# Prediction of Bankruptcy

April. 2020

AeLee Im

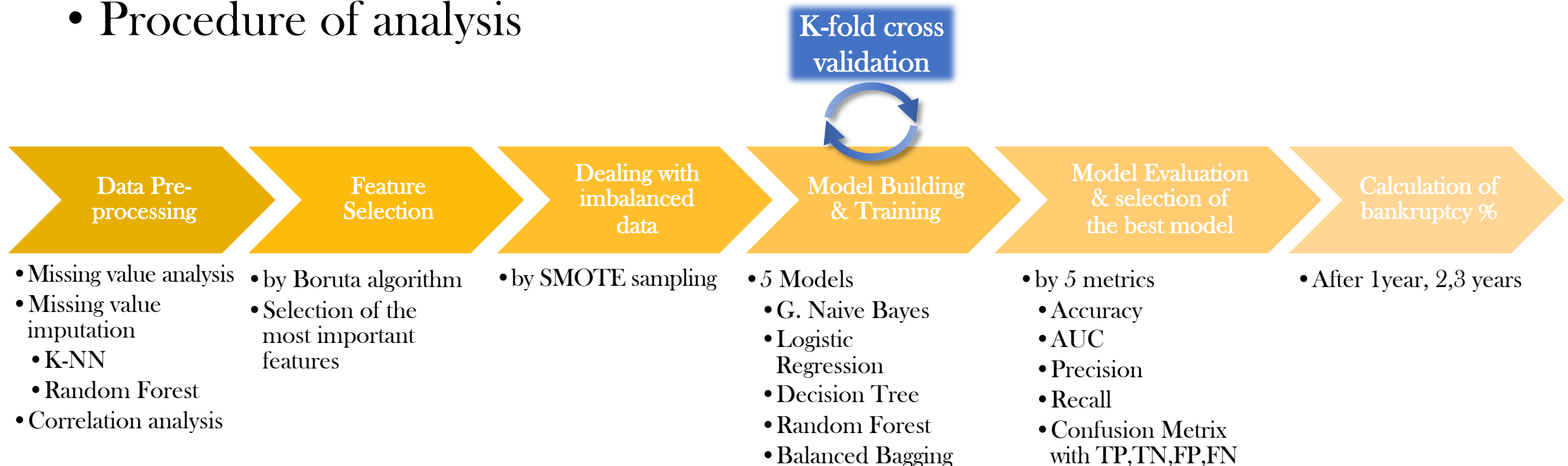# Contents

- **Project overview**
  - Objective
  - Procedure
- **Analysis methods**
- **Analysis results**
- **Next steps**
- **Appendix**

# Project Overview

- Objective of project : Prediction of bankruptcy of company
  - Forecast the likelyhood of bankruptcy of companies
  - What is the probability of bankruptcy after 1 year and 2, 3 years?

- Procedure of analysis

**K-fold cross validation**

| Data Pre-processing | Feature Selection | Dealing with imbalanced data | Model Building & Training | Model Evaluation & selection of the best model | Calculation of bankruptcy % |
|---|---|---|---|---|---|

- Missing value analysis
- Missing value imputation
  - K-NN
  - Random Forest
- Correlation analysis

- by Boruta algorithm
- Selection of the most important features

- by SMOTE sampling

- 5 Models
  - G. Naive Bayes
  - Logistic Regression
  - Decision Tree
  - Random Forest
  - Balanced Bagging

- by 5 metrics
  - Accuracy
  - AUC
  - Precision
  - Recall
  - Confusion Metrix with TP,TN,FP,FN

- After 1year, 2,3 years

# Analysis Methods

### How to deal with missing values?
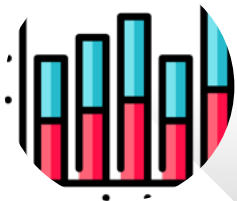✓ Impute missing value with k-NN, Random Forest imputation technique

