

# DMLR DeFi Survival Data Pipeline Example

Hanzhen Qin(qinh2)

12 February 2025

## Survival Data Pipeline

- Project name: DMLR DeFi-LTM

```
source("~/DMLR_DeFi_Survival_Dataset_And_Benchmark/source/survivalData_pipeline/data_preprocessing.R")
source("~/DMLR_DeFi_Survival_Dataset_And_Benchmark/source/survivalData_pipeline/model_evaluation_visualization.R")
source("~/DMLR_DeFi_Survival_Dataset_And_Benchmark/source/survivalData_pipeline/classification_models.R")
source("~/DMLR_DeFi_Survival_Dataset_And_Benchmark/source/survivalData_pipeline/get_classification_cutoff.R")
```

```
# set the indexEvent and outcomeEvent
indexEvent = "borrow"
outcomeEvent = "deposit"

# load the corresponding train and test data
get_train_test_data(indexEvent, outcomeEvent)
```

```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 405 of `x` matches multiple rows in `y`.
## i Row 639 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

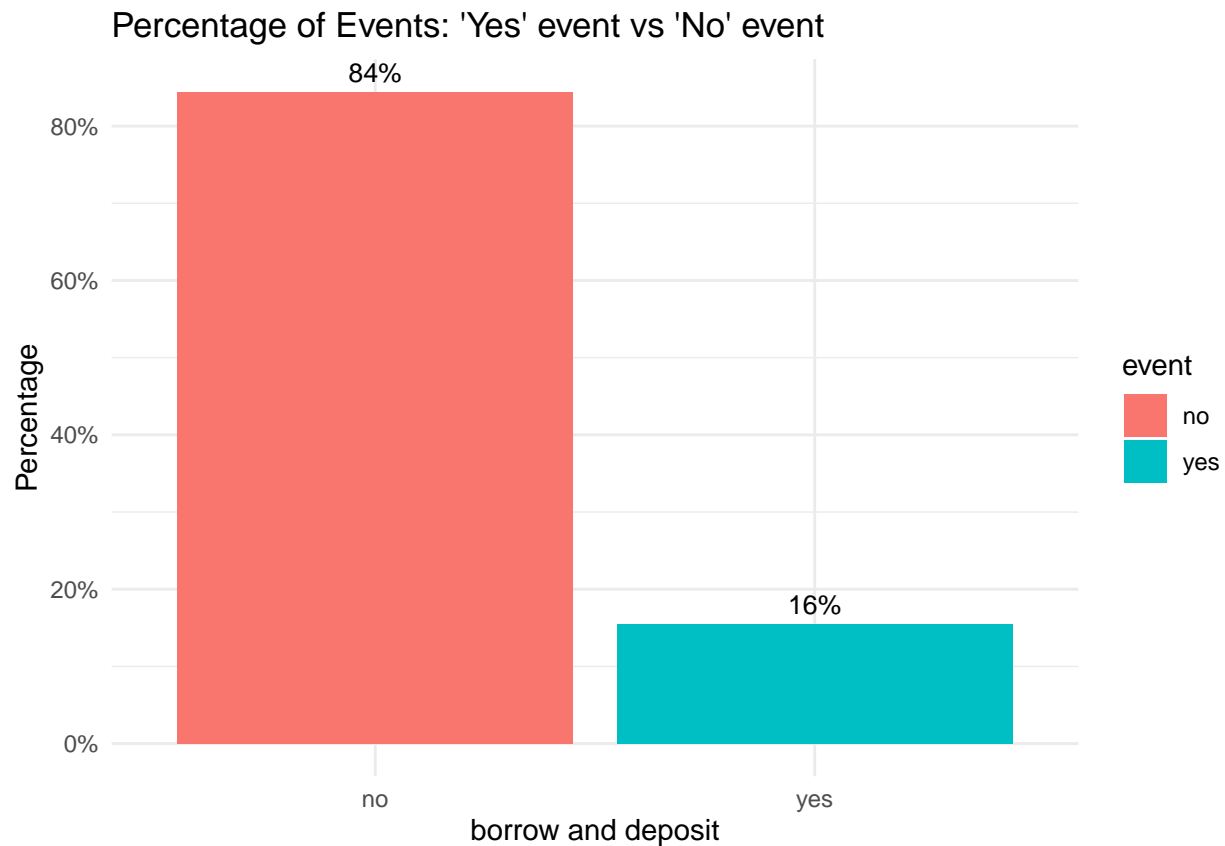
```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 275 of `x` matches multiple rows in `y`.
## i Row 338 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

```
# If you want to check the train and test data, you can run the following codes.
# cat("Train data:\n")
# summary(train)
# cat("Test data:\n")
# summary(test)
```

Using the `get_classification_cutoff` function to get the optimal timeDiff, then we will call the `data_processing` function above to get all the training data and test data.

```
classification_cutoff = get_classification_cutoff(indexEvent, outcomeEvent)
train_data = data_processing(train, classification_cutoff)
test_data = data_processing(test, classification_cutoff)
```

*# If you want to watch the percentages between "Yes" and "No" label, run this code.*  
`get_percentage(train_data, indexEvent, outcomeEvent)`

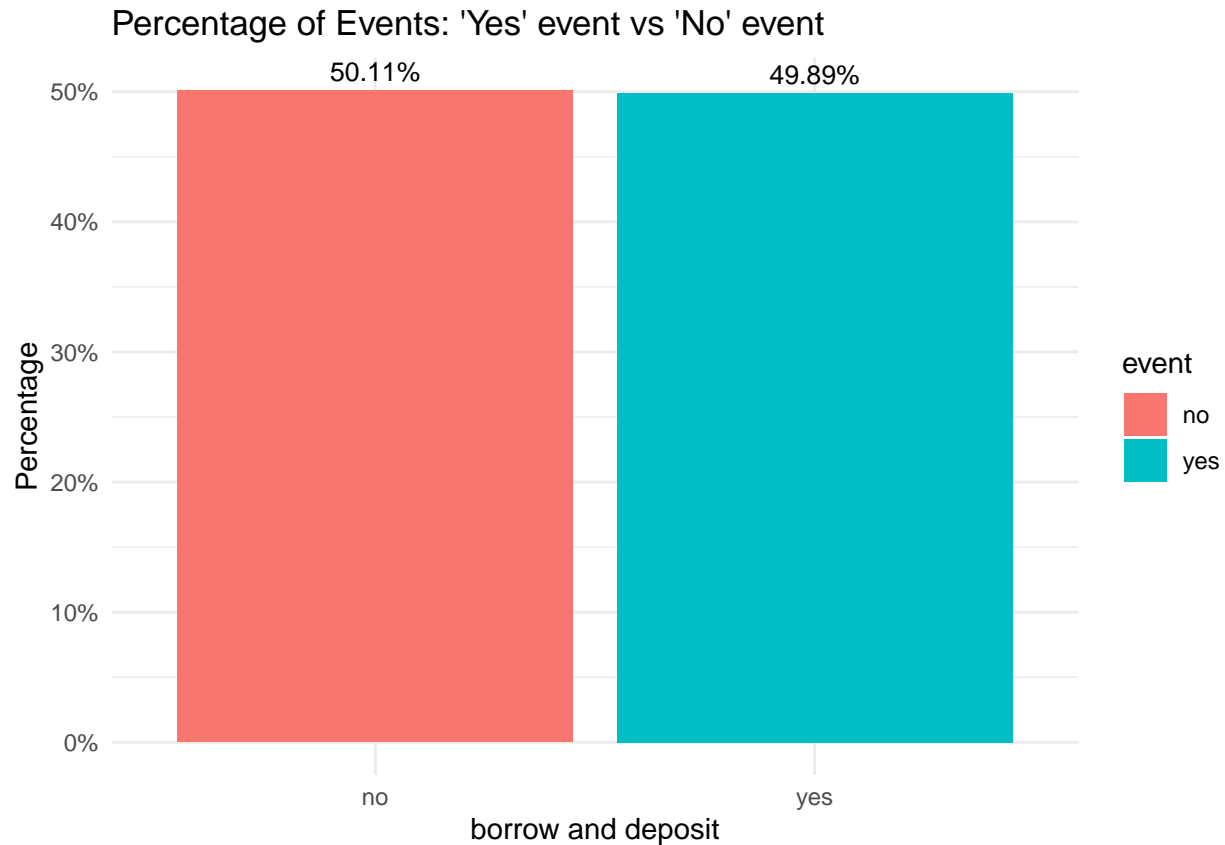


If the ratio of “No” labels to “Yes” labels in the dataset is significantly imbalanced, we can utilize the `smote_data` function to generate a new, more balanced dataset. This balanced dataset ensures that both classes are better represented, helping to mitigate the bias introduced by class imbalance and ultimately improving the accuracy and reliability of our classification model.

```
train_data <- smote_data(train_data)
```

Then you can check the updated balanced version of train data.

