# A Data Symphony of Image Segmentation

"Data Analysis on Image Segmentation Data"

Syed Muhammad Taqi Raza Kazmi
Aamir Khan
Irfan Mahmood

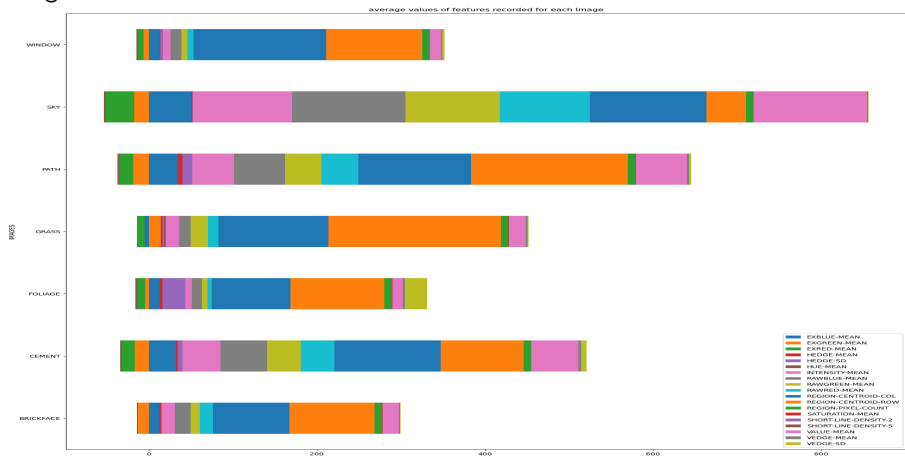Università degli Studi dell'Aquila

July 4, 2023

# Contents

# Data Structure and Visualization

The image segmentation data is the combination of two data sets. One data set consists of 210 rows and 20 columns. Another data set consists of 2100 rows and 20 columns. The plot below show the average values of features recorded for each Image:



average values of features recorded for each Image

# Aim and objective

The image segmentation data is a complex and large data set. The instances were drawn randomly from a database of 7 outdoor images. The images were hand segmented to create a classification for every pixel. Each instance is a 3x3 region. The seven outdoor images are brickface, sky, foliage, cement, window, path, grass. In this project, we applied first implement K-Means clustering to find intrinsic groups within this dataset, getting it misfit, we applied an ML algorithm (SVC) to train an SVM classification model with default options for hyperparameters and computed different classification scores for the training and testing set. Finally, we applied the ensemble algorithm GradientBoostingClassifier also and compared the accuracies.

