# R Notebook

## Parametros:

**Measure =** Accuracy
**Columns =** sampling, weight_space, underbagging, learner
**Performance =** tuning_measure
**Filter keys =** imba.rate
**Filter values =** 0.03

```r
library("scmamp")
library(dplyr)
```

## Tratamento dos dados

Carregando data set compilado

```r
ds = read.csv("/home/rodrigo/Dropbox/UNICAMP/IC/estudo_cost_learning/SummaryResults/summary_compilation
```

```r
ds = filter(ds, learner != "classif.rusboost")
summary(ds)
```

```
##                   learner       weight_space
##  classif.ksvm        :17100   Mode :logical
##  classif.randomForest:17100   FALSE:41040
##  classif.rusboost    :    0   TRUE :10260
##  classif.xgboost     :17100   NA's :0
##
##
##
##                                 measure        sampling      underbagging
##  Accuracy                      :10260   ADASYN:10260   Mode :logical
##  Area under the curve          :10260   FALSE :30780   FALSE:41040
##  F1 measure                    :10260   SMOTE :10260   TRUE :10260
##  G-mean                        :10260                  NA's :0
##  Matthews correlation coefficient:10260
##
##
##  tuning_measure     holdout_measure    holdout_measure_residual
##  Min.   :-0.1277   Min.   :-0.2120   Min.   :-0.4658
##  1st Qu.: 0.6911   1st Qu.: 0.4001   1st Qu.: 0.1994
##  Median : 0.9700   Median : 0.8571   Median : 0.5581
##  Mean   : 0.7903   Mean   : 0.6718   Mean   : 0.5298
##  3rd Qu.: 0.9975   3rd Qu.: 0.9900   3rd Qu.: 0.8755
##  Max.   : 1.0000   Max.   : 1.0000   Max.   : 1.0000
##  NA's   :1077      NA's   :1077      NA's   :1077
##  iteration_count              dataset        imba.rate
##  Min.   :1       abalone           : 900   Min.   :0.0010
##  1st Qu.:1       adult             : 900   1st Qu.:0.0100
##  Median :2       bank              : 900   Median :0.0300
##  Mean   :2       car               : 900   Mean   :0.0286
```

```
## 3rd Qu.:3       cardiotocography-10clases:  900    3rd Qu.:0.0500
## Max.   :3       cardiotocography-3clases :  900    Max.   :0.0500
## NA's   :1077    (Other)                  :45900
```

Filtrando pela metrica

```
ds = filter(ds, measure == params$measure)
```

Filtrando o data set

```
if(params$filter_keys != 'NULL' && !is.null(params$filter_keys)){
  dots = paste0(params$filter_keys," == '",params$filter_values,"'")
  ds = filter_(ds, .dots = dots)
}
```

```
summary(ds)
```

```
##                   learner    weight_space
## classif.ksvm        :990    Mode :logical
## classif.randomForest:990    FALSE:2376
## classif.rusboost    :  0    TRUE :594
## classif.xgboost     :990    NA's :0
##
##
##
##                               measure       sampling    underbagging
## Accuracy                        :2970    ADASYN: 594    Mode :logical
## Area under the curve            :  0     FALSE :1782    FALSE:2376
## F1 measure                      :  0     SMOTE : 594    TRUE :594
## G-mean                          :  0                    NA's :0
## Matthews correlation coefficient:  0
##
##
## tuning_measure     holdout_measure    holdout_measure_residual
## Min.   :0.09041    Min.   :0.02655    Min.   :0.0346
## 1st Qu.:0.96926    1st Qu.:0.96647    1st Qu.:0.3599
## Median :0.98130    Median :0.97619    Median :0.6882
## Mean   :0.95405    Mean   :0.94750    Mean   :0.6478
## 3rd Qu.:0.99560    3rd Qu.:0.99045    3rd Qu.:0.9438
## Max.   :1.00000    Max.   :1.00000    Max.   :1.0000
## NA's   :57         NA's   :57         NA's   :57
## iteration_count           dataset        imba.rate
## Min.   :1        abalone      : 45    Min.   :0.03
## 1st Qu.:1        adult        : 45    1st Qu.:0.03
## Median :2        annealing    : 45    Median :0.03
## Mean   :2        arrhythmia   : 45    Mean   :0.03
## 3rd Qu.:3        balance-scale: 45    3rd Qu.:0.03
## Max.   :3        bank         : 45    Max.   :0.03
## NA's   :57       (Other)      :2700
```

