

Name: *Daniel Velek, Dario Cavada, Joshua Bringhurst*
Course: *BAN615DE 23FA Predictive Business Analytics*
Professor: *Dr. Shaikh*
Assignment: *Final Project*
Due Date: *Saturday, December 16th*

Table of Contents

1) Abstract	1
2) Introduction	2
3) Conclusion	5
4) Appendix	6
4.1) SAS Diagram	
4.2) Decision Tree	
4.3) Decision Tree HP Forest	
4.4) Decision Tree Interactive	

1) Abstract

The text Information Systems Education Journal (ISEDJ) Volume 19, No. 1, February 2021, features a special issue on Fraud Cases. This abstract is focused on the case titled "Can you Predict the Money Laundering Cases?" authored by Richard V. McCarthy, Wendy Ceccucci, Mary McCarthy, and Nirmalkumar Sugumar. My understanding is this case, is now intended for business analytics courses, and revolves around predictive analytics applied to anti-money laundering efforts. There students are tasked with developing and optimizing predictive models using data from People's United Bank. The case explores the challenges of improving upon the bank's baseline model and emphasizes the importance of accurate detection in addressing the critical issue of money laundering. The significant importance of financial crimes, more specifically money laundering, poses an enlarged threat and challenge for banks in complying with regulatory requirements. In order to address this, the case "Can you Predict the Money Laundering Cases?" expands into the realm of business analytics, focusing on predictive analytics as a tool to enhance anti-money laundering efforts. This introduction provides an overview of the case's relevance, the background of money laundering, and the vital role of predictive models in identifying suspicious transactions. This case sets the stage for students to navigate through the steps of hypothesis development, data analysis, and model development, ultimately aiming to surpass the baseline model's effectiveness in detecting potential money laundering cases. Through this case, students will gain practical insights into the application of predictive analytics in a real-world scenario within the banking industry.

2)Introduction

In this Money Laundering case, the prediction of money laundering cases involved a critically thought-out approach to business analytics, specifically leveraging their predictive analytic techniques. The authors utilized a dataset that was provided by People's United Bank,

consisting of 38,515 total transactions sampled from approximately October 2014 to September 2015. The baseline model, created by the bank through six iterations of predictive models, served as the starting point in order to begin analysis. We found the best fit model by using multiple processes including these statistical measures such as the Kolmogorov-Smirnov (KS) and Receiver Operating Characteristic (ROC) statistics, with a focus on optimizing predictive accuracy.

The various predictive modeling techniques, including regression, were utilized to enhance the model's fit. Predictive models were developed and validated using a 70/30 split of the data, where 70% was used for model building, and 30% for validation. The criteria for selecting the best fit model included an emphasis on minimizing misclassification rates and evaluating Type I and Type II errors. In the final model, it aimed to pass the benchmark set by People's United Bank, providing a more defined identifier of transactions that required investigation for potential money laundering. Through this process, the authors sought to demonstrate the practical application of predictive analytics in improving the efficiency and accuracy of anti-money laundering efforts within the banking sector.

The importance of our project is underscored by its direct alignment with our SAS Enterprise Mining diagram (see appendix 1.1). Recognizing that our level of expertise may not reach that of the professionals showcased in the case study, we embarked on our predictive business statistical analysis. Our method encompassed the utilization of diverse metrics, including KS Statistics, ASE, Misclassification Rates, ROC, Decision Trees, and Neural Network nodes. These metrics played a pivotal role in the discovery of key insights, ranging from identifying missing values to evaluating skewed variables and even pinpointing potential outliers.

Through our meticulous analysis, we unearthed that the HPDMForest had emerged as the optimal Neural Network Node, a discovery validated by both the KS Statistic of 0.519 and the ROC index of 0.85. The selection of this node underscored its superior performance within the context of our predictive modeling endeavors. The application of KS Statistics and ASE of 0.123195 facilitated an in-depth assessment of the discriminatory power inherent in our models, while Misclassification Rate of 0.181965 had provided indispensable insights into the precision of our predictions. The ROC index, offering a comprehensive measure of our model's class-distinguishing capability, significantly contributed to strengthening the overall soundness of our analytical framework.

