

"German Credit" scoring data analysis report

Marta Karaś, Jan Idziak

June 11, 2015

Table of content

I	Introduction	4
1	Data analysis context	4
1.1	Motivation	4
1.2	Research questions	4
II	Materials and methods	5
2	Data set	5
2.1	Data set description	5
3	Binning countinuous variables	7
3.1	Weight of Evidence (WoE)	7
3.2	Information Value (IV)	7
3.3	Motivation	8
4	Feature selection	9
4.1	Feature selection algorithms	10
4.1.1	Algorithms for filtering attributes	10
4.1.2	Algorithms for wrapping classifiers and search attribute subset space	12
5	Classification	13
5.1	Classification algorithms	13
5.2	Classification performance metrics	15
5.2.1	Separation measures	16
6	Principal Components Analysis	17
7	Cluster analysis	18
7.1	Introduction	18
7.2	Similarity, Dissimilarity and Distance	19
7.3	Clustering Algorithms	20
7.3.1	UPGMA	20
7.3.2	K-means	20
7.3.3	Diana	21
7.3.4	PAM	21
7.3.5	Clara	21
7.3.6	Fanny	21
7.3.7	SOM	21
7.3.8	Model based clustering	21
7.3.9	SOTA	22
7.4	Evaluation of clustering	22
7.4.1	Internal criterion	22
7.4.2	External criterion	23

7.4.3	Cluster stability	24
III	Results	25
8	Data preprocessing results	26
8.1	Searching for missing, corrupt and invalid data	26
8.2	Creating derived variables	30
8.3	Binning continuous variables	31
8.4	Correcting bins (levels) of categorical variables	32
8.5	Removing variables of 0.0 Information Value	33
8.6	Recoding variables to WoE	34
8.7	Data preprocessing summary	34
9	Classification modelling results	36
10	Cluster analysis results	37
11	Summary	38
IV	Discussion	39

```
## [1] "janek"
```

Part I

Introduction

In this part of the report we provide answers to the following questions about the "German Credit" data analysis we performed.

1. *Why was the study undertaken?*
2. *What was the purpose of the research? What research questions were stated?*

1 Data analysis context

1.1 Motivation

This report presents results of the "German Credit" scoring data analysis which was performed as a project assignment for the "Pozyskiwanie Wiedzy" course, which we attended at Wroclaw University of Technology, Faculty of Fundamental Problems of Technology (W-11), Mathematics program (Master) in the 2014/15 summer semester. The lecturer of the course (both lectures and laboratories) is Ph.D. Adam Zagdański.

The main goal of the project is to make use of the variety of data-mining methods we have become familiar with during the course, in order to perform complete data analysis of selected data set. We also aim to pay attention to the practical applications of some parts of our work.

1.2 Research questions

We stated the following research purposes for our analysis.

1. Find and describe relations in the data (relations between explanatory variables and response variable, relations between explanatory variables).
2. Compare different methods / algorithms to perform exploratory data analysis and predictive data analysis.
3. Provide a summary of the analysis, containing suggestions of practical application and remarks regarding possible further research.

Part II

Materials and methods

In this part of the report we describe the data set we obtained and the methods we use in the analysis.

This section is rather of the decriptional / theoretical character. For a list of actual analysis steps, the outputs of the methods and more, please refer to the III part of this report.

2 Data set

We perform analysis with the use of The (Statlog) German Credit Data we have obtained from the UCI Machine Learning Repository site.

2.1 Data set description

The data set contains data on 20 variables and the classification whether an applicant is considered a Good or a Bad credit risk for 1000 loan applicants. The file provided contains variables with values encoded according to the following schema:

- Attribute 1: (qualitative) Status of existing checking account
A11 : ... < 0 DM
A12 : 0 <= ... < 200 DM
A13 : ... >= 200 DM / salary assignments for at least 1 year
A14 : no checking account
- Attribute 2: (numerical) Duration in month
- Attribute 3: (qualitative) Credit history
A30 : no credits taken/ all credits paid back duly
A31 : all credits at this bank paid back duly
A32 : existing credits paid back duly till now
A33 : delay in paying off in the past
A34 : critical account/ other credits existing (not at this bank)
- Attribute 4: (qualitative) Purpose
A40 : car (new)
A41 : car (used)
A42 : furniture/equipment
A43 : radio/television
A44 : domestic appliances
A45 : repairs
A46 : education
A47 : (vacation - does not exist?)
A48 : retraining
A49 : business
A410 : others

- Attribute 5: (numerical) Credit amount
- Attribute 6: (qualitative) Savings account/bonds
 - A61 : ... < 100 DM
 - A62 : 100 <= ... < 500 DM
 - A63 : 500 <= ... < 1000 DM
 - A64 : .. >= 1000 DM
 - A65 : unknown/ no savings account
- Attribute 7: (qualitative) Present employment since
 - A71 : unemployed
 - A72 : ... < 1 year
 - A73 : 1 <= ... < 4 years
 - A74 : 4 <= ... < 7 years
 - A75 : .. >= 7 years
- Attribute 8: (numerical) Installment rate in percentage of disposable income
- Attribute 9: (qualitative) Personal status and sex
 - A91 : male : divorced/separated
 - A92 : female : divorced/separated/married
 - A93 : male : single
 - A94 : male : married/widowed
 - A95 : female : single
- Attribute 10: (qualitative) Other debtors / guarantors
 - A101 : none
 - A102 : co-applicant
 - A103 : guarantor
- Attribute 11: (numerical) Present residence since
- Attribute 12: (qualitative) Property
 - A121 : real estate
 - A122 : if not A121 : building society savings agreement/ life insurance
 - A123 : if not A121/A122 : car or other, not in attribute 6
 - A124 : unknown / no property
- Attribute 13: (numerical) Age in years
- Attribute 14: (qualitative) Other installment plans
 - A141 : bank
 - A142 : stores
 - A143 : none
- Attribute 15: (qualitative) Housing
 - A151 : rent
 - A152 : own
 - A153 : for free
- Attribute 16: (numerical) Number of existing credits at this bank

- Attribute 17: (qualitative) Job
A171 : unemployed/ unskilled - non-resident
A172 : unskilled - resident
A173 : skilled employee / official
A174 : management/ self-employed/ highly qualified employee/ officer
- Attribute 18: (numerical) Number of people being liable to provide maintenance for
- Attribute 19: (qualitative) Telephone
A191 : none
A192 : yes, registered under the customers name
- Attribute 20: (qualitative) Foreign worker
A201 : yes
A202 : no

The classification variable states whether there was a default case ('bad' client - failed to pay off the credit) or not ('good' client).

- Classification: (qualitative) Default
1 (default)
0 (non-default)

3 Binning continuous variables

In credit scoring, Information Value (IV) is frequently used to compare predictive power among variables. When developing new scorecards using logistic regression, variables are often binned and recoded using WoE concept.

3.1 Weight of Evidence (WoE)

One of our goals when binning variables is to maximize Information Value. Weight of Evidence (WoE) for single bin is defined as:

$$WoE = \left[\ln \left(\frac{\text{Relative Frequency of Goods}}{\text{Relative Frequency of Bads}} \right) \right] \times 100.$$

We can see that value of WoE will be 0 if the odds of Relative Frequency of Goods / Relative Frequency Bads is equal to 1. If the Relative Frequency of Bads in a group is greater than the Relative Frequency of Goods, the odds ratio will be less than 1 and the WoE will be a negative number; if the Relative Frequency of Goods is greater than the Relative Frequency of Bads in a group, the WoE value will be a positive number.

