

Project 2 overdue

Bouruc Petru-Liviu

31/08/2023

Contents

1	Introduction & Project Description	1
2	Data Preprocessing	1
3	Agglomerative	2
4	Mean-Shift	3
5	Comparison with supervised method and random chance	4
6	Conclusions	5

1 Introduction & Project Description

This project was about applying unsupervised methods on dataset and observe the results. For this one, I chose as a dataset english handwritten characters. The dataset contains 3410 images partitioned into 62 classes. For the supervised methods, I chose AgglomerativeClustering and Mean-Shift. The data is not separated, so we will use the build-in function *train_test_split* from sklearn to make our datasets.

2 Data Preprocessing

As for the data preprocessing, I chose to go with 2 variants: simple HSV features from image pixels, and features obtained from getting the image through a VGG network.

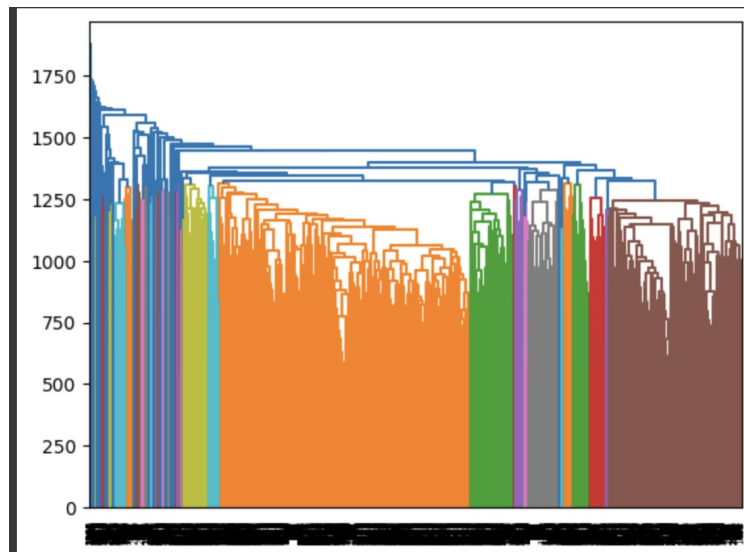
The HSV model is the most accurate representation of how humans perceive color. More than that, after converting to HSV I took the mean value of the pixels and flatten the image in order to have less computations to perform.

As for the second approach, firstly I read the images and simply saved them as arrays. Then I put them through a standard VGG pre-trained model with weights from imagenet, also making use of the build-in keras method *preprocess_input*. For this approach I also tried to use PCA to reduce the dimension of my input and also keep the core information.

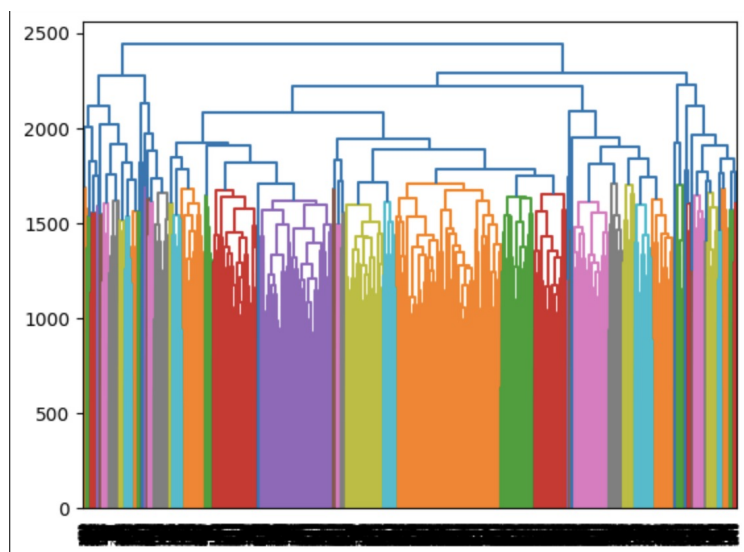
3 Agglomerative

Agglomerative Clustering is a type of clustering in which objects that are close to each other are placed into groups. It works similarly with a tree, such that at the beginning, all objects are individual clusters (leaves) and the algorithm builds the tree from the bottom to the top (root). The similarities between clusters can be computed with different methods, a combination between metrics and linkages. This structure is called a dendrogram, being a common type of hierarchical clustering.