There was a total of 6 variables with a right skewness, starting largest to smallest:

- Wires_Size – 17.0522
- Num_Related – 15.3631
- Num_Trigger – 15.2936
- Wires_Mult – 15.0994
- Num_Acct_Alert – 8.9476
- Num_Tran_Alert – 6.6717

Furthermore, our exploration of Decision Trees (see appendix 2.1 & 2.3) significantly enriched our comprehension of variable importance and relationships within the dataset. This combination of many and various techniques in a predictive analysis highlights our steadfast dedication to withdraw meaningful insights despite the potential limitations in our practical sophistication. This all-inclusive approach instilled confidence in our required transformation and in addressing the potential outliers, thereby reinforcing the reliability of our analytical outcomes.

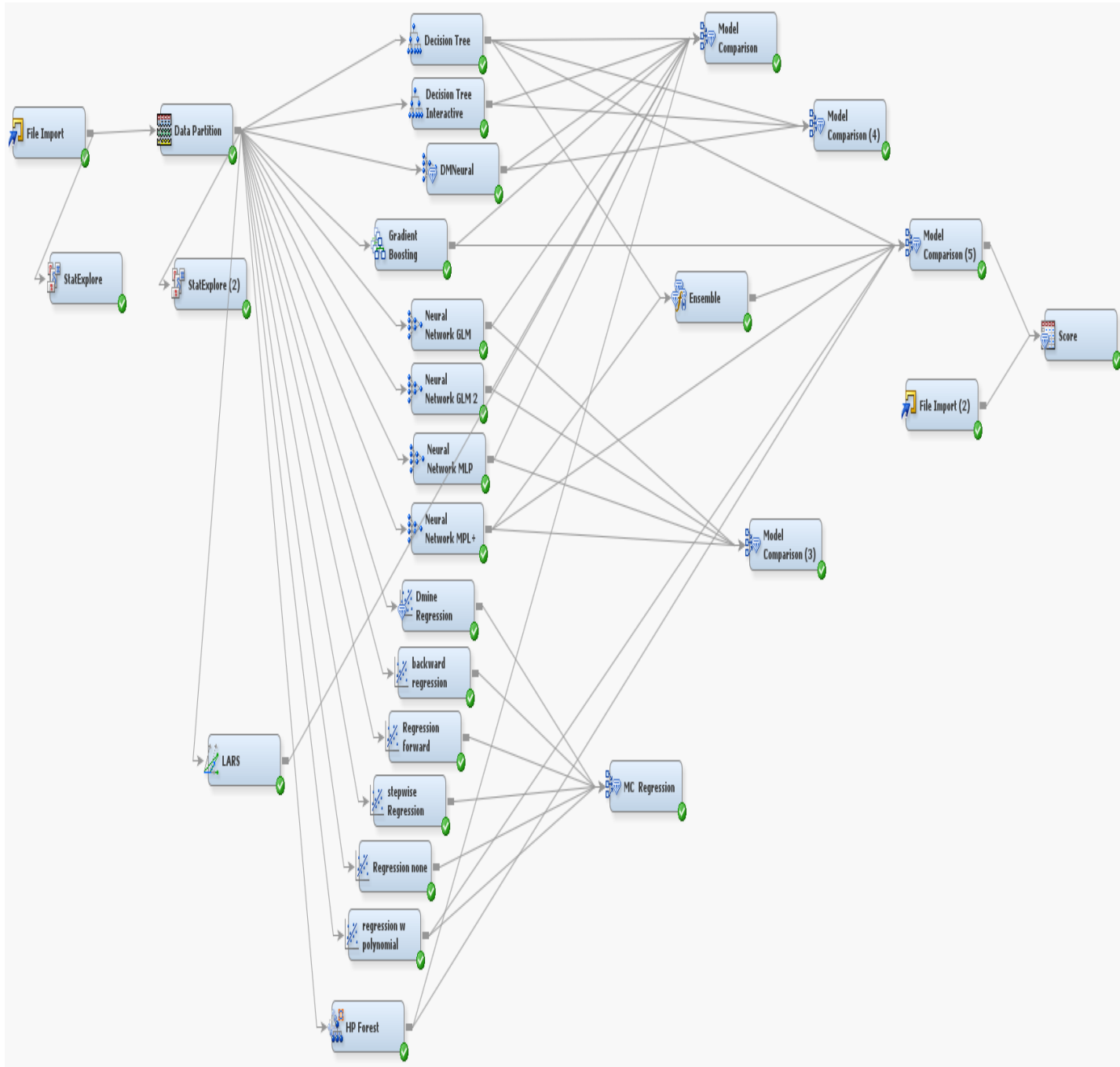
Lastly, our journey through this predictive business statistical analysis has been both enlightening and fruitful. The alignment of our project with the SAS Enterprise Mining diagram provided a solid foundation for our exploration. Acknowledging that our proficiency does not rival that of the professionals who are featured in the case study, we navigated through diverse metrics, including KS Statistics, ASE, Misclassification Rates, ROC, Decision Trees, and Neural Network nodes. What was deemed surprising is the fact that the HPDMForest node emerged as the optimal node, supported by both the KS Statistic and the ROC index, signifies a critical milestone in our predictive modeling efforts. This finding underscores the sturdiness and effectiveness of our chosen approach. The application of KS Statistics and ASE allowed us to rigorously assess the power of our models, while insights from Misclassification Rates provided a tangible measure of accuracy. The ROC index, with its comprehensive evaluation of its capabilities, further solidified the reliability of our analytical work.