*# If you want to watch the percentages between "Yes" and "No" label, run this code.*  
`get_percentage(train_data, indexEvent, outcomeEvent)`



After obtaining the train and test data, we will apply all the classification models to evaluate the relationship between these events.

```
lr_return = logistic_regression(train_data, test_data)
```

```
## [1] "Logistic Regression (Validation) model prediction accuracy:"
## Class accuracy: 67%
## Negative 1 accuracy: 67%
## Balanced accuracy: 67%
## Overall accuracy: 67%
## Precision: 66%
## F1 score: 67%
## [1] "Logistic Regression model prediction accuracy:"
## Class accuracy: 94%
## Negative 1 accuracy: 32%
## Balanced accuracy: 63%
## Overall accuracy: 65%
## Precision: 81%
## F1 score: 46%
```

```
accuracy_lr_dataframe = lr_return$metrics_lr_dataframe
accuracy_lr = lr_return$metrics_lr
```

```
dt_return = decision_tree(train_data, test_data)
```

```
## [1] "Decision Tree (Validation) model prediction accuracy:"
## Class accuracy: 89%
## Negative 1 accuracy: 90%
## Balanced accuracy: 90%
## Overall accuracy: 90%
## Precision: 89%
## F1 score: 89%
## [1] "Decision Tree model prediction accuracy:"
## Class accuracy: 82%
## Negative 1 accuracy: 45%
## Balanced accuracy: 64%
## Overall accuracy: 81%
## Precision: 7%
## F1 score: 12%
```

```
accuracy_dt_dataframe = dt_return$metrics_dt_dataframe
accuracy_dt = dt_return$metrics_dt
```

```
nb_return = naive_bayes(train_data, test_data)
```

```
## [1] "Naive Bayes (Validation) model prediction accuracy:"
## Class accuracy: 80%
## Negative 1 accuracy: 92%
## Balanced accuracy: 86%
## Overall accuracy: 85%
## Precision: 77%
## F1 score: 84%
## [1] "Naive Bayes model prediction accuracy:"
## Class accuracy: 88%
## Negative 1 accuracy: 44%
## Balanced accuracy: 66%
## Overall accuracy: 79%
## Precision: 47%
## F1 score: 45%
```

```
accuracy_nb_dataframe = nb_return$metrics_nb_dataframe
accuracy_nb = nb_return$metrics_nb
```

```
xgb_return = XG_Boost(train_data, test_data)
```

```
## [1] "XGBoost (Validation) model prediction accuracy:"
## Class accuracy: 99%
## Negative 1 accuracy: 98%
## Balanced accuracy: 98%
## Overall accuracy: 98%
## Precision: 99%
## F1 score: 98%
## [1] "XGBoost model prediction accuracy:"
## Class accuracy: 47%
## Negative 1 accuracy: 82%
## Balanced accuracy: 65%
## Overall accuracy: 82%
```

```
## Precision: 100%
## F1 score: 90%
```

```
accuracy_xgb_dataframe = xgb_return$metrics_xgb_dataframe
accuracy_xgb = xgb_return$metrics_xgb
```

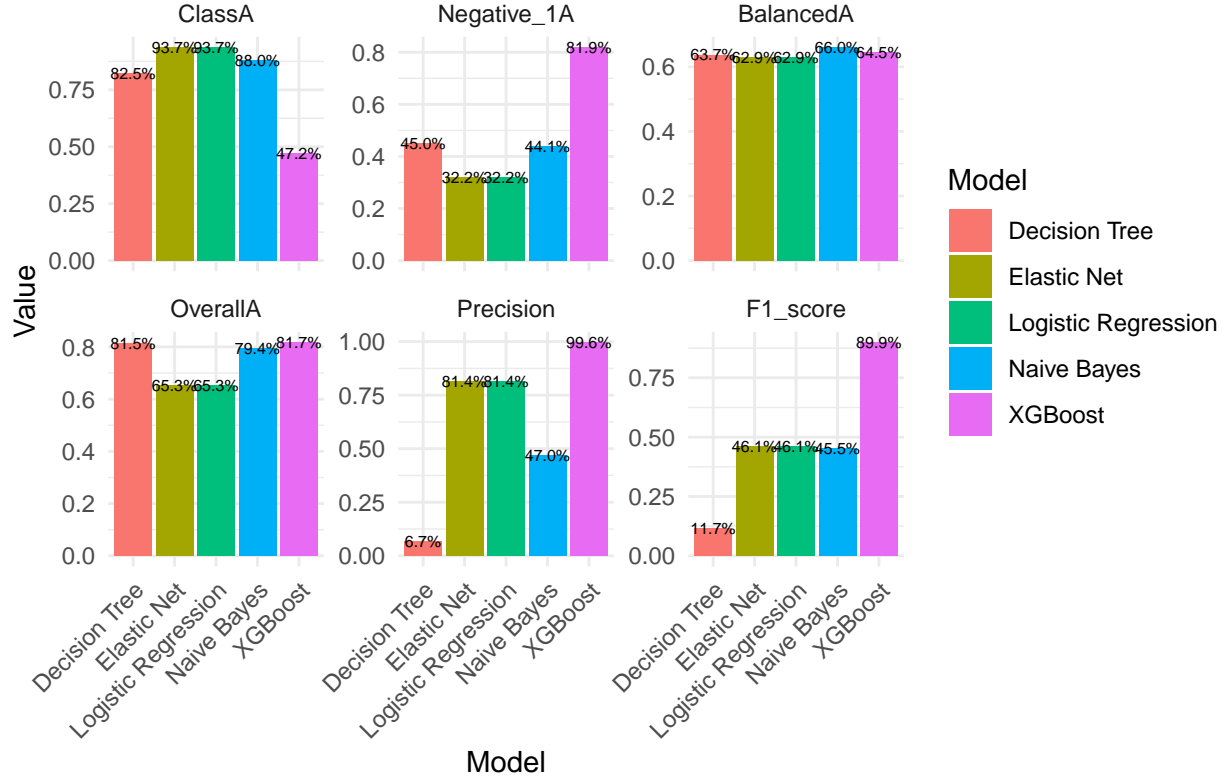
```
en_return = elastic_net(train_data, test_data)
```

```
## [1] "Elastic Net (Validation) model prediction accuracy:"
## Class accuracy: 67%
## Negative 1 accuracy: 67%
## Balanced accuracy: 67%
## Overall accuracy: 67%
## Precision: 66%
## F1 score: 67%
## [1] "Elastic Net model prediction accuracy:"
## Class accuracy: 94%
## Negative 1 accuracy: 32%
## Balanced accuracy: 63%
## Overall accuracy: 65%
## Precision: 81%
## F1 score: 46%
```

```
accuracy_en_dataframe = en_return$metrics_en_dataframe
accuracy_en = en_return$metrics_en
```

```
# compare all the classification models
metrics_list_BD <- list(
  list(accuracy_lr, "Logistic Regression"),
  list(accuracy_dt, "Decision Tree"),
  list(accuracy_nb, "Naive Bayes"),
  list(accuracy_xgb, "XGBoost"),
  list(accuracy_en, "Elastic Net")
)
accuracy_comparison_plot(metrics_list_BD)
```

## Comparison of Accuracy Metrics Across Models



```
# Show the final dataframe for all classification models,
# including the classification model name, accuracy, data combination name.
data_name_BD <- paste(indexEvent, "+", outcomeEvent)
accuracy_dataframe_list_BD <- list(accuracy_lr_dataframe, accuracy_dt_dataframe,
                                   accuracy_nb_dataframe, accuracy_xgb_dataframe,
                                   accuracy_en_dataframe)
combined_results_BD <- combine_classification_results(accuracy_dataframe_list_BD, data_name_BD)

# display the combined dataframe
pander(combined_results_BD, caption = "Classification Model Performance")
```