Dendrogram representation of our dataset with VGG features and reduced with PCA, with average linkage:



Dendrogram representation of our dataset with VGG features and reduced with PCA, with complete linkage:

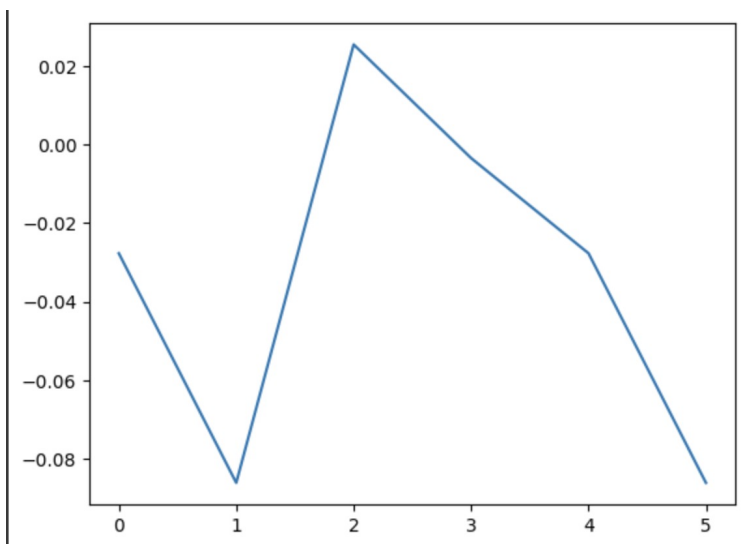


When computing the best combination for Agglomerative Clustering, I took into account

different metrics and linkages. As it can be seen, none of them gave exceptionally results by looking at the silhouettes score.

```
metrics = ["l1", "l2", "manhattan"]
linkages = ["complete", "average"]
s_scores = []
for metric in metrics:
    for linkage in linkages:
        agg_clusters = AgglomerativeClustering(n_clusters=62, metric=metric, linkage=linkage).fit(train_pca_data)
        print(metric, " ", linkage, Counter(agg_clusters.labels_), silhouette_score(train_pca_data, agg_clusters.labels_))
        s_scores.append(silhouette_score(train_pca_data, agg_clusters.labels_))
```

```
l1 complete Counter({0: 1510, 14: 289, 1: 150, 34: 113, 9: 87, 11: 65, 2: 60, 10: 59, 13: 58, 0: 50, 54: 35, 20: 26, 26: 24,
l1 average Counter({0: 2644, 5: 4, 30: 3, 1: 3, 14: 3, 4: 3, 2: 3, 22: 2, 8: 2, 6: 2, 3: 2, 13: 2, 10: 2, 7: 2, 28: 2, 12: 2,
l2 complete Counter({1: 308, 3: 269, 6: 188, 18: 151, 20: 150, 4: 137, 19: 132, 15: 111, 40: 84, 17: 79, 37: 75, 21: 73, 10:
l2 average Counter({1: 1417, 0: 816, 5: 60, 12: 46, 8: 39, 6: 38, 30: 38, 25: 16, 24: 17, 22: 14, 22: 13, 17: 3, 16: 9,
manhattan complete Counter({0: 1510, 14: 289, 1: 150, 34: 113, 9: 87, 11: 65, 2: 60, 10: 59, 13: 58, 0: 50, 54: 35, 20: 26,
manhattan average Counter({0: 2644, 5: 4, 30: 3, 1: 3, 14: 3, 4: 3, 2: 3, 22: 2, 8: 2, 6: 2, 3: 2, 13: 2, 10: 2, 7: 2, 28: 2,
```



One last experiment which I tried with agglomerative was to clusterize all the dataset with and without PCA applied, in order to see if there is a big difference between them. As we can observe, the clusters are approximately the same sizes, which could be an indicator that after the transformation much of the information was being kept.

```
[61] # VGG features, no PCA
agg_clusters = AgglomerativeClustering(n_clusters=62).fit(images_features)
print(sorted(Counter(agg_clusters.labels_).items()))
```

```
[(0, 123), (1, 91), (2, 75), (3, 79), (4, 45), (5, 57), (6, 53), (7, 115), (8, 70), (9, 52), (10, 58), (11, 69), (12, 81), (13,
```

```
# VGG features, with PCA
agg_clusters = AgglomerativeClustering(n_clusters=62).fit(pca_images_features)
print(sorted(Counter(agg_clusters.labels_).items()))
```

