# Exploration of Five Unsupervised Learning Algorithms

**Adedamola Adesoye**
College of Engineering
Northeastern University
Toronto, ON
*adesoye.a@northeastern.edu*

## Abstract

In this report, I explored the performance of unsupervised learning algorithms on two interesting datasets. The five learning algorithms include k-means, Gaussian Mixture Model (GMM), Principal Component Analysis (PCA), Independent Component Analysis (ICA) and Uniform Manifold Approximation and Projection (UMAP). The two datasets were a bank marketing data and adult dataset. I applied all algorithms to both datasets. I then applied a neural network classifier in combination with the unsupervised learning algorithms to run some experiments on the adult dataset.

## 1    Datasets

The two datasets are like the ones used for the supervised learning assignment: bank marketing data and the adult data. The only difference is that the bank marketing data is more balanced with fewer instances than the previously used one. The bank marketing data describes the results of a direct marketing campaign by an unnamed Portuguese bank to convince customers to subscribe to a term deposit, which is a type of deposit account where money is locked up for an agreed time duration. These locked up funds allow banks to invest their customer's money in profitable ventures. The adult data contains demographic information about individuals from a 1994 U.S.A. Census database, and the target column indicates if an individual earns more or less than $50,000 per year, which can help tailor marketing or political campaigns for better conversion rates.

Table 1: The basic feature of both datasets.

|  | Data Set Characteristics | Attribute Characteristics | Associated Tasks | Number of Instances | Number of Attributes |
|---|---|---|---|---|---|
| **Adult** | Multivariate | Real, Categorical | Classification | 32,561 | 14 |
| **Bank Marketing** | Multivariate | Real, Categorical | Classification | 11,163 | 17 |

### 1.1    Data characteristics

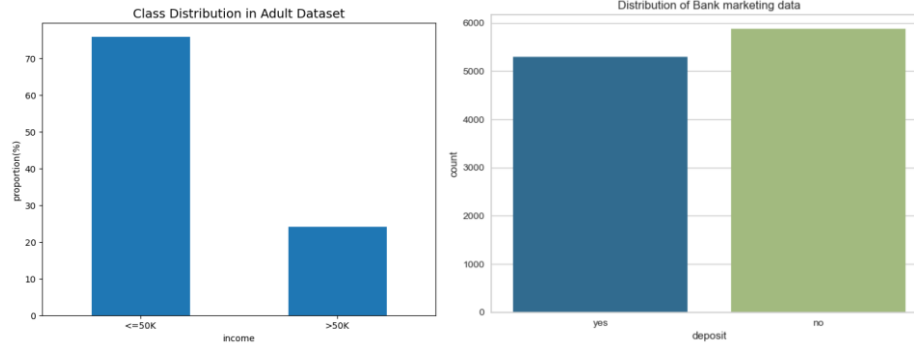The histograms of classes from both datasets are shown in Figure 1. Class unbalance is evident in the adult dataset.

35
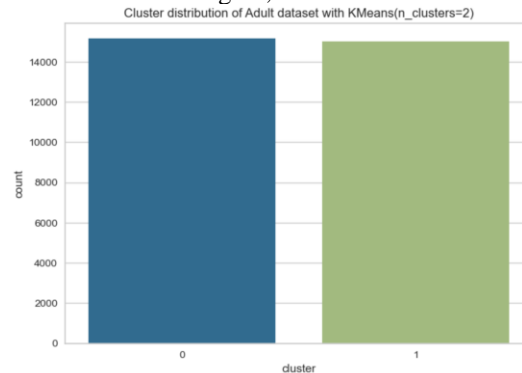36    Figure 1: The class frequency of both datasets.

