

Analysis of the dataset *Credit Card Default*

Gaspare Ferraro (520549)
Javad Khalili (546677)
Mario Matovic (583449)

January 8, 2019



University of Pisa
Exam: Data Mining
Year: 2018/2019
Instructors: Dino Pedreschi, Anna Monreale, Riccardo Guidotti

Contents

1	Introduction	2
2	Data Understanding	2
2.1	Data semantics, distribution and statistics	2
2.2	Assessing data quality	5
2.3	Variables transformations	5
2.4	Correlations and redundant variables	5
3	Clustering	6
3.1	K-means	6
3.1.1	Choice of attributes and distance function	6
3.1.2	Choise of the best value of k	6
3.1.3	Cluster analysis	6
3.2	DBSCAN	8
3.2.1	Choice of attributes and distance function	8
3.2.2	Study of the clustering parameters	8
3.2.3	Characterization and interpretation of the obtained clusters	8
3.3	Hierarchical clustering	8
3.3.1	Choice of attributes and distance function	8
3.3.2	Discussion of dendograms using different algorithms	8
3.4	Evaluation of clustering approaches and comparison of the clustering obtained	8
4	Association Rules Mining	8
4.1	Frequent patterns extraction with different parameters	9
4.2	Discussion of the most interesting frequent patterns	9
4.3	Association rules extraction with different values of confidence	9
4.4	Discussion of the most interesting rules	9
4.5	Use the most meaningful rules to replace missing values	9
4.6	Use the most meaningful rules to predict credit card defaults	9
5	Classification	9
5.1	Choice of attributes for the decision trees	9
5.2	Decision trees interpretation and validation with test and training set	9
5.3	Discussion of the best prediction model	9
6	Conclusion	9

1. Introduction

This report is aimed to illustrate the phases and the results of the analysis that we have conducted regarding the customers default payments in Taiwan. In particular our target is to better understand under which conditions we should consider a client credible or not.

Each customer is modelled by a record in the dataset, which is composed by 24 attributes that describes its personal information and its banking data (like the credit limit, payment amount and others).

The analysis is composed in 4 phases:

- Semantical analysis and data manipulations of each customer (data cleaning, variables transformation, redundant variables, ...)
- Use of 3 clustering algorithms (K-Means, DBSCAN, Hierarchical clustering) to group customers according to similarity properties in order to formulate hypotheses about customers credibility.
- Verification of the hypotheses given in the previous phase and determination of association rules, in order to find co-occurrences between attributes.
- Classification of the customers between who fail to make a payment by time and regular customers.

2. Data Understanding

The dataset is composed by 10000 records. Each record represents a customer, described by 24 different attributes.

2.1 Data semantics, distribution and statistics

In this section we will analyze, for each attribute, its semantic and we will show interesting statistic and plot. We have used two different colors for who went in credit default (green) and not (red) in order to better visualize their distribution among the different attributes.

We have discretized the continuous attributes by using the natural binning method.

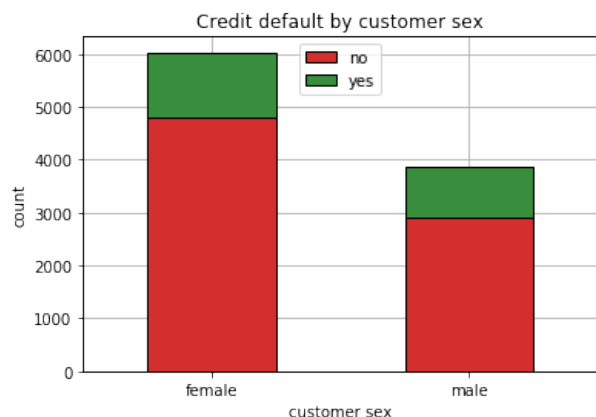
For these attributes the mode of the bin has also been reported as it is more representative.

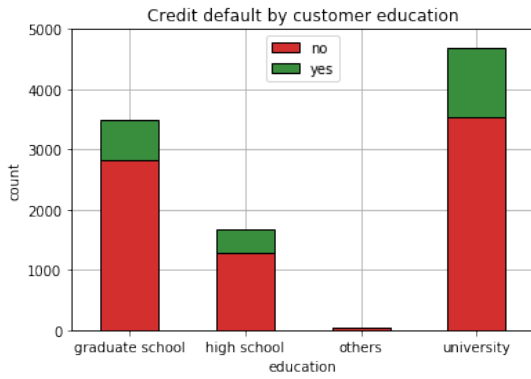
Sex

Gender of the customer.