Computando as médias das iteracoes

```
ds = group_by(ds, learner , weight_space , measure , sampling , underbagging , dataset , imba.rate)
ds = summarise(ds, tuning_measure = mean(tuning_measure), holdout_measure = mean(holdout_measure),
          holdout_measure_residual = mean(holdout_measure_residual))

ds = as.data.frame(ds)
```

Criando dataframe

```r
# Dividindo o ds em n, um para cada técnica
splited_df = ds %>% group_by_at(.vars = params$columns) %>% do(vals = as.data.frame(.)) %>% select(vals)

# Juntando cada uma das partes horizontalmente em um data set
df_tec_wide = do.call("cbind", splited_df)

# Renomeando duplicacao de nomes
colnames(df_tec_wide) = make.unique(colnames(df_tec_wide))

# Selecionando apenas as medidas da performance escolhida
df_tec_wide_residual = select(df_tec_wide, matches(paste("^", params$performance, "$|", params$performan

# Renomeando colunas
new_names = NULL
for(i in (1:length(splited_df))){
  id = toString(sapply(splited_df[[i]][1, params$columns], as.character))
  new_names = c(new_names, id)
}
colnames(df_tec_wide_residual) = new_names

# Verificando a dimensao do df
dim(df_tec_wide_residual)
```

```
## [1] 66 15
```

```r
# Renomeando a variavel
df = df_tec_wide_residual

head(df)
```

```
##   ADASYN, FALSE, FALSE, classif.ksvm
## 1                          0.9595381
## 2                          0.9749599
## 3                          0.9831176
## 4                          0.9241574
## 5                          1.0000000
## 6                          0.9850835
##   ADASYN, FALSE, FALSE, classif.randomForest
## 1                                  0.9566567
## 2                                         NA
## 3                                  0.9926403
## 4                                  0.9970458
## 5                                  1.0000000
## 6                                         NA
##   ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
## 1                             0.9671889                        0.9578372
## 2                             0.9842633                        0.9663640
## 3                             0.9883465                        0.9691089
## 4                             0.9866841                        0.9684141
## 5                             1.0000000                        0.9975170
## 6                             0.9842693                        0.9699278
##   FALSE, FALSE, FALSE, classif.randomForest
## 1                                 0.9697619
## 2                                        NA
```

```
## 3                             0.9798042
## 4                             0.9779849
## 5                             1.0000000
## 6                             0.9695214
##   FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
## 1                            0.9697619                          0.6175442
## 2                            0.9775379                          0.9163519
## 3                            0.9750536                          0.9489037
## 4                            0.9760693                          0.8627718
## 5                            1.0000000                          0.9838255
## 6                            0.9693184                          0.5901602
##   FALSE, FALSE, TRUE, classif.randomForest
## 1                            0.5934099
## 2                            0.7957626
## 3                            0.8300747
## 4                            0.9061451
## 5                            0.9267302
## 6                            0.8092056
##   FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
## 1                            0.6314521                          0.9625227
## 2                            0.8091126                          0.9662495
## 3                            0.8145848                          0.9691089
## 4                            0.9214709                          0.9684141
## 5                            0.9240949                          0.9975170
## 6                            0.7991469                          0.9699278
##   FALSE, TRUE, FALSE, classif.randomForest
## 1                            0.9697619
## 2                            0.9769489
## 3                            0.9786156
## 4                            0.9808670
## 5                            1.0000000
## 6                            0.9697246
##   FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
## 1                            0.9699038                          0.9611646
## 2                            0.9772270                          0.9763133
## 3                            0.9738621                          0.9822298
## 4                            0.9751364                          0.8695467
## 5                            1.0000000                          1.0000000
## 6                            0.9693184                          0.9840940
##   SMOTE, FALSE, FALSE, classif.randomForest
## 1                            0.9550077
## 2                            0.9795664
## 3                            0.9929624
## 4                            0.9950040
## 5                            1.0000000
## 6                            0.9848819
##   SMOTE, FALSE, FALSE, classif.xgboost
## 1                            0.9676123
## 2                            0.9851323
## 3                            0.9880556
## 4                            0.9890114
## 5                            1.0000000
## 6                            0.9849344
```

```r
summary(df)
```

```
##   ADASYN, FALSE, FALSE, classif.ksvm
##   Min.   :0.7764
##   1st Qu.:0.9809
##   Median :0.9942
##   Mean   :0.9811
##   3rd Qu.:0.9977
##   Max.   :1.0000
##   NA's   :2
##   ADASYN, FALSE, FALSE, classif.randomForest
##   Min.   :0.7717
##   1st Qu.:0.9830
##   Median :0.9968
##   Mean   :0.9865
##   3rd Qu.:0.9993
##   Max.   :1.0000
##   NA's   :7
##   ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
##   Min.   :0.7715                         Min.   :0.9569
##   1st Qu.:0.9826                         1st Qu.:0.9697
##   Median :0.9936                         Median :0.9706
##   Mean   :0.9853                         Mean   :0.9752
##   3rd Qu.:0.9983                         3rd Qu.:0.9810
##   Max.   :1.0000                         Max.   :0.9977
##
##   FALSE, FALSE, FALSE, classif.randomForest
##   Min.   :0.9680
##   1st Qu.:0.9705
##   Median :0.9802
##   Mean   :0.9822
##   3rd Qu.:0.9922
##   Max.   :1.0000
##   NA's   :2
##   FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
##   Min.   :0.9656                        Min.   :0.2896
##   1st Qu.:0.9737                        1st Qu.:0.6963
##   Median :0.9819                        Median :0.9189
##   Mean   :0.9833                        Mean   :0.8290
##   3rd Qu.:0.9930                        3rd Qu.:0.9775
##   Max.   :1.0000                        Max.   :0.9973
##
##   FALSE, FALSE, TRUE, classif.randomForest
##   Min.   :0.5074
##   1st Qu.:0.7763
##   Median :0.8945
##   Mean   :0.8529
##   3rd Qu.:0.9666
##   Max.   :1.0000
##   NA's   :2
##   FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
##   Min.   :0.4993                       Min.   :0.9617
##   1st Qu.:0.7990                       1st Qu.:0.9697
##   Median :0.8917                       Median :0.9705
```

```
## Mean   :0.8463              Mean   :0.9751
## 3rd Qu.:0.9540              3rd Qu.:0.9810
## Max.   :1.0000              Max.   :0.9977
##
## FALSE, TRUE, FALSE, classif.randomForest
## Min.   :0.9677
## 1st Qu.:0.9705
## Median :0.9798
## Mean   :0.9823
## 3rd Qu.:0.9928
## Max.   :1.0000
## NA's   :1
## FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
## Min.   :0.9656              Min.   :0.7761
## 1st Qu.:0.9742              1st Qu.:0.9769
## Median :0.9818              Median :0.9924
## Mean   :0.9833              Mean   :0.9808
## 3rd Qu.:0.9929              3rd Qu.:0.9976
## Max.   :1.0000              Max.   :1.0000
##
## SMOTE, FALSE, FALSE, classif.randomForest
## Min.   :0.7750
## 1st Qu.:0.9819
## Median :0.9954
## Mean   :0.9865
## 3rd Qu.:0.9993
## Max.   :1.0000
## NA's   :5
## SMOTE, FALSE, FALSE, classif.xgboost
## Min.   :0.7750
## 1st Qu.:0.9822
## Median :0.9941
## Mean   :0.9860
## 3rd Qu.:0.9981
## Max.   :1.0000
##
```

