

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 [11]: # 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
```

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

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: Just looking at the variation by considering the min and max of the table, we could definitely observe a dimension with capturing high variation of "Fresh", "Milk", and "Grocery". The variance along the principal components play a large role in selecting dimensions. The PCA dimensions tries to mutually orthogonal, maximize variance and has an ordered set of features.

ICA tries to produce components that maximizes the separation between separate classes. From the notes, we see Kurtosis which is the spikiness or peakedness is used to separate the classes. The classes in our case could be the different customers sets we have - based on high volume or low volume spend. Therefore a simple frequency distribution of annual spend 'range' over categories could show any spike or kurtosis we are looking for.

For low volume or family run shops, I would expect to a spike selling fresh and groceries (Fresh, Milk and Grocery) and specially prepared foods, ("Delicatessen").

For high volume customers like Sams club, I would expect them to sell more Frozen and Detergents and Paper products. I don't usually think about Milk and Fresh and Delicatessen from high volume customers.

So, based on the spikes or kurtosis I am expecting from above theories on those two different classes, I would expect two dimensions that maximizes the separation between the class of the customers who are the source of the distribution of the data we have.

PCA

```

In [12]: # 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.explained_variance_ratio_

[[-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]]
[ 0.45961362  0.40517227  0.07003008  0.04402344  0.01502212  0.00613848]

```

In [13]:



```

#!/usr/bin/env python
# coding: utf8
# GistID: 8785541

from __future__ import division

import numpy as np
import matplotlib
import matplotlib.pyplot as plt

# Make a random array and then make it positive-definite
A = np.random.randn(6, 9)
A = np.asmatrix(A) * np.asmatrix(A.T)
U, S, V = np.linalg.svd(A)
eigvals = S**2 / np.cumsum(S)[-1]
eigvals2 = S**2 / np.sum(S)
assert (eigvals == eigvals2).all()

eigvals = pca.explained_variance_ratio_
tot = sum(eigvals)
var_exp = [(i / tot)*100 for i in sorted(eigvals, reverse=True)]
cum_var_exp = np.cumsum(var_exp)

fig = plt.figure(figsize=(8,5))
sing_vals = np.arange(len(eigvals)) + 1
plt.plot(sing_vals, eigvals, 'ro-', linewidth=2)
plt.title('Scree Plot')
plt.xlabel('Principal Component')
plt.ylabel('Eigenvalue')
#I don't like the default legend so I typically make mine like below, e.g.
#with smaller fonts and a bit transparent so I do not cover up data, and make
#it moveable by the viewer in case upper-right is a bad place for it
leg = plt.legend(['Eigenvalues from SVD'], loc='best', borderpad=0.3,
                 shadow=False, prop=matplotlib.font_manager.FontProperties(size='small'),
                 markerscale=0.4)
leg.get_frame().set_alpha(0.4)
leg.draggable(state=True)
plt.show()

#---
fig = plt.figure(figsize=(8,5))
sing_vals = np.arange(len(cum_var_exp)) + 1
plt.plot(sing_vals, cum_var_exp, 'ro-', linewidth=2)
plt.title('Scree Plot')
plt.xlabel('Principal Component')
plt.ylabel('Cumulative Sum of variation')
#I don't like the default legend so I typically make mine like below, e.g.
#with smaller fonts and a bit transparent so I do not cover up data

```

