

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

Hello Student,

You almost get it done, keep it UP!

There is only one minor mistake you need to amend, I believe you already understand most of the concept in this project.

Data Exploration

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

Well done commenting on the different types of establishments the three customers could represent.

A prediction score for the removed feature is accurately reported. Justification is made for whether the removed feature is relevant.

Great analysis to identify that the amount of `Delicatessen` purchased is necessary to identify specific customers!

Student identifies features that are correlated and compares these features to the predicted feature. Student further discusses the data distribution for those features.

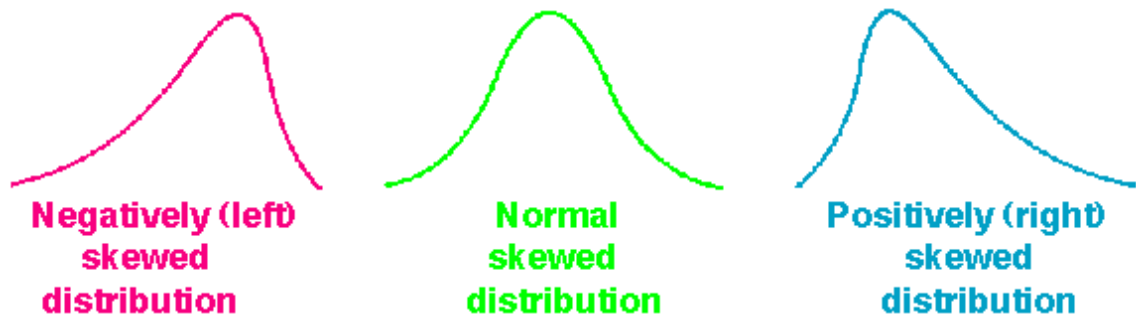
You did great identifying the features that have correlation from the dataset.

And your finding from the previous question on the `Delicatessen` feature nicely aligns with your finding from this section.

And also, well done commenting on the distribution of the dataset.

Suggestions and Comments:

- Kindly note that the distribution is skewed to the right (most points lie to the left of the graph). The graph below interprets how to quickly judge the skewness of a distribution:



- As you can see from the scatter plot, there are a number of outliers for most of the features.
- Also, the median falls below the mean, and there are a large number of data points near 0.

Data Preprocessing

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

Nice work implementing feature scaling :)

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

Great job identifying that 65, 66, 75, 128, 154 are outliers, and giving justifications on why you think they should be removed.

Feature Transformation

The total variance explained for two and four dimensions of the data from PCA is accurately reported. The first four dimensions are interpreted as a representation of customer spending with justification.

Fantastic work in identifying the amount of variance explained by the first 2 and the first 4 dimensions. Also, you did a great job in interpreting the first four dimensions as a representation of customer spending. Please see some suggestions and comments below for more information on this section:

Suggestions and Comments:

These two links helped me a lot in understanding Principal Component Analysis and what it does, hope it helps you too.

<https://onlinecourses.science.psu.edu/stat505/node/54>

<http://setosa.io/ev/principal-component-analysis/>

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

Nice job here!

Clustering

The Gaussian Mixture Model and K-Means algorithms have been compared in detail. Student's choice of algorithm is justified based on the characteristics of the algorithm and data.

Great work elaborating on [GMMs](#) and [K-means](#), and giving a solid reason as to why you chose [GMM] algorithm for this particular problem.

You might want to provide some citations and reference for your work to make it more credible.

Below are some of my comments, feedback, and suggested reading:

Gaussian Mixture Models

- You did well on giving the advantages of Gaussian Mixture Model.

SUGGESTED READING:

- If you feel as going deeper with regards to Gaussian Mixture Models, check out the following links:
 - http://home.deib.polimi.it/matteucc/Clustering/tutorial_html/mixture.html
 - <http://www.nickgillian.com/wiki/pmwiki.php/GRT/GMMClassifier>
 - <http://scikit-learn.org/stable/modules/mixture.html#gmm-classifier>

K-Means Clustering

- Your description on the advantages of K-means is very explicit

SUGGESTED READING:

- Check out these links for even more thorough explanations of K-Means Clustering:
 - <http://playwidtech.blogspot.hk/2013/02/k-means-clustering-advantages-and.html>
 - <https://sites.google.com/site/dataclusteringalgorithms/k-means-clustering-algorithm>
 - http://www.improvedoutcomes.com/docs/WebSiteDocs/Clustering/K-Means_Clustering_Overview.htm
 - <http://stats.stackexchange.com/questions/133656/how-to-understand-the-drawbacks-of-k-means>

Reason for choice of Algorithm:

- Well done here again! Please see the links below for some more information on how to compare the two algorithms:

SUGGESTED READING:

- <https://www.quora.com/What-is-the-difference-between-K-means-and-the-mixture-model-of-Gaussian>

Several silhouette scores are accurately reported, and the optimal number of clusters is chosen based on the best reported score. The cluster visualization provided produces the optimal number of clusters based on the clustering algorithm chosen.

Great job identifying the optimal cluster score as 2, and identifying its associated [silhouette score](#)

The establishments represented by each customer segment are proposed based on the statistical description of the dataset. The inverse transformation and inverse scaling has been properly implemented and applied to the cluster centers in code.

Well done recovering the true centres, and proposing establishments using guidance from the statistical description of the dataset!

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

Nice work here giving detailed comparison of your initial intuition about the sample points and their predicted clusters!

Conclusion

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

Great explanation here! 🙌

The key is to conduct the A/B test on only one segment at the time, rather than on both segments. I like your design of using 2 A/B tests, one per segment.

Suggestions and Comments:

- You can find more information about A/B testing from this [link](#). You can also refer to the [Udacity course](#) to get an overview on how A/B testing is conducted (please note that you don't need to take the whole course, just skimming through the lecture videos and overviews should be alright)

- These links were also really helpful to me when picking up A/B testing. Hope they help you too!
<https://www.quora.com/When-should-A-B-testing-not-be-trusted-to-make-decisions/answer/Edwin-Chen-1>
<http://multithreaded.stitchfix.com/blog/2015/05/26/significant-sample/>
<http://techblog.netflix.com/2016/04/its-all-about-testing-netflix.html>
<https://vwo.com/ab-testing/>

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

Nice try for this section. Below are my comments and feedback:

Required:

- The question actually requires you to use the clustered data in a Supervised Learner, and not some additional data about the customers, but the **clustered data labels**. Please explain how the clustered data labels can be used in a Supervised Learner.

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.

Well done comparing the clusters from your algorithm to the customer 'Channel' data!

 RESUBMIT

 [DOWNLOAD PROJECT](#)



Best practices for your project resubmission

Ben shares 5 helpful tips to get you through revising and resubmitting your project.

[Watch Video](#) (3:01)

Have a question about your review? Email us at review-support@udacity.com and include the link to this review.

RETURN TO PATH

Rate this review

[Student FAQ](#)