #3 clusters
clusters_4 = KMeans(n_clusters=4, random_state=1) #4 clusters
```

```
clusters_2.fit(reduced_data)
clusters_3.fit(reduced_data)
clusters_4.fit(reduced_data)
```

```
# print clusters.labels_
print "2 Clusters: \n {} \n".format(clusters_2)
print "3 Clusters: \n {} \n".format(clusters_3)
print "4 Clusters: \n {} \n".format(clusters_4)
```

2 Clusters:

```
KMeans(copy_x=True, init='k-means++', max_iter=300, n_clusters=2, n_init=10,
      n_jobs=1, precompute_distances='auto', random_state=1, tol=0.0001,
      verbose=0)
```

3 Clusters:

```
KMeans(copy_x=True, init='k-means++', max_iter=300, n_clusters=3, n_init=10,
      n_jobs=1, precompute_distances='auto', random_state=1, tol=0.0001,
      verbose=0)
```

4 Clusters:

```
KMeans(copy_x=True, init='k-means++', max_iter=300, n_clusters=4, n_init=10,
      n_jobs=1, precompute_distances='auto', random_state=1, tol=0.0001,
      verbose=0)
```

```
In [7]: # Plot the decision boundary by building a mesh grid to populate a graph.
x_min, x_max = reduced_data[:, 0].min() - 1, reduced_data[:, 0].max() + 1
y_min, y_max = reduced_data[:, 1].min() - 1, reduced_data[:, 1].max() + 1
hx = (x_max-x_min)/1000.
hy = (y_max-y_min)/1000.
xx, yy = np.meshgrid(np.arange(x_min, x_max, hx), np.arange(y_min, y_max, hy))

# Obtain labels for each point in mesh. Use last trained model.
# Z = clusters.predict(np.c_[xx.ravel(), yy.ravel()])

Z_2 = clusters_2.predict(np.c_[xx.ravel(), yy.ravel()])
Z_3 = clusters_3.predict(np.c_[xx.ravel(), yy.ravel()])
Z_4 = clusters_4.predict(np.c_[xx.ravel(), yy.ravel()])
```

```
In [8]: # TODO: Find the centroids for KMeans or the cluster means for GMM
```

```
centroids_2 = clusters_2.cluster_centers_  
centroids_3 = clusters_3.cluster_centers_  
centroids_4 = clusters_4.cluster_centers_  
  
print "2 Cluster centroids:\n {} \n".format(centroids_2)  
print "3 Cluster centroids:\n {} \n".format(centroids_3)  
print "4 Cluster centroids:\n {} \n".format(centroids_4)
```

```
2 Cluster centroids:  
[[-24088.3328  1218.1794]  
 [ 4175.311   -211.1511]]
```

```
3 Cluster centroids:  
[[ 4114.9538 -3081.0322]  
 [ 1339.4462 25546.4907]  
 [-24220.7119 -4364.4556]]
```

```
4 Cluster centroids:  
[[-23984.5576 -4910.9367]  
 [ 3496.7882 -5024.8081]  
 [-14526.8761 50607.6414]  
 [ 6166.1731 11736.8138]]
```

```

In [9]: # Put the result into a color plot
def plot_clusters(Z, clusters, centroids):
    Z = Z.reshape(xx.shape)
    plt.figure(1)
    plt.clf()
    plt.imshow(Z, interpolation='nearest',
               extent=(xx.min(), xx.max(), yy.min(), yy.max()),
               cmap=plt.cm.Paired,
               aspect='auto', origin='lower')

