# Data mining and visualization for business intelligence

**Assessment**

# Practical Assessment

**Submitted by:**

Johan Sebastian Ramirez Vallejo

# Contents

## Figures

## Tables

# 1. Data Exploration

The dataset contains numeric type data, and it was retrieved by WEKA software, it was required to convert the data to nominal type data to help the algorithm to make sense of the data, in this case, the values represent a scale of 0 – 4, which represents categories where the records are distributed in this scale across to the attributes to determine the most common response across the survey. Therefore, Nominal data is used to label variables without providing any quantitative values. So, the records are segmented in each category. This helps to create a profile and formulate hypotheses according to the problems that are being analysed. For example, the hypothesis who tend to divorce to take prevention actions.

It is generally observed about the attribute values are distributed by the result class where "blue" means no result in divorce, and "Red" means results in divorce, demonstrating that couples who strongly agree (0 on the scale) with the attribute tend to **no divorce**, in contrast, who strongly disagree with the highest category (4 in the scales) the attribute tend to **divorce**.

## 1.1 Attributes Distribution

Each attribute has a numeric domain of 0-4. An attribute value of 0 means the participant strongly agreed with the attribute while 4 means the participant strongly disagreed with the attribute. The class 0 = Blue means- no divorce, and 1 = Red – Divorce. Selection Attributes, a feature in WEKA support to determine the correlation this attribute evaluator was correlation attribute evaluation with standards parameters. The results showed a rank list where the top is the most correlated attributes and the bottom the fewer correlated attributes. It is important to highlight that the attribute evaluator can evaluate nominal and numeric data.

The dataset was modified and created sub-datasets in terms of rank to evaluate the 6 most, least and author-correlated attributes as shown below.

### 1.1.1 Explanation of the 6 most correlated attributes

- 9. I enjoy travelling with my wife.
- 18. My spouse and I have similar ideas about how marriage should be
- 19. My spouse and I have similar ideas about how roles should be in marriage
- 35. I can insult my spouse during our discussions.
- 36. I can be humiliated when we have discussions.
- 40. We're just starting a discussion before I know what is going on.

### 1.1.2 Explanation of the 6 Least correlated attributes

- 6. We do not have time at home as partners.
- 43. I mostly stay silent to calm the environment a little bit.
- 45. I would rather stay silent than discussing with my spouse.
- 46. Even if I am right in the discussion, I stay silent to hurt my spouse.
- 48. I feel right in our discussions.
- 52. I would not hesitate to tell my spouse about her/his inadequacy.

### 1.1.3 Explanation of the 6 most correlated attributes Author

- 2. I know we can ignore our differences, even if things get hard sometimes.
- 6. We do not have time at home as partners.

- 11. I think that one day in the future when I look back, I see that my spouse and I have been in harmony with each other.

- 18. My spouse and I have similar ideas about how marriage should be

- 26. I know my spouse's basic anxieties.

- 40. We're just starting a discussion before I know what is going on.

## 1.2  Correlation Rank

The correlation rank was marked by the select attribute feature in WEKA, where the attributes were listed in a rank and, they were visualized to explore the distribution of each attribute.
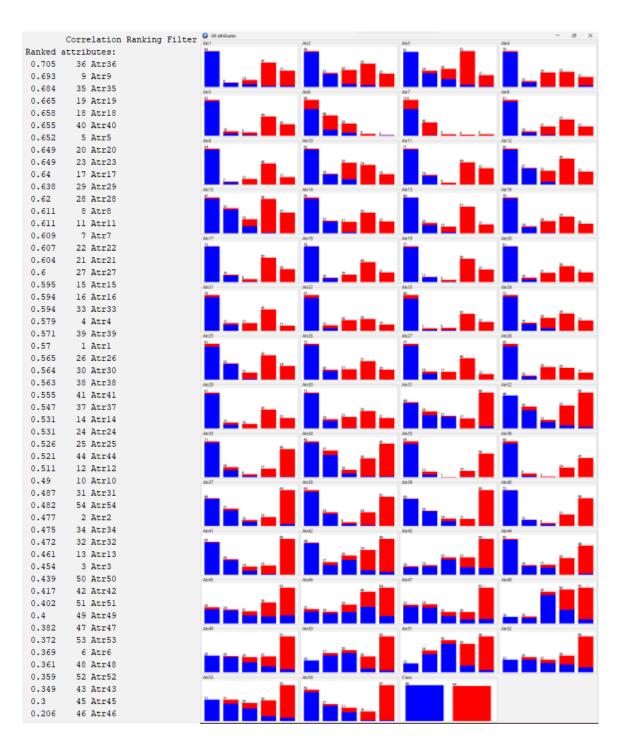