37    **2        Adult data**
38
39    **2.1      Pre-processing**
40    I started by dropping missing values, leaving me with 30,162 instances. I transformed the *income*
41    feature to a numerical variable, making *<=50K* equal to *0* and *>50K* equal to *1*. Although I
42    dropped this feature before building my clustering models, I did this to be able to later evaluate the
43    results of my models with two clusters. To avoid issues with multicollinearity, I dropped
44    *education* feature which provides the same information as *education-num* but in a different format.
45    I one-hot encoded my categorical variables, scaled all my numerical variables (except *education-*
46    *num*) to zero mean and unit variance; I scaled *education-num* (an ordinal variable) to a range
47    between 0 and 1.
48
49    **2.2      K-Means**
50    I started by initializing my K-Means with two clusters. I did this because the initial target feature
51    was a binary variable. The resulting clusters were relatively evenly distributed with *cluster 0*
52    having 15,148 instances and *cluster 1* having 15,014 instances.
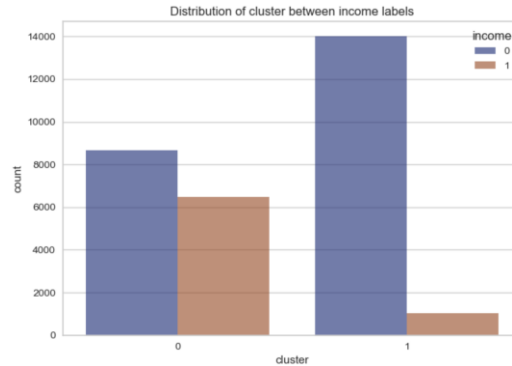


53
54    I then went on to treat the *income* label as the original target feature of a test set and the cluster
55    labels as the predicted labels of the test set. This yielded the classification report shown below.

```
               precision    recall  f1-score   support

           0       0.57      0.38      0.46     22654
           1       0.07      0.14      0.09      7508

    accuracy                           0.32     30162
   macro avg       0.32      0.26      0.27     30162
weighted avg       0.45      0.32      0.37     30162
```
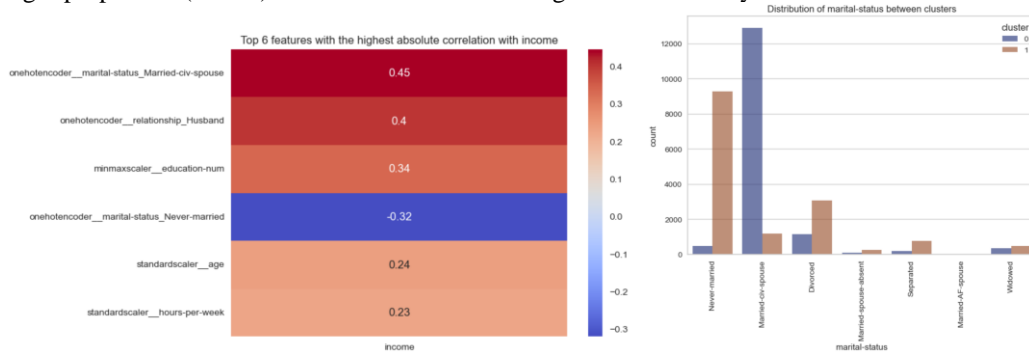
56
57    Because K-Means is unsupervised and *0 & 1* are just ways to differentiate two different clusters, I
58    changed *cluster 0* to *cluster 1* and vice-versa. Again, treating the cluster labels as the prediction
59    labels of a test set, I got the following results.

```
              precision    recall  f1-score   support

           0       0.93      0.62      0.74     22654
           1       0.43      0.86      0.57      7508

    accuracy                           0.68     30162
   macro avg       0.68      0.74      0.66     30162
weighted avg       0.81      0.68      0.70     30162
```
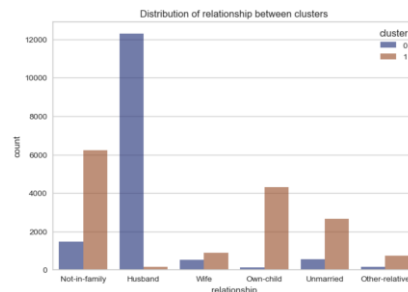
60
61  Every metric improves significantly but this model is still unsuitable for labelling the original
62  manually labelled data, one reason, because the label distributions are unsimilar (3:1 in the
63  original data compared to 1:1 in the clustered data).
64  It can be noted from the count plot below that *cluster 1* is mostly made up of residents with
65  income less than or equal to 50k/year and *cluster 0* contains 86% of the residents whose *income* is
66  greater than 50K/year.