# Data Cleaning

The data frame has no missing values and we observed 224 duplicate rows, which we considered as unhealthy to the dataset so they were removed.
The data statistics are shown below:

```
Uncleaned_Project_Dataset Statistics
Number of variables: 20
Number of observations: 2310
Missing Values: 0
Missing Values : 0.0
Duplicate rows: 446
Duplicate rows (%): 9.697
Variable types:
        int32 2
        int64 1
        float64 16
        object 1
```

```
Cleaned_Project_Dataset Statistics
Number of variables: 20
Number of observations: 1864
Missing Values: 0
Missing Values : 0.0
Duplicate rows: 0
Duplicate rows (%): 0.0
Variable types:
        int32 2
        int64 1
        float64 16
        object 1
```
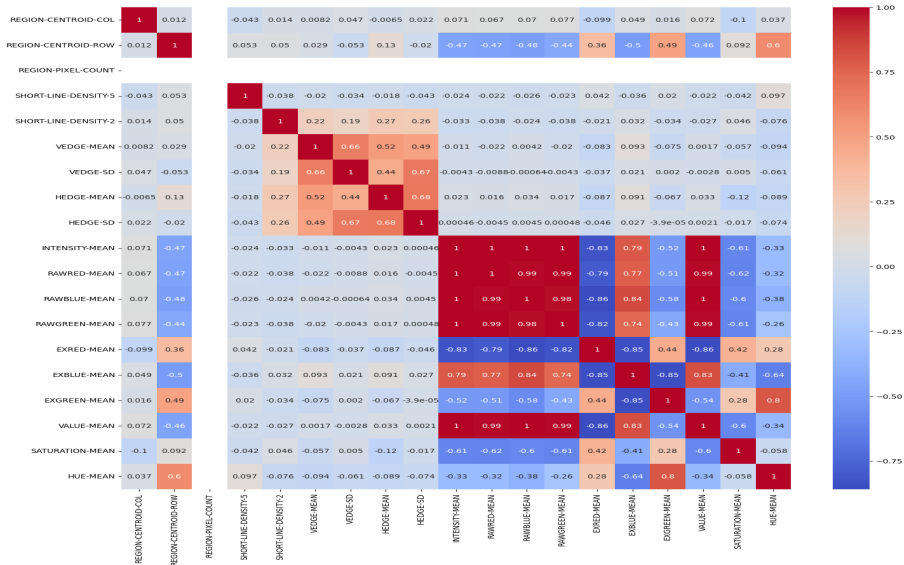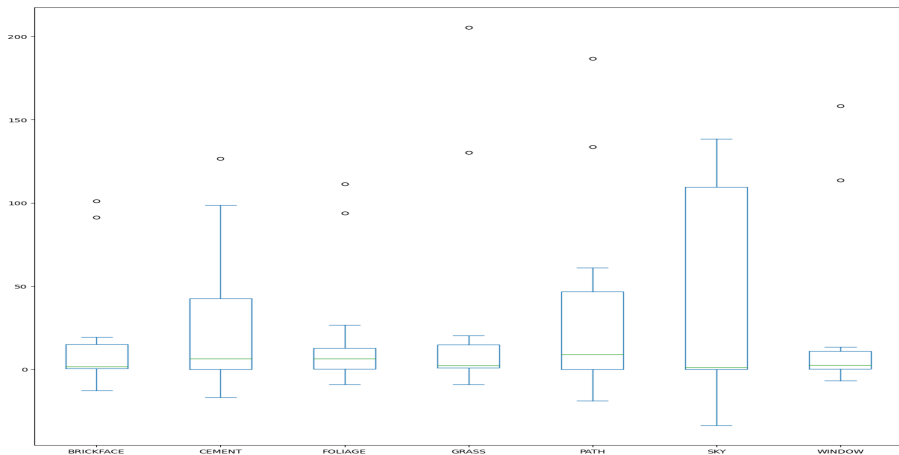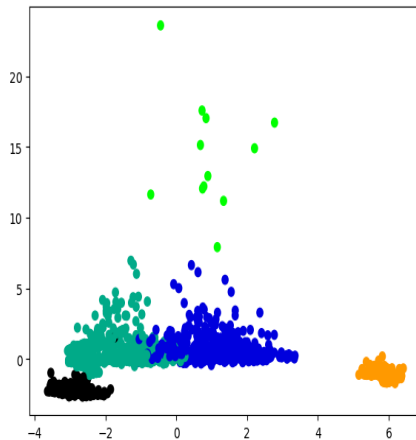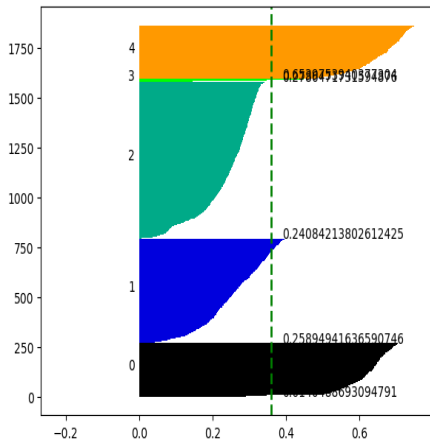
# Correlation matrix

# Outliers



Except for Sky, every image have outliers. Cement has one outlier above the upper whisker. On the other hand, Brickface, Foliage, Grass, Path and Window have several outliers above the upper whisker.

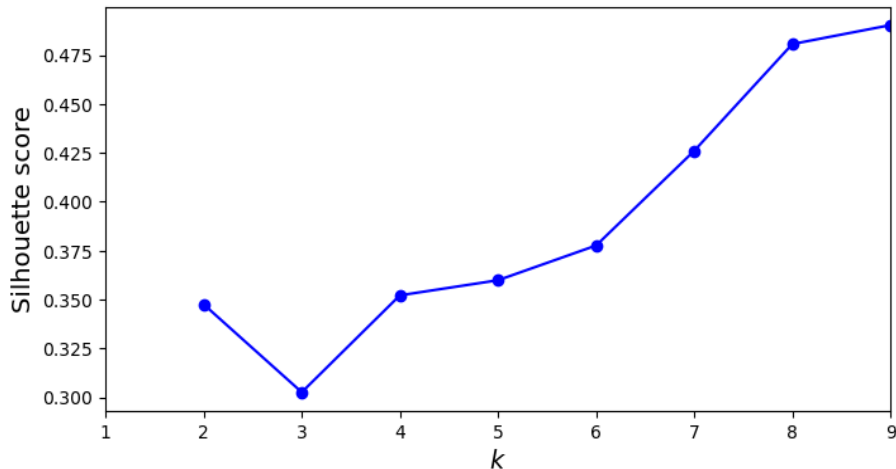# Unsupervised learning: Clustering



Shown is the K-means clustering with K=5
We have repeated the experiment many times with different clusters. Computing the silhouette scores
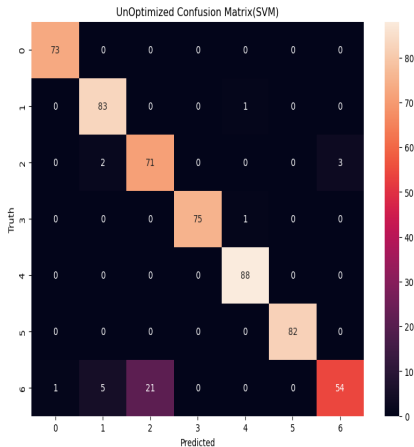
# Silhouette Analysis by Elbow Method



Here we can see that Silhouette's score is increasing with the change in the number of clusters which clearly shows that Clustering is not a better strategy for this data set.