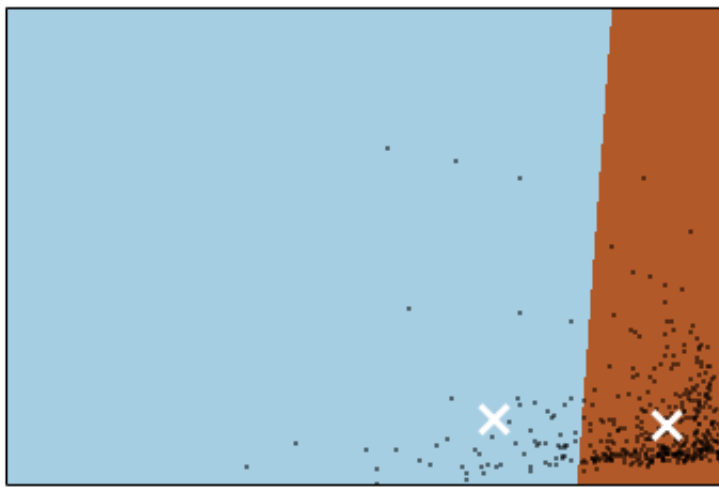
    plt.plot(reduced_data[:, 0], reduced_data[:, 1], 'k.', markersize=2)
    plt.scatter(centroids[:, 0], centroids[:, 1],
               marker='x', s=169, linewidths=3,
               color='w', zorder=10)
    plt.title('Clustering on the wholesale grocery dataset (PCA-reduced data)\n'
             'Centroids are marked with white cross')
    plt.xlim(x_min, x_max)
    plt.ylim(y_min, y_max)
    plt.xticks(())
    plt.yticks(())
    plt.show()

    print "x_min: %d, x_max: %d, y_min: %d, y_max: %d" % (x_min, x_max, y_min, y_max)
    print data[:5]
    print reduced_data[:5]
    print clusters.labels_[:5]

plot_clusters(Z_2, clusters_2, centroids_2)
plot_clusters(Z_3, clusters_3, centroids_3)
plot_clusters(Z_4, clusters_4, centroids_4)

```

Clustering on the wholesale grocery dataset (PCA-reduced data)
Centroids are marked with white cross

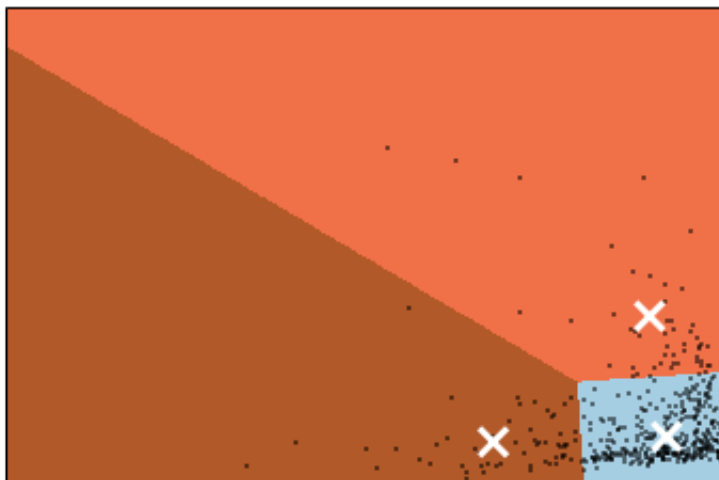


```

x_min: -103864, x_max: 13020, y_min: -14004, y_max: 99227
Fresh Milk Grocery Frozen Detergents_Paper Delicatessen
0 12669 9656 7561 214 2674 1338
1 7057 9810 9568 1762 3293 1776
2 6353 8808 7684 2405 3516 7844
3 13265 1196 4221 6404 507 1788
4 22615 5410 7198 3915 1777 5185
[[-650.0221 1585.5191]
 [ 4426.805 4042.4515]
 [ 4841.9987 2578.7622]
 [-990.3464 -6279.806 ]
 [-10657.9987 -2159.7258]]
[1 1 1 1 0]

```

Clustering on the wholesale grocery dataset (PCA-reduced data)
Centroids are marked with white cross