Distribution of cluster between income labels

67
68  Taking a closer look at some of the features with the highest absolute correlation with *income*
69  (shown in the figure on the left below), there's some relatively high correlation between being
70  married to a civilian spouse and earning more than 50K/year and it's no surprise that most of these
71  residents are grouped in *cluster 0* as shown in the figure on the right below. Also, having a
72  *marital-status* of *never-married*, has some correlation with not earning more than 50K/year and
73  it's no surprise that nearly all these instances are grouped in *cluster 1*, the cluster containing the
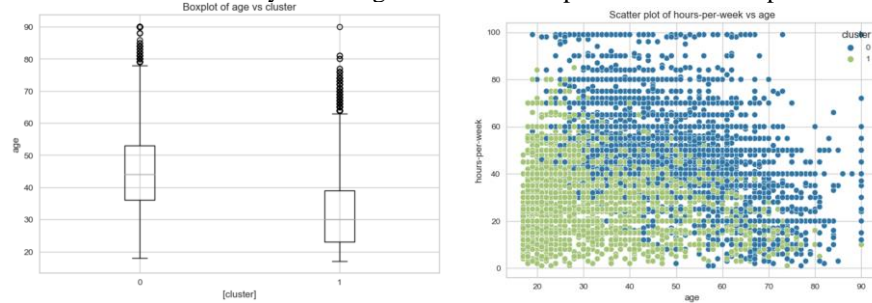74  larger proportion (61.7%) of the residents not earning more than 50K/year.



Top 6 features with the highest absolute correlation with income



Distribution of marital-status between clusters

75
76  As shown below, another relatively strong indicator of earning more than 50K/year, being a
77  *Husband*, is mostly contained in *cluster 0,* the cluster containing 86% of the residents earning
78  more than 50K/year.



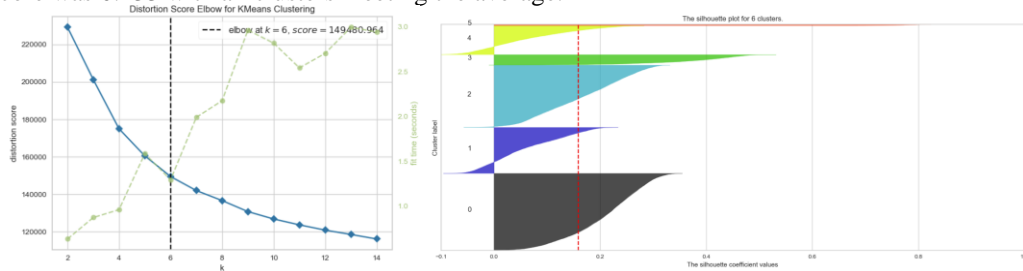Distribution of relationship between clusters

79

80   Other notable mentions are working *hours-per-week* and *age*; the higher the working hours and
81   age, the more likely they are in *cluster 0*. The feature *education-num* displays no visible pattern in
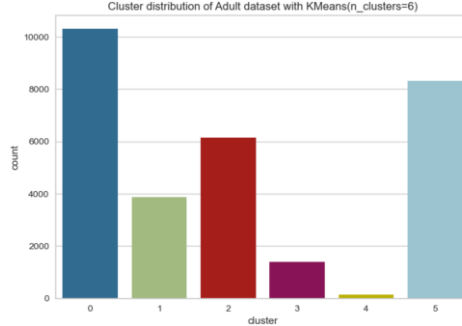82   the data.
83   <u>N.B</u>: The features were scaled before clustering. The original scale is retained for the x and y ticks
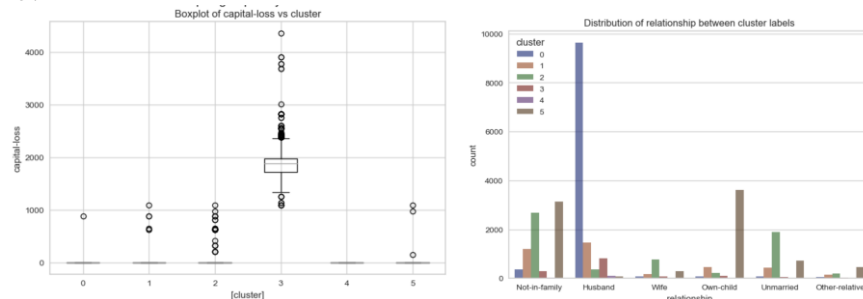84   of the plots below for easier analysis. Using the scaled ticks produces the same pattern in the data.



85
86   Moving on from n_clusters = 2 and looking for the optimal k using elbow and silhouette plots, the
87   optimal k was 6 and 5 respectively. The silhouette average scores were relatively low with the
88   highest (between two and nine clusters) being 0.162. For n_clusters = 2, the average silhouette
89   score was 0.133 with all clusters meeting the average.



