



## PROJECT

## Creating Customer Segments

A part of the Machine Learning Engineer Nanodegree Program

## PROJECT REVIEW

## NOTES

SHARE YOUR ACCOMPLISHMENT!  

## Requires Changes

1 SPECIFICATION REQUIRES CHANGES

## Data Exploration

Three separate samples of the data are chosen and their establishment representations are proposed based on the statistical description of the dataset.

You could also visualize each sample in terms of the underlying features:

```
samples_cpy = samples.copy()
samples_cpy[samples.shape[0]] = data.median()

plt.style.use('ggplot')
samples_cpy.plot(kind='bar')
labels = samples.index.values.tolist()
labels.append("Data Median")
plt.xticks(range(samples.shape[0]), labels)
```

With this plot the Median and Mean points show you the spending of the average customer. Note, however, that you aren't able to compare the spending between different types of spending. For example, suppose you identify a sample with similar Fresh and Milk expenditures. Because the average Milk spending is half that of Fresh, if the values are comparable the Milk spending is **twice** as significant.

Here is another way of looking at the same data. This time, the expenses have been scaled to have zero mean and a standard deviation/variance of 1.

```
import matplotlib.pyplot as plt
import seaborn as sns

samples_for_plot = samples.copy()
samples_for_plot = (samples_for_plot - data.mean())/data.std()

labels = ['Sample 1', 'Sample 2', 'Sample 3']
samples_for_plot.plot(kind='bar')
_ = plt.xticks(range(3), labels)
```

A critical concept that is lying just beneath the surface with this project are the ideas of standardization and normalization.

[https://en.wikipedia.org/wiki/Normalization\\_\(statistics\)](https://en.wikipedia.org/wiki/Normalization_(statistics))

You can start thinking about this idea while looking at the samples, by applying a simple `numpy` transformation to the data.

```
samples_sc = (samples - data.mean())/data.std()
```

Note that this is ignoring the concept of [degrees of freedom](#).

Having performed the transformation, the values of the scaled samples represent the number of standard deviations from the mean of the dataset for each value. This allows you to compare "apples" to "oranges", that is, data collected on difference scales in a meaningful way.

Make sure you are clear that this test reveals that a feature *may* be unnecessary. The removal of a feature may fail to capture non-linear interactions!

```

from sklearn.cross_validation import train_test_split
from sklearn.tree import DecisionTreeRegressor

def calculate_r_2_for_feature(data, feature):
    new_data = data.drop(feature, axis=1)

    X_train, \
    X_test, \
    y_train, \
    y_test = train_test_split(
        new_data, data[feature], test_size=0.25
    )

    regressor = DecisionTreeRegressor()
    regressor.fit(X_train, y_train)

    score = regressor.score(X_test, y_test)
    return score

def r_2_mean(data, feature, runs=200):
    return np.array([calculate_r_2_for_feature(data, feature)
                     for _ in range(200) ]).mean().round(4)

print "{0:17} {1}".format("Fresh: ", r_2_mean(data, 'Fresh'))
print "{0:17} {1}".format("Milk: ", r_2_mean(data, 'Milk'))
print "{0:17} {1}".format("Grocery: ", r_2_mean(data, 'Grocery'))
print "{0:17} {1}".format("Frozen: ", r_2_mean(data, 'Frozen'))
print "{0:17} {1}".format("Detergents_Paper: ", r_2_mean(data, 'Detergents_Paper'))
print "{0:17} {1}".format("Delicatessen: ", r_2_mean(data, 'Delicatessen'))

```

Strong correlations can be connected back to the  $R^2$  from the previous analysis.

I used this to visualize correlations.

```
corr = data.corr()
mask = np.zeros_like(corr)
mask[np.triu_indices_from(mask, 1)] = True
with sns.axes_style("white"):
    ax = sns.heatmap(corr, mask=mask, square=True, annot=True,
                      cmap='RdBu', fmt='+.3f')
    plt.xticks(rotation=45, ha='center')
```

Understanding the distribution of data is critical. A normal distribution is critical for a well behaving principal component analysis.

Feature scaling for both the data and the sample data has been properly implemented in code.