*Figure 1 Correlation Rank and Visualization*

### 1.2.1 Most correlated Attributes



Figure 2 Most correlated Attributes



Figure 3. Visualization of most correlation attributes

### 1.2.2 Least Correlated Attributes

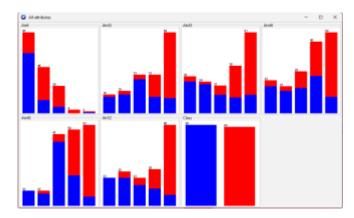| | | |
|---|---|---|
| 0.369 | 6 | Atr6 |
| 0.361 | 48 | Atr48 |
| 0.359 | 52 | Atr52 |
| 0.349 | 43 | Atr43 |
| 0.3 | 45 | Atr45 |
| 0.206 | 46 | Atr46 |

Figure 4. Least Correlation Attributes



Figure 5. Visualization least correlation attributes

### 1.2.3 Comparison 6 most correlation attributes and 6 least correlation attributes

Comparing the most correlated attributes against the least correlated attributes, show that people who strongly agree do not get a divorce (The positive skew), and people who are more disagree with the attribute tend to divorce (The negative skew). Then, the most correlated attributes showed that the subset of attributes is highly correlated to the class, so the segmentation of the class is determined by the extreme values in the 0-4 scale. On the other hand, the lower correlation does not show the correlation to both classes which means that people who agree with the attribute do not get a divorce but there is another proportion of people who disagree re also no get a divorce then cannot be conclusive to any hypothesis or predict the results if we see this uncorrelated attribute. For example, attribute 6 in the 6 least correlated attributes shows a positive skew for both classes that indicates that people who answer the question tend to agree and cannot be classified as who may get a divorce or not. Then, there is no group segmentation to determine the correlation to the class.

### 1.2.4 Author Correlation

```
Ranked attributes:
 0.658   4 Atr18
 0.655   6 Atr40
 0.611   3 Atr11
 0.565   5 Atr26
 0.477   1 Atr2
 0.369   2 Atr6
```
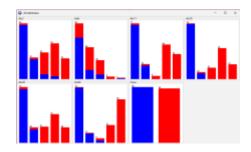
Figure 6. Author Rank attributes



Figure 7. Visualization author attributes

# 2. Model building

The dataset was explored, and it was created new sub-datasets which represent the most correlated attributes, and least correlated attributes given by the selection attribute WEKA feature, where was used the rank method to determine a list that indicates the less to more rank where 0 low correlate and 1 high correlate. The pre-process module in WEKA was selected and filtered the corresponding dataset and saved it. then it was created models with different techniques such as Artificial Neural Network ANN(MLP), Simple Tree J48, and Logistic Regression. Therefore, it was analysed the correlated attributes chosen by the author.
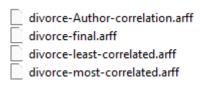
divorce-Author-correlation.arff
divorce-final.arff
divorce-least-correlated.arff
divorce-most-correlated.arff

*Figure 8. Data sets used*

## 2.1 Analysis datasets

It was determined the nominal dataset and sub-datasets, and they were used to create two models using two test options 10-fold cross-validation, and 70:30 train:test. The model measure for comparison was accuracy and kappa which indicate the results of each model.

| Classification Technique | Number of features (Dataset) | Accuracy (10fold cross valid.) | Kappa (10fold cross valid.) | Build Time (secs) | Accuracy (70:30 train:test split) | Kappa (70:30 train:test split) | Build Time (secs) |
|---|---|---|---|---|---|---|---|
| ANN(MLP) | All 54 attributes (divorce _nominal. arff) | 97.6471 % | 0.9529 | 42.32 sec | 96.0784 % | 0.9217 | 40.44 sec |
| J48 | | 95.2941 % | 0.9059 | 0.01 sec | 96.0784 % | 0.9217 | 0 sec |
| Logistic Regression | | 95.2941 % | 0.9059 | 0.17 sec | 96.0784 % | 0.9217 | 0.04 sec |
| ANN(MLP) | 6 most correlated attributes / (divorce _nominal _most6.arff) | 97.6471 % | 0.9529 | 1.43 sec | 94.1176 % | 0.8828 | 1.28 sec |
| J48 | | 95.8824 % | 0.9176 | 0 sec | 96.0784 % | 0.9217 | 0 sec |
| Logistic Regression | | 97.6471 % | 0.9529 | 0.01 sec | 98.0392 % | 0.9607 | 0.01 sec |
| ANN(MLP) | 6 least correlated attributes / (divorce_ Nominal _least6.arff) | 90    % | 0.8001 | 1.57 sec | 80.3922 % | 0.6065 | 1.48 sec |
| J48 | | 89.4118 % | 0.7881 | 0 sec | 74.5098 % | 0.4848 | 0 sec |
| Logistic Regression | | 88.8235 % | 0.7765 | 0.01 sec | 84.3137 % | 0.6852 | 0.01 sec |
| ANN(MLP) | Attributes proposed by the authors (divorce_ Nominal _author6.arff) | 98.2353 % | 0.9647 | 1.48 sec | 98.0392 % | 0.9607 | 1.54 sec |
| J48 | | 95.2941 % | 0.9058 | 0 sec | 96.0784 % | 0.9217 | 0 sec |
| Logistic Regression | | 98.2353 % | 0.9647 | 0.01 sec | 96.0784 % | 0.9217 | 0 sec |

