Optimal Distance Measure for Agglomerative Hierarchical Clustering

CS5100 Semester Project Checkpoint 1 - 4/7/2022

Xichen Liu Jeff Turgeon

Github Repository

Our project code and the datasets that we are currently working with can be found on our Github site https://github.com/PdAlbedo/Al_Final_Project. We will continue to add code and datasets to this site as the project progresses. A separate text file was created for uploading to Gradescope detailing the files in the repository, but we could not submit multiple documents, so the overview is included in the Github Overview section.

Github Overview

The code and datasets identified below can be found on our Github site - https://github.com/PdAlbedo/Al Final Project.

1. agglomerativeClustering.py

This file is the source code for our project. A copy of the file has been uploaded to Gradescope as well so you have the point in time copy. The program currently reads in a sample dataset, and runs the dataset through the sci-kit Agglomerative Clustering algorithm using different linkage methods. The output from each iteration shows the clustering that would occur using the supplied dataset and linkage method.

2. CC General.csv

This dataset originated from Kaggle -

https://www.kaggle.com/datasets/arjunbhasin2013/ccdata. The dataset contains information on credit card accounts that can be used for clustering and analysis. The file contains 8950 records, with 18 different datapoints for each customer record.

3. Mall_Customers.csv

This dataset originated from Kaggle - https://www.kaggle.com/datasets/shwetabh123/mall-customers. The dataset contains information that can be used by marketing teams to target specific customers. The file only has 4 descriptive variables (gender, age, annual income, and spending score), and only 200 records, but it has variety in the data that we beleive may be useful in our analysis.

5. chicago_data_cleaning.ipynb

This is work completed by Xichen to clean the Chicago Crime dataset.

6. readme.txt

This document contains a link to the source Chicago Crime dataset. It is a very large dataset and excedes the space limitations on Github.

Revised Project Schedule

We are relatively close to our initial project schedule. The initial schedule called for us to have our basic clustering algorithm program written, and our data sets cleaned by 4/4/2022. We initially thought that we would need to create some custom heuristics, and our target date for those was 4/15/2022. The remaining focus of the project was on the analysis, which we had targeted 4/25/2022 for completion. The chart below is from our initial proposal.

Date	Xichen Liu	Jeff Turgeon	
4/4	Clustering algorithm Program	Data Cleaning Routine	
4/15	Half of heuristic functions	Half of heuristic functions	
4/25	Analysis		

The revised timeline is based on a better understanding of the project that we are working on. We thought that we would need to create our own custom agglomerative hierarchical clustering algorithm to implement the different distance measures, however, we have since learned that the algorithm provided by Scikit Learn allows us to choose a distance measure, otherwise known as affinity, and it allows us to provide a precomputed distance matrix - https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html. This flexibility allows us to put a larger focus on analyzing the results for our project. We thought it would be easier to find sample datasets to use for analysis, but many of the readily available datasets are small, or are commonly used in other projects like this. We are still working to find datasets that provide the type of

It is difficult to break the project timeline targets apart at this point. Our program itself is not overly large or complicated. The meat of the project is in the analysis phase and with presenting the results of the project to our peers. We will both be working on all aspects of the project as a team so that we both gain a strong understanding of the results and final analysis.

Our new timeline is as follows:

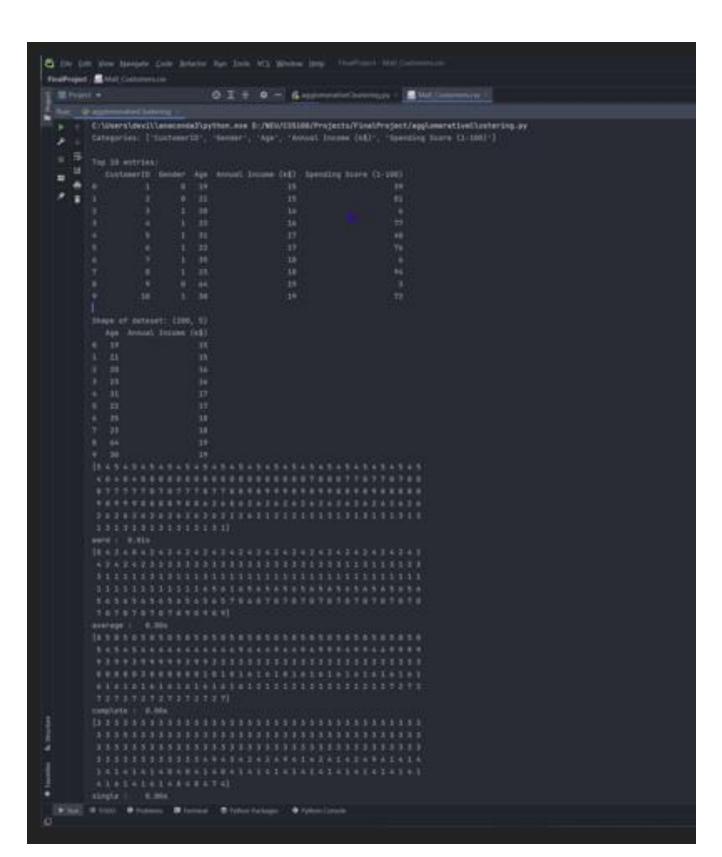
variety that we would like to analyze in our project.

