

# LOAN DEFAULT

**Outlier Group** 

Northeastern University - Toronto

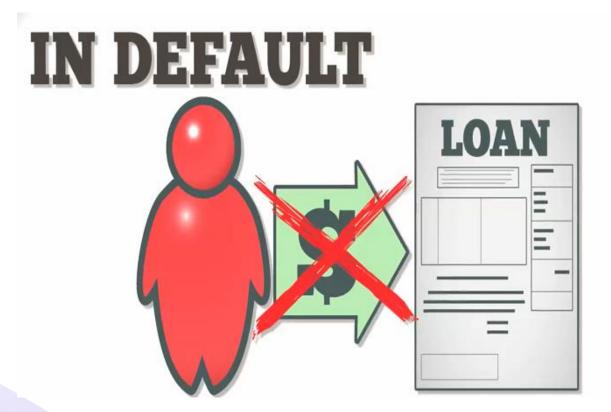
ALY6040: Data Mining





## **INTRODUCTION**







### **Terminology**

Default refers to the ability to repay.

Default: unable to pay back.

Non-default: able to repay.



### Two scenarios of loss

Approve for default customers Reject non-default customers



#### **Research Questions**

- 1. Given current customer information and historical data, would customer be able to repay the loan? Classification problem
- 2. How to deal with imbalanced dataset



## **INTRODUCTION**





Information of current customers with Target variable (TARGET).

**TARGET O: Non-default** 

**TARGET 1: Default** 

✓ 148,670 records and 34 variables



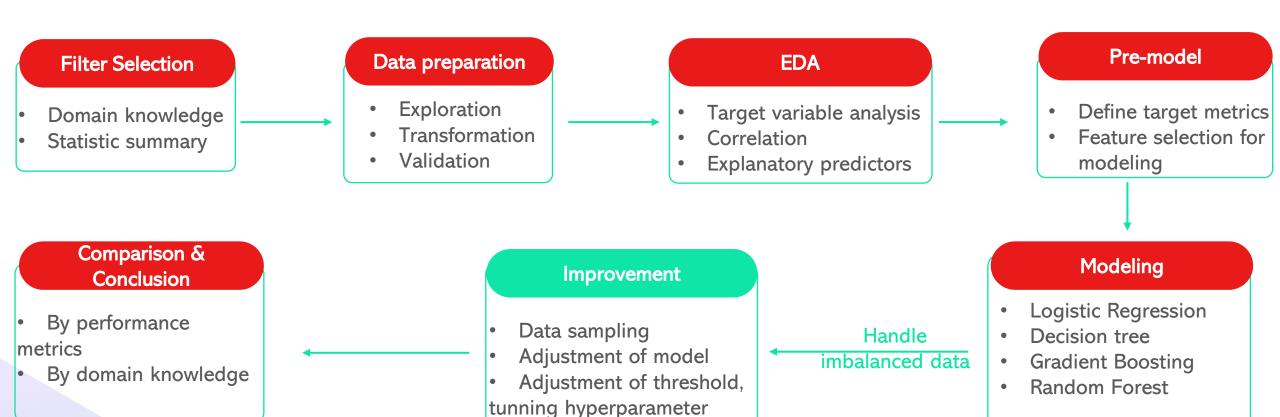
### **Project objective**

- ✓ Master the data munging
- Deploy various methods of data mining for Classification problem
- / Handle imbalanced class of target variable
- Familiarize ourselves with the target industry



## **INTRODUCTION**









#### FINAL DATA SUMMARY

**Final Data set** 

137,336 records x 24 variables

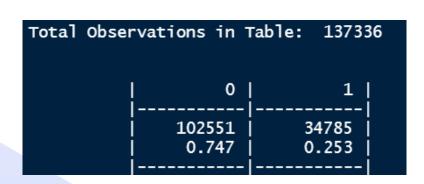
```
str(df)
'data.frame':
               137336 obs. of 24 variables:
$ loan_limit
                           : Factor w/ 3 levels "cf", "Missing", ...: 1 1 1 1
$ Gender
                           : Factor w/ 4 levels "Female", "Joint", ...: 4 3 3
                           : Factor w/ 3 levels "Missing", "nopre",...: 2 2 3
$ approv_in_adv
$ loan_type
                           : Factor w/ 3 levels "type1", "type2", ...: 1 2 1 1
                           : Factor w/ 5 levels "Missing","p1",..: 2 2 2 5
$ loan_purpose
                           : Factor w/ 2 levels "l1","l2": 1 1 1
$ Credit_Worthiness
                           : Factor w/ 2 levels "nopc", "opc": 1 1 1 1 1 1 1
$ open_credit
                           : Factor w/ 2 levels "b/c", "nob/c": 2 1 2 2 2 2
$ business_or_commercial
$ loan_amount
                           : int 116500 206500 406500 456500 696500 706500
$ rate_of_interest
                                 3.99 3.99 4.56 4.25 4
                           : num
$ term
                                  360 360 360 360 360 360 360 360 360
                           : Factor w/ 2 levels "int_only", "not_int": 2 2 2
$ interest_only
                           : Factor w/ 2 levels "lpsm", "not_lpsm": 2 1 2 2
$ lump_sum_payment
$ property_value
                           : num 118000 418000 508000 658000 758000 ...
                           : Factor w/ 3 levels "ir", "pr", "sr": 2 2 2 2 2 2
$ occupancy_type
$ total_units
                                 11111111111...
                           : num
$ income
                                 1740 4980 9480 11880 10440 ...
                           : Factor w/ 4 levels "CIB", "CRIF", "EQUI", ...: 4 3
$ credit_type
$ Credit_Score
                           : int 758 552 834 587 602 864 860 863 580 788
$ co.applicant_credit_type : Factor w/ 2 levels "CIB", "EXP": 1 2 1 1 2 2 2
                           : Factor w/ 7 levels "<25",">74","25-34",...:
$ age
$ submission_of_application: Factor w/ 2 levels "not_inst", "to_inst": 2 2 2
                           : Factor w/ 4 levels "central", "North",...: 4 2 4
$ Region
$ Status
                           : int 1100000000...
```

