

Alex Friedrichsen

Homework 9 Data Mining 395

April 17, 2023

Problem 1

(Comparing Clustering Algorithms)

a) Compare and contrast the performance and suitability of three popular clustering algorithms: K-

Means Clustering, Hierarchical Clustering, and Density-Based Spatial Clustering of Applications with

Noise (DBSCAN). In your report, compare the three algorithms in terms of the following computational aspects:

- Time Complexity: Discuss the time complexity of each algorithm and how it might affect the scalability of the algorithm when dealing with large datasets.
- Robustness: Explain how each algorithm is robust to outliers or noise in the dataset and how it handles such data points.
- Cluster Shape and Size: Describe how each algorithm handles clusters of different shapes and sizes and discuss which algorithm might be more suitable for datasets with irregularly shaped clusters.

Answer Problem 1

Clustering is a an approach to analyzing a data set that involves grouping data points based on their characteristics. It is a type of unsupervised machine learning. It can help find patterns or groups in the data.

The time complexity of each algorithm is dependent on n , the number of data points.

- K-means time complexity: $O(nki)$, where k is the number of clusters, and i is the number of iterations. Worst case: exponential, and therefore risky for large data sets.
- Hierarchical clustering: $O(n^2)$. Slowest.



The robustness is the algorithm's ability to handle noisy data.

- K-means is not very robust as it tends to assign outliers to the nearest cluster, despite the large gap.
- Hierarchical clustering is relatively robust due to its basis in distance metrics that can deal with noise. As the noise grows to great, though, it will eventually cause incorrect clustering.
- DBSCAN is highly robust to outliers because it does not cluster data points that shouldn't belong to a specific cluster. Uses a density based approach to identify clusters. Requires minimum number of data points to form a cluster.

Cluster shape and cluster size refer to the distribution of data points within clusters.

- K-means assume spherical, equal size clusters. It is unsuitable for data sets with irregularly shaped clusters. It also requires the choice of centroids, a central data point for each cluster.
- Hierarchical clustering can make clusters of varying shapes and sizes due to the process of building its tree based structure. This allows it to capture more complex relationships.
- DBSCAN can handle varying shapes and sizes of clusters as well, due to the density-based approach.

▼ Problem 2:

Compare and contrast the performance and suitability of three popular clustering algorithms: K-Means Clustering, Hierarchical Clustering, and Density-Based Spatial Clustering of Applications with Noise (DBSCAN). In your report, compare the three algorithms in terms of the following computational aspects: Time Complexity: Discuss the time complexity of each algorithm and how it might affect the scalability of the algorithm when dealing with large datasets. Robustness: Explain how each algorithm is robust to outliers or noise in the dataset and how it handles such data points. Cluster Shape and Size: Describe how each algorithm handles clusters of different shapes and sizes and discuss which algorithm might be more suitable for datasets with irregularly shaped clusters.

```
1 from sklearn.datasets import make_moons
2 from sklearn.cluster import KMeans, AgglomerativeClustering, DBSCAN
3 from sklearn.metrics import silhouette_score, calinski_harabasz_score, davies_bouldin_score
4 import time
5 import os
6 import psutil
7 import numpy as np
8 import pandas as pd
```

```
8 import pandas as pd
9
10

1 # Define a function to perform hierarchical clustering and save the evaluation scores ar
2 def hierarchical_clustering(X, filename):
3     start_time = time.time()
4     process = psutil.Process(os.getpid())
5     hc = AgglomerativeClustering(n_clusters=2).fit(X)
6     end_time = time.time()
7     execution_time = end_time - start_time
8     memory_usage = process.memory_info().rss / 1024 / 1024 # convert from bytes to MB
9     silhouette = silhouette_score(X, hc.labels_)
10    calinski = calinski_harabasz_score(X, hc.labels_)
11    davies = davies_bouldin_score(X, hc.labels_)
12    with open(filename, 'a') as f:
13        f.write("Hierarchical Clustering Results:\n")
14        f.write("Silhouette Score: {}\n".format(silhouette))
15        f.write("Calinski-Harabasz Index: {}\n".format(calinski))
16        f.write("Davies-Bouldin Index: {}\n".format(davies))
17        f.write("Execution Time: {} seconds\n".format(execution_time))
18        f.write("Memory Usage: {} MB\n".format(memory_usage))

1 # Define a function to perform k-means clustering and save the evaluation scores and ex
2 def kmeans_clustering(X, filename):
3     start_time = time.time()
4     process = psutil.Process(os.getpid())
5     kmeans = KMeans(n_clusters=2, random_state=0).fit(X)
6     end_time = time.time()
7     execution_time = end_time - start_time
8     memory_usage = process.memory_info().rss / 1024 / 1024 # convert from bytes to MB
9     silhouette = silhouette_score(X, kmeans.labels_)
10    calinski = calinski_harabasz_score(X, kmeans.labels_)
11    davies = davies_bouldin_score(X, kmeans.labels_)
12    with open(filename, 'w') as f:
13        f.write("K-Means Clustering Results:\n")
14        f.write("Silhouette Score: {}\n".format(silhouette))
15        f.write("Calinski-Harabasz Index: {}\n".format(calinski))
16        f.write("Davies-Bouldin Index: {}\n".format(davies))
17        f.write("Execution Time: {} seconds\n".format(execution_time))
18        f.write("Memory Usage: {} MB\n".format(memory_usage))

1 def dbscan_clustering(X, filename):
2     start_time = time.time()
3     process = psutil.Process(os.getpid())
4     n_samples = X.shape[0]
5     valid_labels = False
6     while not valid_labels:
7         eps = np.random.uniform(0.1, 0.5) # adjust range as necessary
```

```

8         min_samples = np.random.randint(2, int(n_samples / 2)) # adjust range as necessi
9         try:
10             dbscan = DBSCAN(eps=eps, min_samples=min_samples).fit(X)
11             silhouette = silhouette_score(X, dbscan.labels_)
12             calinski = calinski_harabasz_score(X, dbscan.labels_)
13             davies = davies_bouldin_score(X, dbscan.labels_)
14             valid_labels = len(set(dbscan.labels_)) > 1
15         except:
16             pass
17     end_time = time.time()
18     execution_time = end_time - start_time
19     memory_usage = process.memory_info().rss / 1024 / 1024 # convert from bytes to MB
20     with open(filename, 'a') as f:
21         f.write("DBSCAN Clustering Results:\n")
22         f.write("Eps: {}\n".format(eps))
23         f.write("Min Samples: {}\n".format(min_samples))
24         f.write("Silhouette Score: {}\n".format(silhouette))
25         f.write("Calinski-Harabasz Index: {}\n".format(calinski))
26         f.write("Davies-Bouldin Index: {}\n".format(davies))
27         f.write("Execution Time: {} seconds\n".format(execution_time))
28         f.write("Memory Usage: {} MB\n".format(memory_usage))
29