## Verificando a média de cada coluna selecionada

```
for(i in (1:dim(df)[2])){
  print(paste("Media da coluna ", colnames(df)[i], " = ", mean(df[,i], na.rm = TRUE), sep=""))
}
```

```
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.ksvm = 0.981089172328964"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.randomForest = 0.986507892933772"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.xgboost = 0.985320119782326"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.ksvm = 0.975237303652973"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.randomForest = 0.982232504000285"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.xgboost = 0.983314101281034"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.ksvm = 0.828952284039238"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.randomForest = 0.852911103010748"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.xgboost = 0.846292309755586"
```

```
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.ksvm = 0.975071996792856"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.randomForest = 0.982303264242282"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.xgboost = 0.983275975322345"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.ksvm = 0.980768143223078"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.randomForest = 0.986499616970574"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.xgboost = 0.98596666289117"
```

# Fazendo teste de normalidade

```
plotDensities(data = na.omit(df))
```



# Testando as diferencas

```
friedmanTest(df)
```

```
##
##  Friedman's rank sum test
##
## data:  df
## Friedman's chi-squared = 499.08, df = 14, p-value < 2.2e-16
```

## Testando as diferencas par a par

```
test <- nemenyiTest (df, alpha=0.05)
abs(test$diff.matrix) > test$statistic
```

```
##          ADASYN, FALSE, FALSE, classif.ksvm
##   [1,]                            FALSE
##   [2,]                            FALSE
##   [3,]                            FALSE
##   [4,]                             TRUE
##   [5,]                            FALSE
##   [6,]                            FALSE
##   [7,]                             TRUE
##   [8,]                             TRUE
##   [9,]                             TRUE
##  [10,]                             TRUE
##  [11,]                            FALSE
##  [12,]                            FALSE
##  [13,]                            FALSE
##  [14,]                            FALSE
##  [15,]                            FALSE
##          ADASYN, FALSE, FALSE, classif.randomForest
##   [1,]                                  FALSE
##   [2,]                                  FALSE
##   [3,]                                  FALSE
##   [4,]                                   TRUE
##   [5,]                                   TRUE
##   [6,]                                  FALSE
##   [7,]                                   TRUE
##   [8,]                                   TRUE
##   [9,]                                   TRUE
##  [10,]                                   TRUE
##  [11,]                                   TRUE
##  [12,]                                  FALSE
##  [13,]                                  FALSE
##  [14,]                                  FALSE
##  [15,]                                  FALSE
##          ADASYN, FALSE, FALSE, classif.xgboost
##   [1,]                              FALSE
##   [2,]                              FALSE
##   [3,]                              FALSE
##   [4,]                               TRUE
##   [5,]                               TRUE
##   [6,]                              FALSE
##   [7,]                               TRUE
##   [8,]                               TRUE
##   [9,]                               TRUE
##  [10,]                               TRUE
##  [11,]                               TRUE
##  [12,]                              FALSE
##  [13,]                              FALSE
##  [14,]                              FALSE
##  [15,]                              FALSE
```

```
##         FALSE, FALSE, FALSE, classif.ksvm
##  [1,]                             TRUE
##  [2,]                             TRUE
##  [3,]                             TRUE
##  [4,]                            FALSE
##  [5,]                            FALSE
##  [6,]                             TRUE
##  [7,]                             TRUE
##  [8,]                             TRUE
##  [9,]                             TRUE
## [10,]                            FALSE
## [11,]                            FALSE
## [12,]                             TRUE
## [13,]                             TRUE
## [14,]                             TRUE
## [15,]                             TRUE
##         FALSE, FALSE, FALSE, classif.randomForest
##  [1,]                                    FALSE
##  [2,]                                     TRUE
##  [3,]                                     TRUE
##  [4,]                                    FALSE
##  [5,]                                    FALSE
##  [6,]                                    FALSE
##  [7,]                                     TRUE
##  [8,]                                     TRUE
##  [9,]                                     TRUE
## [10,]                                    FALSE
## [11,]                                    FALSE
## [12,]                                    FALSE
## [13,]                                     TRUE
## [14,]                                     TRUE
## [15,]                                     TRUE
##         FALSE, FALSE, FALSE, classif.xgboost
##  [1,]                               FALSE
##  [2,]                               FALSE
##  [3,]                               FALSE
##  [4,]                                TRUE
##  [5,]                               FALSE
##  [6,]                               FALSE
##  [7,]                                TRUE
##  [8,]                                TRUE
##  [9,]                                TRUE
## [10,]                                TRUE
## [11,]                               FALSE
## [12,]                               FALSE
## [13,]                               FALSE
## [14,]                                TRUE
## [15,]                               FALSE
##         FALSE, FALSE, TRUE, classif.ksvm
##  [1,]                            TRUE
##  [2,]                            TRUE
##  [3,]                            TRUE
##  [4,]                            TRUE
##  [5,]                            TRUE
```