A binary attribute that can assume the values of *male* (3868 of 10000) or *female* (6032 of 10000).

Both of the gender values have a similar credit default rate (25% for males and 20% for females).





Education

Qualification of the customer.

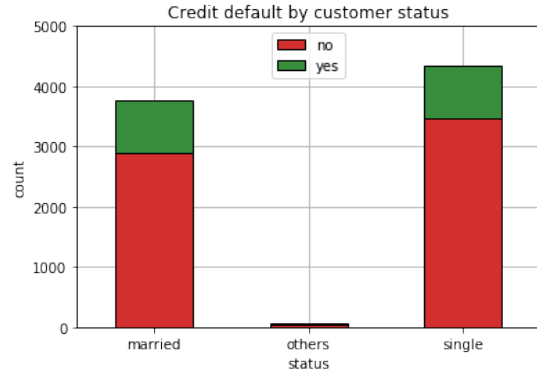
A categorical attribute that can assume the values of *university* (4685 of 10000), *high school* (1672 of 10000), *graduate school* (3480 of 10000) or *others* (36 of 10000). The default rate is again very similar for all the qualifications (around the 20%), except for the *others* which is equal to 5%, but its number of records is very low to make any assumptions.

Status

Marital status of the customer.

A categorical attribute that can assume the values of *married* (4685 of 10000), *single* (3757 of 10000) or *other status* (75 of 10000).

The default rate is very similar for all the status (around the 25%).



Age

Age of the customer.

An attribute that in the dataset assume integer values in [21, 75], the lower limit to 21 is due that in Taiwan the age of majority is 20, on the other hand the upper limit can be as high as humanly possible. The average age is 35.49, the standard deviation is 9.22. The mode is 29 and the median is 34. The 50% of the ages lie in [28, 41].

We decide to set a bin to 5 year as it represent a good trade-off between the size and number of bins. The bin with most elements is 25.

Again the default rate is similar for all the bins (around 25%).

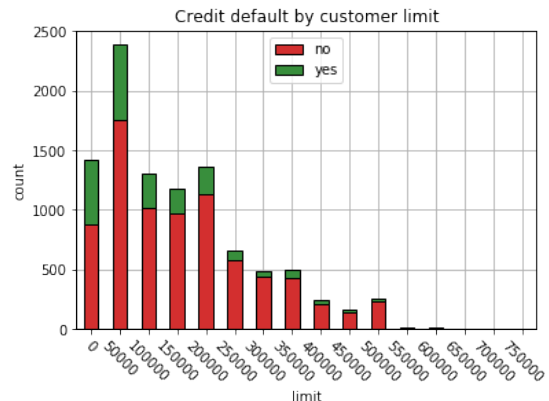
Limit

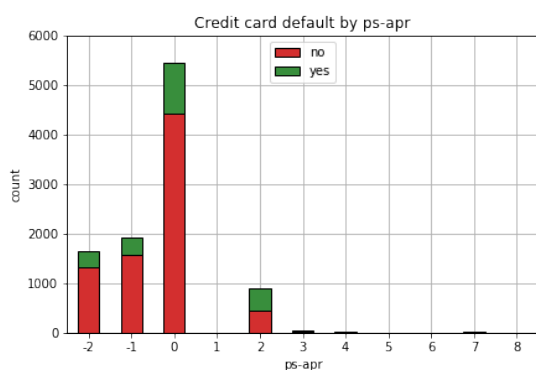
Limit of the credit card (expressed in NT dollar).

It is the maximum amount the credit card company will let borrow on the account, a continuous attribute that can assume values in [10000, 780000] (all values are multiples of 10000).

The average is 167197 and the standard deviation is 128975, 50% of the ages lie in [50000, 240000]. The bin with most elements is 50000.

The default rate in this case is very different for each bin as it decrease with the limit: the first bin has a default rate of 38%, the second one of 26% and around 10% for the last bin.





Payment status

History of past payments.

Six categorical attributes that represent the repayment status, one for each month between April and September. A payment status is an integer number in the range $[-2, 9]$ where:

-2 = no consumption

-1 = paid in full

0 = the use of revolving credit

1 = payment delay for one month

2 = payment delay for two months

...

