

Report

Dan Nash

40217045@napier.ac.uk

Edinburgh Napier University - Module Title (SET09120)

1 Introduction

We have been tasked to clean, prep and analyse a large set of data. This is so we can then identify interesting patterns within the data and give the bank a data driven informed decision, on whom they should give a loan out to.

2 Data Cleaning

The Author used the program "OpenRefine" during the data cleaning stage. When the data has been cleaned, and the relevant data types have been declared either numerical or nominal, then it is ready to be prepped for Weka.

2.1 OpenRefine

When the data file (excel spread sheet) is first downloaded, it must be converted into a CSV file format, instead of it's standard xlsx format. Unfortunately by doing this, any numeric data type (column) is then transformed into a text value. The affected data types can be converted back it's original form within OpenRefine though. Before loading the CSV file into OpenRefine, make sure to place the appropriate headers over the correct columns of the spreadsheet. This is so the user can maintain an understanding of what the data in each column actually represents.

Start OpenRefine and start a new project. Select the CSV file and then click begin. The data will then present it's self in a very similar fashion as if it was Excel. Before we start transforming the relative data columns into it's respective data types; "num_dependents" to numeri-

cal. We must first clean the data in the affected columns.

In "**num_dependents**", select text facet and then you'll see a number of discrepancies located within the data set. The following discrepancies need to be "edited"; "1one" and "one" need to be edited to "1". "two" and "twotwo" need to be edited to "2". This column should now be cleaned and it can be "common transformed" into a numeric data type.

"**Purpose**" column needs to be text facet as well. The following discrepancies need to be fixed; "business" and "busines" to "buisness". "Eduction" to "education". "ather" to "other". You will then see 10 choices left, but when looking at the assignment brief you'll notice that there are 11 entries. "Vacation" is the missing entry, but it's never marked in any of the "case_no" entries for "purpose". So you don't need to worry about it being missing.

In "**job**", select text facet and look at the discrepancy "yes". There are 2 instances of "yes", so select "yes" and then it will show you the data rows for both of those entries. The issue here is we must identify how "yes" is interpreted. The reason being is that there are 4 options for "job" so there are many ways we can identify "yes". After reviewing them, I decided to place them in "unskilled resident", the reason being was that one of them had "<100" savings and was relatively young "26". The other "case_no" was older "37" and had "<100" savings, but his "purpose" was education. Which gave me the impression that he was needing a loan to pay for education, which would lead him to be a "skilled" or "high qualif/self emp/mgmt". So edit "yes" to "unskilled resident".

Common Transform "**credit_amount**" into numeric, then do a "numeric facet" on the same

column. There are 5 records that are over the amount of 10,000,000 which is an obscene amount and has to be a clerical error. I then proceeded to remove all of the zero's for each case. 4 of the 5 data entries will now be fixed and within an appropriate "credit_amount" range, but the 5th one (case_no 432) will still be too high. You must numeric facet the column again and set the search range from between 100,000 to 7,200,000. After doing so, 4 rows will appear, with one of them being "case_no 432" again. Ignore case "432" for now and remove three "0" on the other "credit_amount" cases. I say remove 3 zero's instead of 4 because if you look at their purpose then the amount they need will be above 1000. With case number 432, remove the first digit "1", and that will bring the credit value down to "11328", which would match his profile because his "purpose" is "other" and he's in the "high qualif/self emp/mgmt" for the "job" category.

In the "existing_credits" column, do a text facet, and then proceed to remove any last digit with numbers that contain 2 digits, i.e 11 = 1. Any entities with 3 digits i.e 333, then remove the last two digits. Any numbers that are represented like this; 0.1 , then remove the "0." from the entities. After these edits you should be left with four entities that go as follows: 1, 2, 3, 4.

Common Transform "age" to numeric, then do a numeric facet. Set the range from -40 to 0, after doing so, proceed to edit out any of the "-" values. Do another numeric facet, set the range from 0 to it's maximum, from here you'll see values with decimal points, i.e 0.1. Remove any decimal points and then do another numeric facet setting the range from the first instance "6" to it's maximum range. There will be two instances which seem out of place. These instances are "6" (case_no 26) and "11" (case_no 54). Instead of removing them, I investigated into what their possible age bracket could be. I looked at their "employment" and "credit_history"; both had employment in the "1<=X<4" range and had "existing_paid" within the "credit_history" category. Going by this I believed it was fair to place a "2" in front of their original values, which would place them within the "20's" bracket. Finally do another numeric facet with the range of 80 to 340. Any entity

that's beyond the age of "100", then remove the final digit to place them under the age of "100".