This helped me to see the distribution of data before and after scaling.

```
fig, axes = plt.subplots(2, 3)
axes = axes.flatten()
fig.set_size_inches(18, 6)
fig.suptitle('Distribution of Features')

for i, col in enumerate(data.columns):
    feature = data[col]
    sns.distplot(feature, label=col, ax=axes[i]).set(xlim=(-1000, 20000),)
    axes[i].axvline(feature.mean(),linewidth=1)
    axes[i].axvline(feature.median(),linewidth=1, color='r')
```

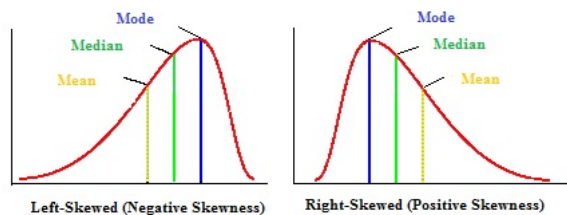
In this plot, we can actually see the distribution. We can see that most of the features are log-normal, that is, they have a distribution such that their logarithm has a normal distribution. This implies that we can achieve a normal distribution for our data simply by taking the logarithm of the data.

After preparing the log data, we can look at the distribution again:

```
fig, axes = plt.subplots(2, 3)
axes = axes.flatten()
fig.set_size_inches(18, 6)
fig.suptitle('Distribution of Features for Log Data')

for i, col in enumerate(log_data.columns):
    feature = log_data[col]
    sns.distplot(feature, label=col, ax=axes[i])
    axes[i].axvline(feature.mean(),linewidth=1)
    axes[i].axvline(feature.median(),linewidth=1, color='r')
```

From these you can think about skewness in data:



To compare the log-transformed feature distributions, you can also plot them on top of each other with a seaborn KDE plot...

```
import matplotlib.pyplot as plt
import seaborn as sns
# set plot style & color scheme
sns.set_style('ticks')
with sns.color_palette("Reds_r"):
    # plot densities of log data
    plt.figure(figsize=(8,4))
    for col in data.columns:
        sns.kdeplot(log_data[col], shade=True)
    plt.legend(loc='best')
```

Perhaps most importantly, by taking the log of the data, we have actually placed the data on similar scales. This will be critical when we perform our principal component analysis.

#### NEXT LEVEL SUGGESTION

Strictly speaking, a principal component analysis assumes that data has a mean of zero and a variance of 1. If you think about it, the PCA is completely based upon the variance and co-variance in the data. This requires that the data can be completely described using the variance and co-variance i.e. that the mean is zero. We need to have the variance of each feature set to 1 so that no feature comes to dominate the PCA. To achieve this you can scale the `log_data`.

```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
scaler.fit(log_data)

log_data = scaler.transform(log_data)
log_samples = scaler.transform(log_samples)
```

If you do this, you will need to apply the following to transform your data back in **Implementation: Data Recovery**.

Student identifies extreme outliers and discusses whether the outliers should be removed. Justification is made for any data points removed.

```
from collections import Counter
```

PCA has been properly implemented and applied to both the scaled data and scaled sample data for the two-dimensional case in code.

## 4/6

The choice of outliers is very important for the silhouette scores of your outliers. What if you ran your analysis again, this time without removing any outliers? How would the number of clusters and silhouette scores vary?

You could also visualize each segment in terms of the underlying features:

Sample points are correctly identified by customer segment, and the predicted cluster for each sample point is discussed.

Student correctly identifies how an A/B test can be performed on customers after a change in the wholesale distributor's service.

You might think about the *power* of the change:

<https://www.quora.com/When-should-A-B-testing-not-be-trusted-to-make-decisions/answer/Edwin-Chen-1>

Student discusses with justification how the clustering data can be used in a supervised learner for new predictions.

Here is an online discussion about using unsupervised learning to generate a target: <https://datascience.stackexchange.com/questions/985/can-i-use-unsupervised-learning-followed-by-supervised-learning>

Comparison is made between customer segments and customer 'Channel' data. Discussion of customer segments being identified by 'Channel' data is provided, including whether this representation is consistent with previous results.

 RESUBMIT

[↓ DOWNLOAD PROJECT](#)

Learn the **best practices** for revising and resubmitting your project.

RETURN TO PATH

Rate this review