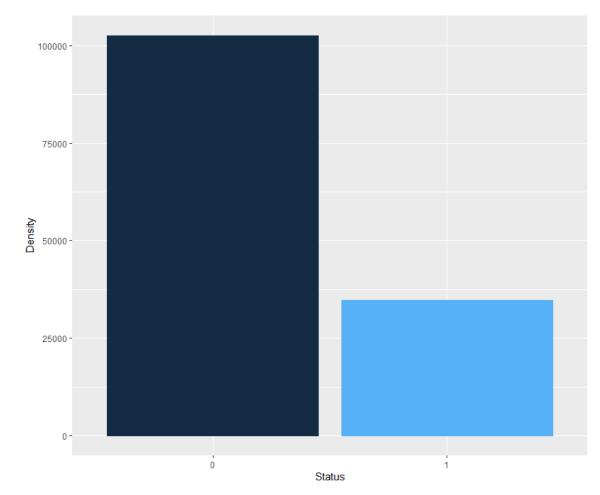




**TARGET VARIABLE** 

Slightly imbalanced
75% non-default 25 default





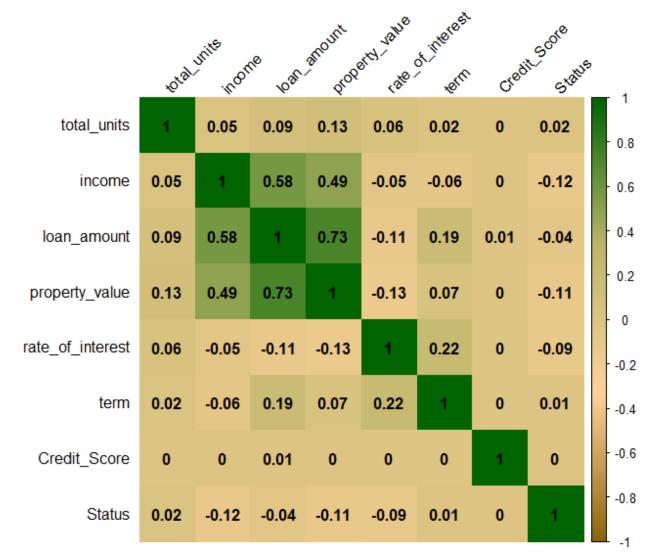


**EDA** 

**CORRELOGRAM** 

To Target variable (Status):

Mostly insignificant
correlations



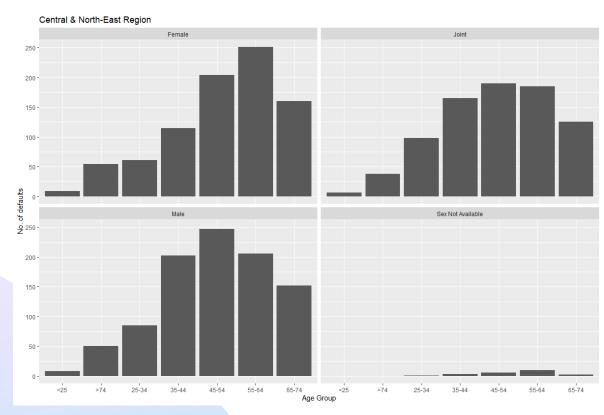


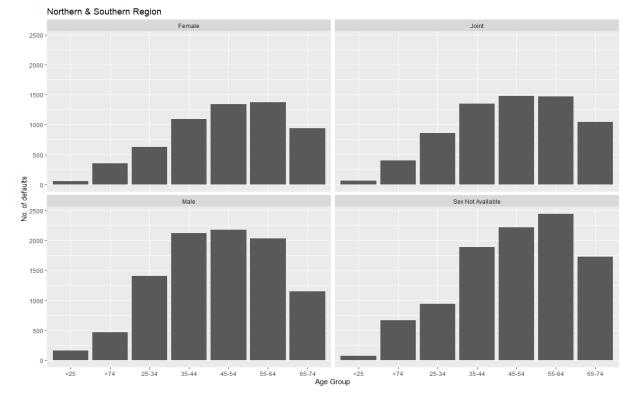


**EXPLANATORY PREDICTORS** 

> table1

central 2271 North North-East 15991 362 south 15926

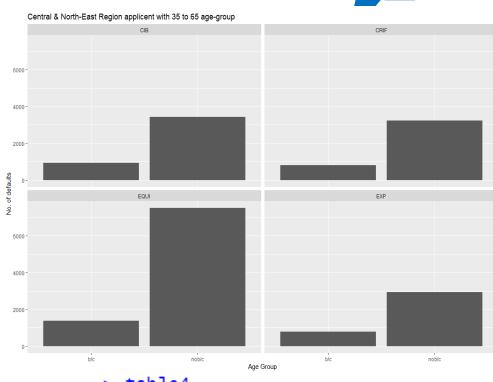


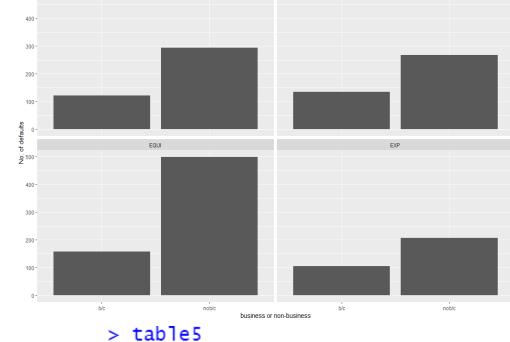






#### **EXPLANATORY PREDICTORS**





> table4

CIB CRIF EQUI EXP b/c 930 815 1373 790 nob/c 3432 3226 7514 2927

Central & North-East Region applicent with 35 to 65 age-group

CIB CRIF EQUI EXP b/c 121 134 158 105 nob/c 294 267 499 206





# **BACKGROUND RATIONALE**



Class bias problem for ML



So how to tackle this?