Date	Xichen Liu	Status	Jeff Turgeon	Status
4/7/2022	Build the framework	The framework	Find or create	We have found
	for running the	has been built and	datasets of varying	a few datasets
	clustering algorithm.	tested with one	sizes and dimensions.	so far. Xichen
		dataset. We have		cleaned the
		been conferring		Chicago Crime
		about the need to		dataset so we
		build our own		could use it.
		algorithm versus		
		using provided		
		library options.		
4/14/2022	Update the program to		Update the program	
	loop through the		to loop through the	
	distance metrics.		distance metrics.	
4/17/2022	Finalize the datasets		Finalize the datasets	
	used for analysis.		used for analysis.	
4/21/2022	Compute the results		Compute the results	
Checkpoint	for analysis.		for analysis.	
2				
4/24/2022	Document the findings.		Document the	
			findings.	
4/28/2022	Project Presentation		Project Presentation	
5/5/2022	Submit Project Paper		Submit Project Paper	

Sample Program Output

Our program was not functioning properly at the time of the oral presentation, so we were advised to include an image of the program output in this document as well.

As you can see from the screen shot on the following page, the Chicago Crime dataset was evaluated using several different linkage methods. Each linkage method provided a different hierarchical clustering, as is evidenced by the different group numbers assigned to the data points in the algorithm output. This is just step one in implementing the algorithm, and is not mean to represent final project output.



Current Source Code

```
Apply agglomerative clustering to the given dataset and find out optimal
distance measure for agglomerative hierarchical clustering
__author__ = "Xichen Liu, Jeff Turgeon"
import string
from time import time
import pandas as pd
from sklearn.cluster import AgglomerativeClustering
import numpy as np
from matplotlib import pyplot as plt
from sklearn import manifold, datasets
digits = datasets.load digits()
X, y = digits.data, digits.target
n_samples, n_features = X.shape
np.random.seed(0)
# -----
# Visualize the clustering
def plot_clustering(X_red, labels, title = None):
  x_min, x_max = np.min(X_red, axis = 0), np.max(X_red, axis = 0)
  X_{red} = (X_{red} - x_{min}) / (x_{max} - x_{min})
  plt.figure(figsize = (6, 4))
  for digit in digits.target_names:
    plt.scatter(
      *X_red[y == digit].T,
      marker = f"${digit}$",
      s = 50,
      c = plt.cm.nipy_spectral(labels[y == digit] / 10),
      alpha = 0.5,
    )
  plt.xticks([])
  plt.yticks([])
  if title is not None:
    plt.title(title, size = 17)
  plt.axis("off")
  plt.tight_layout(rect = [0, 0.03, 1, 0.95])
def plot some():
  # ------
  # 2D embedding of the digits dataset
  print("Computing embedding")
```

```
X_red = manifold.SpectralEmbedding(n_components = 2).fit_transform(X)
  print("Done.")
  for linkage in ("ward", "average", "complete", "single"):
    clustering = AgglomerativeClustering(linkage = linkage, n clusters = 10)
    t0 = time()
    clustering.fit(X red)
    print("%s:\t%.2fs" % (linkage, time() - t0))
    plot_clustering(X_red, clustering.labels_, "%s linkage" % linkage)
  plt.show()
  return
def main():
  input file = 'Mall Customers.csv'
  df = pd.read csv(input file, header = 0)
  # convert gender to numeric data
  df['Gender'] = df['Gender'].map(lambda x: 0 if x == 'Male' else x)
  df['Gender'] = df['Gender'].map(lambda x: 1 if x == 'Female' else x)
  # display some basic info of dataset
  original header = list(df.columns.values)
  print("Categories: ", end = "")
  print(original header)
  print("\nTop 10 entries: ")
  print(df.head(10))
  print("\nShape of dateset: ", end = "")
  print(df.shape)
  df_age_annual_income = df[['Age', 'Annual Income (k$)']].copy()
  print(df age annual income.head(10))
  # print("Computing embedding")
  # embedded = manifold.SpectralEmbedding(n components =
2).fit transform(df age annual income)
  # print("Done.")
  for linkage in ("ward", "average", "complete", "single"):
    clustering = AgglomerativeClustering(linkage = linkage, n_clusters = 10)
    t0 = time()
    clustering.fit(df)
    print(clustering.labels_)
    print("%s:\t%.2fs" % (linkage, time() - t0))
```

```
# plt.figure(figsize = (6, 4))
    # for label in clustering.labels_:
      plt.scatter(
    #
          df_age_annual_income['Age'],
    #
          df_age_annual_income['Annual Income (k$)'],
          marker = f"${label}$",
    #
    #
          s = 50,
    #
          c = plt.cm.nipy_spectral(label * 10),
          alpha = 0.5,
    #
    # )
    plt.xticks([])
    plt.yticks([])
    plt.show()
      plot_clustering(embedded, clustering.labels_, "%s linkage" % linkage)
  # plt.show()
  # numpy_array = df.as_matrix()
  # numeric_headers.reverse()
  # reverse_df = df[numeric_headers]
  # reverse_df.to_excel('path_to_file.xls')
if __name__ == '__main__':
  main()
```