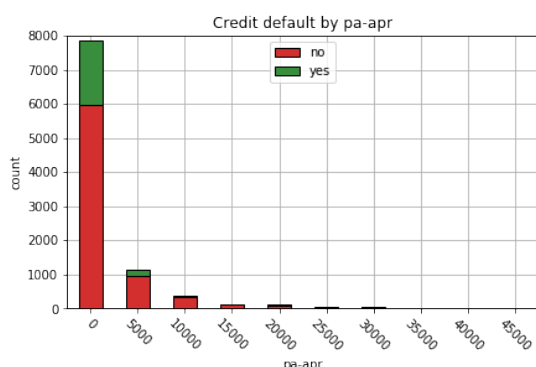
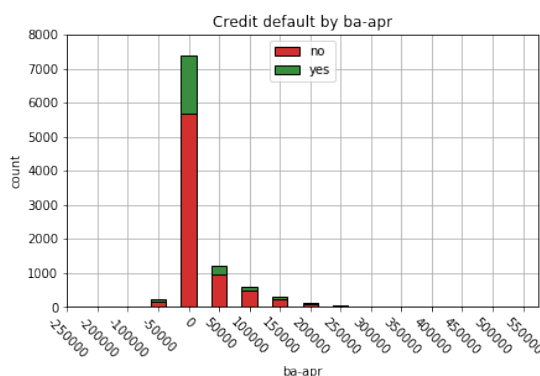
The distribution of the payment status over the six months are all very similar to each other so we only illustrate the first one for simplicity. The most frequent bin is always the 0 and around 90% of the customers always lies between -2 and 0.

Bill Amount

Amount of bill statement (expressed in NT dollar).

Six continuous attributes, again one for each month between April and September.

The bill amount is the sum of all the purchases, payments and other debits and credits made to your credit card account within the billing cycle. The values are in range $[-209051, 616836]$. As in the payment status, the distribution of the bill amount are nearly the same over the six month. The most frequent by is always the 0 (from 0 NTD to 49999 NTD). The mean increases according to the month, it start from 38621 of April up to 51490 of September. On the other hand the credit default rate is always around the 20% for all the bins (except some border case for small bin).



Payment amount

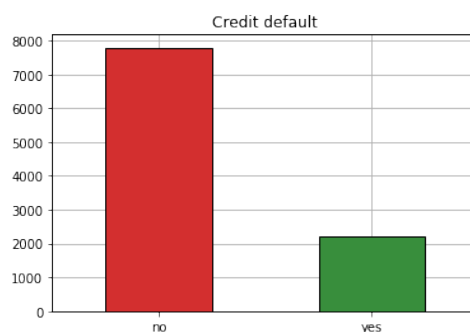
Amount of payment (expressed in NT dollar).

Six continuous attribute, one for each month between April and September. It represent the sum of all the payments made by the customer with the credit card in the last month. The distribution of the payment status over the six months are again all very similar to each other. The most frequent bin is always the 0 (from 0 NTD to 4999 NTD), which always include around the 75% of the customers. The mean is around 4719 and 5973 without any particular distribution other the months.

Credit default

Credit card default.

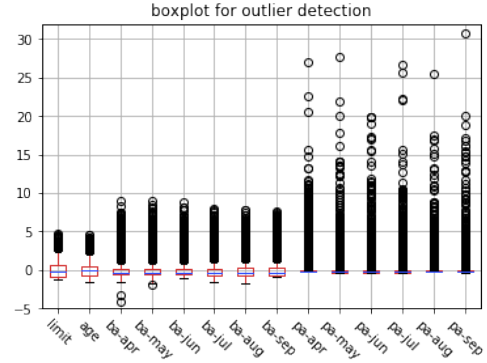
A boolean attribute that represents the credibility of the customer. The credit card default is a status that is applied to the customer when a list of terms are missed (like make the minimum required payment), and it affects for example your possibility to demand for a loan or to get another credit card. In the dataset we have 2212 customer in credit card default and 7788 not.



2.2 Assessing data quality

In this part we have checked for duplicates, missing values and outliers, in order to reduce the quantity of data and to avoid that their presence could negatively affect the analysis.