The data cleaning should now be complete once going through all these steps.

3 Data Preparation

Now it's time to prep the data into the correct formats so that Weka can work with them. First off we need to create a standard data set which will contain both numeric and nominal data. So, make sure to common transform the following data columns (if have not previously been done by now); "case_no", "credit_amount", "age", "existing_credits" and "num_dependent". Then proceed to common transform all the other data columns into text, this is purely just for insurance purposes. There is also no need to text facet or numeric facet any of the other data columns that have not been inspected already. They have already been checked and there is no discrepancies.

Export the project to a Excel document in CSV format, and then proceed to open it in a "Word" document. Proceed to save the document and then open it in a text editor of your choice, I stuck with Notepad. From here you have to reposition the data and add particular value, so that when it's turned into a ".arff" file, Weka can then read it with no errors occurring. Review **appendix 1** to understand how the txt file should look. The main points to understand is that any previous columns is now designated with a "@attribute" and that any "attributes" that deal with "nominal" data must have a their data objects i.e "tv/radio", "0<=X200"; stated in the curly braces, which are located in their relevant attribute fields. Review **appendix 2** to understand. Any data type that's numeric, simply put "real" right after stating the attribute name. Review **appendix 3** to understand. Once this is done convert it into an arff file and Weka will now be able to work with it.

3.1 Numeric to Nominal

We need a ".arff" file with mostly nominal data. Certain algorithms like "Apriori" within "Association" can only work with nominal data. So, go back to the OpenRefine project and start to turn any numeric data into nominal.

"num_dependents" and "existing_credit" can simply just be turned into nominal data by doing Common Transform > To Text.

To turn larger data sets like "age" and "credit_amount" into nominal data, we must start applying the data into clusters, to which we can then express with an appropriate nominal value.

Select "age" and use a numeric facet, select the range 19 to 26 then Edit Cell > Transform the selected range to " $18 \leq X < 26$ ". From here on, repeat the previous steps but work with range increments of 10. For example, 26 to 36 (" $26 \leq X < 36$ "), 36 to 46 (" $36 \leq X < 46$ "). Do this till you get to the end of the range "76". By doing this, the data will be turned into "text" data and will be identifiable clusters of ranged data that we can work with.

Select "credit_amount", and do a numeric facet. We follow almost the same steps that we did for **age** but we break down the range to 2000 increments. I.e. 0 to 2000 will be " $0 \leq X < 2000$ ". We do this all the way up to 20,000. Also, there are no records of data within the range of 16,000 to 18,000. So, do not worry about trying to create this range. You will be left with 9 instance categories of "credit_amount", which can be seen in **appendix 4**.

All data (apart from **case_no**) has been turned into nominal data. Export the file into a excel ".xlsx" file and then open it in word. Save the the file and then open it in your preferred text editor. Similar to what we did previously in (**appendix 1**), we now have to state which "attributes" are "real" (numeric) and which attributes are nominal (" $26 \leq X < 36$ ", " $26 \leq X < 46$ "). **case_no** should be the only "real" attribute there. Look at **appendix 5** to see how the ".txt" file should look. Now save thew file and convert it into a ".arff" file. You now have a mostly nominal ".arff" file which can be interacted within Weka.

4 Weka

The three areas we'll be analysing the data from our newly created ".arff" files are "**Classification**" with the **48** Algorithm, "**Association**", with the **Apriori** Algorithm and finally "**Clusters**" with the **simpleKmeans** algorithm.

4.1 Classification

Load the nominal values only ".arff" file into weka. When using the classification rule set, having a mix of numerics and nominals can muddle the data a bit. Once the file has been loaded, go to the "Select Attribute" tab located at the top right of the Weka interface. Then choose "CorrelationAttributeEval" (evaluates the worth of an attribute by measuring the correlation between it and the class) from the "Attribute Evaluator" terminal. Make sure "(Nom) class" is chosen and then click start. The following results will then be produced; Review Appendix 6

From here, we identify the ones which are ranked poorly (highest **entropy**), "**age**", "**num_dependents**", "**case_no**" and "**job**". This is so we can then remove them from our data set before proceeding to the **Classification** stage. Go back to the "Preprocess" tab and proceed to remove the previously mentioned attributes from the data set.