```

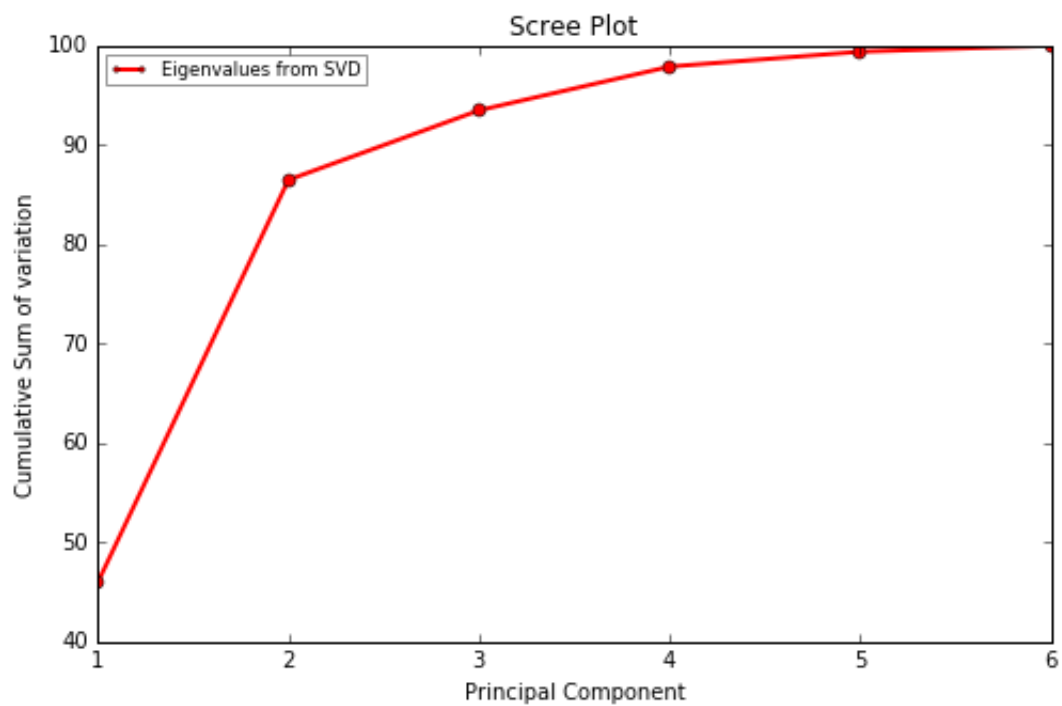
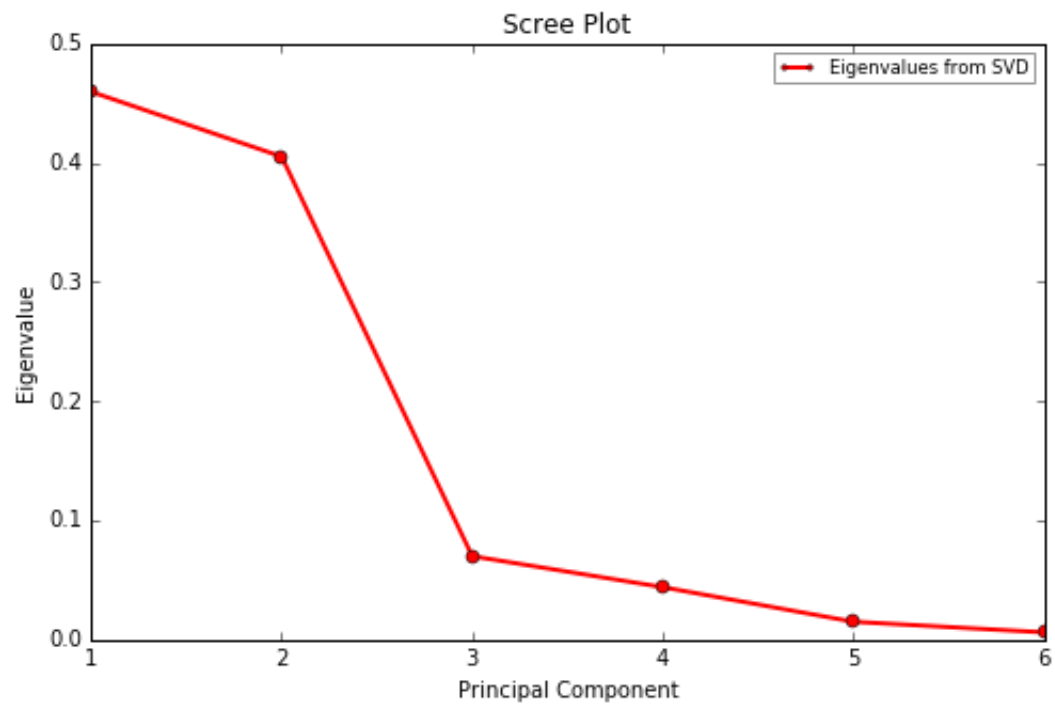
a, and make
#it moveable by the viewer in case upper-right is a bad place for i
t
leg = plt.legend(['Eigenvalues from SVD'], loc='best', borderpad=0.
3,
                shadow=False, prop=matplotlib.font_manager.FontPro
perties(size='small'),
                markerscale=0.4)
leg.get_frame().set_alpha(0.4)
leg.draggable(state=True)
plt.show()

def main():
    """Run main."""