Table 1: Classification Model Performance (continued below)

Model	Class_Accuracy	Negative_1_Accuracy	Balanced_Accuracy
Logistic Regression	94%	32%	63%
Decision Tree	82%	45%	64%
Naive Bayes	88%	44%	66%
XGBoost	47%	82%	65%
Elastic Net	94%	32%	63%

Overall_Accuracy	Precision	F1_Score	Data_Combination
65%	81%	46%	borrow + deposit

Overall_Accuracy	Precision	F1_Score	Data_Combination
81%	7%	12%	borrow + deposit
79%	47%	45%	borrow + deposit
82%	100%	90%	borrow + deposit
65%	81%	46%	borrow + deposit

```
# set the indexEvent and outcomeEvent
```

```
indexEvent = "borrow"
```

```
outcomeEvent = "repay"
```

```
# load the corresponding train and test data
```

```
get_train_test_data(indexEvent, outcomeEvent)
```

```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 56 of `x` matches multiple rows in `y`.
## i Row 10453 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 1858 of `x` matches multiple rows in `y`.
## i Row 6521 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

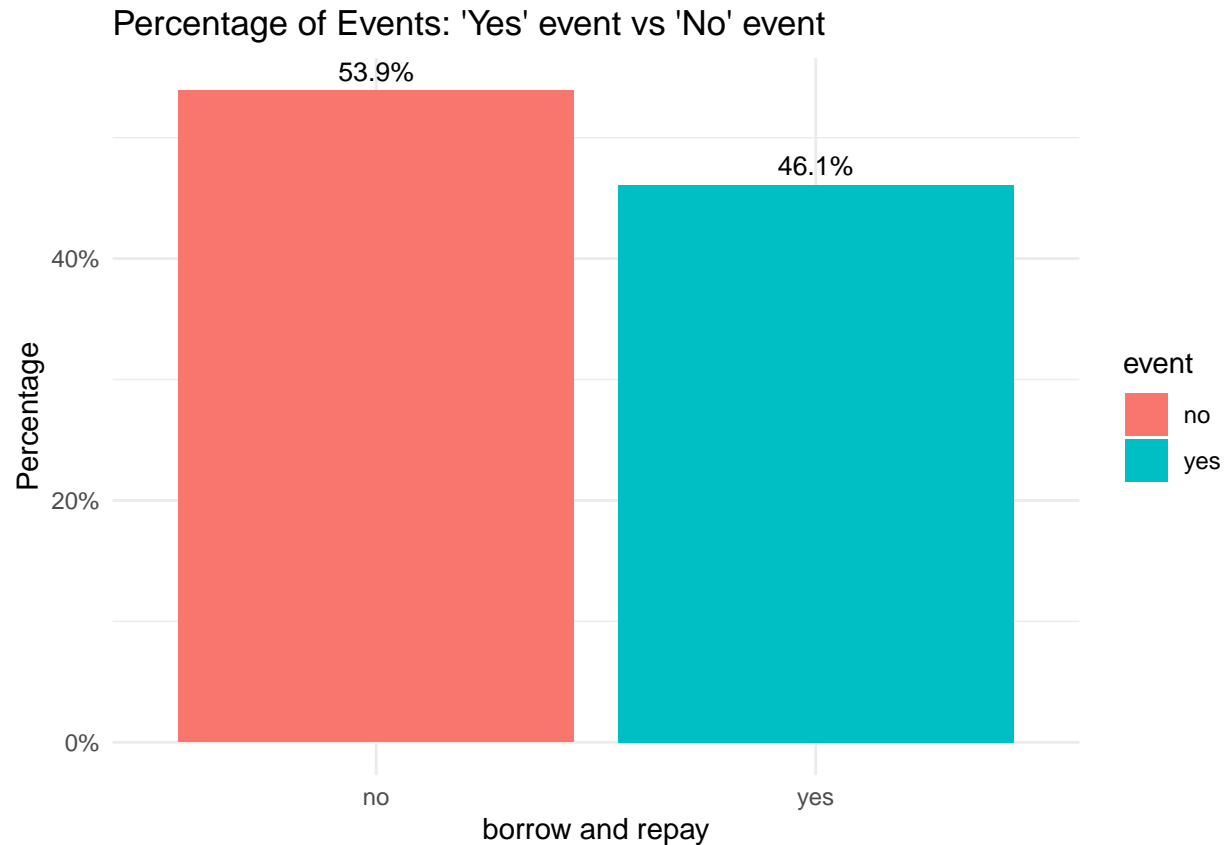
```
classification_cutoff = get_classification_cutoff(indexEvent, outcomeEvent)
```

```
train_data = data_processing(train, classification_cutoff)
```

```
test_data = data_processing(test, classification_cutoff)
```

```
# If you want to watch the percentages between "Yes" and "No" label, run this code.
```

```
get_percentage(train_data, indexEvent, outcomeEvent)
```



```
lr_return = logistic_regression(train_data, test_data)
```

```
## [1] "Logistic Regression (Validation) model prediction accuracy:"  
## Class accuracy: 70%  
## Negative 1 accuracy: 72%  
## Balanced accuracy: 71%  
## Overall accuracy: 71%  
## Precision: 61%  
## F1 score: 66%  
## [1] "Logistic Regression model prediction accuracy:"  
## Class accuracy: 74%  
## Negative 1 accuracy: 64%  
## Balanced accuracy: 69%  
## Overall accuracy: 69%  
## Precision: 68%  
## F1 score: 66%
```

```
accuracy_lr_dataframe = lr_return$metrics_lr_dataframe  
accuracy_lr = lr_return$metrics_lr
```

```
dt_return = decision_tree(train_data, test_data)
```

```
## [1] "Decision Tree (Validation) model prediction accuracy:"  
## Class accuracy: 69%
```



```
## Negative 1 accuracy: 68%
## Balanced accuracy: 69%
## Overall accuracy: 69%
## Precision: 59%
## F1 score: 64%
## [1] "Decision Tree model prediction accuracy:"
## Class accuracy: 76%
## Negative 1 accuracy: 66%
## Balanced accuracy: 71%
## Overall accuracy: 71%
## Precision: 69%
## F1 score: 68%
```

```
accuracy_dt_dataframe = dt_return$metrics_dt_dataframe
accuracy_dt = dt_return$metrics_dt
```

```
nb_return = naive_bayes(train_data, test_data)
```

```
## [1] "Naive Bayes (Validation) model prediction accuracy:"
## Class accuracy: 59%
## Negative 1 accuracy: 81%
## Balanced accuracy: 70%
## Overall accuracy: 61%
## Precision: 21%
## F1 score: 34%
## [1] "Naive Bayes model prediction accuracy:"
## Class accuracy: 65%
## Negative 1 accuracy: 79%
## Balanced accuracy: 72%
## Overall accuracy: 67%
## Precision: 33%
## F1 score: 47%
```

```
accuracy_nb_dataframe = nb_return$metrics_nb_dataframe
accuracy_nb = nb_return$metrics_nb
```

```
xgb_return = XG_Boost(train_data, test_data)
```

```
## [1] "XGBoost (Validation) model prediction accuracy:"
## Class accuracy: 77%
## Negative 1 accuracy: 71%
## Balanced accuracy: 74%
## Overall accuracy: 73%
## Precision: 85%
## F1 score: 77%
## [1] "XGBoost model prediction accuracy:"
## Class accuracy: 64%
## Negative 1 accuracy: 71%
## Balanced accuracy: 67%
## Overall accuracy: 68%
## Precision: 74%
## F1 score: 72%
```