Regarding the duplicates, in the dataset there are only two identical records so we decide to not remove the single duplicate as it is not a big problem. Regarding the missing values we have 100 records with missing sex, 127 records, 1822 records with missing status and 951 records with missing age. We decided to not remove them as they are a big part of the dataset, instead we replace the missing ages with the average age (which is 35.49). We set the missing statuses to *others* as we think that it can include some particular cases (like engaged or cohabitant). We again set the missing educations to *others* and the missing sex to female as they are the more frequents in the dataset. For the outliers we first normalize all the continuous attributes using the z-score algorithm (in this way we can visualize all the attributes in the same scale). Then we plot them in a boxplot and we can see that in the dataset only 3 values are under the lower adjacent value, instead over the upper adjacent value we have 50 limits, 211 limits and an average of 895 for the bill/payment amounts. We decided to not consider them outliers as we think that they represent the small percentage of rich people in the dataset (as they an high bill/payment amounts).



2.3 Variables transformations

We decide to transform the *credit_default* attribute from yes/no to 1/2 in order to better visualize them in the following plots. For the same reason and to limit the dispersion we decided to divide in bin the continuous attributes. In particular we divide ages in bin of 5 years, payment amounts in bin of 1000 NTD, limits and bill amounts in bin of 10000 NTD.

As last transformation we have normalized all the continuous attribute using the min-max scaling.

2.4 Correlations and redundant variables

Analyzing the correlation matrix of all the continuous attribute we can clearly see that all the attributes related to the bill amounts are strongly correlated.

All the pairs of bill amounts have a correlation between 0.8 and 0.95 so we decide to analyze further this 6 attributes.

So we add two new attributes to the dataset called *ba_mean* and *pa_mean* calculated as the mean of all the bill/payment amount for each customer.

We again plot the correlation matrix and we can see that both *ba_mean* and *pa_mean* have a high correlations with their relative attributes. For this reasons we drop all the twelve bill/payment amounts attributes and introduce their mean as new attributes.

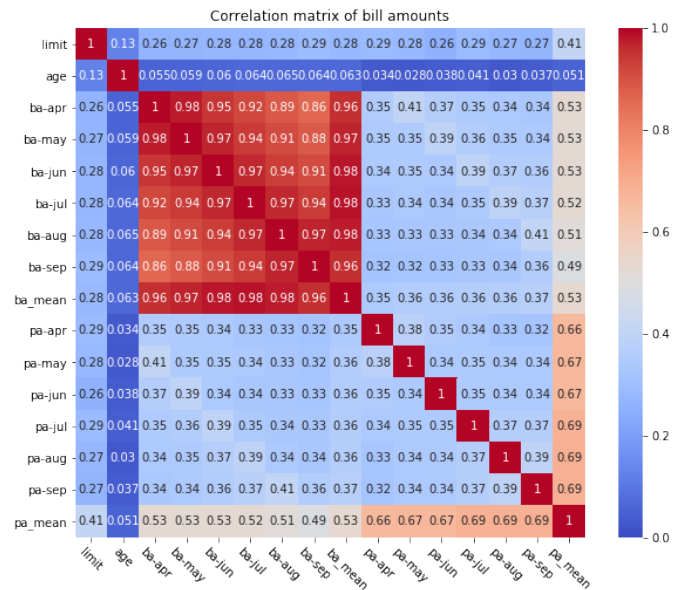


Figure 2.1: correlation matrix with *ba_mean*

3. Clustering

3.1 K-means

3.1.1. Choice of attributes and distance function

We have selected 3 attributes to perform the clustering via the K-Means algorithm. In particular we have ignored the categorical attributes as there isn't a proper metric to define a distance function over these kinds of attributes.

The used attributes are only: *limit*, *ba_mean*, *pa_mean*. Note that we have excluded the age in this list as it is the only one that isn't measured in NTD.

For these attributes we used the euclidean distance as metric function as it is more natural for this kind of attributes.

3.1.2. Choise of the best value of k

To choose the best value of k we have plotted for each K between 2 and 50 the values of the SSE and the silhouette score. Each value is calculated as the best among 10 runs (with the lowest SSE) and with a maximum of 300 iterations of the K-Means algorithm.