90
91   I chose to experiment with k = 6, which had an average silhouette score of 0.159. The clusters
92   were distributed as follows, with *cluster 4* having the lowest number of instances, 148.
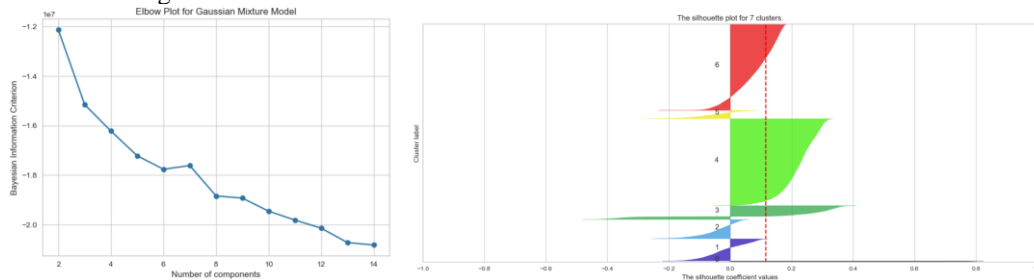


93
94   The most notable patterns were residents in *cluster 3* having *capita-loss* way higher than other
95   clusters, *cluster 0* mostly made up of *Husband* in a *relationship*, and *Own-child* mostly contained
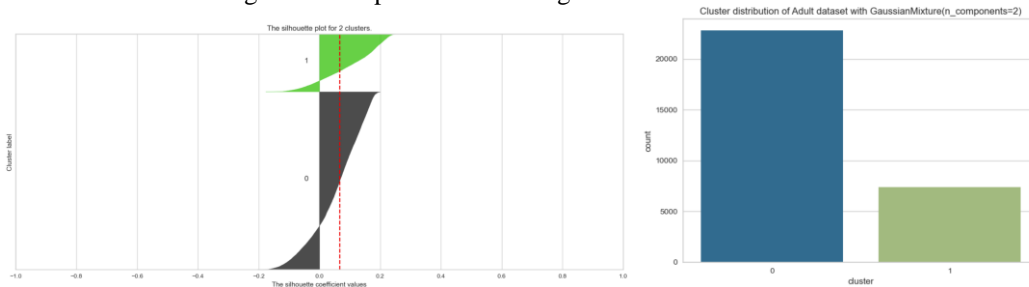96   in *cluster 5*.



97
98
99
100

101 **2.3      Gaussian Mixture Model (GMM)**
102 To search for the optimal number of components for the GMM, at first, I looked at an elbow plot
103 with a range from two to nine components. The Bayesian Information Criterion(BIC) didn't look
104 like levelling off within the selected range, so I widened the range to 14 and there was still no sign
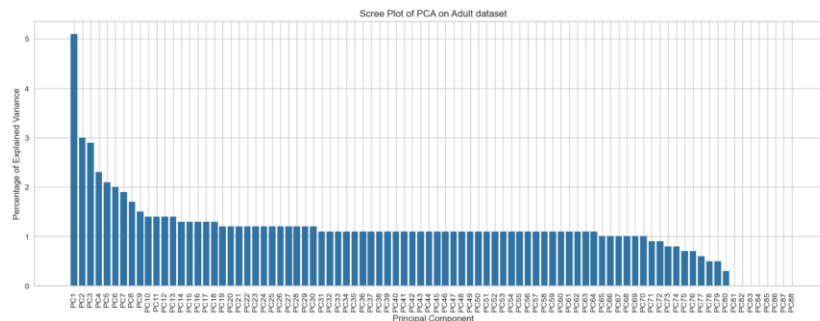105 of BIC levelling off.



106
107 There were similar complications when using the silhouette plot as well. Within the range of two
108 and nine, using seven components had the highest (0.116) silhouette score but as shown in the
109 silhouette plot above, it is a bad pick because at least two clusters do not meet the average
110 silhouette score. Also, the average silhouette score did not seem to have levelled off within the
111 range of two and nine; eight n_components dropped to 0.010 average silhouette score while nine
112 n_components jumped to 0.072. To save time, I decided to experiment with n_components = 2
113 which had the following silhouette plot and an average silhoette score of 0.007.