```
accuracy_xgb_dataframe = xgb_return$metrics_xgb_dataframe
accuracy_xgb = xgb_return$metrics_xgb
```

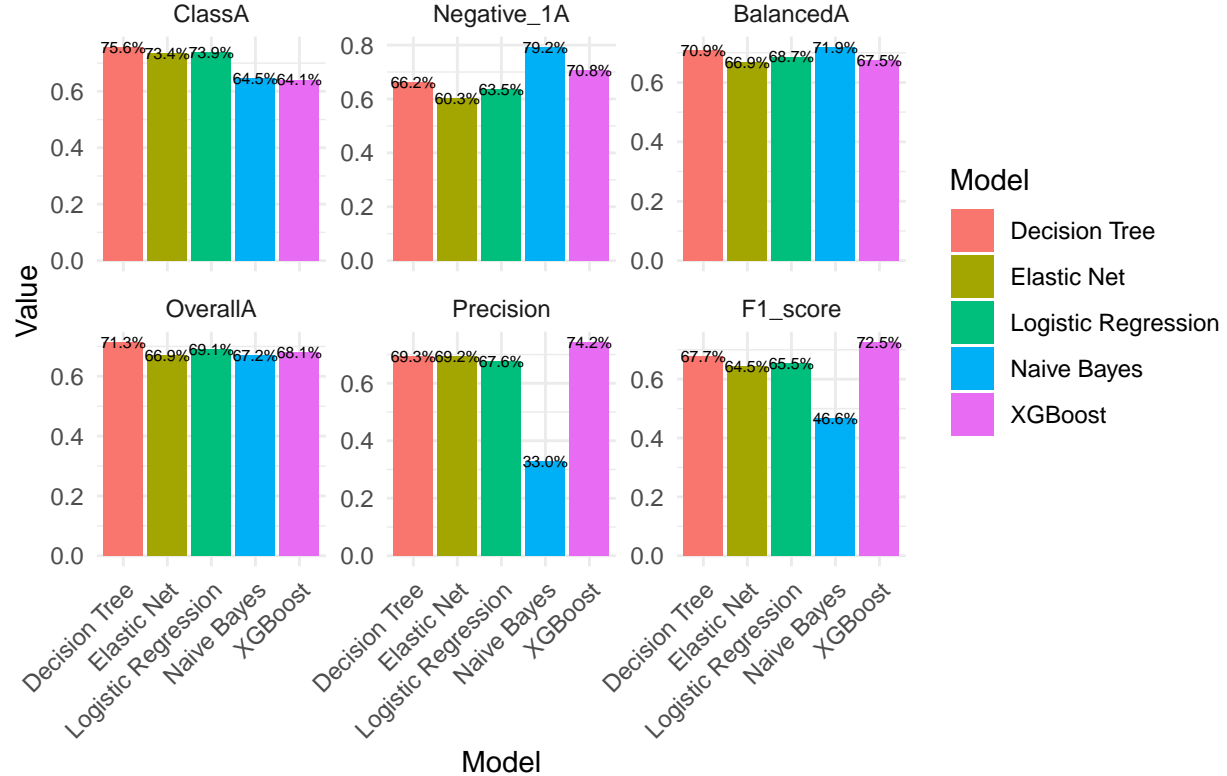
```
en_return = elastic_net(train_data, test_data)
```

```
## [1] "Elastic Net (Validation) model prediction accuracy:"
## Class accuracy: 71%
## Negative 1 accuracy: 72%
## Balanced accuracy: 71%
## Overall accuracy: 71%
## Precision: 61%
## F1 score: 66%
## [1] "Elastic Net model prediction accuracy:"
## Class accuracy: 73%
## Negative 1 accuracy: 60%
## Balanced accuracy: 67%
## Overall accuracy: 67%
## Precision: 69%
## F1 score: 64%
```

```
accuracy_en_dataframe = en_return$metrics_en_dataframe
accuracy_en = en_return$metrics_en
```

```
# compare all the classification models
metrics_list_BR <- list(
  list(accuracy_lr, "Logistic Regression"),
  list(accuracy_dt, "Decision Tree"),
  list(accuracy_nb, "Naive Bayes"),
  list(accuracy_xgb, "XGBoost"),
  list(accuracy_en, "Elastic Net")
)
accuracy_comparison_plot(metrics_list_BR)
```

## Comparison of Accuracy Metrics Across Models



```
# Show the final dataframe for all classification models,
# including the classification model name, accuracy, data combination name.
data_name_BR <- paste(indexEvent, "+", outcomeEvent)
accuracy_dataframe_list_BR <- list(accuracy_lr_dataframe, accuracy_dt_dataframe,
                                   accuracy_nb_dataframe, accuracy_xgb_dataframe,
                                   accuracy_en_dataframe)
combined_results_BR <- combine_classification_results(accuracy_dataframe_list_BR, data_name_BR)

# display the combined dataframe
pander(combined_results_BR, caption = "Classification Model Performance")
```

Table 3: Classification Model Performance (continued below)

Model	Class_Accuracy	Negative_1_Accuracy	Balanced_Accuracy
Logistic Regression	74%	64%	69%
Decision Tree	76%	66%	71%
Naive Bayes	65%	79%	72%
XGBoost	64%	71%	67%
Elastic Net	73%	60%	67%

Overall_Accuracy	Precision	F1_Score	Data_Combination
69%	68%	66%	borrow + repay

Overall_Accuracy	Precision	F1_Score	Data_Combination
71%	69%	68%	borrow + repay
67%	33%	47%	borrow + repay
68%	74%	72%	borrow + repay
67%	69%	64%	borrow + repay

```
# set the indexEvent and outcomeEvent
```

```
indexEvent = "borrow"
```

```
outcomeEvent = "withdraw"
```

```
# load the corresponding train and test data
```

```
get_train_test_data(indexEvent, outcomeEvent)
```

```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 465 of `x` matches multiple rows in `y`.
## i Row 639 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 285 of `x` matches multiple rows in `y`.
## i Row 338 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

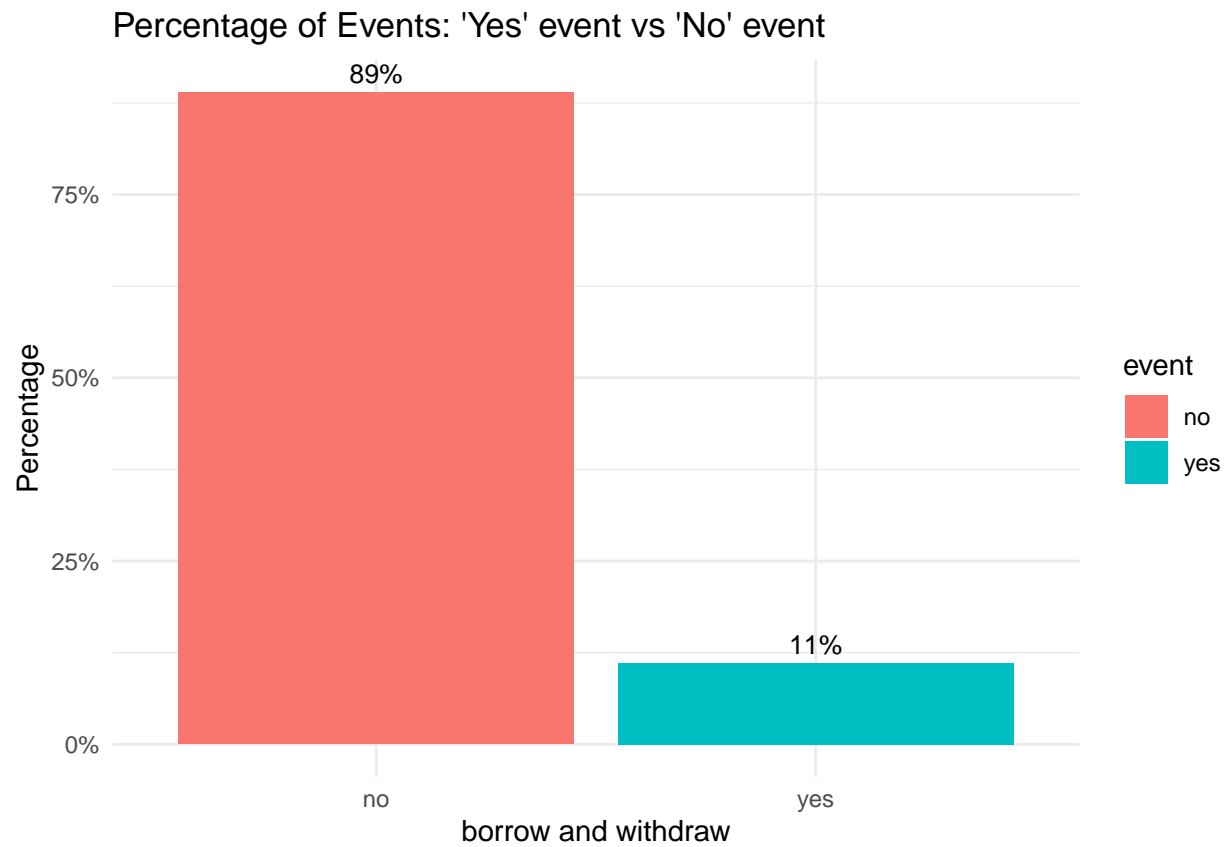
```
classification_cutoff = get_classification_cutoff(indexEvent, outcomeEvent)
```

```
train_data = data_processing(train, classification_cutoff)
```