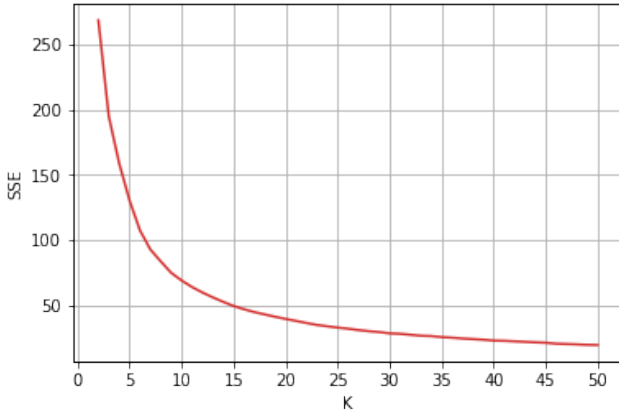


Figure 3.1: SSE value for each K .

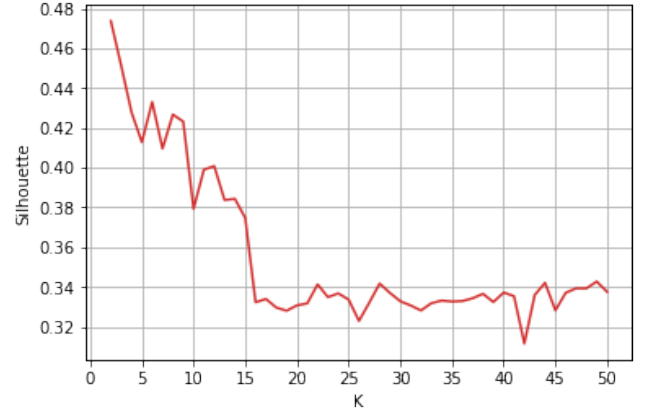


Figure 3.2: silhouette score for each K .

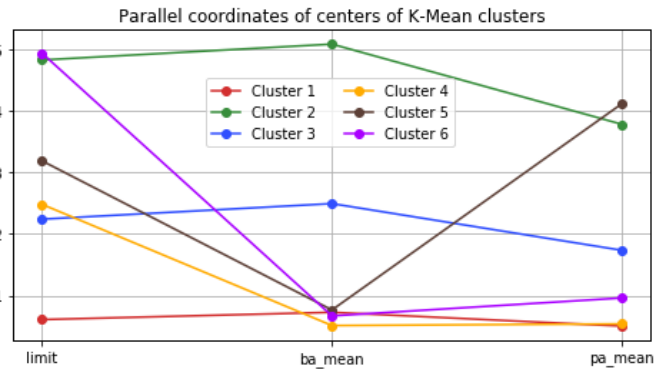
Using the elbow method we decide to set $k = 6$ as it represents a good trade-off between SSE and data interpretability.

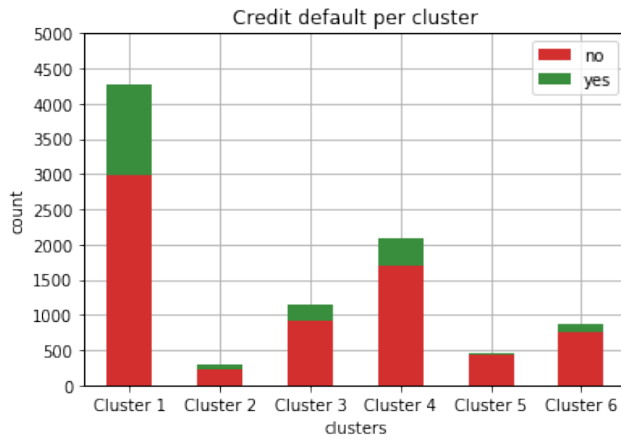
3.1.3. Cluster analysis

To better visualize the obtained clusters we have plotted the centers in parallel coordinates.

From a first point of view we can see that all of the centers are separated in three class (low, medium, high) for all of the three used attributes.