### How to deal with highly correlated & numerous features?
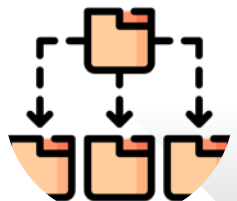• Without removing correlated features, choose the Boruta selection algorithm

### How to deal with imbalanced dataset?
• Make each class balanced by SMOTE oversampling technique

### How to deal with not enough training data?
• Training with every data and iteration with K-fold cross validation

## Benefit of solution

Reduce the loss of data & beneficial imputation for high dimensional data

Considering of correlation between features during selection of important features

No loss of data via oversampling & preparing well balanced dataset for high dimensional data

No bias with certain part of dataset & reliable model training results (No over-/underfitting problem)

# Analysis Results : Model comparision & selection

- The best performed model : Balanced Bagging Model
  - ✓ The most accurate prediction performance with highest metrics score (e.g. AUC 0.94 ~ 0.95 etc.)
  - ✓ Excellent performance for every dataset (1year~5year)
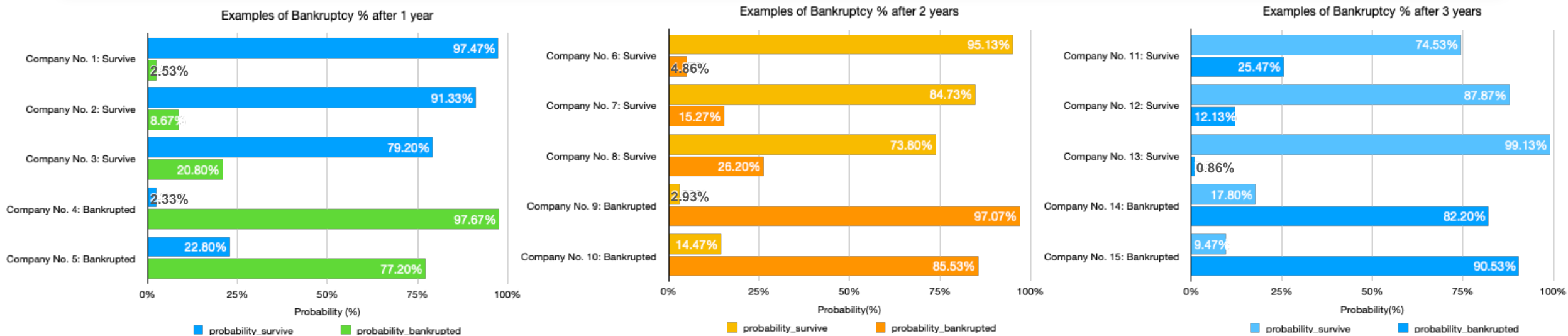  - ✓ No difference between k-NN and Random Forest imputation methods

→ **Final Model Decision : Balanced Bagging Model (trained with Random Forest imputed dataset)**