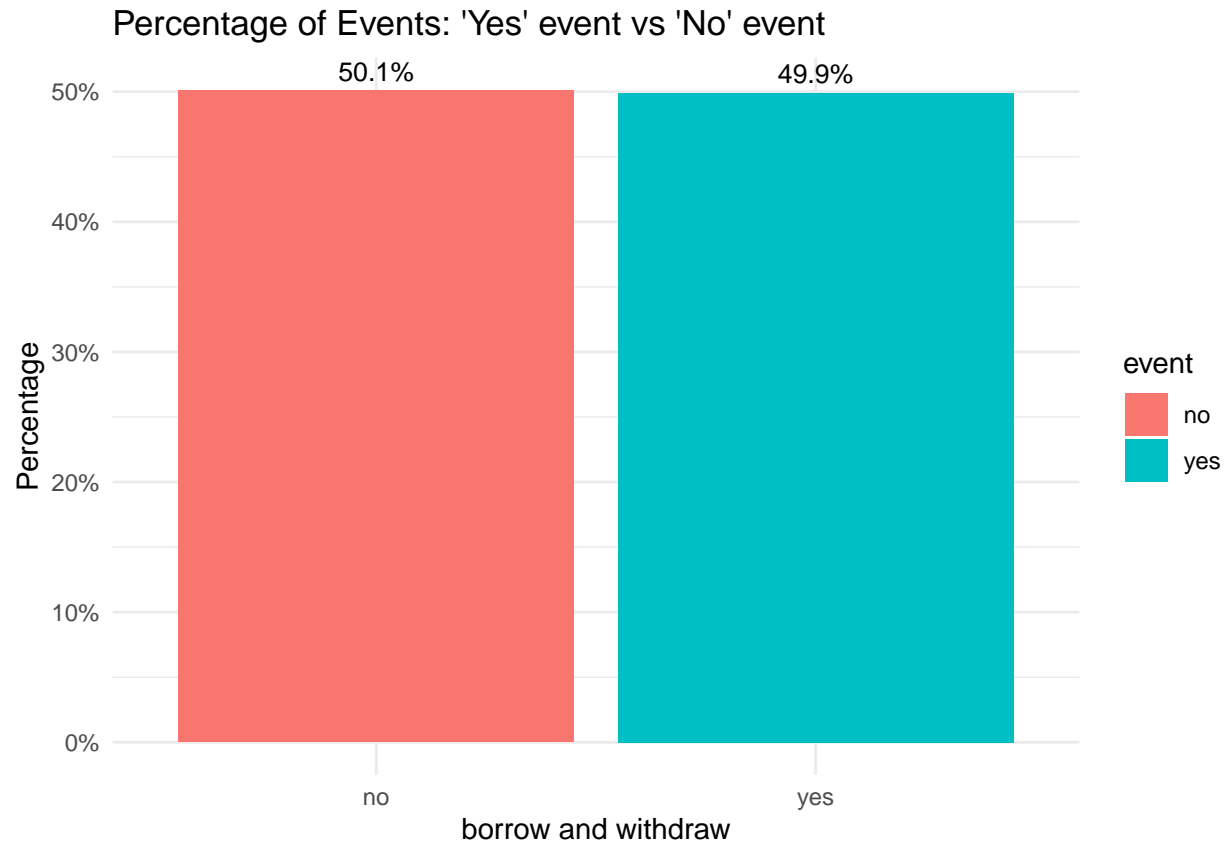
```
test_data = data_processing(test, classification_cutoff)
```

```
# If you want to watch the percentages between "Yes" and "No" label, run this code.
```

```
get_percentage(train_data, indexEvent, outcomeEvent)
```



```
train_data <- smote_data(train_data)
get_percentage(train_data, indexEvent, outcomeEvent)
```



```
lr_return = logistic_regression(train_data, test_data)
```

```
## [1] "Logistic Regression (Validation) model prediction accuracy:"
## Class accuracy: 69%
## Negative 1 accuracy: 70%
## Balanced accuracy: 70%
## Overall accuracy: 70%
## Precision: 69%
## F1 score: 69%
## [1] "Logistic Regression model prediction accuracy:"
## Class accuracy: 95%
## Negative 1 accuracy: 27%
## Balanced accuracy: 61%
## Overall accuracy: 68%
## Precision: 80%
## F1 score: 41%
```

```
accuracy_lr_dataframe = lr_return$metrics_lr_dataframe
accuracy_lr = lr_return$metrics_lr
```

```
dt_return = decision_tree(train_data, test_data)
```

```
## [1] "Decision Tree (Validation) model prediction accuracy:"
## Class accuracy: 91%
```

```
## Negative 1 accuracy: 92%
## Balanced accuracy: 92%
## Overall accuracy: 92%
## Precision: 91%
## F1 score: 92%
## [1] "Decision Tree model prediction accuracy:"
## Class accuracy: 87%
## Negative 1 accuracy: 40%
## Balanced accuracy: 64%
## Overall accuracy: 86%
## Precision: 6%
## F1 score: 11%
```

```
accuracy_dt_dataframe = dt_return$metrics_dt_dataframe
accuracy_dt = dt_return$metrics_dt
```

```
nb_return = naive_bayes(train_data, test_data)
```

```
## [1] "Naive Bayes (Validation) model prediction accuracy:"
## Class accuracy: 83%
## Negative 1 accuracy: 94%
## Balanced accuracy: 88%
## Overall accuracy: 87%
## Precision: 80%
## F1 score: 86%
## [1] "Naive Bayes model prediction accuracy:"
## Class accuracy: 92%
## Negative 1 accuracy: 40%
## Balanced accuracy: 66%
## Overall accuracy: 83%
## Precision: 51%
## F1 score: 45%
```

```
accuracy_nb_dataframe = nb_return$metrics_nb_dataframe
accuracy_nb = nb_return$metrics_nb
```

```
xgb_return = XG_Boost(train_data, test_data)
```

```
## [1] "XGBoost (Validation) model prediction accuracy:"
## Class accuracy: 99%
## Negative 1 accuracy: 99%
## Balanced accuracy: 99%
## Overall accuracy: 99%
## Precision: 99%
## F1 score: 99%
## [1] "XGBoost model prediction accuracy:"
## Class accuracy: 23%
## Negative 1 accuracy: 86%
## Balanced accuracy: 55%
## Overall accuracy: 86%
## Precision: 100%
## F1 score: 92%
```

```
accuracy_xgb_dataframe = xgb_return$metrics_xgb_dataframe
accuracy_xgb = xgb_return$metrics_xgb
```

```
en_return = elastic_net(train_data, test_data)
```