Go to the "**Classify**" tab on the Weka interface and choose the "**J48**" algorithm from "trees" directory within the "Classifier" interface. Click on "J48" and a new menu will appear, look for "Min-NumObj" then change the value from 2 to 10 (**Appendix 7**). After the changes, exit that sub menu and make sure "(Nom) Class" is selected, then click "Start" to begin the calculations.

From output screen there are a number of factors to consider, they are the "**Correctly and Incorrectly Classified Instances**" and the "**Confusion Matrix**". From looking at the correctly and incorrectly classified instances (**Appendix 8**), you can see that 73% instances have been calculated correctly and 27% incorrectly. This percentage ratio is not ideal, but even if you decrease (reducing the "MinNumObj" number) the number of leaves and size of the tree to below the 30's (currently at Leaves:

40, size of Tree: 46). These percentages only changes from 1 to 3%. From looking at the "**Confusion Matrix**" (**Appendix 9**), you can see that 623 instances has been predicted as good and a 107 instances have been predicted classified as bad. But upon further inspection it is stated that in fact 77 instances were identified as bad when in fact they were good, but more alarmingly, there is 193 instances which are stated as good, when they are actually bad. Obviously the assessment of incorrectly identified "good" instances takes more precedence than the incorrectly predicted "bad" instances due to the larger size difference.

Right click on the "Result List" and select visualize tree, from here (**Appendix 10**) you should be able to identify some common rules within the data presented.

Rule 1: If Checking Status = No Checking, **Then** Class = Good. This rule has a coverage of 394 with an accuracy of 88%

Rule 2: If Checking Status = $0 \leq X < 200$ **And** Credit Amount = $0 \leq X < 2000$, **Then** Class = Good. This rule has a coverage of 106 and an accuracy of 62%.

Rule 3: If Checking Status = < 0 **And** Credit History = Existing Paid **And** Purpose = New Car, **Then** class = Bad. This rule has a coverage of 42 with an accuracy of 64%.

Rule 4: If Checking Status = $0 \leq X < 200$ **And** Credit Amount = $2000 \leq X < 4000$, **Then** Class = Good. This rule has a coverage of 77 with an accuracy of 72/

Rule 5: If Checking Status = ≥ 200 , **Then** Class = Good. This rule has a coverage of 63 with an accuracy of 77%.

Rule 6: If Checking Status = < 0 **And** Credit History = Critical/Other Existing Credit, **Then** Class = Good. This rule has a coverage of 67 and an accuracy of 73%.

What can be taken from looking at these rules is that having "No Checking" greatly increases your chance to receive a loan. This is probably due to the fact that if the bank does not already have you as a customer, then they want to make you a customer.

4.2 Association

Association only takes text data, so the ".arff" data file that only contains nominal values must be used. Load that up and proceed to remove "**case_no**", due to it being a numeric value. Also remove "**num_dependents**" and "**existing_credits**" too, because those two attributes seem to manipulate the data in undesirable ways.

Click on the "Associate" tab and choose the "**Apriori**" algorithm, and then click "Start". The algorithm will then produce us the best rules from the attributes given. **Appendix 11** shows these rules printed out, but I will state the best 6 rules from the appendix after this paragraph and it's calculated coverage and confidence values.

Rule 1: If Checking Status = No Checking **And** Purpose = Radio/Tv, **Then** Class = Good. This rule has a coverage of 120 and a confidence/accuracy value at 0.94

Rule 2: If Checking Status = No Checking **And** Credit History = Critical/Other Existing Credit, **Then** Class = Good. This rule has a coverage of 143 and a confidence/accuracy value of 0.93.

Rule 3: If Checking Status = No Checking **And** employment = ≥ 7 . **Then** Class = Good. This rules has a coverage of 107 and a confidence/accuracy value of 0.93.

Rule 4: If Checking Status = No Checking **And** Personal Status = Male Single **And** Job = Skilled, **Then** Class = Good. This rule has a coverage of 139 and a confidence/accuracy value of 0.93.