```
##  [6,]                               TRUE
##  [7,]                              FALSE
##  [8,]                              FALSE
##  [9,]                              FALSE
## [10,]                               TRUE
## [11,]                               TRUE
## [12,]                               TRUE
## [13,]                               TRUE
## [14,]                               TRUE
## [15,]                               TRUE
##        FALSE, FALSE, TRUE, classif.randomForest
##  [1,]                                  TRUE
##  [2,]                                  TRUE
##  [3,]                                  TRUE
##  [4,]                                  TRUE
##  [5,]                                  TRUE
##  [6,]                                  TRUE
##  [7,]                                 FALSE
##  [8,]                                 FALSE
##  [9,]                                 FALSE
## [10,]                                  TRUE
## [11,]                                  TRUE
## [12,]                                  TRUE
## [13,]                                  TRUE
## [14,]                                  TRUE
## [15,]                                  TRUE
##        FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
##  [1,]                                 TRUE                             TRUE
##  [2,]                                 TRUE                             TRUE
##  [3,]                                 TRUE                             TRUE
##  [4,]                                 TRUE                            FALSE
##  [5,]                                 TRUE                            FALSE
##  [6,]                                 TRUE                             TRUE
##  [7,]                                FALSE                             TRUE
##  [8,]                                FALSE                             TRUE
##  [9,]                                FALSE                             TRUE
## [10,]                                 TRUE                            FALSE
## [11,]                                 TRUE                            FALSE
## [12,]                                 TRUE                             TRUE
## [13,]                                 TRUE                             TRUE
## [14,]                                 TRUE                             TRUE
## [15,]                                 TRUE                             TRUE
##        FALSE, TRUE, FALSE, classif.randomForest
##  [1,]                                 FALSE
##  [2,]                                  TRUE
##  [3,]                                  TRUE
##  [4,]                                 FALSE
##  [5,]                                 FALSE
##  [6,]                                 FALSE
##  [7,]                                  TRUE
##  [8,]                                  TRUE
##  [9,]                                  TRUE
## [10,]                                 FALSE
## [11,]                                 FALSE
```

```
## [12,]                                      FALSE
## [13,]                                      FALSE
## [14,]                                       TRUE
## [15,]                                       TRUE
##        FALSE, TRUE, FALSE, classif.xgboost
##  [1,]                                      FALSE
##  [2,]                                      FALSE
##  [3,]                                      FALSE
##  [4,]                                       TRUE
##  [5,]                                      FALSE
##  [6,]                                      FALSE
##  [7,]                                       TRUE
##  [8,]                                       TRUE
##  [9,]                                       TRUE
## [10,]                                       TRUE
## [11,]                                      FALSE
## [12,]                                      FALSE
## [13,]                                      FALSE
## [14,]                                       TRUE
## [15,]                                      FALSE
##        SMOTE, FALSE, FALSE, classif.ksvm
##  [1,]                                      FALSE
##  [2,]                                      FALSE
##  [3,]                                      FALSE
##  [4,]                                       TRUE
##  [5,]                                       TRUE
##  [6,]                                      FALSE
##  [7,]                                       TRUE
##  [8,]                                       TRUE
##  [9,]                                       TRUE
## [10,]                                       TRUE
## [11,]                                      FALSE
## [12,]                                      FALSE
## [13,]                                      FALSE
## [14,]                                      FALSE
## [15,]                                      FALSE
##        SMOTE, FALSE, FALSE, classif.randomForest
##  [1,]                                      FALSE
##  [2,]                                      FALSE
##  [3,]                                      FALSE
##  [4,]                                       TRUE
##  [5,]                                       TRUE
##  [6,]                                       TRUE
##  [7,]                                       TRUE
##  [8,]                                       TRUE
##  [9,]                                       TRUE
## [10,]                                       TRUE
## [11,]                                       TRUE
## [12,]                                       TRUE
## [13,]                                      FALSE
## [14,]                                      FALSE
## [15,]                                      FALSE
##        SMOTE, FALSE, FALSE, classif.xgboost
##  [1,]                                      FALSE
```