    return 0

if __name__ == '__main__':
    main()

```



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: I plotted both the eigen values and their explained variation of each of the PCA dimension using "scree plot" concept as explained in <http://support.minitab.com/en-us/minitab/17/topic-library/modeling-statistics/multivariate/principal-components-and-factor-analysis/what-is-a-scree-plot/> (<http://support.minitab.com/en-us/minitab/17/topic-library/modeling-statistics/multivariate/principal-components-and-factor-analysis/what-is-a-scree-plot/>) and <https://plot.ly/ipython-notebooks/principal-component-analysis/#PCA-Vs.-LDA> (<https://plot.ly/ipython-notebooks/principal-component-analysis/#PCA-Vs.-LDA>)

These are my observations.

There is a steep drop in variation after two dimensions. The first two contribute to nearly 85% of the variation. However, based on the scree plot criteria, we will retain all the dimensions before the point that starts the flat line trend i.e. no more increase in the variation.

Looking at the second graph that shows the cumulative variation by dimension, we see the flat trend starts to happen at the fifth vector. Based on the guideline provided by scree, I will use four dimensions for analysis.

I also tried to see if I can use the Kaiser method as explained in <http://support.minitab.com/en-us/minitab/17/topic-library/modeling-statistics/multivariate/principal-components-and-factor-analysis/number-of-principal-components/> (<http://support.minitab.com/en-us/minitab/17/topic-library/modeling-statistics/multivariate/principal-components-and-factor-analysis/number-of-principal-components/>) but normalization of eigenvalues in scikit eigenvector cannot use the rule to throw out anything having value greater than 1.0.

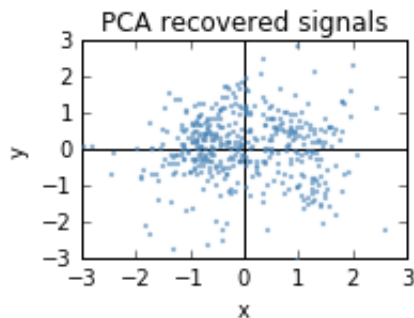
```
In [14]: def plot_samples(S, axis_list=None):
          plt.scatter(S[:, 0], S[:, 1], s=2, marker='o', zorder=10,
                      color='steelblue', alpha=0.5)
          if axis_list is not None:
              colors = ['orange', 'red']
              for color, axis in zip(colors, axis_list):
                  axis /= axis.std()
                  x_axis, y_axis = axis
                  # Trick to get legend to work
                  plt.plot(0.1 * x_axis, 0.1 * y_axis, linewidth=2, color
                           =color)
                  plt.quiver(0, 0, x_axis, y_axis, zorder=11, width=0.01,
                             scale=6,
                             color=color)

          plt.hlines(0, -3, 3)
          plt.vlines(0, -3, 3)
          plt.xlim(-3, 3)
          plt.ylim(-3, 3)
          plt.xlabel('x')
          plt.ylabel('y')
```


In [15]: *#Plotting PCA variance*

```
#data /= data.std(axis=0)  
S_pca_ = pca.fit(np.log(data)).transform(np.log(data))  
plt.subplot(2, 2, 4)  
plot_samples(S_pca_ / np.std(S_pca_, axis=0))  
plt.title('PCA recovered signals')
```

Out[15]: <matplotlib.text.Text at 0x111b81d10>



In [16]: **import pandas as pd**
from sklearn.decomposition import PCA

```
def biplot(df):  
    # Fit on 2 components  
    pca = PCA(n_components=2, whiten=True).fit(df)  
  
    # Plot transformed/projected data  
    ax = pd.DataFrame(  
        pca.transform(df),  
        columns=['PC1', 'PC2']  
    ).plot(kind='scatter', x='PC1', y='PC2', figsize=(10, 8), s=0.  
    8)  
  