Rule 5: If Checking Status = No Checking **And** Credit Amount = $0 \leq X < 2000$ **And** Job = Skilled. **Then** Class = Good. This rule has a coverage of 105 and a confidence/accuracy value of 0.92

Rule 6: If Checking Status = No Checking **And** Credit Amount = $0 \leq X < 2000$, **Then** Class = Good. This rule has a coverage of 156 and a confidence/accuracy value of 0.92.

It can obviously been seen that if the customer has not had any account with the bank, then they are categorized as a good customer. Just

like the **Classification** algorithm showed.

4.3 Clustering

The "**SimpleKMeans**" algorithm in Clustering can handle nominal and numeric values. So load up the ".arff" file which contains both those data types. From here remove the attributes that we had removed from Association; **case_no**, **num_dependents** and **existing_credits**.

In the "**Cluster**" tab, select the algorithm **SimpleKMeans** and click on it again to open a sub menu of it. From here we have to change two values; The number of Clusters the algorithm produces (this is help diverse and spread the data calculation across a wider spread) and the Seed. Change the number of Clusters to 4; **Appendix 12**. And change the number of Seeds to 20; **Appendix 13**.

Proceed to click "**Start**" and a table producing a set of results spread along 4 clusters is shown; **Appendix 14**. At a quick glimpse of the results, we can see that cluster 0 and 3 (1 and 4) are classed as "bad" while cluster 1 and 2 (2 and 3) are classed as "good". Later we'll look into some rules that can be generated from these results. By scrolling down we can see the results for "**Clustered Instances**", **Appendix 15**. Here it shows the number of instances from the total of the data (1000) occurred and placed in which cluster. We can see that cluster 1 and 3 are very similar with a 5% difference between them. Cluster 2 however is the most predominant one within the data set, with it occurring 39% of the time.

If we wanted to visualize the data in correlation to how predominant the properties "**Good**" and "**Bad**" are in relation to which cluster they occur in, then we can do that by right clicking on the results lists and select visualize tree. From here we choose the values of our X and Y range and what property they will output on the graph. For the X coordinate select "X: Cluster (Nom)", for the Y coordinate select "Y: Class (Nom)" and for the property output choose "Colour: class (Nom)". Then move the "Jitter" slider bar to the far right. **Appendix 16** shows the visualisation of the data. Here it's made evident that there are higher cases of "Bad" in cluster 0 and 3 than there is in cluster 1 and 2.

From looking at all the data provided, the following rules can be extracted;

Rule 1: Checking Status = No Checking **And** Class = good, in cluster 1 and 2. But in clusters 0 and 3, Checking Status = to not "No Checking". Essentially, the Checking Status result "**No Checking**" is a major benefactor to making the Class value = "**Good**".

Rule 2: Employment length doesn't seem to effect if the Class = Good. Each Cluster has different employment length values and so no continuity/patter can be established.

Rule 3: Saving Status across all clusters = <100.

Rule 4: The results Male Single from "**Personal Status**" and Skilled from "**Job**" are both apparent in clusters 1 and 2. Which have the class value set at **Good**. They both don't show up together in clusters 0 and 3. Interestingly enough, these results align with rule 4 from **Association** rule set.

Rule 5: It is unclear if "**Purpose**" has a large effect on the outcome of Class. Reason being, in cluster 0 and 1, the purpose is the same but in cluster 2 and 3, the purpose are different.

As can be seen, which is a common pattern that's been established now; If the customer has No Checking status, then they are categorized as good. Also There may be a correlation between what job and personal status the customer has, which dictates what class they end up at. The pattern seems to be:

If job = Skilled **And** Personal Status = Male Single, **Then** Class = Good.

5 Visualisation

What is the relation between Checking Status and the decision to give a loan?

By looking through the rules generated by the previous algorithms, we can see a common trend that having "**No Checking**" stated, will greatly increase the likelihood of the bank giving you a loan. This can directly be seen when comparing Rule 1 from "**Clustering**" and the results from the Weka Cluster Visualizer (**Appendix**

16). With these indicators and many of the rules throughout the various algorithms pointing towards the same conclusion; It is safe to say that if a customer does not have an account with the bank, then they have a high chance of being granted a loan. It is also safe to say that if they do already have an account, unless the customer checking status is not " ≥ 200 " then it doesn't seem to influence the decision to grant a loan or not. If Checking Status does fall within the " ≥ 200 " however, then the likely hood of rejection is increased exponentially. This can all be seen within the Weka Visualizer (**Appendix 17**).