3.2 Information Value (IV)

We define Information Value of the variable as follow:

$$IV = \sum_{i=1}^k \left[(\text{Relative Frequency of Goods}_i - \text{Relative Frequency of Bads}_i) \times \ln \left(\frac{\text{Relative Frequency of Goods}}{\text{Relative Frequency of Bads}} \right) \right],$$

By convention the values of the IV statistic can be interpreted as follows. If the IV statistic is:

- Less than 0.02, then the predictor is not useful for modeling (separating the Goods from the Bads),
- 0.02 to 0.1, then the predictor has only a weak relationship to the Goods/Bads odds ratio,
- 0.1 to 0.3, then the predictor has a medium strength relationship to the Goods/Bads odds ratio,
- 0.3 or higher, then the predictor has a strong relationship to the Goods/Bads odds ratio.

3.3 Motivation

The WoE recoding of predictors is particularly well suited for subsequent modeling using Logistic Regression. Specifically, logistic regression will fit a linear regression equation of predictors (or WoE-recoded continuous predictors) to predict the logit-transformed binary Goods/Bads variable. Therefore, by using WoE-recoded predictors in logistic regression the predictors are all prepared and coded to the same WoE scale, and the parameters in the linear logistic regression equation can be directly compared. ([8])

4 Feature selection

Following [1], feature selection is essentially a task to remove irrelevant and/or redundant features. *Irrelevant features* can be removed without affecting learning performance. *Redundant features* are a type of irrelevant feature. The distinction is that redundant feature implies the co-presence of another feature; individually, each feature is relevant, but the removal of one of them will not affect learning performance.

The selection of features may be achieved in two ways:

1. **Feature ranking.** The idea is to rank features according to some criterion and select the top k features.
2. **Subset selection.** The idea is to select a minimum subset of features without learning performance deterioration.

In other words, subset selection algorithms can automatically determine the number of selected features, while feature ranking algorithms need to rely on some given threshold to select features.

The tree typical feature selection models are:

1. **Filter.** In a filter model, one selects the features firstly and then uses this subset to execute a classification algorithm.
2. **Wrapper.** In a wrapper model, one employs a learning algorithm and uses its performance to determine the quality of selected features.
3. **Embedded.** An embedded model of features selection integrates the selection of features in model building. An example of such model is a decision tree induction algorithm, in which at each branching node, a feature has to be selected.

In literature, various search strategies are proposed, including: forward, backward, floating, branch-and-bound, and randomized. A relevant issue, regarding exhaustive and heuristic searches is whether there is any reason to perform exhaustive searches if time complexity were not a concern. Research shows that exhaustive search can lead to the features that exacerbate data overfitting, while heuristic search is less prone to data overfitting in feature selection, facing small data samples.

The evaluation of feature selection often entails two tasks:

1. One is to compare two cases: before and after feature selection. The goal of this task is to observe if feature selection achieves its intended objectives. The aspects of evaluation may include the number of selected features, time, scalability and learning model's performance.
2. The second task is to compare two feature selection algorithms to see if one is better than other for a certain task.

4.1 Feature selection algorithms

In this subsection we describe methods for feature selection we use in our analysis. In general, we use the FSelector R package exhaustively. This package contains both algorithms for filtering attributes and algorithms for wrapping classifiers and search attribute subset space.

4.1.1 Algorithms for filtering attributes

CFS filter CFS is a correlation-based filter method CFS from [2]. It gives high scores to subsets that include features that are highly correlated to the class attribute but have low correlation to each other. Let *Attribute* be an attribute subset that has k attributes, rcf models the correlation of the attributes to the class attribute, rfc - the intercorrelation between attributes. We define *Attribute* score as:

$$CfsScore(Attribute) = \frac{k rcf}{\sqrt{k + k(k-1)rfc}}.$$

The algorithm from FSelector R package makes use of *Best-first search* for searching the attribute subset space. In *Best-first search*, the algorithm chooses the best node from all already evaluated ones and evaluates it. The selection of the best node is repeated approximately *max.brackets* times in case no better node found.

Chi-squared filter The algorithm evaluates the worth of an attribute by computing the value of the chi-squared statistic with respect to the class.

Information Gain filter One of the entropy-based filters. Algorithm evaluates the worth of an attribute by measuring the information gain with respect to the class.

$$InfoGain(Class, Attribute) = H(Class) + H(Attribute) - H(Class|Attribute),$$

where H is the information entropy.

Gain Ratio filter One of the entropy-based filters. Algorithm evaluates the worth of an attribute by measuring the gain ratio with respect to the class.

$$GainR(Class, Attribute) = \frac{H(Class) + H(Attribute) - H(Class|Attribute)}{H(Attribute)},$$

where H is the information entropy.

Symmetrical Uncertainty filter One of the entropy-based filters. Algorithm evaluates the worth of a set attributes by measuring the symmetrical uncertainty with respect to another set of attributes.

$$SymmU(Class, Attribute) = 2 \frac{H(Class) + H(Attribute) - H(Class|Attribute)}{H(Attribute) + H(Class)},$$

where H is the information entropy.

Linear Correlation filter The algorithm finds weights of continuous attributes basing on their Pearson's correlation with continuous class attribute.

Rank Correlation filter The algorithm finds weights of continuous attributes basing on their Spearman's correlation with continuous class attribute.

OneR algorithm The algorithms find weights of discrete attributes basing on very simple association rules involving only one attribute in condition part. In other words, it uses the minimum-error attribute for prediction, discretizing numeric attributes. For more information, see [4].

RReliefF filter The algorithm evaluates the worth of an attribute by repeatedly sampling an instance and considering the value of the given attribute for the nearest instance of the same and different class. Considering that result, it evaluates weights of attributes. Can operate on both discrete and continuous class data. For more information see [5,6,7].

Consistency-based filter Evaluates the worth of a subset of attributes by the level of consistency in the class values when the training instances are projected onto the subset of attributes. Consistency of any subset can never be lower than that of the full set of attributes, hence the usual practice is to use this subset evaluator in conjunction with a Random or Exhaustive search which looks for the smallest subset with consistency equal to that of the full set of attributes. The FSelector R package implementation makes use of *Best-first search* for searching the attribute subset space. Works for continuous and discrete data.

RandomForest filter It is a wrapper for variable importance measure produced by randomForest algorithm. The FSelector R package implementation allows for two types of importance measure:

1. mean decrease in accuracy,
2. mean decrease in node impurity.

The first measure is computed from permuting OOB (out-of-bound) data: For each tree, the prediction error on the out-of-bag portion of the data is recorded (error rate for classification, MSE for regression). Then the same is done after permuting each predictor variable. The difference between the two are then averaged over all trees, and normalized by the standard deviation of the differences. If the standard deviation of the differences is equal to 0 for a variable, the division is not done (but the average is almost always equal to 0 in that case).

The second measure is the total decrease in node impurities from splitting on the variable, averaged over all trees. For classification, the node impurity is measured by the Gini index. For regression, it is measured by residual sum of squares.

4.1.2 Algorithms for wrapping classifiers and search attribute subset space

In general, the wrapper approach depends on the so called *evaluation function* that is used to return a numeric value (a score) indicating how important a given subset of features is. Typically, one uses the classification-accuracy (usually based on cross-validation) as the score for the subset.

Below we provide a brief description of the algorithms for searching attribute subset space.

Greedy search At first, greedy search algorithms expand starting node, evaluate its children and choose the best one which becomes a new starting node. This process goes only in one direction. *Forward search* starts from an empty and *backward search* from a full set of attributes.

Best-first search The algorithm is similar to *Forward search* besides the fact that it chooses the best node from all already evaluated ones and evaluates it. In the FSelector R package implementation, the selection of the best node is repeated approximately *max.brackets* times in case no better node found.

Hill climbing search The algorithm starts with a random attribute set. Then it evaluates all its neighbours and chooses the best one. It might be susceptible to local maximum.

Exhaustive search The algorithm searches the whole attribute subset space in breadth-first order.

5 Classification

5.1 Classification algorithms

kNN k- nearest neighbours Method is used for modeling in problem of regression or classification. It is simple algorithm using lazy learning. There is no actual model so all the computation is done while classification. In the problem of classification the result for every single observation is a class for which in k closest neighbours from the training set is the most popular.