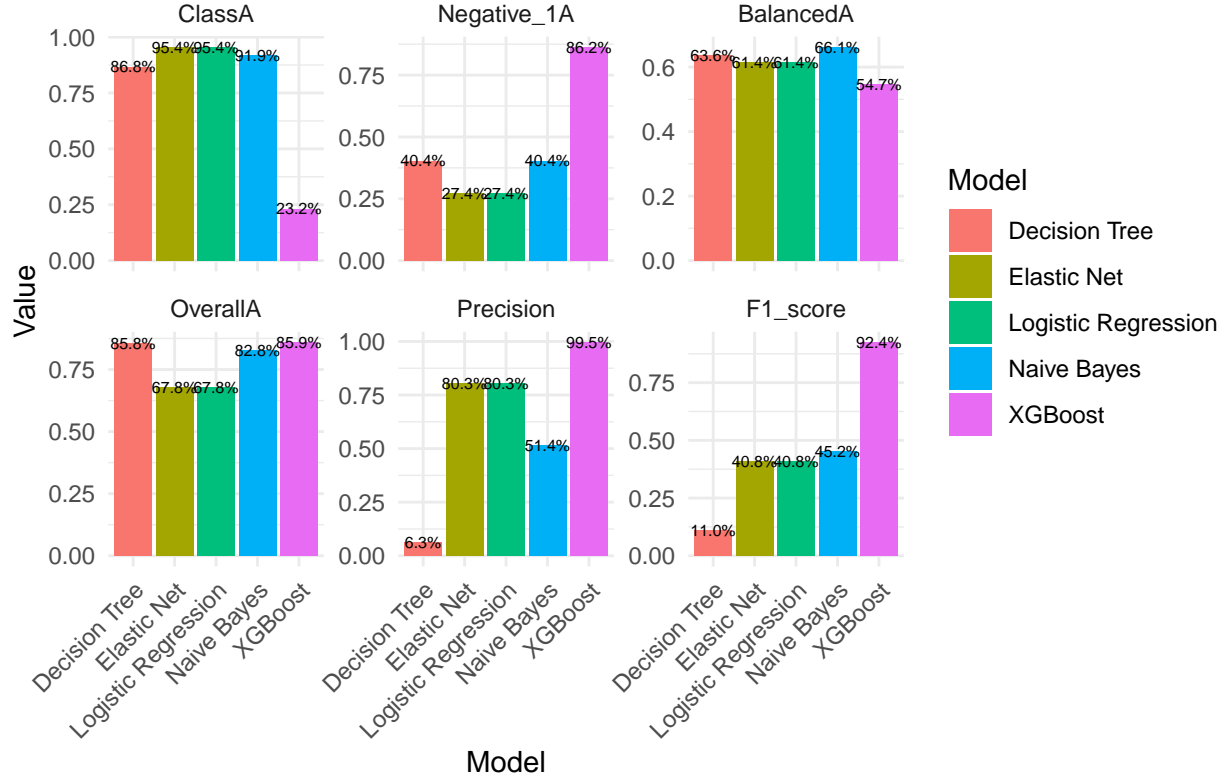
```
## [1] "Elastic Net (Validation) model prediction accuracy:"
## Class accuracy: 69%
## Negative 1 accuracy: 70%
## Balanced accuracy: 70%
## Overall accuracy: 70%
## Precision: 69%
## F1 score: 69%
## [1] "Elastic Net model prediction accuracy:"
## Class accuracy: 95%
## Negative 1 accuracy: 27%
## Balanced accuracy: 61%
## Overall accuracy: 68%
## Precision: 80%
## F1 score: 41%
```

```
accuracy_en_dataframe = en_return$metrics_en_dataframe
accuracy_en = en_return$metrics_en
```

```
# compare all the classification models
metrics_list_BW <- list(
  list(accuracy_lr, "Logistic Regression"),
  list(accuracy_dt, "Decision Tree"),
  list(accuracy_nb, "Naive Bayes"),
  list(accuracy_xgb, "XGBoost"),
  list(accuracy_en, "Elastic Net")
)
accuracy_comparison_plot(metrics_list_BW)
```



## Comparison of Accuracy Metrics Across Models



```
# Show the final dataframe for all classification models,
# including the classification model name, accuracy, data combination name.
data_name_BW <- paste(indexEvent, "+", outcomeEvent)
accuracy_dataframe_list_BW <- list(accuracy_lr_dataframe, accuracy_dt_dataframe,
                                   accuracy_nb_dataframe, accuracy_xgb_dataframe,
                                   accuracy_en_dataframe)
combined_results_BW <- combine_classification_results(accuracy_dataframe_list_BW, data_name_BW)

# display the combined dataframe
pander(combined_results_BW, caption = "Classification Model Performance")
```

Table 5: Classification Model Performance (continued below)

Model	Class_Accuracy	Negative_1_Accuracy	Balanced_Accuracy
Logistic Regression	95%	27%	61%
Decision Tree	87%	40%	64%
Naive Bayes	92%	40%	66%
XGBoost	23%	86%	55%
Elastic Net	95%	27%	61%

Overall_Accuracy	Precision	F1_Score	Data_Combination
68%	80%	41%	borrow + withdraw

Overall_Accuracy	Precision	F1_Score	Data_Combination
86%	6%	11%	borrow + withdraw
83%	51%	45%	borrow + withdraw
86%	100%	92%	borrow + withdraw
68%	80%	41%	borrow + withdraw

```
# set the indexEvent and outcomeEvent
```

```
indexEvent = "borrow"
```

```
outcomeEvent = "account liquidated"
```

```
# load the corresponding train and test data
```

```
get_train_test_data(indexEvent, outcomeEvent)
```

```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 570 of `x` matches multiple rows in `y`.
## i Row 637 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

```
## Warning in inner_join(y, X, by = "id"): Detected an unexpected many-to-many relationship between `x`
## i Row 82 of `x` matches multiple rows in `y`.
## i Row 82 of `y` matches multiple rows in `x`.
## i If a many-to-many relationship is expected, set `relationship =
##   "many-to-many"` to silence this warning.
```

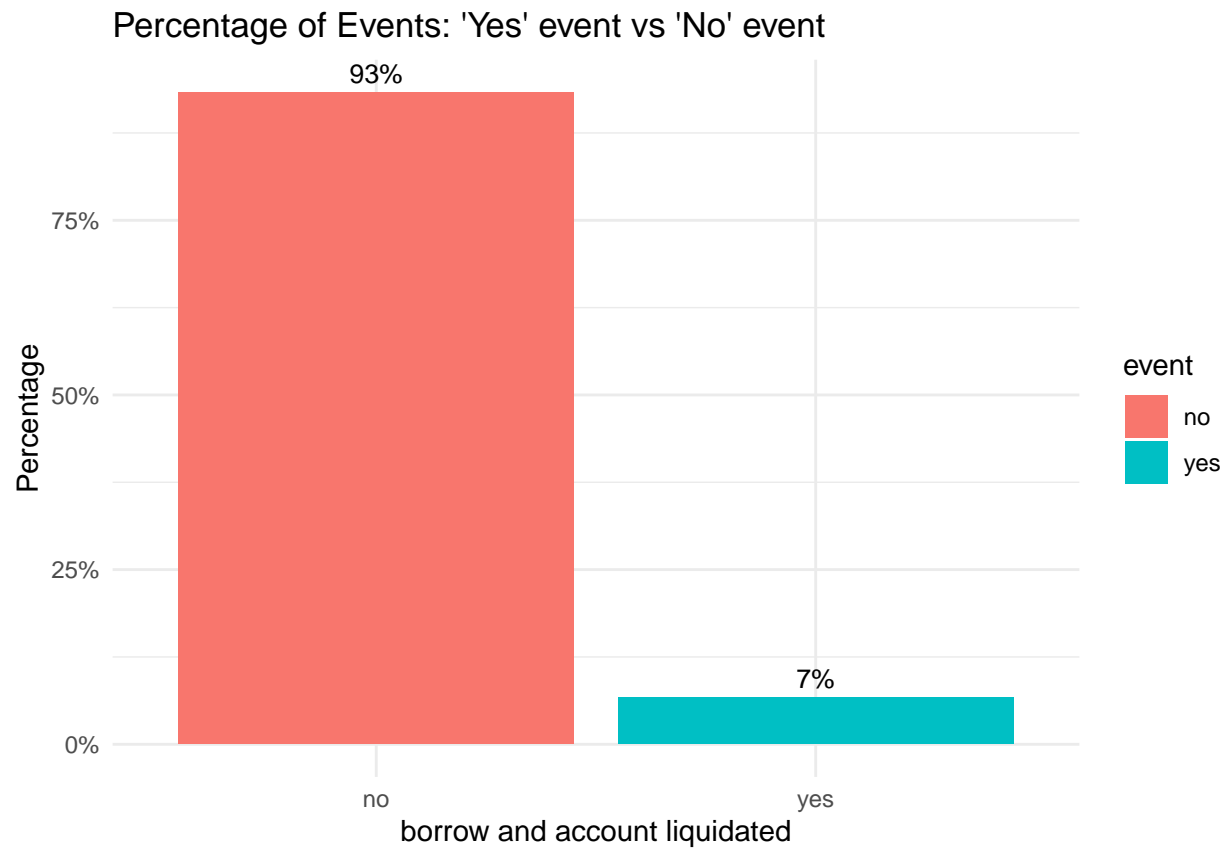
```
classification_cutoff = get_classification_cutoff(indexEvent, outcomeEvent)
```