| Model performance with k-NN imputed dataset | | | | | | Model performance with Random Forest(missForest) imputed dataset | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Imputer** | **Model** | **Data set(Year)** | **AUC** | **Accuracy** | **Precision** | **Recalls** | | | | |
| k-NN | Balanced Bagging | 1year | 0.9559 | 0.9560 | 0.9718 | 0.9391 | | | | |
| | | 2year | 0.9494 | 0.9494 | 0.9670 | 0.9305 | | | | |
| | | 3year | 0.9466 | 0.9466 | 0.9659 | 0.9258 | | | | |
| | | 4year | 0.9459 | 0.9459 | 0.9752 | 0.9152 | | | | |
| | | 5year | 0.9476 | 0.9476 | 0.9689 | 0.9252 | | | | |
| | Random Forest | 1year | 0.9256 | 0.9257 | 0.9383 | 0.9115 | | | | |
| | | 2year | 0.9148 | 0.9148 | 0.9289 | 0.8982 | | | | |
| | | 3year | 0.9173 | 0.9173 | 0.9324 | 0.8998 | | | | |
| | | 4year | 0.9168 | 0.9168 | 0.9395 | 0.8909 | | | | |
| | | 5year | 0.9217 | 0.9217 | 0.9322 | 0.9097 | | | | |
| | Decision Tree | 1year | 0.8950 | 0.8950 | 0.9078 | 0.8792 | | | | |
| | | 2year | 0.8781 | 0.8781 | 0.8904 | 0.8623 | | | | |
| | | 3year | 0.8806 | 0.8806 | 0.8894 | 0.8692 | | | | |
| | | 4year | 0.8818 | 0.8818 | 0.8953 | 0.8648 | | | | |
| | | 5year | 0.8971 | 0.8971 | 0.9124 | 0.8786 | | | | |
| | Gaussian Naive Bayes | 1year | 0.5036 | 0.5036 | 0.5908 | 0.0232 | | | | |
| | | 2year | 0.5081 | 0.5081 | 0.6684 | 0.0326 | | | | |
| | | 3year | 0.5137 | 0.5137 | 0.7406 | 0.0432 | | | | |
| | | 4year | 0.5078 | 0.5078 | 0.5663 | 0.0661 | | | | |
| | | 5year | 0.5189 | 0.5187 | 0.6077 | 0.2350 | | | | |
| | Logistic Regression | 1year | 0.6174 | 0.6154 | 0.6174 | 0.7386 | | | | |
| | | 2year | 0.5771 | 0.5765 | 0.6523 | 0.5194 | | | | |
| | | 3year | 0.6380 | 0.6389 | 0.7129 | 0.5833 | | | | |
| | | 4year | 0.6477 | 0.6479 | 0.6179 | 0.7873 | | | | |
| | | 5year | 0.6960 | 0.6960 | 0.7622 | 0.5697 | | | | |

| **Imputer** | **Model** | **Data set(Year)** | **AUC** | **Accuracy** | **Precision** | **Recalls** |
|---|---|---|---|---|---|---|
| missForest | Balanced Bagging | 1year | 0.9657 | 0.9657 | 0.9806 | 0.9502 |
| | | 2year | 0.9628 | 0.9628 | 0.9793 | 0.9455 |
| | | 3year | 0.9448 | 0.9448 | 0.9685 | 0.9196 |
| | | 4year | 0.9464 | 0.9464 | 0.9755 | 0.9159 |
| | | 5year | 0.9439 | 0.9439 | 0.9614 | 0.9251 |
| | Random Forest | 1year | 0.9454 | 0.9454 | 0.9610 | 0.9283 |
| | | 2year | 0.9283 | 0.9283 | 0.9411 | 0.9137 |
| | | 3year | 0.9138 | 0.9138 | 0.9278 | 0.8974 |
| | | 4year | 0.9173 | 0.9173 | 0.9395 | 0.8921 |
| | | 5year | 0.9206 | 0.9206 | 0.9353 | 0.9038 |
| | Decision Tree | 1year | 0.9234 | 0.9235 | 0.9330 | 0.9122 |
| | | 2year | 0.8869 | 0.8870 | 0.8986 | 0.8725 |
| | | 3year | 0.8762 | 0.8761 | 0.8869 | 0.8623 |
| | | 4year | 0.8806 | 0.8806 | 0.8946 | 0.8629 |
| | | 5year | 0.8874 | 0.8873 | 0.9016 | 0.8698 |
| | Gaussian Naive Bayes | 1year | 0.5112 | 0.5112 | 0.6812 | 0.0422 |
| | | 2year | 0.5193 | 0.5193 | 0.7405 | 0.0584 |
| | | 3year | 0.5225 | 0.5225 | 0.7787 | 0.0632 |
| | | 4year | 0.5146 | 0.5147 | 0.6015 | 0.0879 |
| | | 5year | 0.5325 | 0.5315 | 0.5695 | 0.6195 |
| | Logistic Regression | 1year | 0.6548 | 0.6541 | 0.7596 | 0.4556 |
| | | 2year | 0.5084 | 0.5097 | 0.5138 | 0.5665 |
| | | 3year | 0.5968 | 0.5968 | 0.5628 | 0.8666 |
| | | 4year | 0.6621 | 0.6620 | 0.6424 | 0.7330 |
| | | 5year | 0.7250 | 0.7251 | 0.7638 | 0.6528 |