114
115 Evidently, there is a high imbalance (3:1) very similar to the proportion of the original target
116 labels. Comparing the cluster labels with the original labels shows that 95% of the residents who
117 earned more than 50K/year were contained in the first cluster.
118
119 **2.4      Principal Component Analysis (PCA)**
120 Before applying PCA to my dataset, I scaled every feature using scikit-learn's StandardScaler to
121 ensure each feature was on the same zero mean and unit variance scale and to avoid PCA being
122 biased towards any one feature. I applied PCA to the scaled data and got the same number of
123 components as original features, 88. The scree plot below shows how much variation each
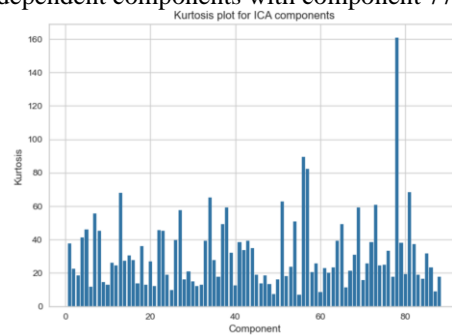124 component accounts for.



125
126 The last eight components, together, account for less than 1% of variation in the dataset. The
127 clustering result of retaining all the principal components was the same with the clustering result
128 of the original dataset.
129 It seems the initialization of the PCA influences how much variation the resulting components will
130 account for because when I initialised PCA with n_components = 10, I noticed that the resulting

131    ten components accounted for a slightly higher variation (0.35%) than did the first ten components
132    of the PCA with n_components set to the default "None". On the other hand, when I initialized
133    PCA with n_components = 0.90, I got the same 67 components with the first 67 components of the
134    PCA with n_components set to the default "None". Using only these 67 components and
135    clustering my data with GMM(n_component=2), I got the same cluster labels as when I kept all
136    components. Using GMM(n_components =3) yielded different results; only 12,343 instances had
137    the same labels. I also checked and confirmed that this was not a case of *cluster 0* in one result
138    being *cluster 2* in another result; the data distribution was different among clusters for both results.
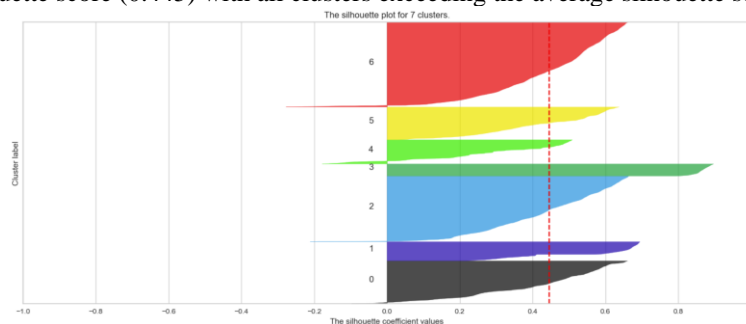139

140 **2.5      Independent Component Analysis (ICA)**
141    Fitting the FastICA with its default parameters yielded a warning with the description, "FastICA
142    did not converge. Consider increasing tolerance or the maximum number of iterations." I
143    increased the maximum iterations to 1000 and the tolerance to 0.001 but I got the same warning.
144    Adding *whiten = False* also yielded the same warning. I compared GMM(n_components=2)
145    clustering of FastICA(max_iter=1000, tol=0.001, random_state = 26, whiten=False) fitted and
146    transformed data with FastICA() fitted and transformed data. The results were different. The
147    former yielded the same clustering results as the original data. The following plot the distribution
148    of kurtosis values of the independent components with component 77 having the highest value.
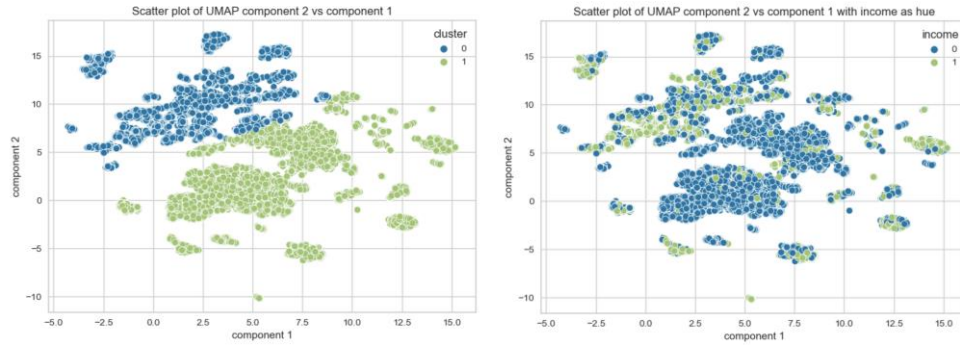


Kurtosis plot for ICA components

149
150 **2.6      Uniform Manifold Approximation and Projection (UMAP)**
151    After reducing my data's dimensions to two using UMAP, I combined silhouette plot and GMM
152    to find the optimal number of components betweeen two and nine. Seven clusters had the highest
153    average silhouette score (0.445) with all clusters exceeding the average silhouette score.