*Table 1 Analysis table divorce data set*

## 2.2 Discussion

### Algorithms used

**Artificial neuron Network ANN(MLP):** This is a technique that was created to simulate how a human brain works in a simple way. These behave like neurons are interconnected to transmit the inputs generated by activation functions from sensory organs to the core. According to the distribution of the dataset, the values are automatically updated until achieving the target outputs values based on the learning rules. The network created after the training process can classify the data given for the test (Yöntem et al., 2019). The model accuracy is based on the number of neurons that best fit the data.

**Simple Tree J48:** It is a supervised learning approach and predictive decision support tool that create representation from observations to possible effects, characterized by nodes useful for decision-making, and at the end of every node is the outcome of making those decisions (Ngai et al., 2011). Some advantages stand in that requires little data preparation which means can handle

missing values. It can handle categorical and numerical data and can represent multi-output problems. the decision tree by its selves do not perform predictively compare to other algorithms, but, using it as a base of a complex algorithm such as random forest is useful (Hanafy & Ming, 2021).

**Logistic regression:** is a classical technique that removes the variables which have no relations with the claimed model. Therefore, the technique is helpful to predict the expected outcome of a binary dependent variable, and it is defined by a given set of predictor variables in this case a linear combination of risk factors that is relevant to the probability of observing an event. (Pesantez-Narvaez et al., 2019). The main characteristics are based that the maximum likelihood estimates are easily found, and the predictions are given 0 o 1 which is easily interpreted as the event probability event. Some disadvantages – it may lead to overfitting - The major limitation of Logistic Regression is the assumption of linearity between the dependent variable and the independent variables.

a.  **Compare your results for Accuracy (10fold cross valid.) and Accuracy (70:30 train:test split). Discuss the difference and justify the suitability of 10-fold cross-validation.**

Cross-validation is a standard evaluation technique, which consists in running in a systematic way a percentage split, for example, the dataset is split into 10 pieces or folds and uses each split as a testing dataset against the remaining 9 pieces. This process gives 10 evaluation results and computes the average; at the 11th time, it is invoked the data set to obtain the final model. Therefore, this process takes more time to build a model, the results indicate that ANN (MLP) algorithm takes more time than the 70:30 train: Test validation. This indicates that building the model takes 42.32 sec, but this takes more time to make the cross-validation against the model, cross-validation does this repetitive time until it reaches the k factor to get the final test. Then, the total time obtained is approximately 7 min to build the model and evaluate the model. Hence, 70:30 train:test split validation is faster in terms that this validation method split the dataset in 70% for training and 30% for testing and builds the final model in 40.44 sec. The final result indicates that the model with cross-validation obtained an accuracy of 97.6471 %, slightly superior to 70:30 train:test which got 96.0784 %. This shows that cross-validation has better performance, but this does not work well for all models, for example, the performance on J48 and logistic regression is less accurate than the 70:30 train:test. The purpose of cross-validation is to evaluate the predictive ability of the model with new data, and it offers the option of evaluating the data with different data sets to obtain more accurate results.

b.  **From your table, analyses the results of 6 most correlated attributes and attributes proposed by the authors. Can you justify the difference, if any? [5 marks]**

Comparing the results obtained between the 6 most correlated attributes selected in the selection method by WEKA, and the 6 attributes selected by the author, it is conclusive that the attributes observed by an expert can improve considerably the accuracy of the model even using different algorithms. The analysis table 1 shows that the accuracy of the algorithm ANN(MLP), J48 and logistic Regression are 97.6471 %, 95.8824 %, and 97.6471 % respectively, they are slightly less than the accuracy of the author choose, are 98.2353 %, 95.2941 %, 98.2353 % of the same algorithms. This makes more effective attributes that were used by the author. Hence, the algorithm ANN(MLP) and logistic regression have better performance on the attributes selected by the author correlation feature, their attributes have a common meaning and failed attempts to repair, love map and negative conflict behaviours (Yöntem et al., 2019).