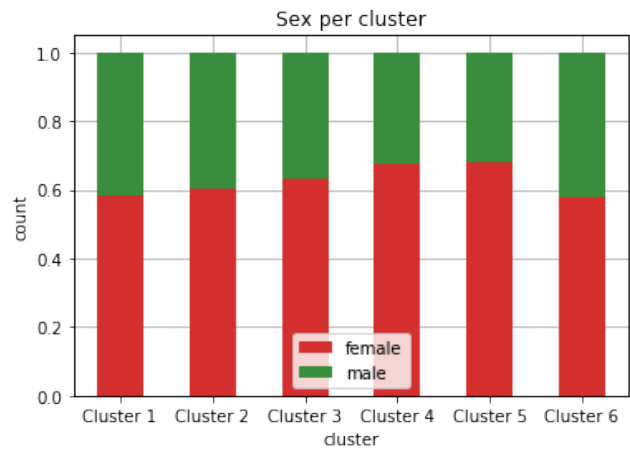
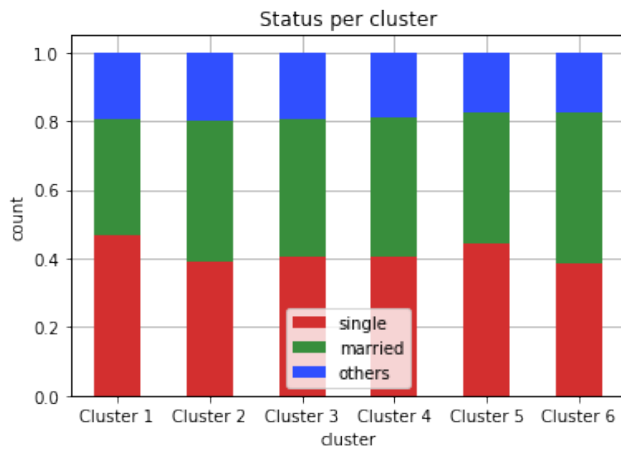
We want to go further in order to understand the different kind of customers and if there are some interesting correlation.



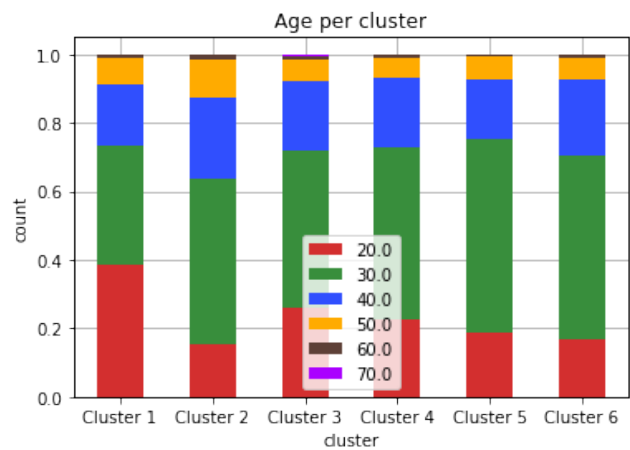
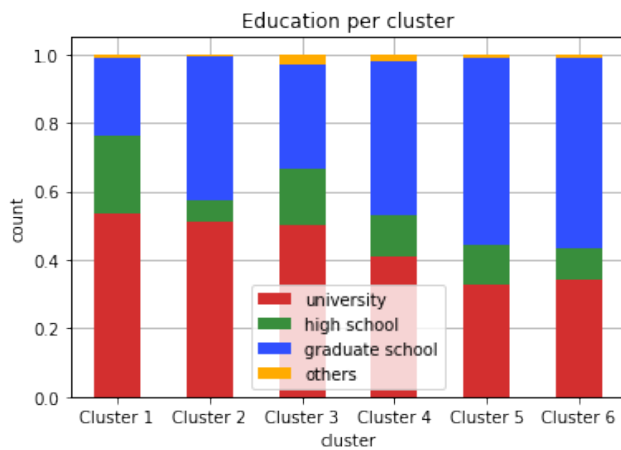


We have also plotted the number of credit defaults for each cluster to trying to understand the distribution of the customers. The cluster with the most credit default is the first one, with a ratio of 30% and the cluster with the lowest credit default is the fifth, with a ratio of 6%.

From the figures below we can see that the distributions of sex and status are the same for all the clusters, in particular their distributions are the same in the whole datasets so we can't use these attributes to make any assumptions over the obtained clusters.



More useful are instead the distributions of the attributes relative to the educations and ages:



We can discuss each cluster individually to see the most interesting properties.

Cluster 1 is the biggest cluster, with 4276 total elements include nearly half of the customers, and it is also the cluster with the highest credit default ratio (40%). From the parallel coordinates plot we can see that the customers in this class have the credit cards with the lowest limit and also have the lowest bill/payment amounts over the six months. The lowest limit could be that this cluster is the one with most of young people (40% are below the 30), but also is the cluster with most customer with an education in university and in high school. The high positive credit default ratio could be a combination of the distribution of these attributes (younger and university/high school).