    # Plot arrows and labels  
    for i, (pc1, pc2) in enumerate(zip(pca.components_[0], pca.comp  
    onents_[1])):  
        ax.arrow(0, 0, pc1, pc2, width=0.001, fc='orange', ec='oran  
    ge')  
        ax.annotate(df.columns[i], (pc1, pc2), size=12)  
  
    return ax
```

In [113]: *#ax = biplot(data)*
Play around with the ranges for scaling the plot
#ax.set_xlim([-1.5, 1])
#ax.set_ylim([-1, 1])

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

Answer:

The dimensions represent how best they could maximize the variation along the orthogonal axis and therefore minimize the loss of information. In case of 'curse of dimensionality', we could employ PCA to do a 'dimensionality reduction' while preserving information as best as possible. As a result, the computation becomes faster and efficient.

Looking at the matrix below, Dimension 5 - Detergent Paper and Frozen. - High Volume

Dimension 1 - Fresh

```
[[-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]]
```

Working backwards, I theorize that a high volume customer like Sams Club or Costco would sell (and therefore buy) more Frozen, Detergents and Paper. So, a high variation of these components could indicate this type of 'not-so' mom and pop shops. With the same thought, a mom and pop, gas station or convenience store for example, would sell (and therefore buy) more Milk, Delicatessen.

Using the above theory,

I could say the following fall to 'High volume customers'

- Dimension No. 6 has more variation of Detergents and Paper followed by Frozen and Delicatessen
- Dimension No. 2 has more variation of Frozen followed by Milk and Detergent and Paper

The following fall to 'Low volume customers'

- Dimension No. 5 has more variation of Delicatessen followed by almost same variation of Milk, Frozen, Grocery and Detergents and Paper
- Dimension No. 1 has more variation of Fresh (very high) followed by Milk and Grocery. So, it could indicate a small grocery or produce store or a mom and pop store which doesn't focus a lot on frozen .
- Dimension No. 3 and No. 4 has more variation of Groceries and Milk and less of Fresh but less selling of Detergents and Delicatessen and Fresh - could indicate a grocery store.

ICA

```
In [25]: # 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)
data /= data.std(axis=0)
ica.fit(data)

# Print the independent components
print ica.components_
```

```
[[-0.05032637  0.00666892  0.00764541  0.00325243 -0.01166535  0.0
0262688]
 [ 0.00494119  0.00162137  0.00559176  0.00235742 -0.00240617 -0.0
5095765]
 [-0.00326085  0.01831886  0.11010883 -0.0072403  -0.13372387 -0.0
1605528]
 [-0.00267131  0.01297231 -0.06806923 -0.00151385  0.0121041  0.0
0514807]
 [-0.01087121 -0.00106729  0.00735949  0.05412038 -0.00258322 -0.0
169443 ]
 [-0.00207458 -0.07249322  0.05596866  0.00170598 -0.01683229  0.0
1699047]]
```

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:

ica attempts to decompose a multivariate signal into independent non-gaussian signals even though it cannot identify the actual number of source signals. [taken from wikipedia]

what we are seeing is the unmixing matrix which decomposes the linearly mixed data into a sum of temporally independent and spatially fixed components. The columns give the relative projection strengths of the respective features at each of the component or dimension.

By setting number of components = 6, i am parameterizing the fast ica to give me six independent dimensions corresponding to six independent customer types based on their spending (or buying) patterns.

directionality and magnitude - since we are trying to source out the non-gaussian source signals, we can consider positive direction as a projection of super-gaussian and negative direction as a projection of sub-gaussian. zero kurtosis would indicate that there are gaussian signals. in other words, super-gaussian signals have pdfs that are more peaky than that of gaussian signals and the sub-gaussian signals have pdfs that are less peaky than that of a gaussian signals. [referring to mit.it/informatica/pdf/info878.pdf]

with the above theory becoming clear, looking at the first component, [-0.00369099 0.01138811 0.12559625 -0.00671944 -0.13374133 -0.01575071]

it is evident that we have weight of grocery and milk projected more while detergents and paper are projected less. it seems the source could be a grocery store and definitely not a very high volume customer stocking up on detergents and paper.

looking at the fourth component, [0.00434314 0.05466034 -0.02307661 -0.0003578 0.01884458

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: K-means is a. computationally efficient as it uses hard clustering - that is, our case, each customer will be assigned to only one cluster based on the eucladian distance. b. runs usually fast - converges using local minimum c. cluster scatter is calculated as a measure of eucladian distance

GMM is

- a. it is the fastest algorithm for learning mixture models.
- b. It can also draw confidence ellipsoids for multivariate models,
- c. and compute the Bayesian Information Criterion to assess the number of clusters in the data.
- d. uses soft clustering - a probability score - for assigning objects to clusters. Objects, customers in our case, could be assigned to multiple clusters of buyer or spending types. So, computationally, GMM is doing more than a straight forward eucladian based k-means clustering.

Confidence Ellipsoid - defines the region where certain percent of all the samples that could be drawn from the underlying Gaussian distribution.

Bayesian Information Criterion - BIC - is a criterion for model selection among a finite set of models; the model with the lowest BIC is preferred. BIC limits the parameters used in the model by penalizing while trying to increase the likelihood.

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