Decision trees Decision tree is a method that perform recursive partition of the set for every predictor. In each step there is chosen split that separates the set between classes the most according to one of the measures. The most popular measures are Information GAIN or GINI. For continuous data it is desired to partition variable into categorical (It could cause loss of the information). Result is highly correlated with the learning set. Nonetheless it is easy to interpret, and attractive visually. Another plus is that decision trees do not have any assumptions about distribution of the data and algorithms works fast.

Random forest It is currently one of the most popular method in machine learning. Its popularity grows thanks to good performance and small assumptions. However it performs well, it is hard to interpret the results, as long as model is complicated and consists many decision trees. For this method in each step of decision tree creation there is taken random subset of the features and then one of them is taken for split. This is done until appropriate settled level. For Random forests the computation time is much higher than for decision trees. Mostly it is because not only one tree is fitted but usually much more. One of the biggest disadvantages of this model is hard interpretation of the output. Even though the subset of predictors is only taken it shows much better results than other regular methods for different data sets.

Logistic regression The most popular method among application in banks, insurance companies and the industries for modeling binary data (It could serve also for prediction multiclass data). It owes popularity to its simplicity, easy open form and straight interpretation. It is subject to produce Score Card. Method is a particular type of generalized linear model where link function has logit form $\text{logit}(p) = \log(\frac{p}{1-p})$. It means that probability of occurrence particular event, is modeled indirectly, as a appropriate transformation.

$$\text{logit}(p_j) = \log\left(\frac{p_j}{1-p_j}\right) = \sum_{i=1}^n \beta_i X_{i,j}$$

Where p_j is estimated probability, β_i is factor for $X_{i,j}$ and X represents the features of observation.

Linear discriminant analysis It is another linear method. Under the assumption of normality and equality of covariance matrices within classes.

$$Pr(C = k|X = x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l}$$

where $C = k$ represents particular class affiliation, x is observation vector and $f_k(x)$, has appropriate Gaussian distribution with the mean μ_k and covariance matrix variance Σ and π_k is a-priori

classes probability. It is enough to compare numerator as long as denominator for all classes would be the same.

Quadric discriminant analysis It is similar method to the linear discriminant analysis. It keeps assumption about normality, but in this case covariance matrices could differ. Appropriate probability function keep its form:

$$Pr(C = k|X = x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l}$$

As before it is enough to compare numerators for all classes.

Naive Bayes Another method that uses Bayesian rule. It is called Naive Bayes as long as it has a naive assumption about loss of correlation between predictors. However this model has easy form, it also perform well in many applications.

$$P(Y = k|X = x) = \frac{P(X = x|Y = y)}{P(X = x)} = \frac{P(X_1 = x_1, \dots, X_n = x_n|Y = y)}{P(X = x)}$$

In this case probabilities are just taken as an empirical realisation of the data. It could also fall into problem of zero class frequencies. To omit this situation it is recommended to use one of the smoothing methods.

5.2 Classification performance metrics

Confusion matrix This is one of simple method to grade quality of the classification. It serves to compare actual class of an observation to one predicted by model.

		Predicted Class	
		True	False
Actual Class	True	True Positive (TP)	False Negative (FN)
	False	False Positive (FP)	True Negative (TN)

Using confusion matrix it is easier to calculate many of the goodness of fit measures for the models such as sensitivity, specificity or many more.

Sensitivity One of the simple measures called also Recall or True Positive Rate. It measures proportion of predicted as true and actual true. Using confusion matrix it could be given as:

$$TPR = \frac{TP}{TP + FN}$$

Specificity This factor is also called True Negative Rate. It measures proportion predicted correctly as false and actual number of false observations. It is given by:

$$SPC = \frac{TN}{TN + FP}$$

Precision Popular measure in information retrieval. It represents fraction of documents relevant to retrieved. In binary classification it is defined as:

$$PPV = \frac{TP}{TP + FP}$$

False discovery rate It is complementary to the Precision measure that is getting more popular thanks to growing dimensionality of data sets. In many applications it is of an interest of scientist to control this factor. The formula for this coefficient is subsequent:

$$FDR = \frac{FP}{TP + FP} = 1 - PPV$$

Accuracy It is simple measure that could be taken as a good indicator for model performance. It takes proportion of all positive classified to total number of observations. For not equally distributed observations between groups (for example in spam detection where there is many spam files classifier that predicts everything as a spam would have high Accuracy)

$$ACC = \frac{TP + TN}{TP + FP + TN + FN}$$

F-measure F-measure is defined as a combination of Precision and Recall. Both of earlier described indexes gives some information about model, but using them separately can effect with falling into missclassification for specific types of data.

$$F = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

5.2.1 Separation measures

Kolmogorow Smirnov statistics and distributions Kolmogorov Smirnov statistics is the biggest distance between distribution functions of scores (probabilities) for actual groups of True and False. Distributions shows just simple cumulative distribution functions for classes.

ROC curve and GINI, AUROC indexes ROC or Receiver Operating Characteristic is a graphical illustration that represents performance of the binary classifier. It is being used in medicine, radiology, biometrics and many more applications. It shows proportion of the True Positive rate (on the vertical axis) and False Positive rate (on the horizontal axis) at various thresholds. AUC is factor strongly connected with the ROC curve. Abbreviation stands for area under curve. It could be calculated as follows:

$$AUC = \int_{-\infty}^{\infty} TPR(x)FPR(x)dx$$

Histogram Good vs Bad It is just histogram showing distribution of scores (probabilities) within classes. For good models there suppose to be visible difference between height of the scores for different classes.

6 Principal Components Analysis

According to *Wikipedia* [10], Principal component analysis (PCA) is a statistical procedure that uses an orthogonal transformation to convert a set of observations of (possibly) correlated variables into a set of values of linearly uncorrelated variables called *principal components*. This transformation is defined in such a way that the first principal component is chosen to account for as much variance in plot dispersion from the centroid as possible; the second is chosen by the same criterion but subject to the constraint of being orthogonal the first, and so on. The principal components are the *eigenvectors of the covariance matrix* of distribution of the variables that a set of observations is distributed from.

The desired outcome of the principal component analysis is to project a feature space onto a smaller subspace that represents our data "well". Depending on how successful we are at reducing the data set, we can seek patterns among the distribution of plots in ordination space, and explore possible environmental correlates with these.

The results of a PCA are complex; we may be interested in knowing:

- variance explained by each eigenvector,
- cumulative variance explained by subsequent eigenvectors,
- *loadings* for each column; the loadings are the contribution of the column vector to each of the eigenvectors. A large positive component means that that column is positively correlated with that eigenvector; a large negative value is negative correlation; and small values mean that the species is unrelated to that eigenvector;
- *scores* for each row in the matrix or dataframe.

Generally, the vast majority of the variance is described on the first few eigenvectors, and we can save space by only calculating the scores for only the first few eigenvectors.

7 Cluster analysis

Clustering is the unsupervised classification of patterns (observations, data items, or feature vectors) into groups (clusters). In this section we include some basic informations about cluster analysis algorithms and cluster analysis performance metrics. The notes below are based mainly on 3 different sources:

- "Data Clustering: A Review" article [9],
- notes from Laboratory for Dynamic Synthetic Vegetation Phenology (LabDSV), the University of California [11],
- `clValid` (R package for cluster validation) vignette [12].

7.1 Introduction

Cluster analysis is the organization of a collection of patterns (usually represented as a vector of measurements, or a point in a multidimensional space) into clusters based on similarity.

Typical pattern clustering activity involves the following steps:

1. pattern representation (optionally including feature extraction and/or selection),
2. definition of a pattern proximity measure appropriate to the data domain,
3. clustering or grouping,
4. data abstraction (if needed), and
5. assessment of output (if needed).