```
train_data = data_processing(train, classification_cutoff)
```

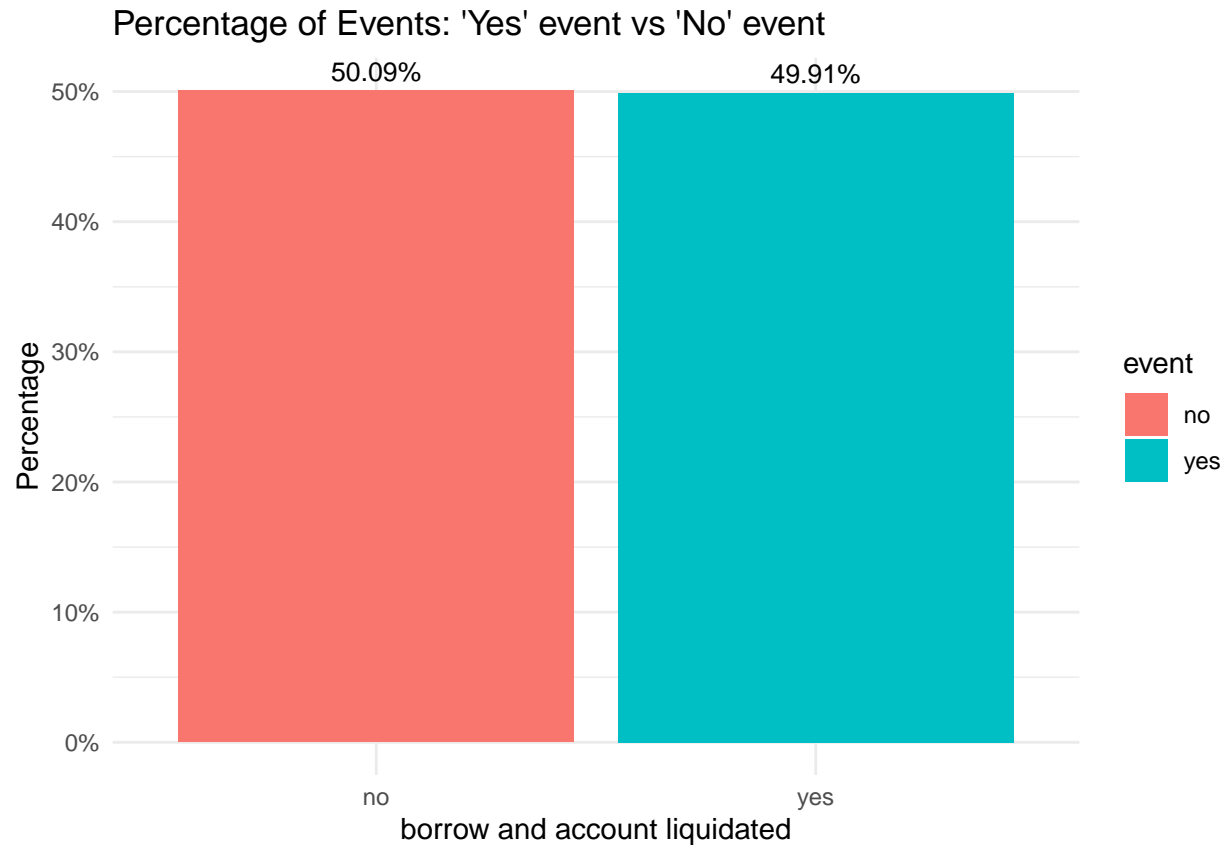
```
test_data = data_processing(test, classification_cutoff)
```

```
# If you want to watch the percentages between "Yes" and "No" label, run this code.
```

```
get_percentage(train_data, indexEvent, outcomeEvent)
```



```
train_data <- smote_data(train_data)
get_percentage(train_data, indexEvent, outcomeEvent)
```



```
lr_return = logistic_regression(train_data, test_data)
```

```
## [1] "Logistic Regression (Validation) model prediction accuracy:"  
## Class accuracy: 71%  
## Negative 1 accuracy: 73%  
## Balanced accuracy: 72%  
## Overall accuracy: 72%  
## Precision: 69%  
## F1 score: 71%  
## [1] "Logistic Regression model prediction accuracy:"  
## Class accuracy: 99%  
## Negative 1 accuracy: 2%  
## Balanced accuracy: 50%  
## Overall accuracy: 60%  
## Precision: 57%  
## F1 score: 4%
```

```
accuracy_lr_dataframe = lr_return$metrics_lr_dataframe  
accuracy_lr = lr_return$metrics_lr
```

```
dt_return = decision_tree(train_data, test_data)
```

```
## [1] "Decision Tree (Validation) model prediction accuracy:"  
## Class accuracy: 100%
```

```
## Negative 1 accuracy: 98%
## Balanced accuracy: 99%
## Overall accuracy: 99%
## Precision: 100%
## F1 score: 99%
## [1] "Decision Tree model prediction accuracy:"
## Class accuracy: 100%
## Negative 1 accuracy: 2%
## Balanced accuracy: 51%
## Overall accuracy: 19%
## Precision: 99%
## F1 score: 3%
```

```
accuracy_dt_dataframe = dt_return$metrics_dt_dataframe
accuracy_dt = dt_return$metrics_dt
```

```
nb_return = naive_bayes(train_data, test_data)
```

```
## [1] "Naive Bayes (Validation) model prediction accuracy:"
## Class accuracy: 98%
## Negative 1 accuracy: 100%
## Balanced accuracy: 99%
## Overall accuracy: 99%
## Precision: 98%
## F1 score: 99%
## [1] "Naive Bayes model prediction accuracy:"
## Class accuracy: 100%
## Negative 1 accuracy: 2%
## Balanced accuracy: 51%
## Overall accuracy: 21%
## Precision: 97%
## F1 score: 3%
```

```
accuracy_nb_dataframe = nb_return$metrics_nb_dataframe
accuracy_nb = nb_return$metrics_nb
```

```
xgb_return = XG_Boost(train_data, test_data)
```

```
## [1] "XGBoost (Validation) model prediction accuracy:"
## Class accuracy: 100%
## Negative 1 accuracy: 100%
## Balanced accuracy: 100%
## Overall accuracy: 100%
## Precision: 100%
## F1 score: 100%
## [1] "XGBoost model prediction accuracy:"
## Class accuracy: 0%
## Negative 1 accuracy: 99%
## Balanced accuracy: 49%
## Overall accuracy: 98%
## Precision: 99%
## F1 score: 99%
```

