# R Notebook

## Parametros:

**Measure =** Accuracy
**Columns =** sampling, weight_space, underbagging, learner
**Performance =** holdout_measure_residual
**Filter keys =** imba.rate
**Filter values =** 0.05

```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          :1230  Mode :logical
## classif.randomForest:1230  FALSE:2952
## classif.rusboost    :   0  TRUE :738
## classif.xgboost     :1230  NA's :0
##
##
##
##                              measure        sampling     underbagging
## Accuracy                       :3690  ADASYN: 738  Mode :logical
## Area under the curve           :   0  FALSE :2214  FALSE:2952
## F1 measure                     :   0  SMOTE : 738  TRUE :738
## G-mean                         :   0              NA's :0
## Matthews correlation coefficient:   0
##
##
## tuning_measure    holdout_measure    holdout_measure_residual
## Min.   :0.2470   Min.   :0.04739   Min.   :0.0367
## 1st Qu.:0.9494   1st Qu.:0.94505   1st Qu.:0.3902
## Median :0.9688   Median :0.96078   Median :0.7223
## Mean   :0.9425   Mean   :0.93413   Mean   :0.6602
## 3rd Qu.:0.9908   3rd Qu.:0.98413   3rd Qu.:0.9315
## Max.   :1.0000   Max.   :1.00000   Max.   :1.0000
## NA's   :42       NA's   :42        NA's   :42
## iteration_count        dataset        imba.rate
## Min.   :1       abalone      : 45  Min.   :0.05
## 1st Qu.:1       adult        : 45  1st Qu.:0.05
## Median :2       annealing    : 45  Median :0.05
## Mean   :2       arrhythmia   : 45  Mean   :0.05
## 3rd Qu.:3       balance-scale: 45  3rd Qu.:0.05
## Max.   :3       bank         : 45  Max.   :0.05
## NA's   :42      (Other)      :3420
```