3) Conclusion

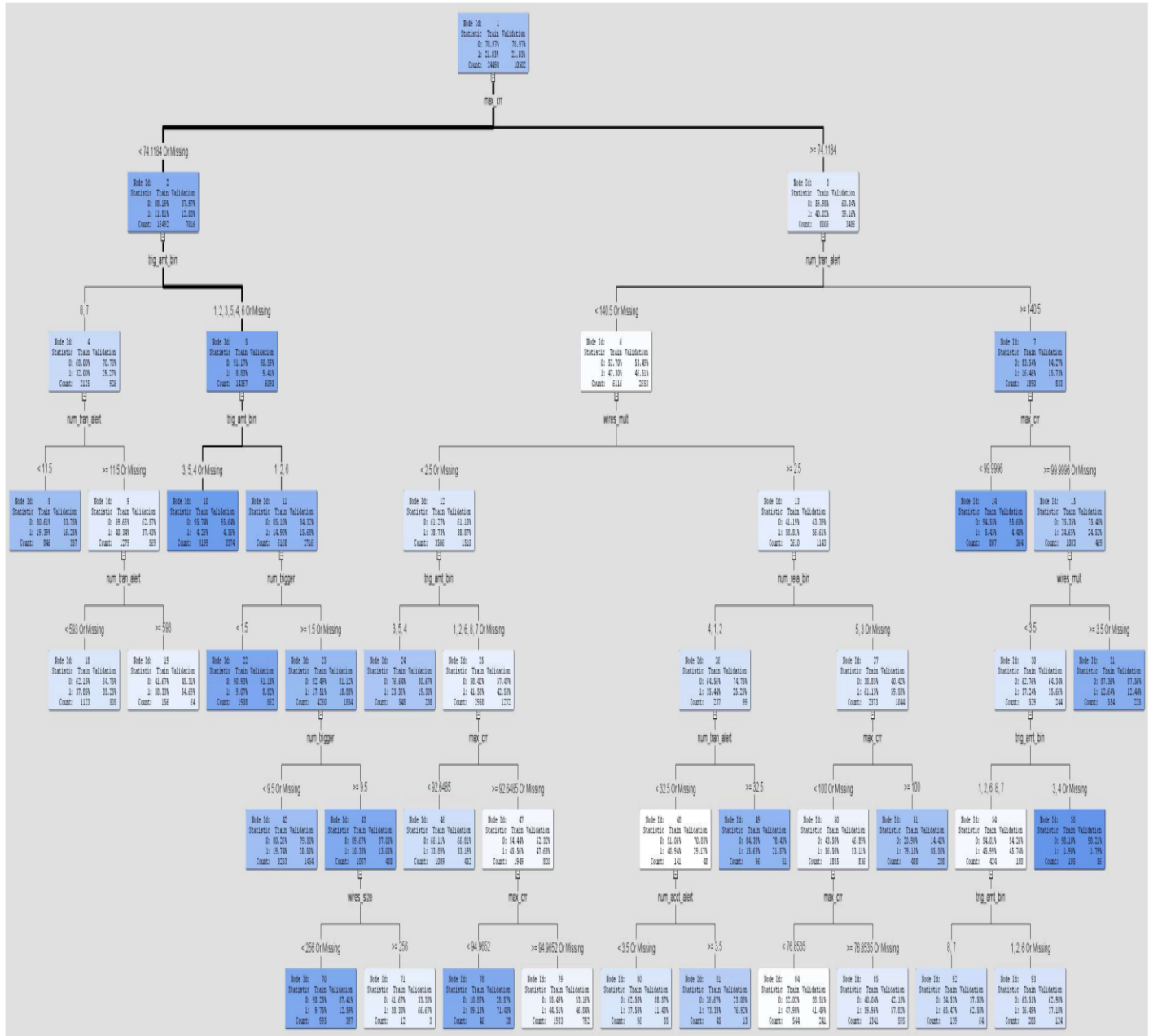
In conclusion, our journey through predictive analytics has significantly enhanced our understanding of all variable importance and its relationships within the dataset, adding depth to the overall analysis. With these techniques, our group has successfully conducted a thorough predictive analysis, showcasing our commitment to extracting meaningful insights, even when faced with limitations in hi-tech sophistication. This extensive approach instilled confidence in the integrity of us and our data, affirming that no transformational adjustments were required, and all outliers were effectively addressed.