#### Risk for finance

Predict default customer (TARGET = 1) is more crucial than for non-default customer (TARGET = 0).

# False Negative (FN) is a

dangerous error

#### Data process

- Upsample
- Downsample
- · SMOTE
- · ROSE

#### Modeling

- Tune hyperparamet ers
- · Update model
- Pruning

#### Threshold

- · Cut-off
- · Cp















Stepwise Logistics

#### **Training set**

#### Full vs Stepwise

```
> anova (log.full,log.stepwise,test="LRT")
Analysis of Deviance Table
Model 1: Status ~ loan_limit + Gender + approv_in_adv + loan_type + loan_purpose +
   Credit_Worthiness + open_credit + business_or_commercial +
   loan_amount + rate_of_interest + term + interest_only + lump_sum_payment +
   property_value + occupancy_type + total_units + income +
   credit_type + Credit_Score + co.applicant_credit_type + age +
   submission_of_application + Region
Model 2: Status ~ loan_limit + Gender + approv_in_adv + loan_type + loan_purpose +
   Credit_Worthiness + open_credit + loan_amount + rate_of_interest +
   term + interest_only + lump_sum_payment + property_value +
   occupancy_type + total_units + income + credit_type + Credit_Score +
   co.applicant_credit_type + age + submission_of_application +
   Region
 Resid. Df Resid. Dev Df Deviance Pr(>Chi)
    109827
                 77867
    109827
                 77867 0
```

```
confusionMatrix(t(conf_matrix),positive = "1")
Confusion Matrix and Statistics
    Actual
Pred
   0 20368 3589
      166 3345
              Accuracy: 0.8633
                95% CI: (0.8592, 0.8673)
    No Information Rate: 0.7476
    P-Value [Acc > NIR] : < 2.2e-16
                 Kappa: 0.567
 Mcnemar's Test P-Value : < 2.2e-16
           Sensitivity : 0.4824
           Specificity: 0.9919
         Pos Pred Value: 0.9527
        Neg Pred Value: 0.8502
            Prevalence: 0.2524
        Detection Rate: 0.1218
   Detection Prevalence: 0.1278
      Balanced Accuracy: 0.7372
       'Positive' Class: 1
```







#### Create weight

#### Train with weight

```
summary(log.weighted)
Call:
glm(formula = Status ~ loan_limit + Gender + approv_in_adv +
    loan_type + loan_purpose + Credit_Worthiness + open_credit +
   loan_amount + rate_of_interest + term + interest_only + lump_sum_payment +
   property_value + occupancy_type + total_units + income +
   credit_type + Credit_Score + co.applicant_credit_type + age +
   submission_of_application + Region, family = binomial(link = logit),
   data = training_set, weights = model_weights)
Deviance Residuals:
             10 Median
-3.5227 -0.7788 -0.5804
                         0.0027 4.0350
Coefficients:
                                  Estimate Std. Error z value Pr(>|z|)
                                5.446e+00 3.718e-01 14.649 < 2e-16 ***
(Intercept)
loan_limitMissing
                                2.112e-01 4.978e-02 4.243 2.21e-05 ***
                                7.115e-01 3.177e-02 22.398 < 2e-16 ***
loan_limitncf
GenderJoint
                                4.403e-02 3.170e-02 1.389 0.164934
GenderMale
                                1.584e-01 2.265e-02
                                                      6.994 2.68e-12
GenderSex Not Available
                                                      0.505 0.613569
                                 1.464e-02 2.898e-02
                                -1.113e-01 9.299e-02 -1.197 0.231334
approv_in_advnopre
approv_in_advpre
                                -4.623e-01 9.490e-02 -4.871 1.11e-06 ***
loan_typetype2
                                6.200e-01 2.263e-02 27.401 < 2e-16 ***
                                -1.217e-01 3.252e-02 -3.741 0.000183 ***
loan_typetype3
loan_purposep1
                                4.428e-01 3.218e-01 1.376 0.168804
                                1.162e+00 3.252e-01 3.572 0.000354 ***
loan_purposep2
loan_purposep3
                                5.112e-01 3.218e-01
                                 5.372e-01 3.218e-01
loan_purposep4
                                                      1.669 0.095072
Credit Worthiness12
                                 5.272e-01 3.692e-02 14.278 < 2e-16 ***
open_creditopc
                                8.452e-01 1.785e-01 4.736 2.18e-06 ***
loan_amount
rate_of_interest
                                -9.390e-01 2.103e-02 -44.660 < 2e-16 ***
                                1.218e-03 1.428e-04 8.531 < 2e-16 ***
interest_onlynot_int
                                -1.148e-01 3.683e-02 -3.116 0.001835 **
lump_sum_paymentnot_lpsm
                                -2.418e+00 5.971e-02 -40.498 < 2e-16 ***
```