```

```

1 def clustering_pipeline(X, filename):
2     # Perform K-Means clustering
3     kmeans_clustering(X, filename)
4
5     # Perform Hierarchical clustering
6     hierarchical_clustering(X, filename)
7
8     # Perform DBSCAN clustering
9     dbscan_clustering(X, filename)
10

```

```

1
2 # # Load the make_moons data set
3 # X, y = make_moons(n_samples=1000, noise=0.05, random_state=0)
4
5
6 # clustering_pipeline(X, "make_moons_results.txt")
7

```

```

1 # # Load the smartphone data set
2 # # Load data from CSV file
3 # data = pd.read_csv("smartphone/train.csv")
4
5
6 # # Extract X and y data
7 # X = data.iloc[:, :-1].values # assuming last column is the target variable

```

```

7 # x = data.iloc[:, :-1].values # assuming last column is the target variable
8 # y = data.iloc[:, -1].values
9 # clustering_pipeline(X, "smartphone_results.txt")

1 def find_optimal_dbscan_params(X, eps_range, min_samples_range):
2     best_silhouette = -1
3     best_eps = 0
4     best_min_samples = 0
5
6     for eps in eps_range:
7         for min_samples in min_samples_range:
8             dbscan = DBSCAN(eps=eps, min_samples=min_samples)
9             labels = dbscan.fit_predict(X)
10            if len(set(labels)) > 1:
11                silhouette = silhouette_score(X, labels)
12                if silhouette > best_silhouette:
13                    best_silhouette = silhouette
14                    best_eps = eps
15                    best_min_samples = min_samples
16
17    return best_eps, best_min_samples, best_silhouette

1 # Load the dataset from file
2 data = pd.read_csv("shuttle-landing-control.data", header=None)
3
4 # Replace asterisks with NaN
5 data = data.replace("*", np.nan)
6
7 # Perform one-hot encoding on categorical variables
8 data = pd.get_dummies(data, columns=[0, 1, 2, 3, 4, 5])
9
10 # Drop any rows with NaN values
11 data = data.dropna()
12
13 # Convert data to a NumPy array
14 X = data.to_numpy()
15
16
17 eps_range = np.arange(0.1, 1.0, 0.1)
18 min_samples_range = range(1, 10)
19
20 # Find optimal values for eps and min_samples
21 best_eps, best_min_samples, best_silhouette = find_optimal_dbscan_params(X, eps_range, r
22
23
24 # Print the optimal hyperparameters and clustering quality
25 print("Optimal eps:", best_eps)
26 print("Optimal min_samples:", best_min_samples)
27 print("Best silhouette score:", best_silhouette)
28 # clustering_pipeline(X, 'shuttle landing results.txt')

