

Implementation of K-means Clustering in Banknote Authentication Dataset

Chissanu Kittipakorn

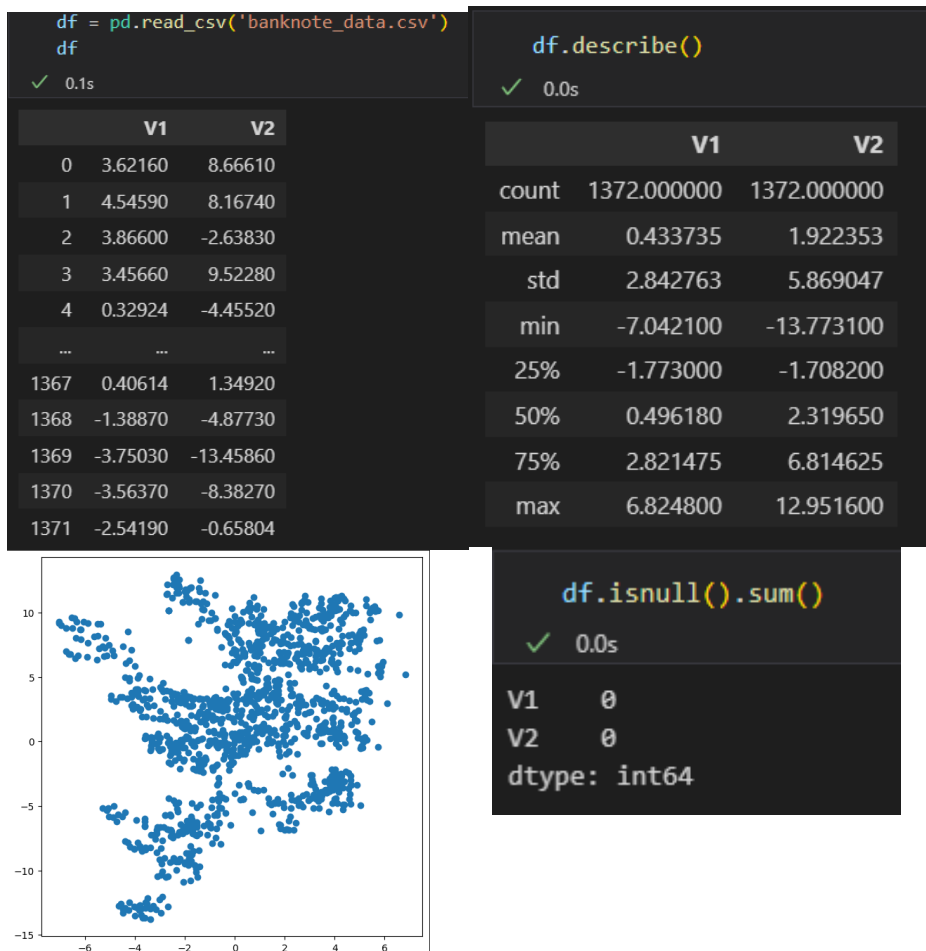
64011728

a. Dataset Overview

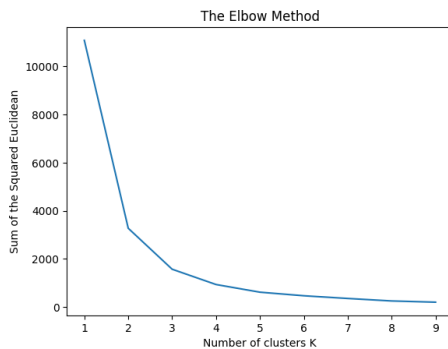
This dataset contains 1,372 data that has 2 columns which are V1 (Variance of Wavelet Transformed Image) and V2 (Skewness of Wavelet Transformed Image). The data range from negative to positive floats which I think the positive means that the banknote is likely to be authentic

b. Data Processing and Exploratory Data Analysis

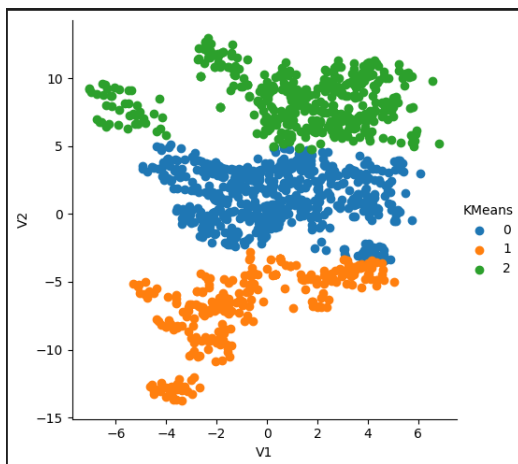
Use Pandas and Matplotlib to extract the content and visualize for better understanding as well as check if the data contain invalid inputs.



c. Training Models



Use Elbow method to find best number of clusters for this case I think either 2 or 3 would be the best choice and I will go with 3 number of clusters



I use seaborn to make a group of clusters and the result are this image.

d. Model Evaluation

I use the Silhouette score to calculate with 3 clusters and the reason the percentage is low because I think there are other factor that affect the prediction due to bad data

```
KMeans
KMeans(n_clusters=3)

silhouette_score(data, model.labels_)
✓ 0.0s
0.43362017621560883
```

```
from sklearn.metrics import davies_bouldin_score
davies_bouldin_score(data, model.labels_)
✓ 0.0s
0.8181686203187953
```

For the second evaluation use Davies boudin score which gives better results.

e. Conclusion

After testing with different number of clusters I found out that 2 clusters give the best result which increase from mine about 1-5% which is significant. I also learned that it is quite hard to train a model with limited number of features and unclean data.