```
accuracy_xgb_dataframe = xgb_return$metrics_xgb_dataframe
accuracy_xgb = xgb_return$metrics_xgb
```

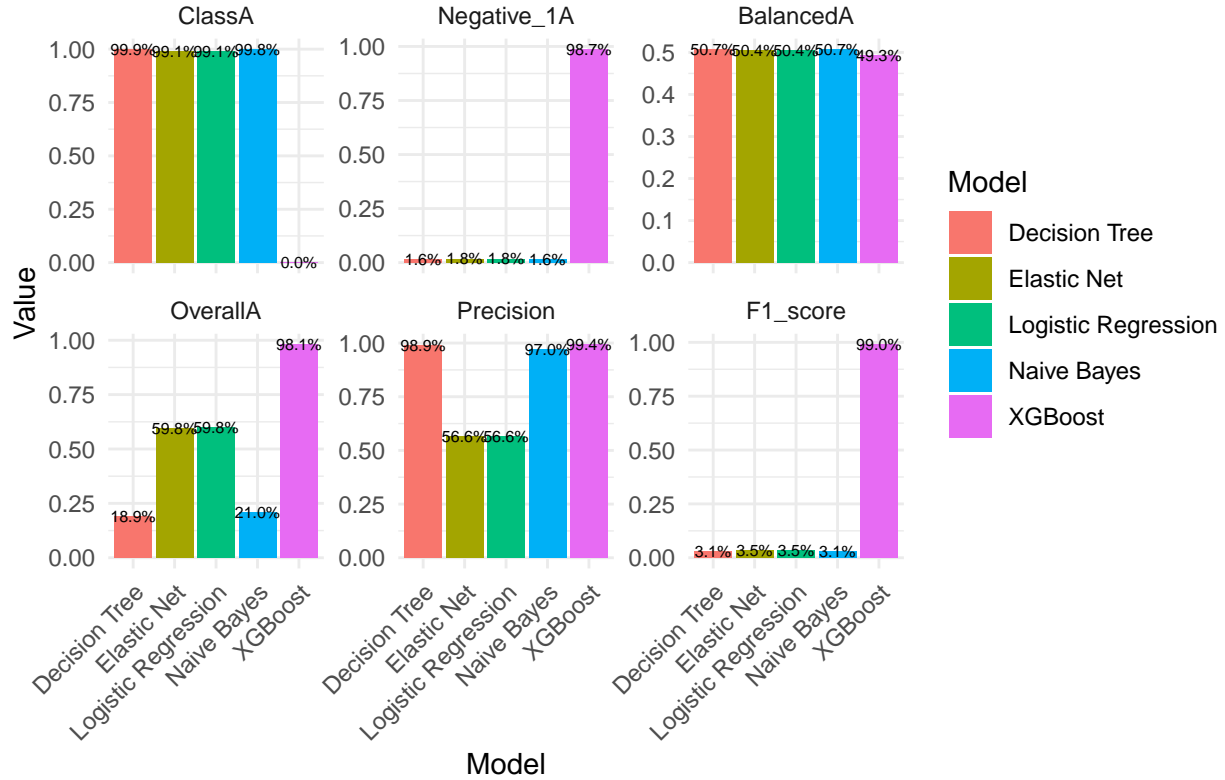
```
en_return = elastic_net(train_data, test_data)
```

```
## [1] "Elastic Net (Validation) model prediction accuracy:"
## Class accuracy: 71%
## Negative 1 accuracy: 73%
## Balanced accuracy: 72%
## Overall accuracy: 72%
## Precision: 69%
## F1 score: 71%
## [1] "Elastic Net model prediction accuracy:"
## Class accuracy: 99%
## Negative 1 accuracy: 2%
## Balanced accuracy: 50%
## Overall accuracy: 60%
## Precision: 57%
## F1 score: 4%
```

```
accuracy_en_dataframe = en_return$metrics_en_dataframe
accuracy_en = en_return$metrics_en
```

```
# compare all the classification models
metrics_list_BAL <- list(
  list(accuracy_lr, "Logistic Regression"),
  list(accuracy_dt, "Decision Tree"),
  list(accuracy_nb, "Naive Bayes"),
  list(accuracy_xgb, "XGBoost"),
  list(accuracy_en, "Elastic Net")
)
accuracy_comparison_plot(metrics_list_BAL)
```

## Comparison of Accuracy Metrics Across Models



```
# Show the final dataframe for all classification models,
# including the classification model name, accuracy, data combination name.
data_name_BAL <- paste(indexEvent, "+", outcomeEvent)
accuracy_dataframe_list_BAL <- list(accuracy_lr_dataframe, accuracy_dt_dataframe,
                                   accuracy_nb_dataframe, accuracy_xgb_dataframe,
                                   accuracy_en_dataframe)
combined_results_BAL <- combine_classification_results(accuracy_dataframe_list_BAL, data_name_BAL)

# display the combined dataframe
pander(combined_results_BAL, caption = "Classification Model Performance")
```

Table 7: Classification Model Performance (continued below)

Model	Class_Accuracy	Negative_1_Accuracy	Balanced_Accuracy
Logistic Regression	99%	2%	50%
Decision Tree	100%	2%	51%
Naive Bayes	100%	2%	51%
XGBoost	0%	99%	49%
Elastic Net	99%	2%	50%

Overall_Accuracy	Precision	F1_Score	Data_Combination
60%	57%	4%	borrow + account liquidated

Overall_Accuracy	Precision	F1_Score	Data_Combination
19%	99%	3%	borrow + account liquidated
21%	97%	3%	borrow + account liquidated
98%	99%	99%	borrow + account liquidated
60%	57%	4%	borrow + account liquidated

## Classification Model Performance For All Data Combinations

After we run all the data combinations, we can use the `combine_accuracy_dataframes` to combine all the classification models' performance into one dataframe.

```
combined_classification_results <- combine_accuracy_dataframes(
  list(combined_results_BAL, combined_results_BD, combined_results_BR, combined_results_BW))
pander(combined_classification_results, caption = "Classification Model Performance for all data")
```

Table 9: Classification Model Performance for all data (continued below)

Model	Class_Accuracy	Negative_1_Accuracy	Balanced_Accuracy
Logistic Regression	99%	2%	50%
Decision Tree	100%	2%	51%
Naive Bayes	100%	2%	51%
XGBoost	0%	99%	49%
Elastic Net	99%	2%	50%
Logistic Regression	94%	32%	63%
Decision Tree	82%	45%	64%
Naive Bayes	88%	44%	66%
XGBoost	47%	82%	65%
Elastic Net	94%	32%	63%
Logistic Regression	74%	64%	69%
Decision Tree	76%	66%	71%
Naive Bayes	65%	79%	72%
XGBoost	64%	71%	67%
Elastic Net	73%	60%	67%
Logistic Regression	95%	27%	61%
Decision Tree	87%	40%	64%
Naive Bayes	92%	40%	66%
XGBoost	23%	86%	55%
Elastic Net	95%	27%	61%

Overall_Accuracy	Precision	F1_Score	Data_Combination
60%	57%	4%	borrow + account liquidated
19%	99%	3%	borrow + account liquidated
21%	97%	3%	borrow + account liquidated
98%	99%	99%	borrow + account liquidated
60%	57%	4%	borrow + account liquidated
65%	81%	46%	borrow + deposit
81%	7%	12%	borrow + deposit
79%	47%	45%	borrow + deposit



Overall_Accuracy	Precision	F1_Score	Data_Combination
82%	100%	90%	borrow + deposit
65%	81%	46%	borrow + deposit
69%	68%	66%	borrow + repay
71%	69%	68%	borrow + repay
67%	33%	47%	borrow + repay
68%	74%	72%	borrow + repay
67%	69%	64%	borrow + repay
68%	80%	41%	borrow + withdraw
86%	6%	11%	borrow + withdraw
83%	51%	45%	borrow + withdraw
86%	100%	92%	borrow + withdraw
68%	80%	41%	borrow + withdraw

## Generating Dataframe For Specified Accuracy

This section is only for a special need, not required for the whole pipeline workflow!!!

In this section, the final output is a combined data frame that consolidates performance metrics for multiple classification models across different data scenarios. Each row represents a specific scenario (e.g., “borrow + withdraw” or “borrow + repay”), while the columns display the selected performance metric (e.g., “balanced\_accuracy”) and the corresponding values for each classification model (e.g., Logistic Regression, Decision Tree).

```
ba_accuracy_dataframe_BAL <- specific_accuracy_statistics(data_name_BAL, "balanced_accuracy",
                                                         metrics_list_BAL)
ba_accuracy_dataframe_BD <- specific_accuracy_statistics(data_name_BD, "balanced_accuracy",
                                                         metrics_list_BD)
ba_accuracy_dataframe_BR <- specific_accuracy_statistics(data_name_BR, "balanced_accuracy",
                                                         metrics_list_BR)
ba_accuracy_dataframe_BW <- specific_accuracy_statistics(data_name_BW, "balanced_accuracy",
                                                         metrics_list_BW)
combined_accuracy_dataframe <- combine_accuracy_dataframes(
  list(ba_accuracy_dataframe_BAL, ba_accuracy_dataframe_BD, ba_accuracy_dataframe_BR,
       ba_accuracy_dataframe_BW))
pander(combined_accuracy_dataframe, caption = "Combined accuracy dataframe")
```

Table 11: Combined accuracy dataframe (continued below)

balanced_accuracy	Logistic.Regression	Decision.Tree
borrow + account liquidated	50.4	50.7
borrow + deposit	62.9	63.7
borrow + repay	68.7	70.9
borrow + withdraw	61.4	63.6

Naive.Bayes	XGBoost	Elastic.Net
50.7	49.3	50.4
66	64.5	62.9
71.9	67.5	66.9

Naive.Bayes	XGBoost	Elastic.Net
66.1	54.7	61.4