Pattern representation process that one usually performs is to gather facts and conjectures about the data, optionally perform feature selection and extraction, and design the subsequent elements of the clustering system. A careful investigation of the available features and any available transformations (even simple ones) can yield significantly improved clustering results. A good pattern representation can often yield a simple and easily understood clustering; a poor pattern representation may yield a complex clustering whose true structure is difficult or impossible to discern.

Since similarity is fundamental to the definition of a cluster, a *measure of the similarity* between two patterns drawn from the same feature space is essential to most clustering procedures. Because of the variety of feature types and scales, the distance measure (or measures) must be chosen carefully. It is most common to calculate the *dissimilarity* between two patterns using a distance measure defined on the feature space.

The *grouping* step can be performed in a number of ways. For example, we can distinguish between *hard* output clustering (a partition of the data into groups) or *fuzzy* (where each pattern has a variable degree of membership in each of the output clusters). We may also distinguish between *hierarchical* clustering algorithms that produce a nested series of partitions based on a criterion for merging or splitting clusters based on similarity and *partitional* clustering algorithms that identify the partition that optimizes (usually locally) a clustering criterion.

Data abstraction is the process of extracting a simple and compact representation of a data set. In the clustering context, a typical data abstraction is a compact description of each cluster, usually in terms of cluster prototypes or representative patterns such as the centroid.

Eventually, validity assessments are performed to determine whether the output is meaningful. A clustering structure is valid if it cannot reasonably have occurred by chance or as an artifact of a clustering algorithm. There are three types of validation studies. An *external* assessment of validity compares the recovered structure to an a priori structure. An *internal* examination of validity tries to determine if the structure is internally appropriate. A *relative* test compares two structures and measures their relative merit.

7.2 Similarity, Dissimilarity and Distance

Similarity is a characterization of the ratio of the number of attributes two objects share in common compared to the total list of attributes between them. Objects which have everything in common are identical, and have a similarity of 1.0. Objects which have nothing in common have a similarity of 0.0. There is a large number of similarity indices proposed and employed.

Dissimilarity is the complement of similarity, and is a characterization of the number of attributes two objects have uniquely compared to the total list of attributes between them. In general, dissimilarity can be calculated as $1 - \text{similarity}$.

Distance is a geometric conception of the proximity of objects in a high dimensional space defined by measurements on the attributes. R calculates distances with functions from at least a few packages. Among them there are:

- `stats::dist` - computes and returns the distance matrix computed by using the specified distance measure to compute the distances between the rows, x and y (p -elements vectors), of a data matrix; the distance measure to be used is one of:
 - "euclidean":
usual distance between the two vectors (L_2 norm): $\sqrt{\sum_i^p (x_i - y_i)^2}$,
 - "maximum":
maximum distance between two components of x and y (supremum norm),
 - "manhattan":
absolute distance between the two vectors (L_1 norm),
 - "canberra":
 $\sum_i^p (|x_i - y_i| / |x_i + y_i|)$; terms with zero numerator and denominator are omitted from the sum and treated as if the values were missing,
 - "binary":
(*asymmetric binary*) the vectors are regarded as binary bits, so non-zero elements are "on" and zero elements are "off"; the distance is the proportion of bits in which only one is on amongst those in which at least one is on,
 - "minkowski":
the p norm, the p th root of the sum of the p th powers of the differences of the components,

- `labdsv::dsvdis` - calculates dissimilarity or distance between rows of a matrix of observations according to a specific index; returns an object of class "dist", equivalent to that from `stats::dist`.

Three indices convert the data to presence/absence automatically. In contingency table notation, they are:

- "steinhaus": $1 - a / (a + b + c)$,
- "sorensen": $1 - 2a / (2a + b + c)$,
- "ochiai": $1 - a / \sqrt{(a + b)(a + c)}$.

In practice, distances and dissimilarities are sometimes used interchangeably. They have quite distinct properties; e.g., dissimilarities are bounded within $[0, 1]$ whereas distances are unbounded on the upper end.

7.3 Clustering Algorithms

R has wide variety of clustering algorithms available. We make use of 9 algorithms from the base distribution (`stats`) and add-on packages. A brief description of each clustering method and its availability is given below.

7.3.1 UPGMA

Unweighted Pair Group Method with Arithmetic Mean is probably the most frequently used clustering algorithm. It is an agglomerative, hierarchical clustering algorithm that yields a dendrogram which can be cut at a chosen height to produce the desired number of clusters. Each observation is initially placed in its own cluster, and the clusters are successively joined together in order of their "closeness". The closeness of any two clusters is determined by a dissimilarity matrix, and can be based on a variety of agglomeration methods.

UPGMA is included with the base distribution of R in function `hclust()`, and is also implemented in the `agnes()` function in package `cluster`.

7.3.2 K-means

K-means is an iterative method which minimizes the within-class sum of squares for a given number of clusters. The algorithm starts with an initial guess for the cluster centers, and each observation is placed in the cluster to which it is closest. The cluster centers are then updated, and the entire process is repeated until the cluster centers no longer move. Often another clustering algorithm (e.g., UPGMA) is run initially to determine starting points for the cluster centers.

K-means is implemented in the function `kmeans()`, included with the base distribution of R.

7.3.3 Diana

Diana is a divisive hierarchical algorithm that initially starts with all observations in a single cluster, and successively divides the clusters until each cluster contains a single observation. Along with SOTA, Diana is one of a few representatives of the divisive hierarchical approach to clustering.

Diana is available in function `diana()` in package `cluster`.

7.3.4 PAM

Partitioning around medoids (PAM) is similar to K-means, but is considered more robust because it admits the use of other dissimilarities besides Euclidean distance. Like K-means, the number of clusters is fixed in advance, and an initial set of cluster centers is required to start the algorithm.

PAM is available in the `cluster` package as function `pam()`.

7.3.5 Clara

Clara is a sampling-based algorithm which implements PAM on a number of sub-datasets. This allows for faster running times when a number of observations is relatively large.

Clara is also available in package `cluster` as function `clara()`.

7.3.6 Fanny

This algorithm performs fuzzy clustering, where each observation can have partial membership in each cluster. Thus, each observation has a vector which gives the partial membership to each of the clusters. A hard cluster can be produced by assigning each observation to the cluster where it has the highest membership.

Fanny is available in the `cluster` package (function `fanny()`).

7.3.7 SOM

Self-organizing maps is an unsupervised learning technique that is popular among computational biologists and machine learning researchers. SOM is based on neural networks, and is highly regarded for its ability to map and visualize high-dimensional data in two dimensions.

SOM is available as the `som()` function in package `kohonen`.

7.3.8 Model based clustering

Under this approach, a statistical model consisting of a finite mixture of Gaussian distributions is fit to the data. Each mixture component represents a cluster, and the mixture components and group memberships are estimated using maximum likelihood (EM algorithm).

The function `Mclust()` in package `mclust` implements model based clustering.

7.3.9 SOTA

Self-organizing tree algorithm (SOTA) is an unsupervised network with a divisive hierarchical binary tree structure. It was originally proposed by Dopazo and Carazo (1997) for phylogenetic reconstruction, and later applied to cluster microarray gene expression data in (Herrero et al., 2001). It uses a fast algorithm and hence is suitable for clustering a large number of objects.

SOTA is included with the `clValid` package as function `sota()`.

7.4 Evaluation of clustering

In this analysis we focus on two types of clustering evaluation measurements: *internal* and *external* criterions.

7.4.1 Internal criterion

Typical objective functions in clustering formalize the goal of attaining:

- high intra-cluster similarity (objects within a cluster are similar) and
- low inter-cluster similarity (objects from different clusters are dissimilar).

This is an *internal criterion* for the quality of a clustering.

One drawback of using internal criteria in cluster evaluation is that high scores on an internal measure do not necessarily result in effective information retrieval applications. Additionally, this evaluation is biased towards algorithms that use the same cluster model. For example k-Means clustering naturally optimizes object distances, and a distance-based internal criterion will likely overrate the resulting clustering. Therefore, the internal evaluation measures are best suited to get some insight into situations where one algorithm performs better than another, but this shall not imply that one algorithm produces more valid results than another.