# Unoptimized SVM Classifier: Confusion matric with Classification Report



UnOptimized Confusion Matrix(SVM)

```
0.9392857142857143
              precision    recall  f1-score   support

           0       0.99      1.00      0.99        73
           1       0.92      0.99      0.95        84
           2       0.77      0.93      0.85        76
           3       1.00      0.99      0.99        76
           4       0.98      1.00      0.99        88
           5       1.00      1.00      1.00        82
           6       0.95      0.67      0.78        81

    accuracy                           0.94       560
   macro avg       0.94      0.94      0.94       560
weighted avg       0.94      0.94      0.94       560
```
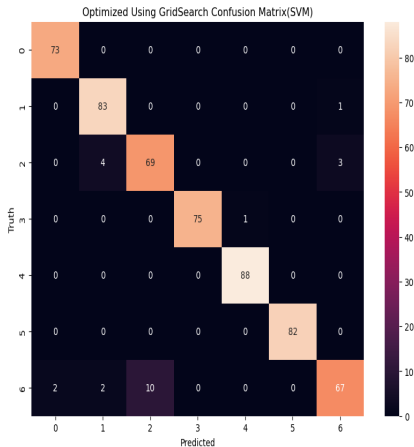
Prediction accuracy is 94 percent

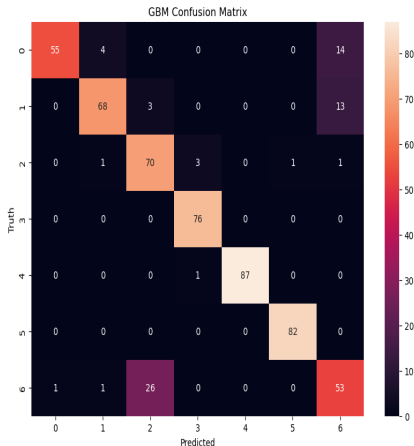# Optimized SVM Classifier: Confusion matric with Classification Report



Optimized Using GridSearch Confusion Matrix(SVM)

0.9589285714285715

|           | precision | recall | f1-score | support |
|-----------|-----------|--------|----------|---------|
| 0         | 0.97      | 1.00   | 0.99     | 73      |
| 1         | 0.93      | 0.99   | 0.96     | 84      |
| 2         | 0.87      | 0.91   | 0.89     | 76      |
| 3         | 1.00      | 0.99   | 0.99     | 76      |
| 4         | 0.99      | 1.00   | 0.99     | 88      |
| 5         | 1.00      | 1.00   | 1.00     | 82      |
| 6         | 0.94      | 0.83   | 0.88     | 81      |
| accuracy  |           |        | 0.96     | 560     |
| macro avg | 0.96      | 0.96   | 0.96     | 560     |
| weighted avg | 0.96   | 0.96   | 0.96     | 560     |

Prediction accuracy is 96 percent

# Gradient Boosting Classifier



GBM Confusion Matrix

0.8767857142857143

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.98 | 0.75 | 0.85 | 73 |
| 1 | 0.92 | 0.82 | 0.87 | 84 |
| 2 | 0.70 | 0.92 | 0.80 | 76 |
| 3 | 0.95 | 1.00 | 0.97 | 76 |
| 4 | 1.00 | 0.99 | 0.99 | 88 |
| 5 | 0.99 | 1.00 | 0.99 | 82 |
| 6 | 0.66 | 0.64 | 0.65 | 81 |
| accuracy |  |  | 0.88 | 560 |
| macro avg | 0.89 | 0.88 | 0.88 | 560 |
| weighted avg | 0.89 | 0.88 | 0.88 | 560 |

Overall, the prediction accuracy is above 80 percent, which is considered a good performance. The results suggested that F1-score is high in "Window" and "Path" features (Class 4 and 5), which had the largest sample size. The opposite is true for "grass" (Class 6).

# Learning Curve for GBC Model and SVM Model



The learning curve suggests that the SVM model achieves higher accuracy than GBC Model and is ready to predict new datasets.

# Results

In this exercise, we have built a successfully optimized Support Vector Machine classifier, of which we were able to improve the accuracy slightly from 0.93 to 0.95 by tuning the model's parameters. The optimum set of parameters for this SVM model is 'C': 100, 'kernel': 'rbf'. We also compared it with the accuracy of the Gradient Boosting Classifier and found that it has accuracy of 0.875, which less than that of the SVM Classifier

# Thank you!