#### Test performance

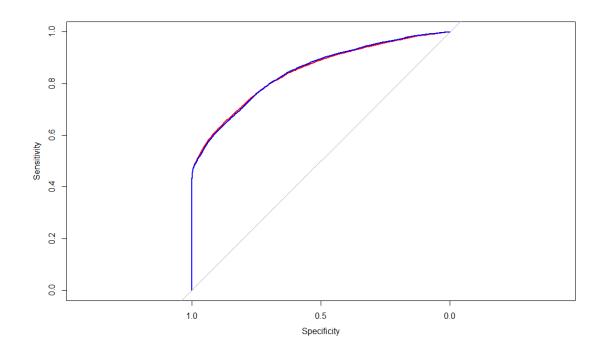
```
confusionMatrix(t(conf_matrix2),positive = "1")
Confusion Matrix and Statistics
   Actual
Pred
           2489
  0 17997
  1 2537
           4445
              Accuracy: 0.817
                95% CI : (0.8124, 0.8216)
   No Information Rate: 0.7476
   P-Value [Acc > NIR] : <2e-16
                 Kappa: 0.5163
 Mcnemar's Test P-Value: 0.5074
           Sensitivity: 0.6410 + 18\%
           Specificity: 0.8764
        Pos Pred Value: 0.6366
        Neg Pred Value: 0.8785
            Prevalence: 0.2524
        Detection Rate: 0.1618
  Detection Prevalence: 0.2542
     Balanced Accuracy: 0.7587
       'Positive' Class: 1
```







Model	Accuracy	Sensitivity	FN	F1	AUC
Log.stepwise	86.33%	48.24%	3589	0.6404	0.8453
Log.weighted	81.70%	64.10%	2489	0.6388	0.8461

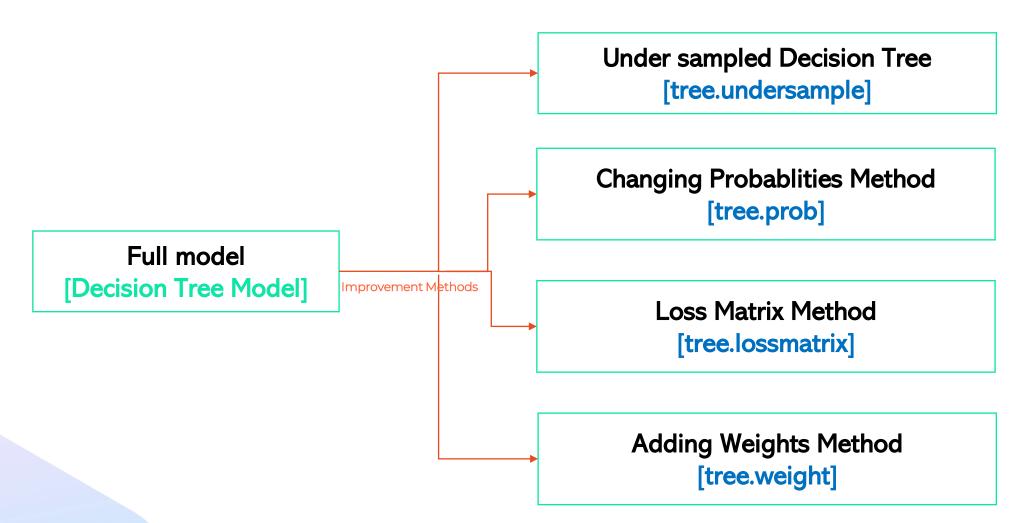






**DECISION TREE** 









## **DECISION TREE**

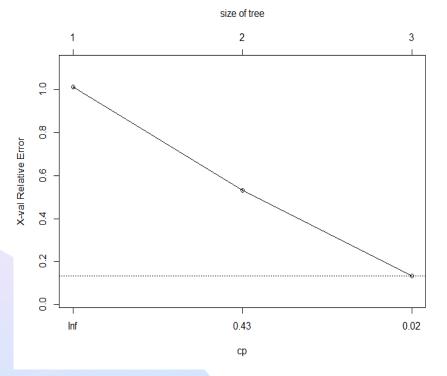


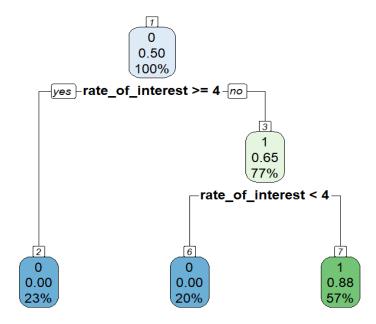
**Under sampled Model** 

#### Training set

# Training Sec







```
> confusionMatrix(t(confmat_undersample),positive =
Confusion Matrix and Statistics
    Actual
   0 17754 0
  1 2780 6934
              Accuracy: 0.8988
                95% CI: (0.8952, 0.9023)
    No Information Rate: 0.7476
    P-Value [Acc > NIR] : < 2.2e-16
                 Kappa : 0.7633
 Mcnemar's Test P-Value : < 2.2e-16
           Sensitivity: 1.0000
           Specificity: 0.8646
        Pos Pred Value: 0.7138
        Neg Pred Value : 1.0000
            Prevalence: 0.2524
        Detection Rate: 0.2524
   Detection Prevalence: 0.3536
      Balanced Accuracy: 0.9323
       'Positive' Class: 1
```