Evaluation metrics

# Analysis Results : Computation of Bankruptcy %

- Test computation of chance of bankruptcy (%) after 1year, 2, 3 years
  - ✓ By applying selected prediction model : Pre-trained model with Balanced Bagging algorithm
  - ✓ Shows the probability(%) of survive and bankruptcy as an evidence of bankruptcy prediction

**Q. Will Company be bankrupted after 1 year? Or 2 years? Or 3 years?**



According to the prediction analysis result,
→ Company No. 1 will be survived with 97.47% probability after 1 year.
→ However, Company No. 10 will be bankrupted with 85.53% probability after 2 years!

# Next steps

## Diversify

- **Additional prediction of accurate date/month/year of bankruptcy**
  - Applying Time series data to predict the precise moment of bankruptcy
  - Q. When this company will be bankrupted?

## Advancing

- **Try another Ensemble model**
  - Boosting method (e.g. Xgboost)
  - Beneficial for high-bias data
  - Compare with our current best model (i.e. by Bagging method) for further model enhancement

## Accelerate

- **Real time monitoring available**
  - Generate auto-updating dashboard
  - Deployment of developed model with implementation of database
  - No waiting time for reporting

## Big data Acceptable

- **Convert to Cloud based analysis**
  - With big data analysis systems (e.g. Hadoop, HIVE, Apache)
  - Deep learning approach available due to increase of data volume

# Appendix.

**Results of data processing**

1. Missing value analysis

2. Correlation analysis

3. Data balancing analysis

4. Feature selection

# Appendix. 1. Missing value analysis

- Original dataset has high volume of missing values (53.99%)
  - ✓ Each dataset has around 50% of missing values (48.71~59.82%)
  - ✓ Most of missing values are concentrated on certain attributes
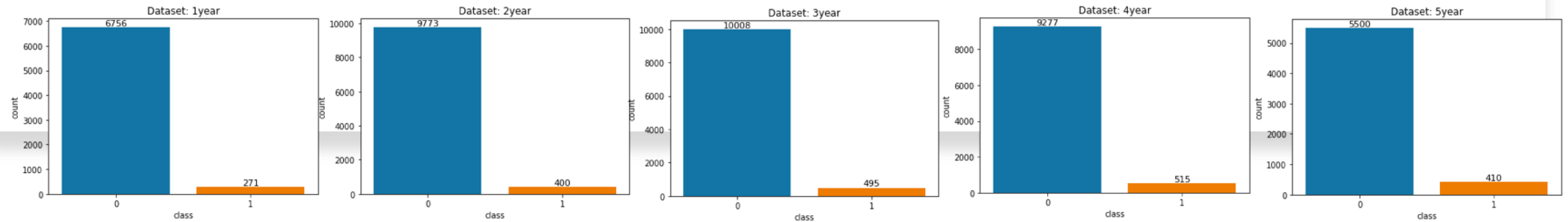    (e.g. Attr 37 has 80.99% of missing value among the total missing value of whole dataset)

# Appendix. 2. Correlation analysis

- How to know correlation from this heat map?
  - Dark area : No correlated cells
  - Red/Light Red : Positive correlated cells
  - Blue/Light Blue : Negative correlated cells
- How's the correlation of our dataset?
  - 39 out of 64 features(60.93% of data) shows high correlation over +/- 0.95
  - Highly correlated data(high bias)



k-NN imputed dataframes, Dataset: 1year

# Appendix. 3. Data balancing analysis

<Before : Highly imbalanced dataset>



| Shape | 6756 | 271 | 9773 | 400 | 10008 | 495 | 9277 | 515 | 5500 | 410 |
|-------|------|-----|------|-----|-------|-----|------|-----|------|-----|