The methods that be used to assess the quality of clustering algorithms based on internal criterion include:

- *Silhouette coefficient*:

The silhouette coefficient contrasts the average distance to elements in the same cluster with the average distance to elements in other clusters. Objects with a high silhouette value are considered well clustered, objects with a low value may be outliers. This index works well with k-means clustering, and is also used to determine the optimal number of clusters.

Implementation of this measurement may be found in the `fpc` package - it is included in the results of `cluster.stats()` function (see: `clus.avg.silwidths` and `avg.silwidth` values).

- *Dunn index*:

The Dunn index aims to identify dense and well-separated clusters. It is defined as the ratio

between the minimal inter-cluster distance to maximal intra-cluster distance (in other words: minimum separation / maximum diameter). For each cluster partition, the Dunn index can be calculated by the following formula:

$$D = \frac{\min_{1 \leq i < j \leq n} d(i, j)}{\max_{1 \leq k \leq n} d'(k)}$$

where $d(i, j)$ represents the distance between clusters i and j , and $d'(k)$ measures the intra-cluster distance of cluster k .

Since internal criterion seek clusters with high intra-cluster similarity and low inter-cluster similarity, algorithms that produce clusters with high Dunn index are more desirable.

Implementation of this measurement may be found in the `fpc` package - it is included in the results of `cluster.stats()` function (see: `dunn` value).

- *Davies-Bouldin index:*

The Davies–Bouldin index can be calculated by the following formula:

$$DB = \frac{1}{n} \sum_{i=1}^n \max_{j \neq i} \left(\frac{\sigma_i + \sigma_j}{d(c_i, c_j)} \right),$$

where n is the number of clusters, c_x is the centroid of cluster x , σ_x is the average distance of all elements in cluster x to centroid c_x , and $d(c_i, c_j)$ is the distance between centroids c_i and c_j .

Since algorithms that produce clusters with low intra-cluster distances (high intra-cluster similarity) and high inter-cluster distances (low inter-cluster similarity) will have a low Davies–Bouldin index, the clustering algorithm that produces a collection of clusters with the smallest Davies–Bouldin index is considered the best algorithm based on this criterion.

The function `index.DB()` in package `clusterSim` implements Davies-Bouldin index.

7.4.2 External criterion

An alternative to internal criteria is direct evaluation in the application of interest. In external evaluation, clustering results are evaluated based on data that was not used for clustering, such as known class labels. These types of evaluation methods measure how close the clustering is to the predetermined classes.

Some of the measures of quality of a cluster algorithm using external criterion include:

- *Rand index:*

Given a set of n elements $S = \{o_1, \dots, o_n\}$ and two partitions of S to compare, $X = \{X_1, \dots, X_r\}$ which is a partition of S into r subsets, and $Y = \{Y_1, \dots, Y_s\}$ which is a partition of S into s subsets, define the following:

- a , the number of pairs of elements in S that are in the same set in X and in the same set in Y ,
- b , the number of pairs of elements in S that are in different sets in X and in different sets in Y ,
- c , the number of pairs of elements in S that are in the same set in X and in different sets in Y
- d , the number of pairs of elements in S that are in different sets in X and in the same set in Y

The Rand index, R , is:

$$R = \frac{a + b}{a + b + c + d} = \frac{a + b}{\binom{n}{2}}.$$

Intuitively, $a + b$ can be considered as the number of agreements between X and Y and $c + d$ as the number of disagreements between X and Y .

The Rand index has a value between 0 and 1, with 0 indicating that the two data clusters do not agree on any pair of points and 1 indicating that the data clusters are exactly the same.

The function `RRand()` in package `phyclust` implements this index.

7.4.3 Cluster stability

Studies of cluster stability may be considered as another type of cluster result validation. Method to perform this study is readily available in the `fpc` R package as `clusterboot()`.

`clusterboot()` function allows for the assessment of the clusterwise stability of a clustering of data. The data is resampled using several schemes (bootstrap, subsetting, jittering, replacement of points by noise) and the Jaccard similarities of the original clusters to the most similar clusters in the resampled data are computed. The mean over these similarities is used as an index of the stability of a cluster (other statistics can be computed as well).

The `fpc` package documentation provides some guidelines for the interpretation when bootstrap resampled scheme is used:

- There is some theoretical justification to consider a Jaccard similarity value smaller or equal to 0.5 as an indication of a "dissolved cluster", see Hennig (Hennig, C. (2008) Dissolution point and isolation robustness: robustness criteria for general cluster analysis methods. *Journal of Multivariate Analysis* 99, 1154-1176).
- Generally, a valid, stable cluster should yield a mean Jaccard similarity value of 0.75 or more.
- Between 0.6 and 0.75, clusters may be considered as indicating patterns in the data, but which points exactly should belong to these clusters is highly doubtful.
- Below average Jaccard values of 0.6, clusters should not be trusted.
- "Highly stable" clusters should yield average Jaccard similarities of 0.85 and above.

Part III
Results

8 Data preprocessing results

In this section we present the results of data preprocessing we performed. The important parts of this process are: creating derived variables, binning continuous variables and correcting bins (factor levels) if it seems reasonable.

The procedures and methods used to perform this part of the analysis include:

- visual inspection of density estimator plots and fitting distribution from different distribution families to numeric variables data,
- comparison of optimal discretization method (*smbinning* package) and equal frequency discretization with Information Value as criterion,
- using WoE criterion to define factor levels "similarity".

8.1 Searching for missing, corrupt and invalid data

We started the preprocessing in searching for missing, corrupt and invalid data. In our dataset most of the variables are factor variables (*PURPOSE*, *EMPLOYMENT* etc.) Some of them are numeric but consist only of a few unique values and thus may be seen rather like ordered factors (*NUM_OF_MAINTAINED_PEOPLE*, *RESIDENCE* etc.) We investigated frequency tables and did not notice anything unusual in the values.

Three variables in the data set are "truly" numeric: *DURATION*, *AMOUNT* and *AGE*. We did not find anything particularly unusual in the values of these variables. To satisfy our curiosity, we tried to fit a probabilistic distribution to the values. We used *MASS* :: *fitdistr* function to perform maximum-likelihood fitting of univariate distribution from selected distribution families. In each case we tried to fit parameters for three distributions families: *Gamma*, *Log-normal* and *Weibull*.

On the graphs below we can see kernel density estimates of variable density (black line) and curves representing density of the distribution fitted. We was not able to fit *Gamma* distribution to the *AMOUNT* variable values (*Error in stats::optim(x = c(1169, 5951, 2096, 7882, 4870, 9055, 2835, : non-finite finite-difference value [1])*). It seems that *Log-normal* distribution fits quite well in each three cases.

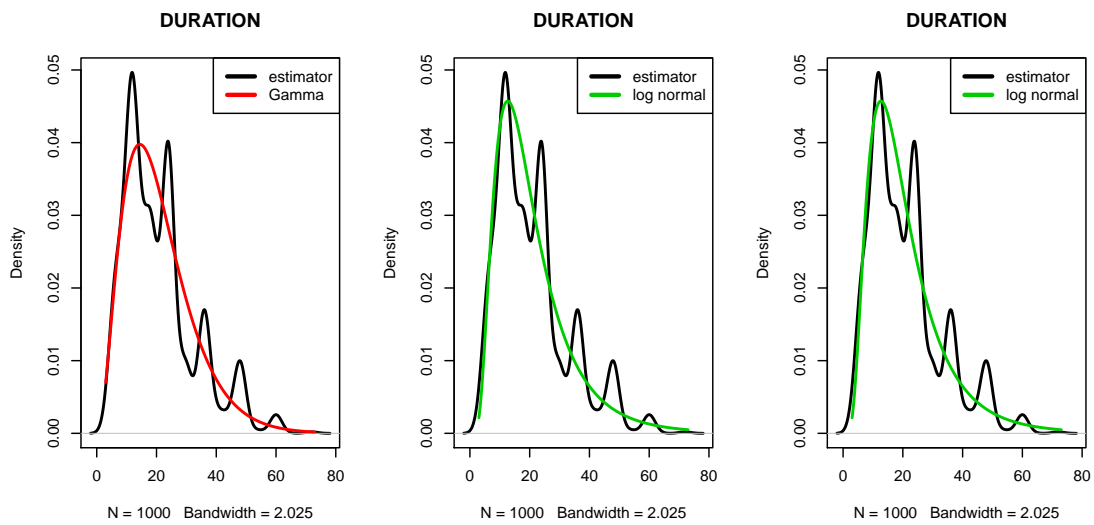


Figure 1: Kernel density estimate of *DURATION* variable density (black line) and curves representing density of the distribution fitted.