### **DECISION TREE**

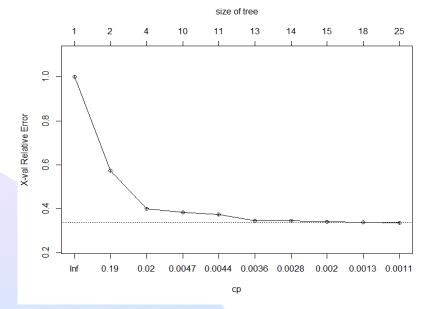


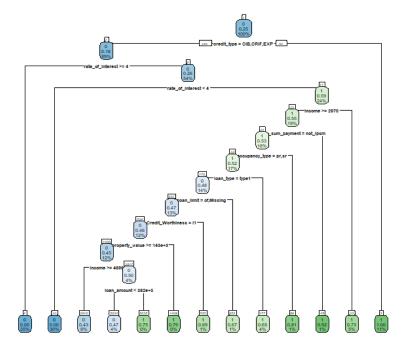
**Changing Probabilities** 

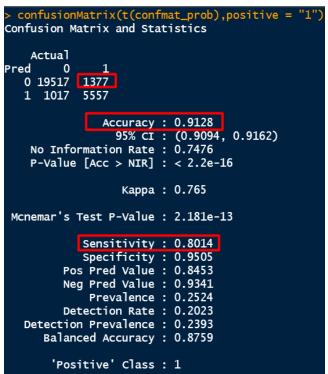
#### Training set

#### 

#### Probabilities: Non-default to 75% and Defaults to 25%











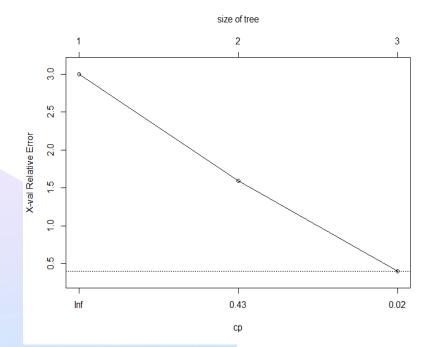
### **DECISION TREE**

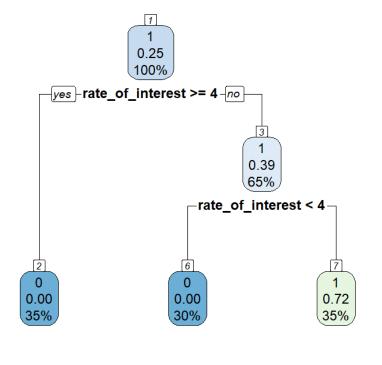


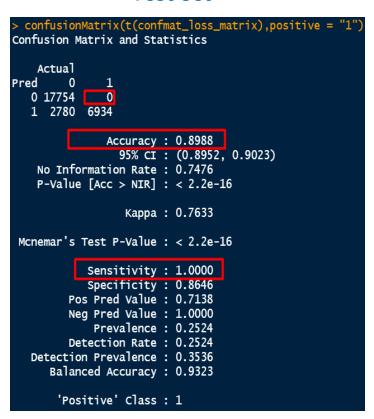
#### Training set

#### 

#### Loss Matrix penalizes FN more strongly than FP









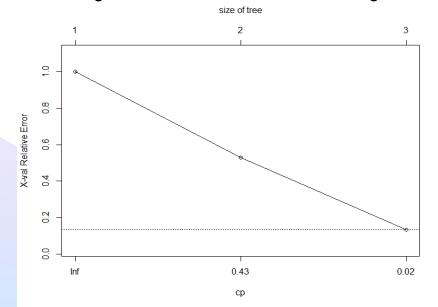


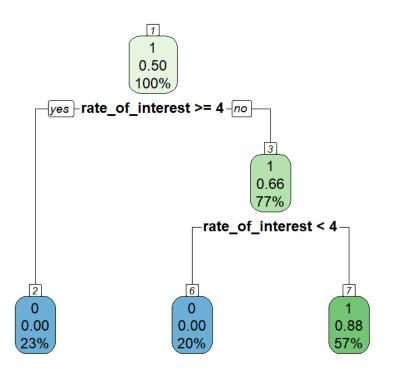
### **DECISION TREE**

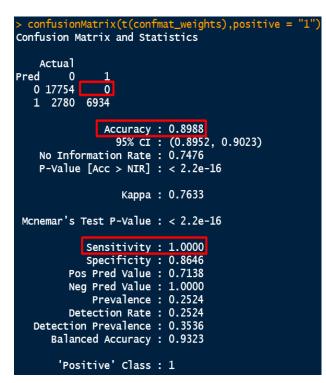


#### Training set

This vector contains weights of 1 for the non-default and weights of 3 for defaults in the training set.









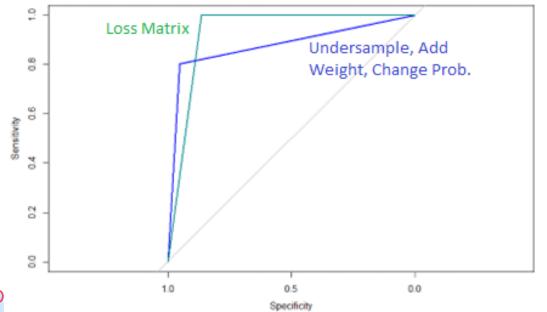


## **DECISION TREE**



## Model recap

Model	Accuracy	Sensitivity	FN	F1	AUC
Under sample	89.88%	100.00%	0	0.8330	0.9323
Change Prob	91.28%	80.14%	1377	0.8223	0.8759
Loss Matrix	89.88%	100.00%	0	0.8330	0.9323
Add Weights	89.88%	100.00%	0	0.8330	0.9323





1st, 3rd and 4th are the same, while the 2nd model is somewhat different. Thus, in the final comparison, author may use only 2 models from the table to evaluate each method performance.