Cluster 2 with a total of 283 customers is the smallest cluster. It has a credit default ratio of 21% and it is the cluster with the older customers, nearly 40% of the customers have an age higher that 40. From the parallel coordinates we can see that is the only cluster with high limit and high bill amount and payment amount mean.

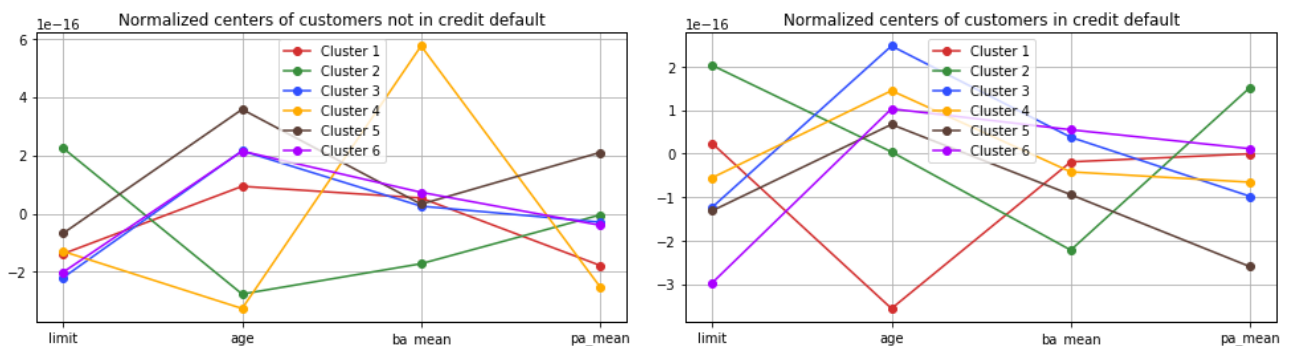
Cluster 3 has a total of 1143 customers and a positive credit default ratio of 19%. In this cluster we have customers with a medium value of all the attributes.

Cluster 4 with a total of 2083 customers is the second cluster for size and is has a positive credit default ratio of 21%. In this cluster we have customers with a medium credit card limit and with the lowest bill/payment amount mean.

Cluster 5 has a total of 457 customers and is the cluster with the lowest positive credit default ratio (only 6%). In this cluster we have customers with a credit card of medium limit, with a low bill amount mean but with the highest payment amount mean. From a point of view of the education in this cluster (and in the last cluster) we have the lowest number of education in university but the highest number of education in graduate school.

Cluster 6 has a total of 879 customers. It is the cluster with the highest limit in the credit cards but it has a low bill/payment amount mean.

We decided to analyze further the obtained cluster, so we separate the customers with a positive credit default from the customers with a negative credit default and plot all the normalized centers of the continuous attributes in two separated parallel coordinates plots:



The most interesting cluster are **Cluster 1** and **Cluster 5**, which are the ones with the highest and lowest positive credit default ratio, so we decided to analyze only these two clusters.

For the **Cluster 1** the customers in credit defaults have an higher credit card limit and payment amount mean but a lower average age (compared to the customers not in credit default from the same cluster).

For the **Cluster 5** the customers in credit defaults have an higher credit card limit and payment amount mean but a lower average age (compared to the customers not in credit default from the same cluster).

3.2 DBSCAN

3.2.1. Choice of attributes and distance function

3.2.2. Study of the clustering parameters

3.2.3. Characterization and interpretation of the obtained clusters

3.3 Hierarchical clustering

3.3.1. Choice of attributes and distance function

3.3.2. Discussion of dendograms using different algorithms

3.4 Evaluation of clustering approaches and comparison of the clustering obtained

4. Association Rules Mining

4.1 Frequent patterns extraction with different parameters

Cose

4.2 Discussion of the most interesting frequent patterns

Cose

4.3 Association rules extraction with different values of confidence

Cose

4.4 Discussion of the most interesting rules

Cose

4.5 Use the most meaningful rules to replace missing values

Cose

4.6 Use the most meaningful rules to predict credit card defaults

Cose

5. Classification

Classificazione

5.1 Choice of attributes for the decision trees

cose

5.2 Decision trees interpretation and validation with test and training set

cose

5.3 Discussion of the best prediction model

cose

6. Conclusion

Conclusione dove parlo di cose.