```

29

30

```

-----
ValueError                                Traceback (most recent call last)
c:\Users\alexp\Documents\GitHub\cs395_data_mining\HW9.ipynb Cell 10 in <module>
      1 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
      2 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=17'>18</a> min_samples_range =
      3 range(1, 10)
      4 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
      5 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=19'>20</a> # Find optimal values
      6 for eps and min_samples
      7 ---> <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
      8 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=20'>21</a> best_eps,
      9 best_min_samples, best_silhouette = find_optimal_dbscan_params(X, eps_range,
     10 min_samples_range)
     11 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
     12 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=23'>24</a> # Print the optimal
     13 hyperparameters and clustering quality
     14 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
     15 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=24'>25</a> print("Optimal eps:",
     16 best_eps)

c:\Users\alexp\Documents\GitHub\cs395_data_mining\HW9.ipynb Cell 10 in
find_optimal_dbscan_params(X, eps_range, min_samples_range)
      1 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
      2 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=8'>9</a> labels =
      3 dbscan.fit_predict(X)
      4 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
      5 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=9'>10</a> if len(set(labels)) >
      6 1:
      7 ---> <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
      8 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=10'>11</a> silhouette =
      9 silhouette_score(X, labels)
     10 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
     11 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=11'>12</a> if silhouette >
     12 best_silhouette:
     13 <a href='vscode-notebook-cell:/c%3A/Users/alexp/Documents/GitHub
     14 /cs395_data_mining/HW9.ipynb#X13sZmlsZQ%3D%3D?line=12'>13</a>
     15 best_silhouette = silhouette

File c:\Users\alexp\AppData\Local\Programs\Python\Python39\lib\site-packages\sklearn
\metrics\cluster\_unsupervised.py:117, in silhouette_score(X, labels, metric,
sample_size, random_state, **kwargs)
    115     else:
    116         X, labels = X[indices], labels[indices]
--> 117 return np.mean(silhouette_samples(X, labels, metric=metric, **kwargs))

File c:\Users\alexp\AppData\Local\Programs\Python\Python39\lib\site-packages\sklearn
\metrics\cluster\_unsupervised.py:231, in silhouette_samples(X, labels, metric,
...

```

Results:

make_moons:

K-Means Clustering Results: Silhouette Score: 0.48975082695298533 Calinski-Harabasz Index: 1481.4753424968233 Davies-Bouldin Index: 0.7817364605843775 Execution Time: 0.04900717735290527 seconds Memory Usage: 146.68359375 MB Hierarchical Clustering Results: Silhouette Score: 0.4391867809034311 Calinski-Harabasz Index: 1087.7713369249986 Davies-Bouldin Index: 0.837284087751716 Execution Time: 0.029980182647705078 seconds Memory Usage: 146.69140625 MB DBSCAN Clustering Results: Eps: 0.4503519041566185 Min Samples: 147 Silhouette Score: 0.3060017383266752 Calinski-Harabasz Index: 361.7530597734725 Davies-Bouldin Index: 2.8816508031309773 Execution Time: 0.24503803253173828 seconds Memory Usage: 147.47265625 MB

smartphone:

K-Means Clustering Results: Silhouette Score: 0.37196991401903456 Calinski-Harabasz Index: 5934.114761241196 Davies-Bouldin Index: 1.074520597574652 Execution Time: 0.3592417240142822 seconds Memory Usage: 311.51953125 MB Hierarchical Clustering Results: Silhouette Score: 0.3735621062678581 Calinski-Harabasz Index: 5765.523209652622 Davies-Bouldin Index: 1.0735359209035171 Execution Time: 8.188697576522827 seconds Memory Usage: 318.9765625 MB

shuttle-landing-control:

K-Means Clustering Results: Silhouette Score: 0.20685539935851324 Calinski-Harabasz Index: 5.1140495867768605 Davies-Bouldin Index: 1.518699240339859 Execution Time: 0.1279003620147705 seconds Memory Usage: 204.6953125 MB Hierarchical Clustering Results: Silhouette Score: 0.19435098907732273 Calinski-Harabasz Index: 4.71151515151515 Davies-Bouldin Index: 1.6027259104482954 Execution Time: 0.1157076358795166 seconds Memory Usage: 205.296875 MB

Discussion

For the make_moons data, k-means had the best silhouette score by a small margin, though all 3 algorithms had positive scores indicating a decent clustering. K-means outperformed in the other metrics as well, and both hierarchical and k-means outperformed DBSCAN, except for in memory usage. Smartphone data had the hierarchical clustering slightly ahead of k-means in silhouette score but not enough to be significant. Hierarchical also took more memory and time. For the shuttle-landing-control dataset, the clustering scores were relatively low at around 0.2, indicating poor clustering.

poor clustering.

DBSCAN did not run for either the shuttle landing control or smartphone datasets. The smartphone dataset was much larger, and I let DBSCAN iterate randomly over eps values between 0.1 and 0.5, in combination with min_sample sizes between 2 and $n/2$. I let the code run over night for 635 minutes and it still hadn't found a valid clustering for the data. For the shuttle-landing-control data, I believe it was too noisy for DBSCAN to create meaningful clusters. I got the error "ValueError: Number of labels is 15. Valid values are 2 to $n_samples - 1$ (inclusive)" which indicates it was clustering into 15 different classes, which is clearly too many for the data.

[Colab paid products](#) - [Cancel contracts here](#)