The silhouette plot for 7 clusters.

154
155    I used these number of clusters in running experiments with the neural network.
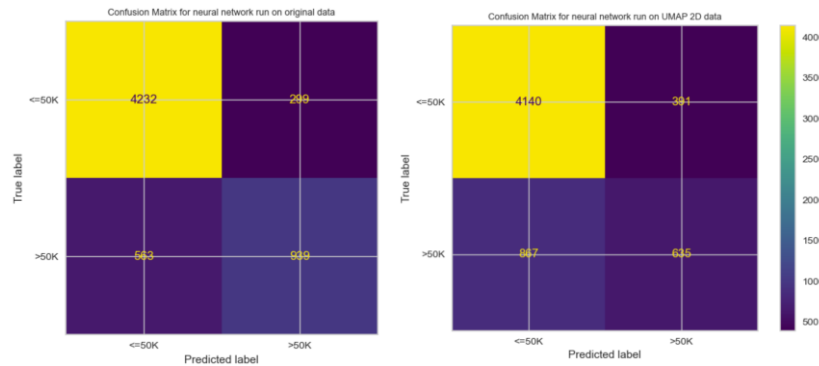156
157    The plot on the left below shows the result of using GMM to group 2D UMAP data into two
158    clusters. The plot on the right shows the same 2D data with the original target variable as the hue.

159
160
### 2.7    Neural Network

161
162
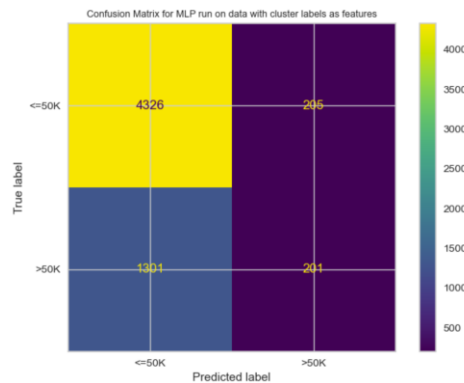#### 2.7.1    Results of Neural Network on UMAP 2D data

163
164 An MLP with these parameters ({'solver': 'adam', 'hidden_layer_sizes': (100,), 'alpha': 0.1,
165 'activation': 'relu'}) yielded an f1-score of 0.69 for the original data and 0.50 on the UMAP 2-D
166 data. The following confusion matrix shows the classification results on both datasets.

167



168
#### 2.7.2    Result of Neural Network run on data with cluster labels as features

169
170 I used the result of GMM grouping the UMAP 2D data into seven clusters for this part of the
171 experiment. I one hot encoded the cluster labels and got data frame shown in the picture on the left
172 below.



173
174 Running the neural network learner on the data yielded an f1-score of 0.21. The results of the
175 classification is shown in the confusion matrix above.
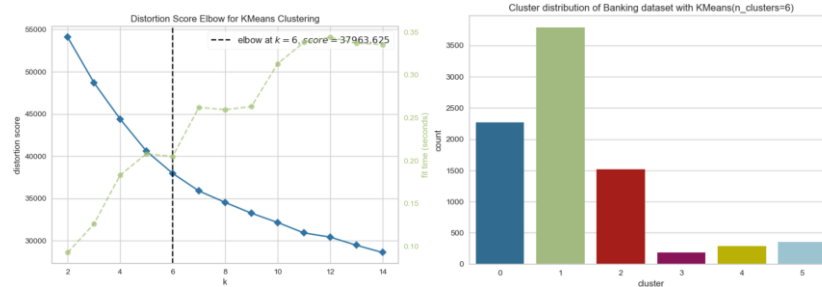
176
177
178
179
180