Simply put; If you have no account with the bank, you are attractive. If you are at the threshold of their Checking Status range, then you are not attractive.

5.1 Sketches

The three types of graphs made in quick sketches are Grouped Bar Chart (**Appendix 18**), Stacked Bar Chart (**Appendix 19**) and finally a Bubble Chart (**Appendix 20**).

5.2 Grouped Bar Chart

The different Checking Status will be located along the X coordinate plain and the number of Instances that they occur will be shown along the y coordinate plain. The "**Good**" class value will be identified with the colour green and the "**Bad**" class value will be red. The Grouped Bar Chart will separate the "Good" and "Bad" values within the same "Checking Status" category and so visually the contrast of the number of instances of these values occurring, will be easily distinguished by the length (**encoding**) of the grouped graphs. **Appendix 18**.

5.3 Stacked Bar Chart

Same as before; The different Checking Status will be located along the X coordinate plain and the number of Instances that they occur will be shown along the y coordinate plain. The colour coding will also be identical too. The main difference with this kind of graph is that instead of running the "**Good**" and **Bad** number of instances side by side, we stack them on top of one another. Stacked Bar Charts use the same

encodings as Grouped Bar Charts and they are good for visually distinguishing the amount each value contributes to the total of the x coordinate value; i.e "No Checking" in Checking Status. **Appendix 19**.

5.4 Bubble Graph

The different Checking Status values will be along the x coordinate plain and on the y coordinate plain, the values of Class will be displayed. Since the encodings of this graph work with size and areas, instead of length (like the previous two graphs), the number of instances of what checking status value occurs is recorded and visually displayed via the size of the bubble. So the height/length is not a factor in this graph, hence why the y coordinates doesn't record the number of "**Instances**".

5.5 Choice for R

I have decided to use the "**Grouped Bar Chart**" for my R diagram. The main reason being that it is easy to compare the number of instances in each class value to the different Checking Status values. And it also labels/categorises all the data the reader needs to know, in a very concise and orderly fashion.

6 Appendix

```

%relation germanfinancedata
@attribute case_no real
@attribute checking_status {'<0', '0<=X<200', '>=200', 'no checking'}
@attribute credit_history {'no credits/all paid', 'all paid', 'existing paid', 'delayed previously', 'critical/other'}
@attribute purpose {'new car', 'used car', 'furniture/equipment', 'radio/tv', 'domestic appliance', 'repairs', 'education'}
@attribute credit_amount real
@attribute saving_status {'<100', '100<=X<500', '500<=X<1000', '>=1000', 'no known savings'}
@attribute employment {'unemployed', '<1', '1<=X<4', '4<=X<7', '>=7'}
@attribute personal_status {'male div/sep', 'female div/dep/mar', 'male single', 'male mar/wid', 'female single'}
@attribute age real
@attribute existing_credits real
@attribute job {'unemp/unskilled non res', 'unskilled resident', 'skilled', 'high qualif/self emp/mgmt'}
@attribute num_dependents real
@attribute class {'good', 'bad'}

```

Figure 1: **Appendix 1** - Data Prep, Nominal and Numeric

```

@attribute checking_status {'<0', '0<=X<200', '>=200', 'no checking'}