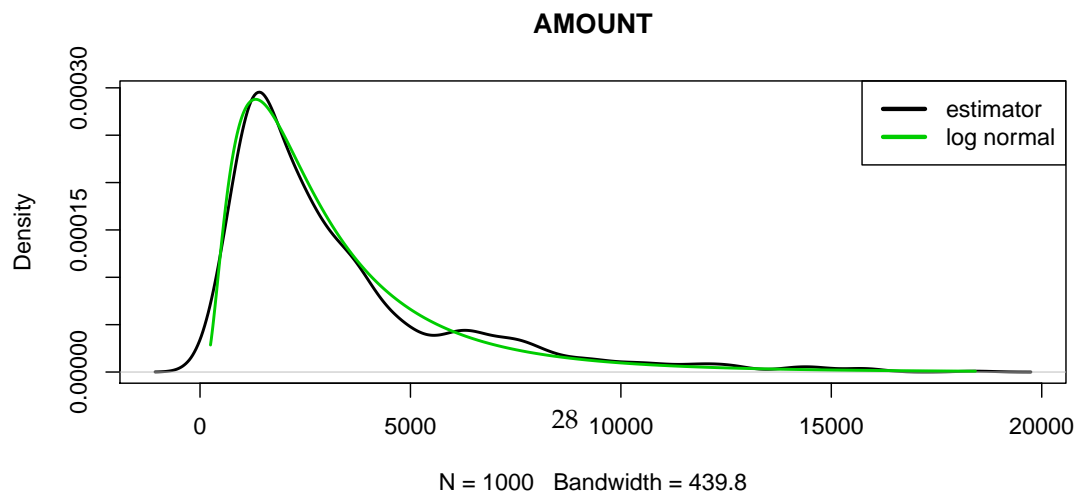
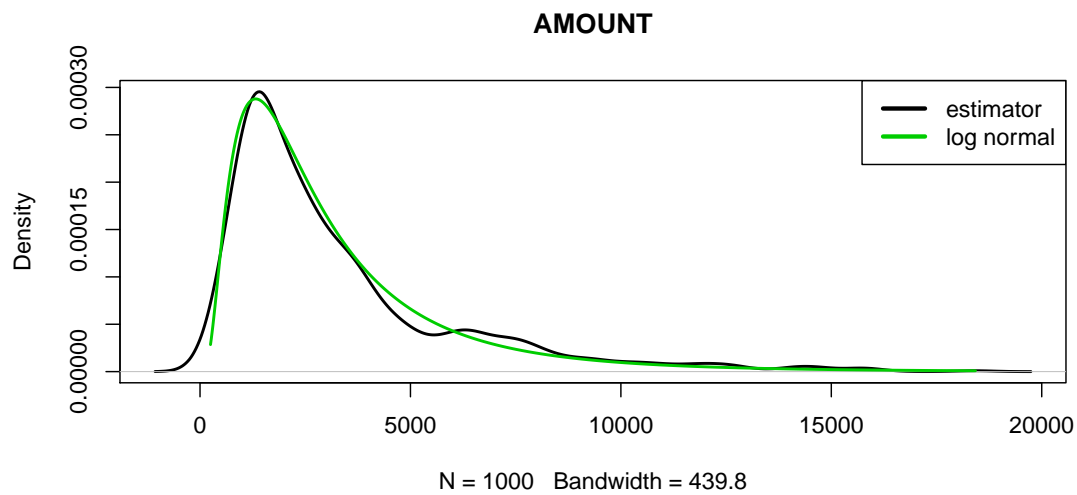
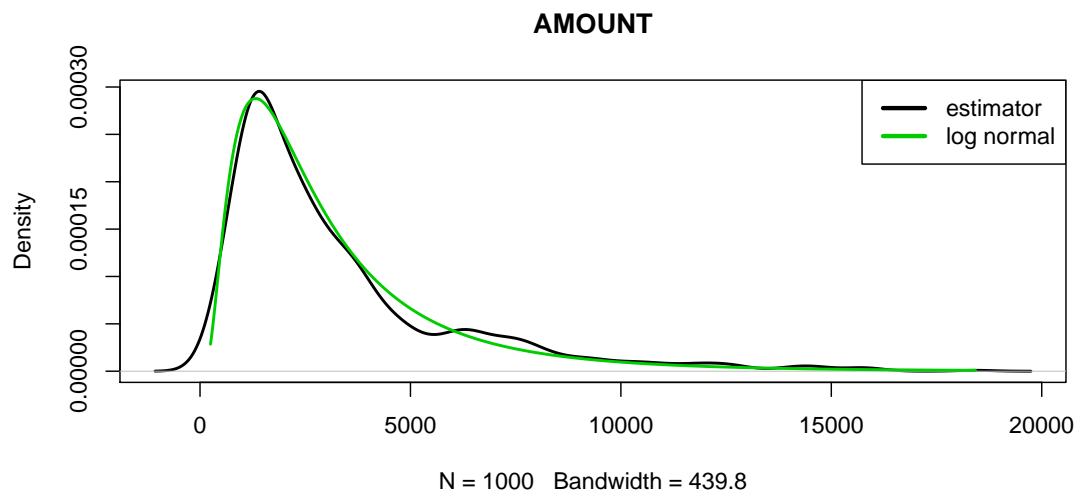


Figure 2: Kernel density estimate of *AMOUNT* variable density (black line) and curves representing density of the distribution fitted.

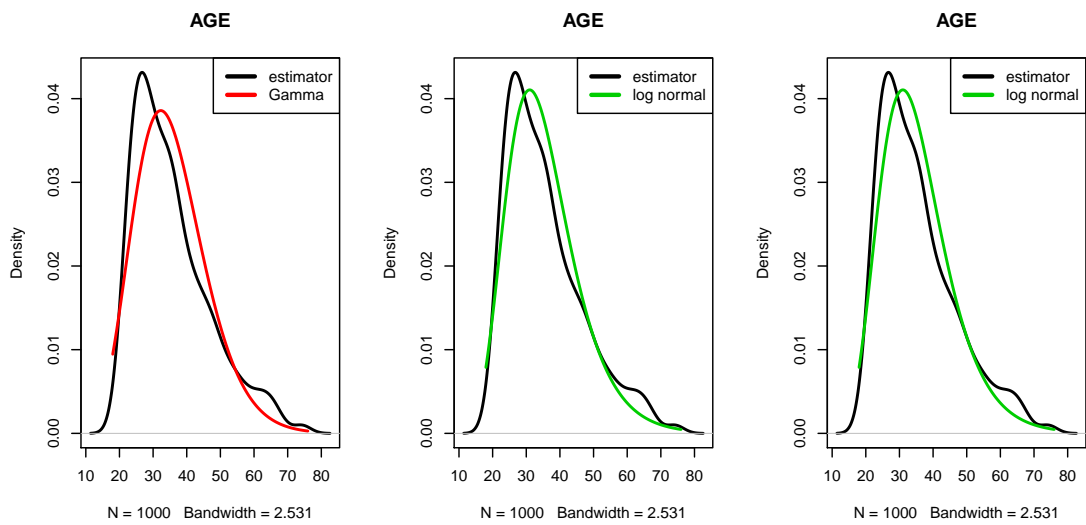


Figure 3: Kernel density estimate of AGE variable density (black line) and curves representing density of the distribution fitted.