## 3  Banking data

### 3.1  Preprocessing
As in the adult dataset, I dropped the instances with unknown values and was left with 8,393 instances. I one-hot encoded the categorical variables, scaled the numerical variables with StandardScaler and used MinMaxScaler for the the one ordinal variable contained in the data. Before applying PCA, I scaled all the features using StandardScaler.
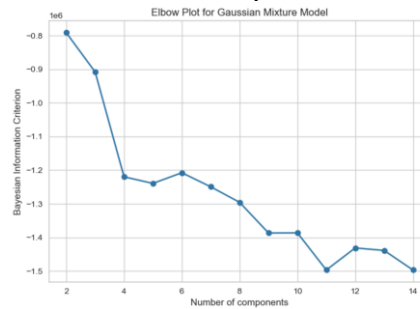
### 3.2  K-means
Running K-Means (n_clusters=2) on the data yielded clusters distribution very different (73:27) from the original data's target label distribution (54:46).
Based on the result of the elbow plot in the diagram on the left below, I ran a K-Means(n_components = 6) on the dataset which yielded the cluster distribution in the diagram on the right below.
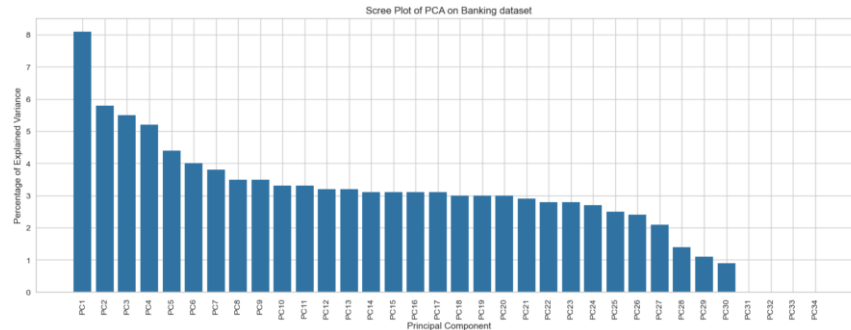


### 3.3  GMM
A GMM(n_components=2) on the dataset yielded a (65:35) data distribution between the clusters. Like in the case of the adult dataset, the BIC for the elbow plot didn't seem to level off in the range of k equals two and 14 as shown in the elbow plot below.
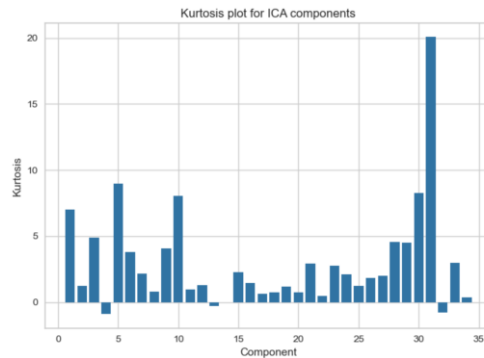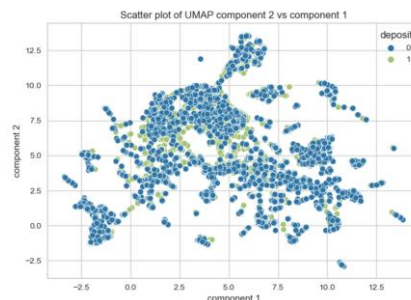


### 3.4  PCA
Applying PCA to the dataset yielded the following plot showing the components and the amount of variation they account for. The last five components account for less than 1% variation.

209 **3.5    ICA**
210 Applying ICA to the data yielded the following kurtosis for each component. Notably, they result
211 in a lower average kurtosis than those derived with the adult dataset.


Kurtosis plot for ICA components

212
213 **3.5    UMAP**
214 Reducing the data to two dimensions and visualizing it like in the diagram below show no clear
215 pattern in the data.


Scatter plot of UMAP component 2 vs component 1

216
217

218 **4    Conclusion**

219 In this report I have applied five unsupervised machine learning algorithms
220 to two datasets.

221 **Acknowledgments**

222 The learning code was produced with multiple discussions with OpenAI's ChatGPT.

223 **References**

224  [1] Chen, J. (2022, March 20). Time Deposit: Definition, How It's Used, Rates, and How to Invest

225 Investopedia.
226 https://www.investopedia.com/terms/t/termdeposit.asp#:~:text=What%20Is%20a%20Term%20Deposit,level
227 s%20of%20required%20minimum%20deposits.

228