<After : Balanced dataset after SMOTE sampling>



| Shape | 6756 | 6756 | 9773 | 9773 | 10008 | 10008 | 9277 | 9277 | 5500 | 5500 |
|-------|------|------|------|------|-------|-------|------|------|------|------|

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| FALSE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE | TRUE | FALSE | FALSE | TRUE |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
| TRUE | FALSE | FALSE | TRUE | FALSE | FALSE | TRUE | FALSE | FALSE | FALSE | TRUE | TRUE |
| 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
| FALSE | TRUE | TRUE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE | FALSE |
| 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 |
| FALSE | TRUE | FALSE | FALSE | TRUE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE |
| 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| 61 | 62 | 63 | 64 | | | | | | | | |
| FALSE | FALSE | FALSE | FALSE | | | | | | | | |

# Appendix 4-1. Feature selection

- Dataset 1year
- TRUE : Feature selected (17)
- FALSE : Not selected(47)

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|----|----|----|
| FALSE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
| FALSE | TRUE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE | TRUE |
| 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
| FALSE | TRUE | TRUE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE |
| 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 |
| TRUE | TRUE | TRUE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | TRUE | FALSE | FALSE |
| 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 |
| FALSE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE | FALSE |
| 61 | 62 | 63 | 64 | | | | | | | | |
| FALSE | FALSE | FALSE | FALSE | | | | | | | | |

# Appendix 4-2. Feature selection

- Dataset 2year
- TRUE : Feature selected (18)
- FALSE : Not selected (46)

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| FALSE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
| TRUE | FALSE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE | TRUE | TRUE |
| 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
| TRUE | TRUE | TRUE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE |
| 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 |
| FALSE | TRUE | FALSE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE | FALSE |
| 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 |
| FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE | FALSE |
| 61 | 62 | 63 | 64 | | | | | | | | |
| FALSE | FALSE | FALSE | FALSE | | | | | | | | |

# Appendix 4-3. Feature selection

- Dataset 3year
- TRUE : Feature selected (17)
- FALSE : Not selected (47)

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|----|----|----|
| FALSE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
| TRUE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE | TRUE |
| 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
| TRUE | TRUE | TRUE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE |
| 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 |
| FALSE | TRUE | TRUE | FALSE | TRUE | TRUE | FALSE | FALSE | TRUE | TRUE | FALSE | TRUE |
| 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 |
| FALSE | FALSE | TRUE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE |
| 61 | 62 | 63 | 64 | | | | | | | | |
| FALSE | FALSE | FALSE | FALSE | | | | | | | | |

# Appendix 4-4. Feature selection

- Dataset 4year
- TRUE : Feature selected (25)
- FALSE : Not selected (39)

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| TRUE | FALSE | FALSE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | FALSE | TRUE |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
| TRUE | TRUE | TRUE | TRUE | FALSE | TRUE | FALSE | FALSE | TRUE | TRUE | FALSE | TRUE |
| 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
| TRUE | TRUE | TRUE | FALSE | TRUE | FALSE | FALSE | FALSE | FALSE | FALSE | TRUE | FALSE |
| 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 |
| FALSE | TRUE | TRUE | FALSE | TRUE | TRUE | TRUE | FALSE | TRUE | TRUE | FALSE | FALSE |
| 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 |
| FALSE | FALSE | TRUE | FALSE | FALSE | TRUE | TRUE | FALSE | FALSE | FALSE | FALSE | TRUE |
| 61 | 62 | 63 | 64 | | | | | | | | |
| FALSE | FALSE | FALSE | FALSE | | | | | | | | |

# Appendix 4-5. Feature selection

- Dataset 5year
- TRUE : Feature selected (28)
- FALSE : Not selected (36)