```

x_min: -103864, x_max: 13020, y_min: -14004, y_max: 99227
Fresh Milk Grocery Frozen Detergents_Paper Delicatessen
0 12669 9656 7561 214 2674 1338
1 7057 9810 9568 1762 3293 1776
2 6353 8808 7684 2405 3516 7844
3 13265 1196 4221 6404 507 1788
4 22615 5410 7198 3915 1777 5185
[[ -650.0221 1585.5191]
 [ 4426.805 4042.4515]
 [ 4841.9987 2578.7622]
 [ -990.3464 -6279.806 ]
 [-10657.9987 -2159.7258]]
[0 0 0 0 2]

```

Clustering on the wholesale grocery dataset (PCA-reduced data)
Centroids are marked with white cross



```

x_min: -103864, x_max: 13020, y_min: -14004, y_max: 99227
Fresh Milk Grocery Frozen Detergents_Paper Delicatessen
0 12669 9656 7561 214 2674 1338
1 7057 9810 9568 1762 3293 1776
2 6353 8808 7684 2405 3516 7844
3 13265 1196 4221 6404 507 1788
4 22615 5410 7198 3915 1777 5185
[[ -650.0221 1585.5191]
 [ 4426.805 4042.4515]
 [ 4841.9987 2578.7622]
 [ -990.3464 -6279.806 ]
 [-10657.9987 -2159.7258]]
[1 3 1 1 0]

```

7) What are the central objects in each cluster? Describe them as customers.

Answer:

Using n=3 for number of clusters, the resulting centroids in the reduced space will be

Centroid 1: [4114.95375632 -3081.03219608]

Centroid 2: [1339.44615464 25546.49074629]

Centroid 3: [-24220.71188261 -4364.45560022]

With PC1 focusing mainly on 'Fresh' (in reverse, because the first coefficient is negative), and PC2 positively correlated with 'Grocery', 'Milk', and 'Detergent_Papers' in that order.

Cluster 1 (blue), having the biggest population, can be described as customers who, relatively speaking, have less 'Fresh' and also less 'Grocery, Milk, and Detergents_Paper' in their annual basket.

Cluster 2 (red), can be described as customers who, relatively speaking, have less 'Fresh' but more of 'Grocery, Milk, and Detergents_Paper' in their annual basket.

Cluster 3 (brown), can be described as customers who, relatively speaking, have much more 'Fresh' but less of 'Grocery, Milk, and Detergents_Paper' in their annual basket.

As a side note, for better visualization and analysis of the clusters, it would be worthwhile to consider additional interactive 3D plots and Silhouette plots, which is out of scope of this project.

###Conclusions

8) Which of these techniques did you feel gave you the most insight into the data?

Answer: PCA was very helpful in providing insights into customer segmentation. Specifically:

- it suggested that ~86% of the data variance could be explained by the first two principal components, hence the data could be reduced into 2 dimensions for further analysis and visualization.

- it helped in understanding which features those two principal components comprise of and most correlate with, hence giving insights into what set of features the future customer research and data collection should be focused on.

ICA was also helpful in understanding the potential set of independent customer sets/types, that in aggregate form the overall annual sales. As such, future research could be guided to focus on these specific and independent customer types with specific market baskets to understand their needs separately.

9) How would you use that technique to help the company design new experiments?

Answer:

With a better understanding of customer clusters and different customer types from PCA and ICA analysis, the company perform A/B testing (e.g. effect of morning vs. evening delivery on different customer segments) on new initiatives and gauge the effect of those initiatives on different customer types. This will help us understand which customers would be most affected by the change and how, and separate efforts can be taken to address the customer needs separately.

Another thing worth consideration is collect and track a higher number of features on each customer to get higher variance and differentiation, and in turn better segmentation and needs assessment.

10) How would you use that data to help you predict future customer needs?

Answer: By knowing the different customer clusters, each potentially with different set of needs, the business can better tailor its services to individual customer segments without necessarily lumping them all to the same group and regarding their needs as one.

In addition, supervised learning techniques such as classification could be leveraged to assign new customers to the already known clusters to better understand their needs and provide better targeted the services to them.