```

Figure 2: **Appendix 2** - Data Prep, Nominal attribute value

@attribute age real

Figure 3: **Appendix 3** - Data Prep, Numeric attribute value

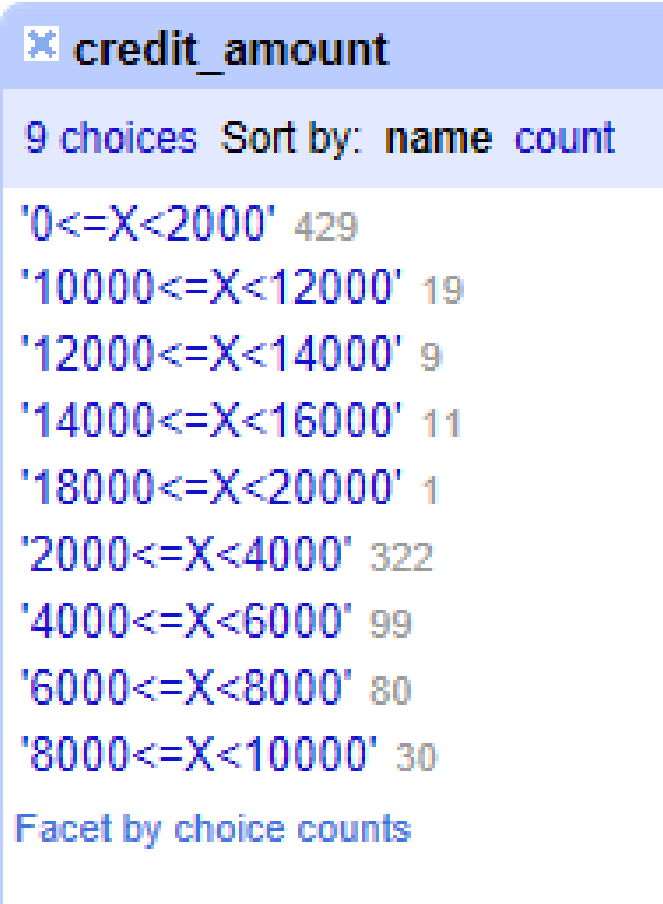


Figure 4: **Appendix 4** - Data Prep, Credit Amount instances



Figure 5: **Appendix 5** - Data Prep, Nominal txt file

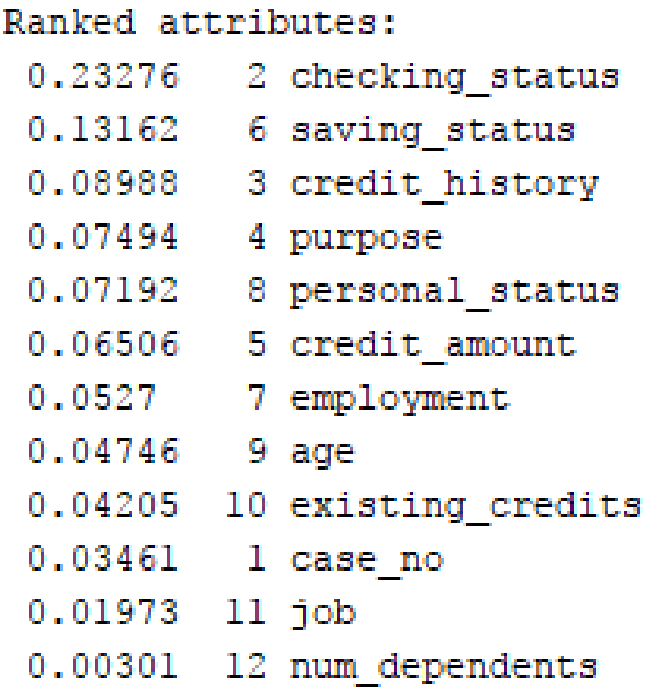


Figure 6: **Appendix 6** - Classification, Ranked Attributes

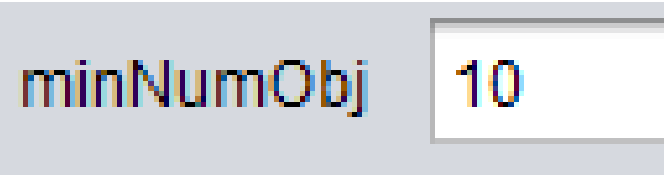


Figure 7: **Appendix 7** - Classification, MinNumObj

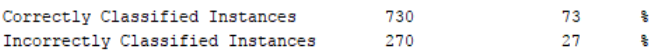


Figure 8: **Appendix 8** - Classification, Instances

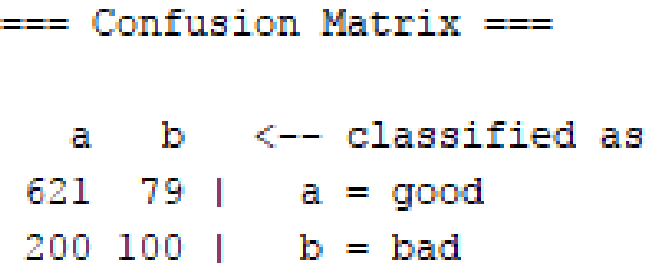


Figure 9: **Appendix 9** - Classification, Confusion Matrix

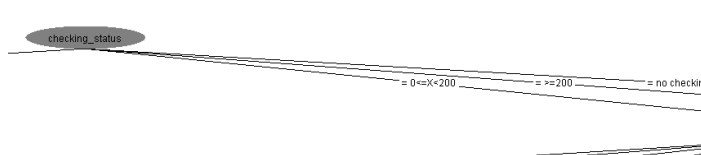


Figure 10: **Appendix 10** - Classification, Visualize Tree

Best rules found:

1. checking_status=0 checking_purpose=radio/tv 127 ==> class=good 120 <conf:(0.94)> lift:(1.35) lev:(0.03) [31] covr:(4.74)
2. checking_status=0 checking_credit_history=critical/other existing credit 153 ==> class=good 143 <conf:(0.93)> lift:(1.34) lev:(0.04) [35] covr:(4.17)
3. checking_status=0 checking_employment=>7 115 ==> class=good 107 <conf:(0.93)> lift:(1.33) lev:(0.03) [24] covr:(3.83)
4. checking_status=0 checking_personal_status=male single job=skilled 150 ==> class=good 139 <conf:(0.93)> lift:(1.32) lev:(0.03) [34] covr:(3.75)
5. checking_status=0 checking_credit_amount=0<X<2000 job=skilled 114 ==> class=good 105 <conf:(0.92)> lift:(1.32) lev:(0.03) [25] covr:(3.42)
6. checking_status=0 checking_credit_amount=0<X<2000 170 ==> class=good 156 <conf:(0.92)> lift:(1.31) lev:(0.04) [37] covr:(3.4)

Figure 11: **Appendix 11** - Association, Best Rules

numClusters

Figure 12: **Appendix 12** - Clustering, Number of Clusters

seed

Figure 13: **Appendix 13** - Clustering, Number of Seeds

| Cluster# | 0 | 1 | 2 | 3 |
|------------------------------|-----------------|---------------|--------------------|---------------------|
| | (146.0) | (255.0) | (386.0) | (213.0) |
| <0 | no checking | no checking | no checking | 0<=X<200 |
| existing paid critical/other | existing credit | existing paid | existing paid | existing paid |
| new car | new car | radio/tv | radio/tv | furniture/equipment |
| 4093.0137 | 3158.8588 | 2854.7513 | 3674.3333 | |
| <100 | <100 | <100 | <100 | <100 |
| 4<=X<7 | >7 | 1<=X<4 | <1 | <1 |
| male single | male single | male single | female div/dep/mar | |
| 37.7534 | 43.2157 | 32.6528 | 30 | |
| unskilled resident | skilled | skilled | skilled | |
| bad | good | good | bad | |

Figure 14: **Appendix 14** - Clustering, Cluster Results

Clustered Instances

| | | |
|---|-----|--------|
| 0 | 146 | (15%) |
| 1 | 255 | (26%) |
| 2 | 386 | (39%) |
| 3 | 213 | (21%) |