4) Appendix

4.1 SAS Diagram



4.2 Decision Tree



4.3 Decision Tree HP Forest

Results - Node: HP Forest Diagram: final

File Edit View Window

Score Rankings Overlay: Prod_ind

Cumulative Lift

Depth

— TRAIN — VALIDATE

Leaf Plot

Number of Leaves

Number of Trees

■ Base ■ Increment

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation
Prod_ind	Prod_ind	_ASE_	Average Sq...	0.118724	0.1231
Prod_ind	Prod_ind	_DIV_	Divisor for A...	48996	210
Prod_ind	Prod_ind	_MAX_	Maximum A...	0.969569	0.9752
Prod_ind	Prod_ind	_NOBS_	Sum of Fre...	24498	105
Prod_ind	Prod_ind	_RASE_	Root Avera...	0.344563	0.3505
Prod_ind	Prod_ind	_SSE_	Sum of Squ...	5816.995	2587.5
Prod_ind	Prod_ind	_DISF_	Frequency...	24498	105
Prod_ind	Prod_ind	_MISC_	Misclassific...	0.171402	0.1815
Prod_ind	Prod_ind	_WPRMNC_	Number of	1100	10

Variable Importance

Variable Name	Number of Splitting Rules	Train: Gini Reduction	Train: Margin Reduction	OOB: Gini Reduction	OOB: Margin Reduction
max_crr	1112	0.021923	0.043845	0.020645	0.0423
num_tran_...	693	0.004587	0.009174	0.002959	0.0074
trig_amt_bin	666	0.019146	0.038292	0.017925	0.0371
num_tran_L	495	0.007447	0.014895	0.005448	0.0130
num_trigger	479	0.002465	0.004930	0.001665	0.0041
num_acct_...	464	0.002049	0.004097	0.001249	0.0032
wires_size	443	0.002801	0.005601	0.001716	0.0044
num_acct	441	0.007893	0.015785	0.007494	0.0145

Iteration Plot

Misclassification Rate

Misclassification Rate

Number of Trees

— Train — Out of Bag — Validate

Output

```

1  *-----*
2  User:           dvelek
3  Date:           December 11, 2023
4  Time:           15:21:13
5  *-----*
6  * Training Output
7  *-----*
8
9
10
11
    
```

Leaf Statistics

Frequency

Number of Leaves

Iteration History

Number of Trees	Number of Leaves	Average Square Error (Train)	Average Square Error (Out of Bag)	Average Square Error (Validate)	Misclassification Rate (Train)	Misclassification Rate (Out of Bag)	Misclassification Rate (Validate)	Log Loss (Train)	Log Loss (Out of Bag)	Log Loss (Validate)
1	62	0.133	0.134	0.137	0.187	0.188	0.196	0.418	0.434	0.421
2	125	0.125	0.131	0.129	0.182	0.186	0.191	0.390	0.422	0.402
3	201	0.122	0.129	0.127	0.176	0.185	0.186	0.384	0.412	0.398
4	264	0.122	0.128	0.126	0.177	0.185	0.187	0.383	0.410	0.395
5	342	0.121	0.127	0.126	0.173	0.183	0.184	0.381	0.405	0.393
6	411	0.120	0.127	0.125	0.172	0.182	0.183	0.380	0.401	0.393
7	476	0.121	0.126	0.125	0.173	0.182	0.184	0.380	0.399	0.393
8	554	0.120	0.126	0.125	0.173	0.182	0.184	0.379	0.396	0.392
9	609	0.120	0.126	0.125	0.173	0.182	0.184	0.379	0.396	0.392

4.4 Decision Tree Interactive