8.2 Creating derived variables

We were considering possibilities of creating derived variables. We came up with propositions of the following formulas:

$$AMOUNT_TO_DURATION = AMOUNT / DURATION,$$

$$DURATION_TO_AGE = DURATION / AGE,$$

$$AMOUNT_TO_AGE = AMOUNT / AGE.$$

On the *Figure 4*, we can see boxplots of *DURATION*, *AGE* and *AMOUNT* across two levels of response variable *RES*. On the *Figure 5*, we can see boxplots of derived variables. This comparison can give us intuition how well our new variables separate good and bad bank clients.

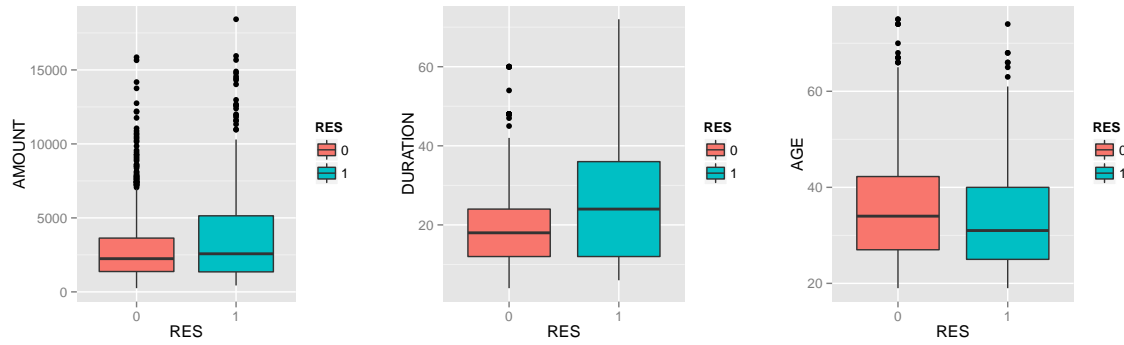


Figure 4: Boxplots of *AMOUNT*, *DURATION* and *AGE* variables across two levels of response variable *RES*.

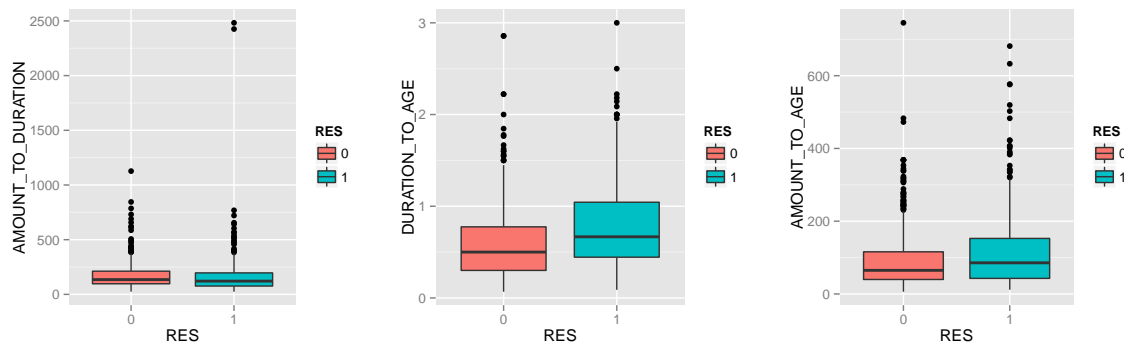


Figure 5: Boxplots of derived variables: *AMOUNT_TO_DURATION*, *DURATION_TO_AGE* and *AMOUNT_TO_AGE* across two levels of response variable *RES*.

The boxplots on the *Figure 5*. agree with our intuition: we expect higher *DURATION_TO_AGE* and *AMOUNT_TO_AGE* values for those clients who defaulted (*RES* = 1). On the other hand, we would expect the same for the *AMOUNT_TO_DURATION* variable whereas the plot shows something slightly opposite. However, the value differences on the *AMOUNT_TO_DURATION* boxplot are so fine that we suppose this variable is going to turn out to be of low information value and will not be consider as important one.

8.3 Binning continuous variables

In our analysis we compared three approaches of dividing continuous variables into categories:

1. equal frequency discretization (resulted bins are of equal number of observations),
2. supervised discretization which utilizes Recursive Partitioning to categorize the numeric characteristic and compute cutpoints based on Conditional Inference Trees algorithm (*smbinning* package),
3. categorize variable with simple tree model (from *rpart* package, with the use of default tree building parameters).

The IV comparison is presented on the *Figure 6*. below. Note that equal frequency discretization results are presented for:

- the same number of bins as in variable resulted from *smbinning* method; signature: *equal_nbins*,
- the same number of bins as in variable resulted from *rpart* method; signature: *equal_nrparts*.

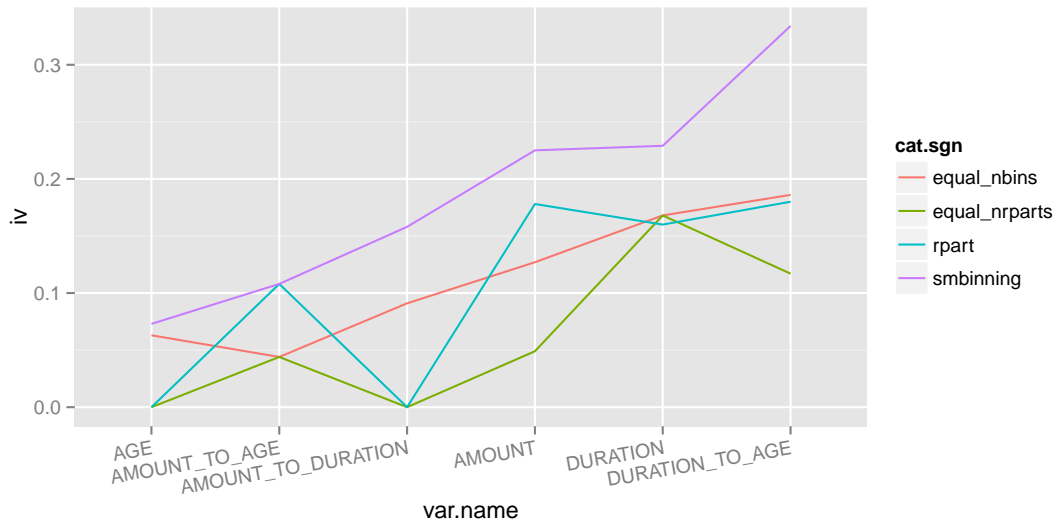


Figure 6: Comparison of Information Values of variables resulted from binning performed by means of 3 differend approaches.

We can see a few interesting things from the plot above:

- The *smbinning* function from the *smbinning* seems to beat other methods in terms of IV of resulted binned variable.
- Derived variable *DURATION_TO_AGE* has a strong relationship to the Goods/Bads odds ratio (over 0.33 IV value), whereas two other derived variables (*AMOUNT_TO_AGE*, *AMOUNT_TO_DURATION*) seem to have not.

8.4 Correcting bins (levels) of categorical variables

The last element of data pre-processing we performed is correcting bins (levels) of categorical variables. In the domain of particular variable's levels, we searched for levels which:

- reasonably similar in terms of their meaning (e.g. for *PURPOSE* variable we have logically similar levels: *car (new)* and *car (old)*,
- do not vary in terms of *GOOD* / *BAD* observations fraction.

To assess whether or not levels vary in terms of *GOOD* / *BAD* observations fraction we wrote a function which compute Information Value for non-changed variable and the same variable with a pair of levels joined, for all possible pair of levels. Below we present the head table output of this function and WoE plot for each of the levels for a particular variable.