In the document it is clear that the author used values of significance to choose the attributes given by the correlated-base feature selection method, but, first It is difficult to say if the author use the correlation with numeric or nominal values because although the correlation feature can handle nominal and numeric data (Hall, 1999), which the author do not indicate how was used. On the other hand, the correlation feature ranks the attributes indicating the correlation coefficient and the author indicates that used the values of significance obtained by applying the correlation base feature from WEKA. Then, it is not comparable to the method, but effectively slightly increases the

performance of the model. Therefore. It is conclusive to say that the author used their expertise in the subject to determine the attributes combined with an uncertain method

The author attributes show correlation with the class except attribute 6 which does not clearly show the data segmentation or correlation to a specific class. it could be concluded that eliminating attribute 6 and building the model again would expect to obtain the same rate of accuracy

**c. From your table, compare the results of the 6 most correlated attributes to All 54 attributes. What do they tell you about ANN and logistic regression? [5 marks]**

Methods were applied to all 54 attributes and 6 most correlated attributes. The results obtained with direct application of methods with all 54 attributes showed that the highest success rate was 97.6471 % by ANN using 10-fold cross-validation while the logistic regression method obtained 96.0784 %. In contrast, the 6 most correlated attributes were 97.6471 % for ANN and 98.0392 % for logistic regression which is the highest in this comparison. the most successful result is obtained with the logistic regression model applied together with the correlation and using the test evaluation 70:30 train:test. However, it is recognised that the cross-validation purpose of cross-validation is to evaluate the predictive ability of the model with new data, then, the model built with cross-validation is more reliable. In addition, logistic regression is less prone to overfitting than ANN due to that it is simpler compared to the complex structure created by the ANN technique. This also is shown during the build model time where ANN take 42.36 secs when are 54 attributes and 1.40 secs when there are 6 attributes while logistic regression takes 0.01 secs to build the models in both cases. then again, logistic regression overperformed the ANN technique to predict and classify accurately faster.

**d. Given that accuracy of predictions is a priority, for each of the 4 configurations of datasets, identify and justify which model you will choose.**

The table above shows different attempts to build a model using the divorce dataset and correlation-based feature to evaluate the techniques against subsets that show differences in the correlation results. The study demonstrates the effectiveness of using the correlation feature to identify the critical attributes that can determine the resulting class, in this case, divorce or not divorce. The analysis is also used to evaluate evaluation where cross-validation remains as effective to the capability to evaluate new data and increase the accuracy of a model. Therefore, the model suggested is the logic regression that used cross-validation on the 6 most effective attributes chosen by the author understanding that those attributes choosing by the experts. The accuracy rate obtained was 98.2353% which is the same obtained by the ANN technique, but this is prone to overfitting for the complex structure, where Logistic regression is less and is simplest. Another reason is that Logistic regression requires less computational resources than the others. then, the best model to predict the outcome is Logistic Regression.

# 3. References

Hall, M. A. (1999). Correlation-based feature selection for machine learning.

Hanafy, M., & Ming, R. (2021). Machine learning approaches for auto insurance big data. *Risks, 9*(2), 42.

Ngai, E. W., Hu, Y., Wong, Y. H., Chen, Y., & Sun, X. (2011). The application of data mining techniques in financial fraud detection: A classification framework and an academic review of literature. *Decision support systems, 50*(3), 559-569.

Pesantez-Narvaez, J., Guillen, M., & Alcañiz, M. (2019). Predicting motor insurance claims using telematics data—XGBoost versus logistic regression. *Risks, 7*(2), 70.

Yöntem, M. K., Kemal, A., Ilhan, T., & KILIÇARSLAN, S. (2019). Divorce prediction using correlation based feature selection and artificial neural networks. *Nevşehir Hacı Bektaş Veli Üniversitesi SBE Dergisi, 9*(1), 259-273.

# 4. Appendix A Results

The system used in the study for the application Intel(R) Core (TM) i7-8550U CPU @ 1.80GHz 1.99 GHz processor, 8.00 GB DDR3 RAM.

| Classification Technique | Number of features (Dataset) | Model (10fold cross valid.) | Model 70:30 train:test split) |
|---|---|---|---|
| ANN(MLP) | All 54 attributes (di.vorce_nominal.arff) | | |
| J48 | | | |
| Logistic Regression | | | |
| ANN(MLP) | 6 most correlated attributes / (divorce_nominal_most6.arff) | | |

| | | | |
|---|---|---|---|
| J48 | | | |
| Logistic Regression | | | |
| ANN(MLP) | 6 least correlated attributes / (divorce_nominal_least6.arff) | | |
| J48 | | | |

| | | | |
|---|---|---|---|
| Logistic Regression | | | |
| ANN(MLP) | | | |
| J48 | Attributes proposed by the authors (divorce_nominal_author6.arff) | | |
| Logistic Regression | | | |