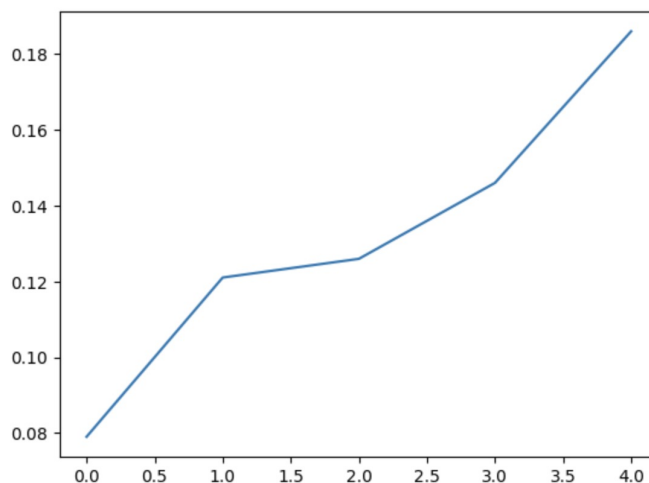
```
[(0, 80), (1, 77), (2, 57), (3, 82), (4, 48), (5, 60), (6, 50), (7, 52), (8, 55), (9, 63), (10, 50), (11, 39), (12, 52), (13,
```

4 Mean-Shift

The second approach was Mean-Shift, which is a clustering algorithm that is based on the idea that it finds high density in a region and move the other points towards that. The shift is done with the help of a kernel, in which the distance of which the mean is calculated is represented by the bandwidth. Sklearn offers us a build-in nearest-neighbor analysis with the help of *estimate_bandwidth* function. Using this, on my dataset it approximated a bandwidth of around 1200, so I computed the Mean-Shift with values from that range.

It was a heavily computational task, which took some time to complete. In the end, an improvement could be seen on terms of silhouette scores.

Using the HSV features, the Mean-Shift algorithm could not make the clustering in reasonable time, as the number of features proved to be too big and it's runtime exceeding one hour. So, for this approach, I could only use the second set of features, those based on VGG and then reduced with PCA.

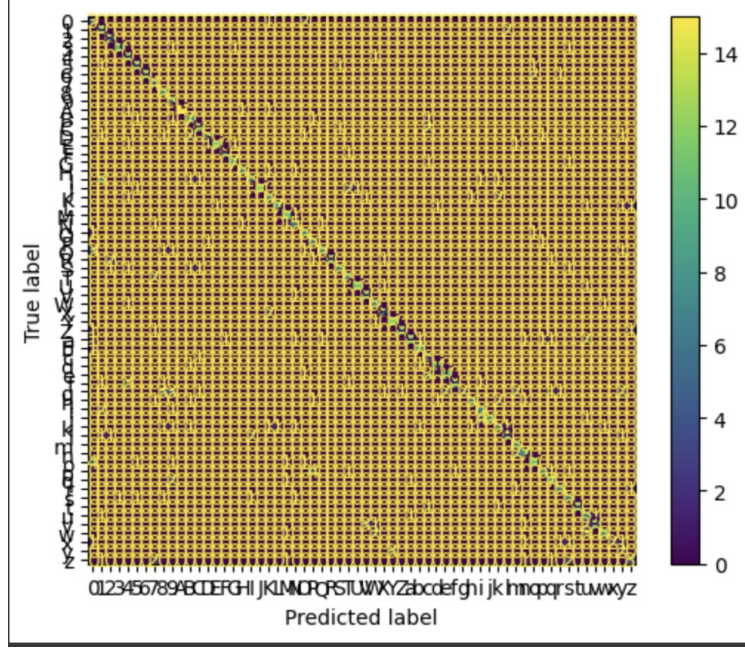


5 Comparison with supervised method and random chance

For the supervised method, I chose to work with RandomForest classifier, and as for the dataset, I used the VGG-features after PCA transformation. I tried a grid search to find the best parameters in order to maximize the accuracy.

The results were:

	n=100	n=200	n=500	n=1000
kernel=rbf	0.7	0.74	0.77	0.78
kernel=poly	0.64	0.69	0.76	0.76



As for the random chance, by assigning random labels to each data point from the possible values of the dataset labels, the accuracy that it managed to obtain was 0.017. It is so low because of the great number of classes.

6 Conclusions

In conclusion, it is clear that the supervised methods outperform the unsupervised one for this task, but with the right improvements, they can give a good analysis of the dataset.