```
In [87]: # TODO: First we reduce the data to two dimensions using PCA to capture variation  
pca = PCA(n_components=2);  
reduced_data = pca.fit_transform(data)  
print reduced_data[:10] # print upto 10 elements
```

```
[[-0.19307077  0.30475306]  
 [-0.43392596  0.32803921]  
 [-0.81022096 -0.81416893]  
 [ 0.7777625  -0.65201155]  
 [-0.16609819 -1.26998809]  
 [ 0.15599237  0.29480541]  
 [ 0.33490718  0.52440632]  
 [-0.14042659  0.23073005]  
 [ 0.51673134  0.65861312]  
 [-1.59029884  0.74016879]]
```

```
In [100]: # 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'
```

```
kmeans_clusters = KMeans(init='k-means++', n_clusters=4, n_init=10)  
#?  
kmeans_clusters.fit(reduced_data)  
print kmeans_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=None, tol=  
0.0001,  
      verbose=0)
```

```

In [101]: lowest_bic = np.infty
bic = []
n_components_range = range(1, 7)
cv_types = ['spherical', 'tied', 'diag', 'full']
for cv_type in cv_types:
    for n_components in n_components_range:
        # Fit a mixture of Gaussians with EM
        gmm = GMM(n_components=n_components, covariance_type=cv_type)

        gmm.fit(reduced_data)
        bic.append(gmm.bic(reduced_data))
        if bic[-1] < lowest_bic:
            lowest_bic = bic[-1]
            best_gmm = gmm
gmm_clusters=best_gmm
print gmm_clusters

GMM(covariance_type='full', init_params='wmc', min_covar=0.001,
     n_components=5, n_init=1, n_iter=100, params='wmc', random_state
     =None,
     thresh=None, tol=0.001, verbose=0)

```

```

In [102]: # 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.
print Z

[3 3 3 ..., 1 1 1]

```

```

In [107]: # TODO: Find the centroids for KMeans or the cluster means for GMM

centroids = kmeans_clusters.cluster_centers_
#centroids = clusters.means_
print centroids

[[-1.29980866  0.73073732]
 [ 0.68390594 -0.15602615]
 [-7.04226309  0.63562646]
 [-4.50795382 -10.02928697]]

```

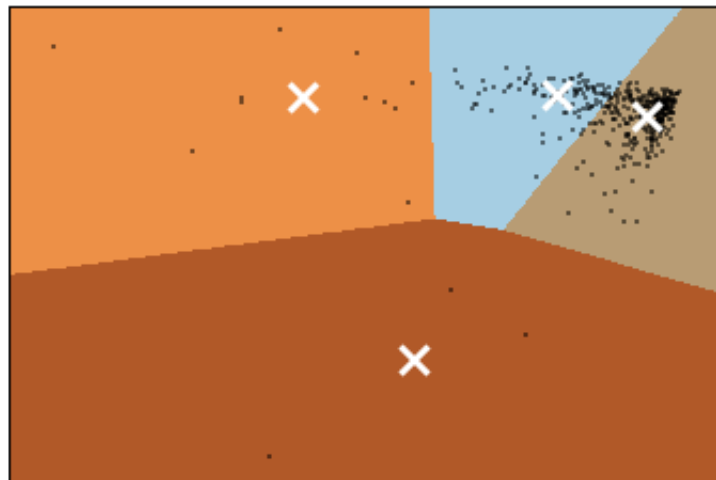
```