11

```
##  [2,]                               FALSE
##  [3,]                               FALSE
##  [4,]                                TRUE
##  [5,]                                TRUE
##  [6,]                               FALSE
##  [7,]                                TRUE
##  [8,]                                TRUE
##  [9,]                                TRUE
## [10,]                                TRUE
## [11,]                                TRUE
## [12,]                               FALSE
## [13,]                               FALSE
## [14,]                               FALSE
## [15,]                               FALSE
```

# Plotando os ranks

```r
print(colMeans(rankMatrix(df)))
```

```
##           ADASYN, FALSE, FALSE, classif.ksvm
##                                     5.659091
## ADASYN, FALSE, FALSE, classif.randomForest
##                                     4.803030
##        ADASYN, FALSE, FALSE, classif.xgboost
##                                     4.727273
##            FALSE, FALSE, FALSE, classif.ksvm
##                                    10.121212
##  FALSE, FALSE, FALSE, classif.randomForest
##                                     7.954545
##        FALSE, FALSE, FALSE, classif.xgboost
##                                     7.143939
##             FALSE, FALSE, TRUE, classif.ksvm
##                                    12.893939
##   FALSE, FALSE, TRUE, classif.randomForest
##                                    13.272727
##        FALSE, FALSE, TRUE, classif.xgboost
##                                    13.674242
##             FALSE, TRUE, FALSE, classif.ksvm
##                                    10.196970
##   FALSE, TRUE, FALSE, classif.randomForest
##                                     7.787879
##        FALSE, TRUE, FALSE, classif.xgboost
##                                     7.234848
##            SMOTE, FALSE, FALSE, classif.ksvm
##                                     5.250000
##  SMOTE, FALSE, FALSE, classif.randomForest
##                                     4.492424
##        SMOTE, FALSE, FALSE, classif.xgboost
##                                     4.787879
```

# Plotando grafico de Critical Diference

```
result = tryCatch({
    plotCD(df, alpha=0.05, cex = 0.35)
}, error = function(e) {})
```