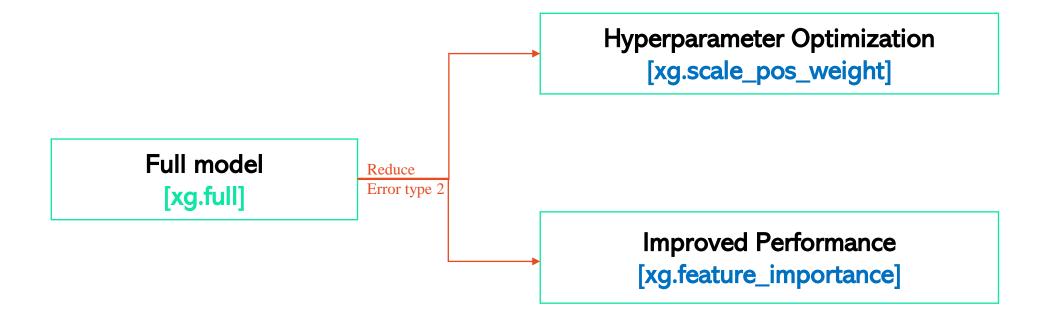
When taking account to overall performance (AUC) and Sensitivity, Loss Matrix Method outperforms with AUC of 0.9323





**XG BOOSTING** 









### **XG BOOSTING**

```
Default - Full model
```

```
sparse.model.matrix(Status ~ .-1, data = trainset)
head(trainm)
train label <- trainset[, 'Status']</pre>
train matrix <- xgb.DMatrix(data = as.matrix(trainm), label = train labe</pre>
```

```
<u>rr[errStest mlogloss == min(errStest mlogloss),</u>]
iter train mlogloss test mlogloss
 66
            0.246145
                            0.26591
```

- One-Hot encoding performed.
- Full-Model All features were selected.
- Sensitivity for the model is 69.4%.
- Minimal test error at Iteration 66.
- Overfitting needs to be addressed.

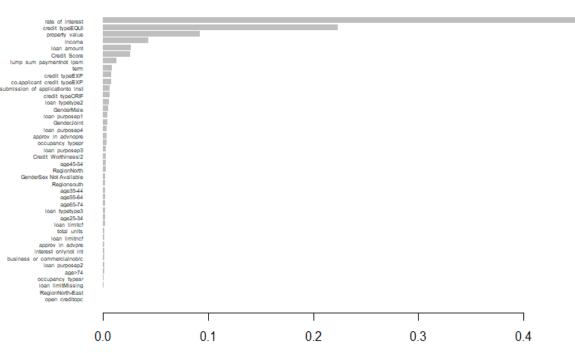
```
##### xgb.Booster
raw: 813.4 Kb
  xgb.train(params = xgb params, data = train matrix, nrounds = 100,
    watchlist = watchlist)
params (as set within xqb.train):
  objective = "multi:softprob", eval metric = "mlogloss", num class = "2"
alidate parameters = "TRUE"
xgb.attributes:
  niter
callbacks:
  cb.print.evaluation(period = print every n)
  cb.evaluation.log()
 of features: 42
                                                                          Confusion Matrix and Statistics
nfeatures: 42
 valuation log:
    iter train mlogloss test mlogloss
               0.494340
                             0.494301
               0.384170
                             0.384138
      99
               0.114238
                             0.135762
     100
               0.113959
                             0.135785
```

```
Reference
Prediction
         0 20114 2108
            327 4784
              Accuracy: 0.9109
                 95% CI: (0.9075, 0.9143)
    No Information Rate: 0.7479
    P-Value [Acc > NIR] : < 2.2e-16
                 Kappa : 0.7417
 Mcnemar's Test P-Value : < 2.2e-16
           Sensitivity: 0.6941
           Specificity: 0.9840
         Pos Pred Value: 0.9360
        Neg Pred Value: 0.9051
             Prevalence: 0.2521
         Detection Rate: 0.1750
   Detection Prevalence: 0.1870
      Balanced Accuracy: 0.8391
       'Positive' Class: 1
```





### **XG BOOSTING**





```
Feature Importance
  feat imp <- xgb.importance(colnames(train matrix), model = xgb model)</pre>
                              Feature
                                               Gain
                                                           Cover
 1:
                    rate of interest 0.4281560671 0.0680954844
 2:
                     credit typeEQUI 0.2354519833 0.0459696883
 3:
                      property value 0.1040538091 0.0965805785
 4:
                               income 0.0400439743 0.1132670056
 5:
                        Credit Score 0.0357485215 0.0839445554
 6:
                          loan amount 0.0303445304 0.1056444167
            lump sum paymentnot lpsm 0.0138323287 0.0300740218
 7:
 8:
                       credit typeEXP 0.0091163524 0.0126709181
 9:
                                 term 0.0090800566 0.0322754752
10:
                      credit typeCRIF 0.0072846171 0.0121248303
11: submission of application to inst 0.0066846194 0.0362462541
         co.applicant credit typeEXP 0.0065434217 0.0232617458
12:
```

- Feature Importance Rate Of Interest, Credit Type,
   Property value, Income.
- ✓ Hyperparameter tuning eta, max\_depth.

```
> xgb_model
##### xgb.Booster
raw: 18.1 Mb
call:
    xgb.train(params = xgb_params, data = train_matrix, prounds = 66,
    watchlist = watchlist, eta = 0.025, max_depth = 15, subsample = 0.5,
    colsample_bytree = 0.5, set.seed = 567)
params (as set within xgb.train):
```

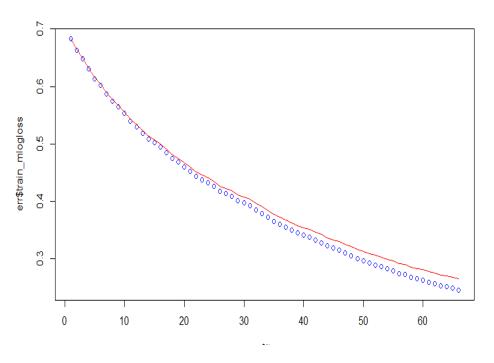


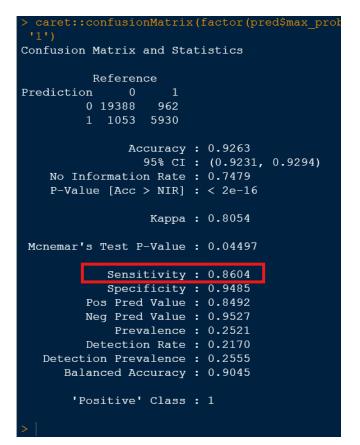


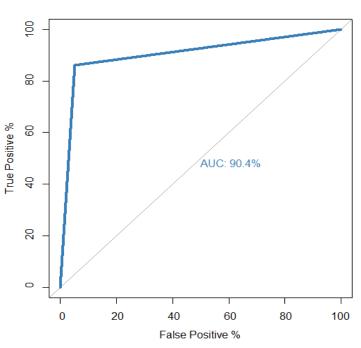
### **XG BOOSTING**



## Default – Hyperparameter Tuning







- ✓ Comparing the Train error and test error.
- ✓ Model's accuracy is 92.63 percent
- ✓ AUC is 90.4 percent.

