







## CONCLUSION

### Results:

The **adjusted Standard Random Forest** was selected as the final model based to its superior overall accuracy. According to the confusion matrix, **Class 2** has the best accuracy, while **Class 1** has the lowest and is incorrectly labeled as Class 3.

It reflects the scatterplot seen above. Class 1 is **intermixed** with Class 3.

As previously stated, **class imbalance** contributed to bias, and Class1, with the smallest size, has the lowest accuracy and F1 score.

### Recommendations

To tackled the class imbalance issue, techniques such as Balanced Random Forest(BRF), SMOTE (Synthetic Minority Oversampling Technique) using Standard Random Forest, can be explored. These techniques can enhance the performance of the minority class without compromising the model's overall accuracy, as discussed in the article "Surviving in a Random Forest with Imbalanced Datasets" (SFU CSPMP, 2024).

### Caution

While the methods mentioned above may boost the F1 score of the minority class, it is important to also consider maintaining overall accuracy and minimizing misclassification of the majority class.

## DATA SOURCES

This data is provided by the professor for the imPlementation of the project CSP571.

## SOURCE CODE

### Listings

- Jupyter Notebook (.ipynb & .pdf)
- Dataset (.csv/as a .zip) & parquet)
- ONNX File (pipeline.onnx)

To use the trained model in runtime 3. ONNX Model Diagram (.pdf)  
The graph for the ONNX Model is attached as separate pdf file, due to issues in the diagram appropriately in the notebook pdf. The model's architecture and conversion process are detailed, ensuring reproducibility and ease of deployment. 4. Requirements (.txt)

### Dependencies:

The project relies on several open-source libraries, primarily:

- Python 3.10+
- Pandas for data manipulation.
- ONNXRuntime 1.18+ for inference, to run the IR 10 model.
- Seaborn and Matplotlib for data visualization.
- Scikit-Learn 1.2+ for building and evaluating models.
- Graphviz for visualizing model architecture.
- ONNX 1.16+ for model conversion, which creates a model in the ONNX IR version of 10.
- Numpy for numerical computations.
- PyArrow for efficient data storage in Parquet format.

### Installation Instructions

- Install the required dependencies.
- Install the graphviz also into the system using the link provided in the references.
- Open the Jupyter Notebook and run the cells sequentially.

### Reproducibility

The project is developed with reproducibility in mind, offering detailed instructions and all essential files to allow others to replicate and expand upon the work.

## 10. Bibliography

(Chicago style – AMS/AIP)

- ChatGPT. "Assistance with Project Report on the Abstract, Overview & Bibliography sections." OpenAI, accessed November 15, 2024.

<https://chatgpt.com/?model=gpt-4o>.

- Stack Overflow. "Pandas Concat Increases Number of Rows." Accessed November 19, 2024.

<https://stackoverflow.com/questions/50368145/pandas-concat-increases-number-of-rows>.

- Analytics Vidhya. "Evaluating a Random Forest Model." Medium, accessed November 19, 2024.

<https://medium.com/analytics-vidhya/evaluating-a-random-forest-model-9d165695ad56>.

- ONNX Runtime. "Compatibility." Accessed November 19, 2024.

<https://onnxruntime.ai/docs/reference/compatibility.html>.

- Graphviz. "Graphviz." Accessed November 19, 2024.

<https://pypl.org/project/graphviz/>.

- SFU CSPMP. "Surviving in a Random Forest with Imbalanced Datasets." Medium, accessed November 19, 2024.

<https://medium.com/sfu-csmpmp/surviving-in-a-random-forest-with-imbalanced-datasets-b98b963d52eb>.