In [108]: # Put the result into a color plot
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()

```

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




```
In [105]: # Bringing the centroid back to Original Space.
centers = pca.inverse_transform(kmeans_clusters.cluster_centers_)
print(centers)

[[ 0.61880172  1.43313755  1.6963725   0.25263027  1.50399413  0.49550998]
 [ 1.00188312  0.42552842  0.41774756  0.69314539  0.18931713  0.44930829]
 [ 0.91527296  4.5713648   5.00883114  0.60471755  4.63026   1.97151049]
 [ 6.43694216  4.0768422   1.9827986   6.9942129   0.51779656  6.71859793]]
```

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

Answer: There are five clusters and they are visible clearly. The centroids basically denotes the 'average' customer in each cluster. What I mean by average is that the PCA fit data basically identifies a customer model (a vector) whose feature values are the mean of each feature for customers represented in that cluster. The same is applicable for the set of clusters the model produced.

```
[[ 0.61880172  1.43313755  1.6963725   0.25263027  1.50399413  0.49550998] [ 1.00188312  0.42552842
 0.41774756  0.69314539  0.18931713  0.44930829] [ 0.91527296  4.5713648   5.00883114  0.60471755
 4.63026   1.97151049] [ 6.43694216  4.0768422   1.9827986   6.9942129   0.51779656  6.71859793]]
```

Looking at the first centroid, It seems the clusters are having customers that have more spending on Groceries, Detergent and Paper followed by Milk, Fresh, and at last Delicatessen and Frozen. The cluster is trying to pull those who have the spending habits of heavy on grocery store and detergents and paper.

```
In [ ]: Conclusions
8) Which of these techniques did you feel gave you the most insight
into the data?
```

Answer: GMM gave the most insight into the data. The main advantage we have is the ability to determine the number of components in an efficient way using BIC Criterion which I also implemented and figured out the best cluster parameters to use.

It is a probabilistic model that assumes all the data points are generated from a mixture of a finite number of Gaussian distributions with unknown parameters. One can think of mixture models as generalizing k-means clustering to incorporate information about the covariance structure of the data as well as the centers of the latent Gaussians. Taken from <http://scikit-learn.org/stable/modules/mixture.html#pros> (<http://scikit-learn.org/stable/modules/mixture.html#pros>)

PCA and ICA helps us to look at more than one dimension at a time. PCA helps us cluster the customers through maximizing the variation of their spending while ICA working through maximizing independents on their spending habits. Both the methods helps us separate the signal from the noise. Even though ICA cannot exactly pin point the number of independent signals, research into maximizing the parameters of ICA and Clustering techniques could improve the overall process.

The thing that I learnt the most from this project is about the unmixing matrix in ICA and covariance matrix in PCA and helped look at the problem holistically and not let one dimension drive the analysis.

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

Answer: A/B testing allows individuals, teams, and companies to make careful changes to their user experiences while collecting data on the results.

Using the above technique, we can now have five different clusters or segments of customers. Instead of implementing a completely different delivery option on all customers, we will treat the segments independently and test the delivery changes using A/B testing method on each group.

Depends on the response for each control group, we will decide whether to rollout the delivery options to the entire segment of the customers. We will repeat the process for all five customer segments.

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

Answer: Having decided that we are going to treat the identified segments independently, we can build supervised techniques like randomForest, DecisionTree to build a model with the target being the number of items sold in each category (derived from total spend).

Once we have built the models, we can split them as train and test sets, fit a model for every category of food (milk, fresh, etc) over every single segment then use the test to predict the demand.

Once we have built models with sufficient accuracy, we can then group a set of new customers as a test set, identify what cluster they will belong in (unsupervised) and then apply the correct supervised trained model to make the prediction. The key lesson is that we will have different supervised models for each segmentation.