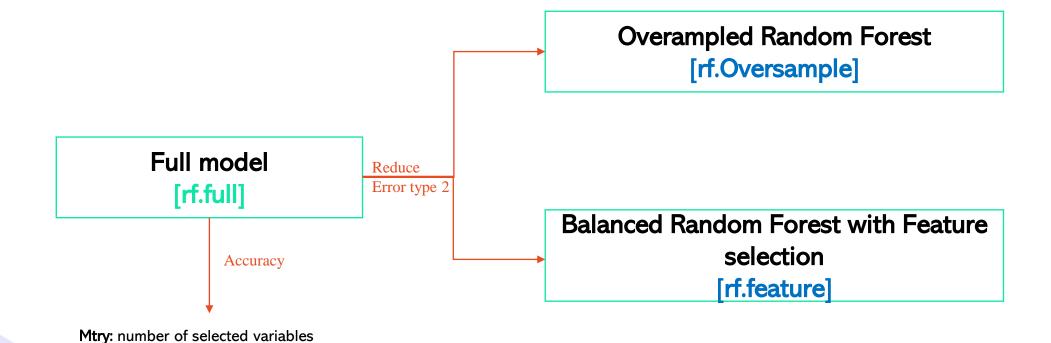


Ntree: number of trees

## **MODELING**







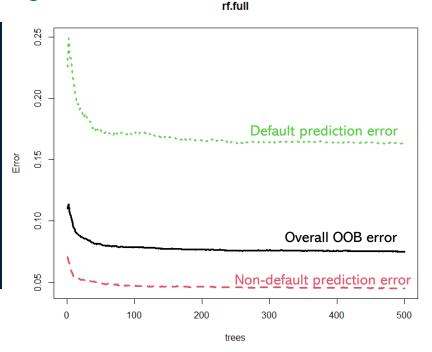




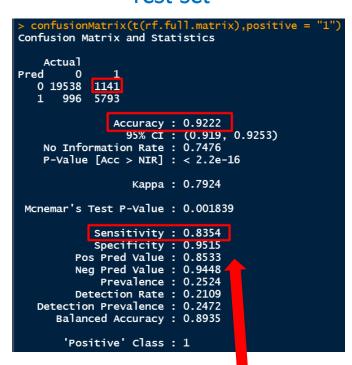
## **RANDOM FOREST**



#### Training set



#### Test set



FN is 4 times higher FP





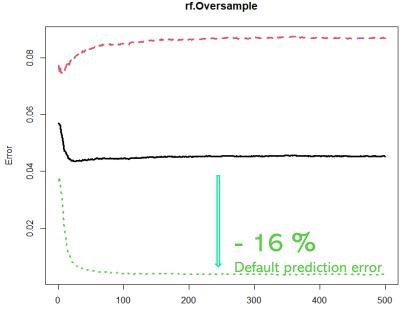


#### Oversample: to balance the set

```
> table(oversampled_training_set$Status)

0 1
82017 81997
```

#### Training set



```
confusionMatrix(t(rf.Oversample.matrix),
Confusion Matrix and Statistics
                      Actual
rf.Oversample.pred.test
                             440
                     0 18701
                     1 1833 6494
              Accuracy: 0.9172
                95% CI : (0.9139, 0.9205)
   No Information Rate: 0.7476
   P-Value [Acc > NIR] : < 2.2e-16
                 Kappa: 0.7944
Mcnemar's Test P-Value : < 2.2e-16
           Sensitivity : 0.9365 +10%
           Specificity: 0.9107
        Pos Pred Value: 0.7799
        Neg Pred Value: 0.9770
            Prevalence: 0.2524
        Detection Rate: 0.2364
  Detection Prevalence: 0.3032
     Balanced Accuracy: 0.9236
       'Positive' Class: 1
```





## **RANDOM FOREST**



### Balanced model with feature selection

#### Training set

#### Call: randomForest(formula = Status ~ credit\_type + rate\_of\_inter property\_value + lump\_sum\_payment + income + co.a pplicant\_credit\_type + submission\_of\_application + loan \_type + business\_or\_commercial + g\_set, ntree = 100, mtry = 4, strata = training\_set\$Status, sampsize = rep(sum(training\_set\$Status == 1), 2)) Type of random forest: classification Number of trees: 100 No. of variables tried at each split: 4 OOB estimate of error rate: 10% Confusion matrix: 0 1 class.error 71055 10962 0.1336552178 26 27825 0.0009335392

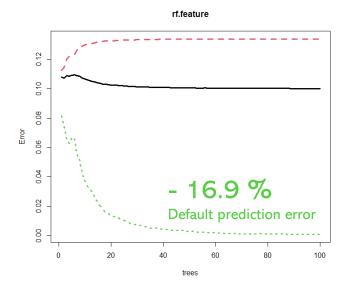
Feature selection: only powerful predictors

-> more precise



Strata: Used for stratified sampling. Sampsize: Sample size.