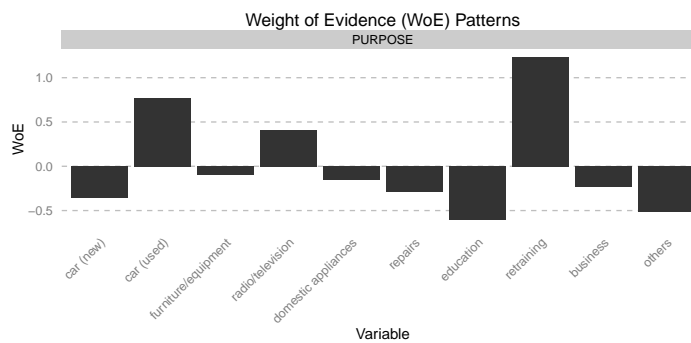
We include results for the two variables in case of which we decided to perform bin correcting.

1. *PURPOSE*

Below we can see what is the *PURPOSE* variable IV (3rd column) after joining a pair of levels (1st and 2nd column).

	lvl1	lvl2	iv
46	NONE	NONE	0.1692
21	furniture/equipment	domestic appliances	0.1692
37	business	repairs	0.1691
36	business	domestic appliances	0.1691
16	education	others	0.1691
28	car (new)	repairs	0.1691

We can also investigate the WoE values plot:



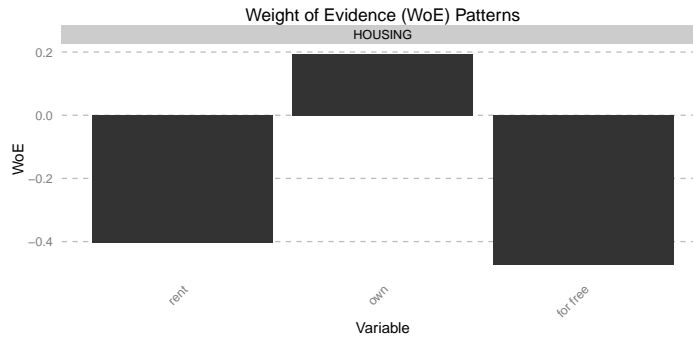
Judging from the above study we concluded that it is reasonable to join *furniture/equipment* and *domesticappliances* levels and we performed it.

2. HOUSING

Below we can see what is the *HOUSING* variable IV (3rd column) after joining a pair of levels (1st and 2nd column).

	lv11	lv12	iv
4	NONE	NONE	0.0833
3	for free	rent	0.0830
1	own	for free	0.0389
2	own	rent	0.0296

We can also investigate the WoE values plot:



Judging from the above study we concluded that it is reasonable to join *for free* and *rent* levels and we performed it.

8.5 Removing variables of 0.0 Information Value

After the transformations described above, we checked what are IV for each of the resulting variable. The values are presented below.

Variable	InformationValue	Bins	ZeroBins	Strength
CHK_ACCT	0.6660	4	0	Very strong
DURATION_TO_AGE	0.3338	3	0	Strong
HISTORY	0.2932	5	0	Strong
DURATION	0.2293	3	0	Strong
AMOUNT	0.2251	4	0	Strong
SAVINGS_ACCT	0.1960	5	0	Average
PURPOSE	0.1692	9	0	Average
AMOUNT_TO_DURATION	0.1575	5	0	Average
PROPERTY	0.1126	4	0	Average
AMOUNT_TO_AGE	0.1082	2	0	Average
EMLOYMENT	0.0864	5	0	Weak
HOUSING	0.0830	2	0	Weak
AGE	0.0732	2	0	Weak
OTHER_INSTALLMENT_PLANS	0.0576	3	0	Weak
STATUS_AND_SEX	0.0447	4	0	Weak
IS_FOREIGN_WORKER	0.0439	2	0	Weak
OTHER_DEBTORS	0.0320	3	0	Weak
RATE_TO_DISP_INCOME	0.0239	2	0	Weak
NUM_OF_THIS_BANK_CREDITS	0.0101	2	0	Wery weak
JOB	0.0081	3	0	Wery weak
TELEPHONE	0.0064	2	0	Wery weak
NUM_OF_MAINTAINED_PEOPLE	0.0000	1	0	Wery weak
RESIDENCE	0.0000	1	0	Wery weak

We can see that two variables: *NUM_OF_MAINTAINED_PEOPLE* and *RESIDENCE* have 0.0 information value. We decided to remove them from the data set as we assume they are kind of useless data.

8.6 Recoding variables to WoE

After the operations described above, we made a copy of our data and converted all of the categorical variables in this copy into their WoE representatives. Since then we have been working on two data sets:

1. *german_data_cat.txt* - categorized data set,
2. *german_data_woe.txt* - categorized and recoded to WoE data set (besides the response variable, this data set contains numeric variables only).

8.7 Data preprocessing summary

To conclude, the data preprocessing we performed above results in:

- creating derived variables and including them into the data set: *AMOUNT_TO_DURATION*, *DURATION_TO_AGE* and *AMOUNT_TO_AGE*,
- binning continuous variables with the use of *smbinning* method and keeping them in those binned form in the data set; related to: *AGE*, *AMOUNT*, *DURATION*, *AMOUNT_TO_DURATION*, *DURATION_TO_AGE* and *AMOUNT_TO_AGE*,

- correcting bins (levels) of categorical variables and keeping them in those transformed form in the data set: *PURPOSE* and *HOUSING*,
- removing variables of 0.0 information value: *NUM_OF_MAINTAINED_PEOPLE* and *RESIDENCE*,
- creating separated data set which contains all the categorical variables in the recoded to WoE form.

9 Classification modelling results

10 Cluster analysis results

11 Summary

Part IV
Discussion

References

- [1] Computational Methods of Feature Selection (Chapman & Hall/CRC Data Mining and Knowledge Discovery Series), Huan Liu, Hiroshi Motoda, 2007, ISBN-13: 978-1584888789
- [2] Hall, M. A., Smith, L. A. (1998). Practical feature subset selection for machine learning. Australian Computer Science Conference. Springer. 181-191.
- [3] Liu, H. and Setiono, R., Chi2: Feature selection and discretization of numeric attributes, Proc. IEEE 7th International Conference on Tools with Artificial Intelligence, 338-391, 1995
- [4] R.C. Holte (1993). Very simple classification rules perform well on most commonly used datasets. Machine Learning. 11:63-91.
- [5] Kenji Kira, Larry A. Rendell: A Practical Approach to Feature Selection. In: Ninth International Workshop on Machine Learning, 249-256, 1992.
- [6] Igor Kononenko: Estimating Attributes: Analysis and Extensions of RELIEF. In: European Conference on Machine Learning, 171-182, 1994.
- [7] Marko Robnik-Sikonja, Igor Kononenko: An adaptation of Relief for attribute estimation in regression. In: Fourteenth International Conference on Machine Learning, 296-304, 1997.
- [8] StatSoft, Formula Guide Weight of Evidence Module <http://documentation.statsoft.com/portals/0/formula>
- [9] Data Clustering: A Review, A.K. Jain, M.N. Murty, P.J. Flynn, 1999, <https://www.cs.rutgers.edu/mlittman/courses/lightai03/jain99data.pdf>
- [10] Principal component analysis From Wikipedia, the free encyclopedia, http://en.wikipedia.org/wiki/Principal_component_analysis
- [11] notes from Laboratory for Dynamic Synthetic Vegetation Phenology (LabDSV), the University of California, <http://ecology.msu.montana.edu/labdsv/R/labs/>
- [12] clValid, an R package for cluster validation, Guy Brock, Vasyl Pihur, Susmita Datta, and Somnath Datta, Department of Bioinformatics and Biostatistics, University of Louisville, <http://cran.r-project.org/web/packages/clValid/vignettes/clValid.pdf>