Figure 15: **Appendix 15** - Clustering, Cluster Instances

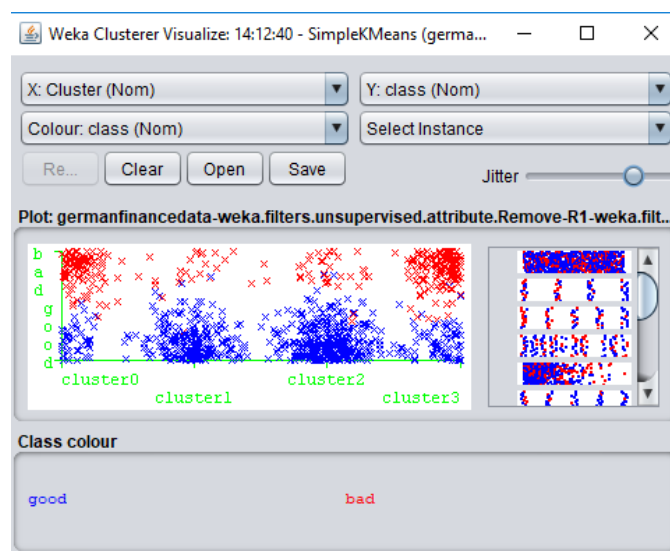


Figure 16: **Appendix 16** - Clustering, Visualize Data

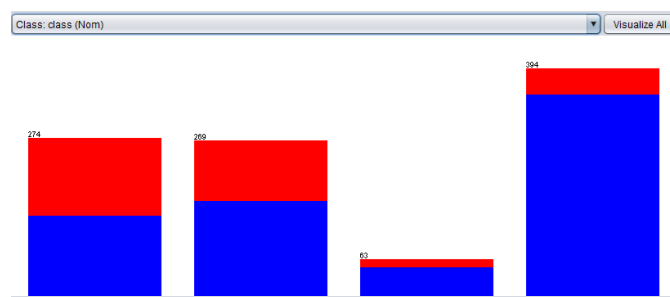


Figure 17: **Appendix 17** - Weka Visualizer, Stacked Graph

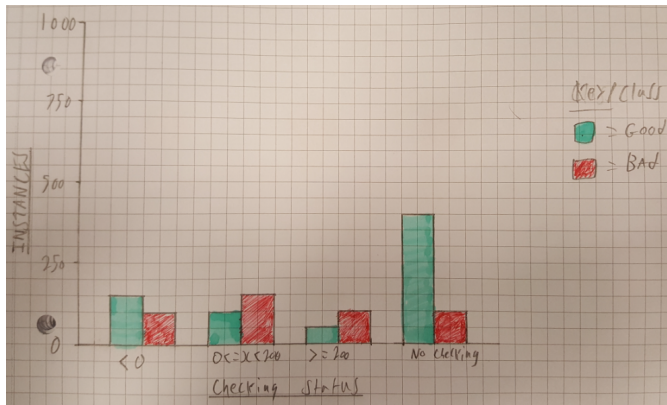


Figure 18: **Appendix 18** - Sketches, Grouped Bar Chart

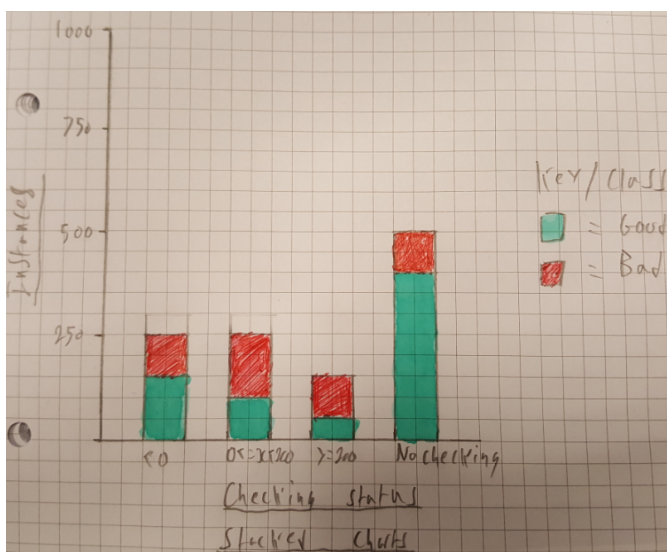


Figure 19: **Appendix 19** - Sketches, Stacked Bar Chart

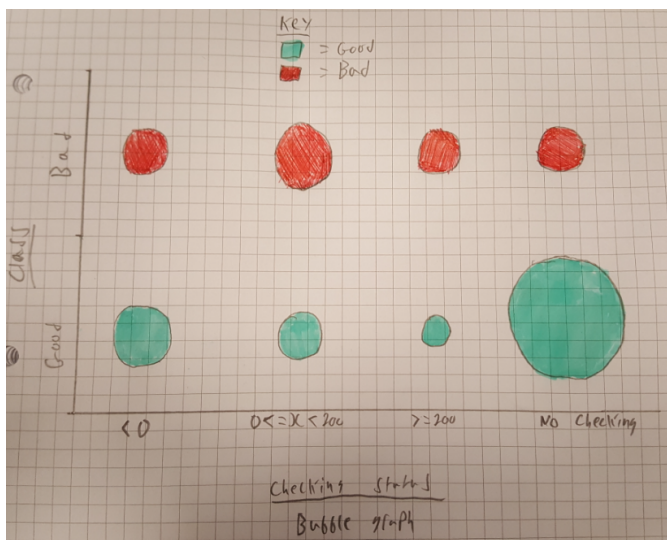


Figure 20: **Appendix 20** - Sketches, Bubble Graph