-> enable the balanced bootstrapping



```
> confusionMatrix(t(rf.feature.matrix),po
Confusion Matrix and Statistics
                   Actual
rf.feature.pred.test
                  0 17756
                  1 2778
                           6934
              Accuracy: 0.8989
                95% CI : (0.8952, 0.9024)
    No Information Rate: 0.7476
    P-Value [Acc > NIR] : < 2.2e-16
                 Kappa: 0.7634
 Mcnemar's Test P-Value : < 2.2e-16
           Sensitivity: 1.0000 + 17\%
           Specificity: 0.8647
        Pos Pred Value: 0.7140
        Neg Pred Value: 1.0000
             Prevalence: 0.2524
         Detection Rate: 0.2524
  Detection Prevalence: 0.3536
     Balanced Accuracy: 0.9324
       'Positive' Class: 1
```

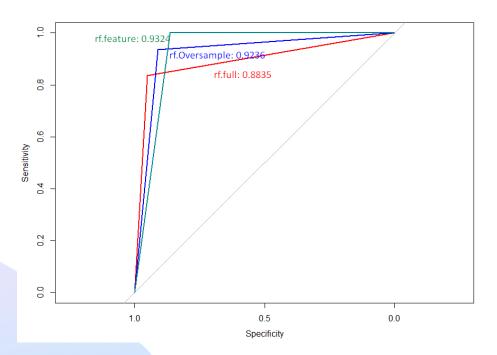


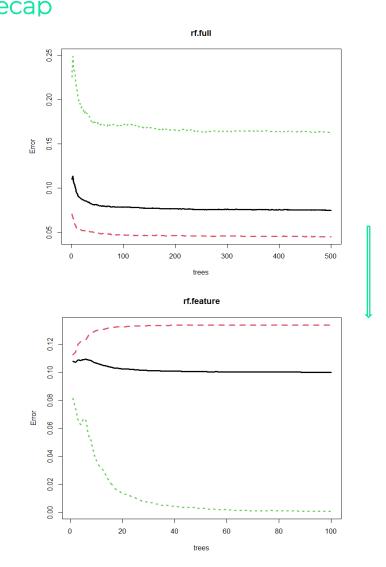


# **RANDOM FOREST**



Model	Model Accuracy		FN	F1	AUC
Rf.full	92.22%	83.54%	1141	0.8442	0.8935
Rf.Oversample	91.72%	93.65%	440	0.8510	0.9236
Rf.feature	89.89%	100.00%	0	0.8331	0.9324





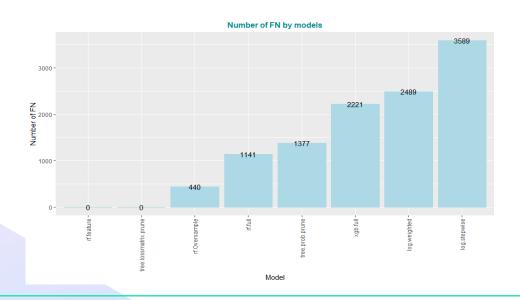


## **COMPARISON - SELECTION**



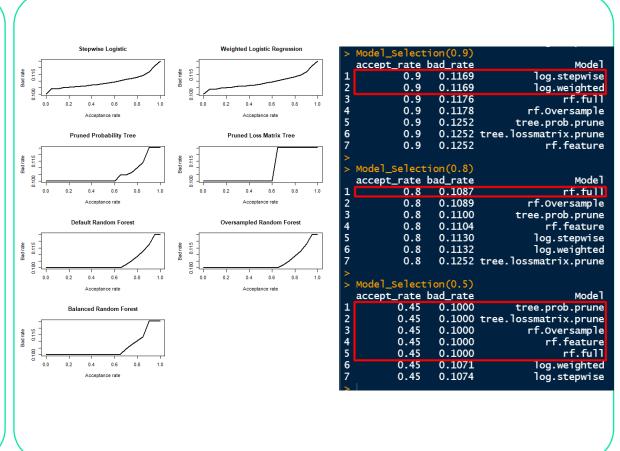
### Selection by **Model Performance metrics**

Туре	Model	Accuracy	Sensitivity	FN	F1	AUC
Logistic Regression	log.stepwise	86.33%	48.24%	3589	0.640	0.845
Logistic Regression	log.weighted	81.70%	64.10%	2489	0.639	0.846
Decision Tree	tree.prob.prune	91.28%	80.14%	1377	0.822	0.876
Decision Tree	tree.lossmatrix.prune	89.88%	100.00%	0	0.833	0.932
XGBoost	xgb.full	91.09%	69.41%	2108	0.850	0.839
XGBoost	xgb.feature	92.63%	86.04%	962	0.874	0.904
Random Forest	rf.full	92.22%	83.54%	1141	0.844	0.894
Random Forest	rf.Oversample	91.72%	93.65%	440	0.851	0.924
Random Forest	rf.feature	89.89%	100.00%	0	0.833	0.932





### Selection by **Domain Knowledge: Bad rate strategy**



### **CONCLUSION**

1. Accuracy is not always the prioritized metrics

- 2. F1 score is also misleading, the right metrics is Sensitivity
- 3. Consider **various solutions** when dealing with Imbalanced data set
- 4. Beside metrics, **domain knowledge** is also another imperative choice when comparing models
- 5. For this data set, overall, Random Forest turns out to be the most powerful models as it uses boosting methods

#### Data process

- Upsample
- $\cdot {\sf Downsample}$
- ·SMOTE
- · ROSE

#### Modeling

- ·Tune hyperparameters
- ·Update model
- · Pruning

#### Threshold

·Cut-off

·Cp



## **LOAN DEFAULT**

**Group Project - OUTLIERS** 

# THANK YOU