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)
```

2

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] 82 15
```

```r
# Renomeando a variavel
df = df_tec_wide_residual

head(df)
```

```
##    ADASYN, FALSE, FALSE, classif.ksvm
## 1                           0.3933596
## 2                           0.5230756
## 3                           0.7788618
## 4                           0.9917695
## 5                           0.3276515
## 6                           0.6973873
##    ADASYN, FALSE, FALSE, classif.randomForest
## 1                                   0.4003001
## 2                                          NA
## 3                                   0.9154472
## 4                                   0.9862826
## 5                                   0.2765152
## 6                                   0.7536606
##    ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
## 1                             0.3674733                          0.3513412
## 2                             0.6378860                          0.4414169
## 3                             0.9154472                          0.7414634
## 4                             0.9670782                          0.9917695
## 5                             0.3011364                          0.3418561
## 6                             0.7444732                          0.6962389
##    FALSE, FALSE, FALSE, classif.randomForest
## 1                                  0.3230163
## 2                                  0.5904008
```

```
## 3                              0.8975610
## 4                              0.9670782
## 5                              0.3011364
## 6                              0.6962389
##   FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
## 1                            0.3232039                        0.6178953
## 2                            0.5825954                        0.5486490
## 3                            0.8601626                        0.7528455
## 4                            0.9821674                        0.8888889
## 5                            0.3011364                        0.3106061
## 6                            0.7137525                        0.6342234
##   FALSE, FALSE, TRUE, classif.randomForest
## 1                              0.6614144
## 2                                     NA
## 3                              0.8943089
## 4                              0.8971193
## 5                              0.3617424
## 6                              0.7912719
##   FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
## 1                           0.6631026                       0.3425249
## 2                           0.8341564                       0.4757891
## 3                           0.8585366                       0.7772358
## 4                           0.9108368                       0.9917695
## 5                           0.4204545                       0.3418561
## 6                           0.8079242                       0.6962389
##   FALSE, TRUE, FALSE, classif.randomForest
## 1                              0.3230163
## 2                              0.5923308
## 3                              0.9154472
## 4                              0.9766804
## 5                              0.2954545
## 6                              0.6962389
##   FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
## 1                           0.3230163                       0.4019884
## 2                           0.5813465                       0.5280143
## 3                           0.8634146                       0.7853659
## 4                           0.9739369                       0.9917695
## 5                           0.3011364                       0.2755682
## 6                           0.7074361                       0.7034166
##   SMOTE, FALSE, FALSE, classif.randomForest
## 1                              0.3869818
## 2                              0.6434208
## 3                              0.9284553
## 4                              0.9835391
## 5                              0.2831439
## 6                              0.7519380
##   SMOTE, FALSE, FALSE, classif.xgboost
## 1                           0.3627837
## 2                           0.6342813
## 3                           0.9073171
## 4                           0.9753086
## 5                           0.2831439
## 6                           0.7482056
```

```
summary(df)
```

```
##   ADASYN, FALSE, FALSE, classif.ksvm
##   Min.    :0.03682
##   1st Qu.:0.34645
##   Median :0.52308
##   Mean    :0.60643
##   3rd Qu.:0.92785
##   Max.    :0.99985
##   NA's    :1
##   ADASYN, FALSE, FALSE, classif.randomForest
##   Min.    :0.03983
##   1st Qu.:0.43067
##   Median :0.73860
##   Mean    :0.68415
##   3rd Qu.:0.93951
##   Max.    :0.99987
##   NA's    :4
##   ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
##   Min.    :0.04563                      Min.    :0.0367
##   1st Qu.:0.44801                       1st Qu.:0.3242
##   Median :0.75505                       Median :0.4871
##   Mean    :0.69543                      Mean    :0.5914
##   3rd Qu.:0.92881                       3rd Qu.:0.9278
##   Max.    :0.99985                      Max.    :0.9999
##
##   FALSE, FALSE, FALSE, classif.randomForest
##   Min.    :0.1302
##   1st Qu.:0.3636
##   Median :0.7077
##   Mean    :0.6472
##   3rd Qu.:0.9321
##   Max.    :0.9999
##   NA's    :1
##   FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
##   Min.    :0.03977                     Min.    :0.1657
##   1st Qu.:0.36236                      1st Qu.:0.4411
##   Median :0.72511                      Median :0.7075
##   Mean    :0.65416                     Mean    :0.6629
##   3rd Qu.:0.94552                      3rd Qu.:0.8559
##   Max.    :0.99986                     Max.    :0.9993
##
##   FALSE, FALSE, TRUE, classif.randomForest
##   Min.    :0.2376
##   1st Qu.:0.6493
##   Median :0.7958
##   Mean    :0.7529
##   3rd Qu.:0.9186
##   Max.    :0.9998
##   NA's    :3
##   FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
##   Min.    :0.2192                     Min.    :0.0367
##   1st Qu.:0.6255                      1st Qu.:0.3242
##   Median :0.7911                      Median :0.4888
```

```
## Mean    :0.7414                    Mean    :0.5872
## 3rd Qu.:0.9140                    3rd Qu.:0.9242
## Max.    :0.9998                    Max.    :0.9999
##
## FALSE, TRUE, FALSE, classif.randomForest
## Min.    :0.1010
## 1st Qu.:0.3537
## Median :0.7094
## Mean    :0.6454
## 3rd Qu.:0.9265
## Max.    :1.0000
## NA's    :1
## FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
## Min.    :0.04244                   Min.    :0.03682
## 1st Qu.:0.36122                   1st Qu.:0.34071
## Median :0.70828                   Median :0.49006
## Mean    :0.65223                   Mean    :0.60576
## 3rd Qu.:0.94398                   3rd Qu.:0.93463
## Max.    :1.00000                   Max.    :0.99971
##
## SMOTE, FALSE, FALSE, classif.randomForest
## Min.    :0.04093
## 1st Qu.:0.42732
## Median :0.74822
## Mean    :0.68263
## 3rd Qu.:0.94508
## Max.    :0.99985
## NA's    :4
## SMOTE, FALSE, FALSE, classif.xgboost
## Min.    :0.04523
## 1st Qu.:0.45109
## Median :0.74671
## Mean    :0.69786
## 3rd Qu.:0.92855
## Max.    :0.99986
##
```

## Verificando a média de cada coluna selecionada

```r
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 = 0.606430443124862"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.randomForest = 0.684150717203129"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.xgboost = 0.695427881082246"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.ksvm = 0.591383456632022"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.randomForest = 0.647159253696668"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.xgboost = 0.654157297459171"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.ksvm = 0.66292252931822"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.randomForest = 0.752869168914993"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.xgboost = 0.741430717787393"
```
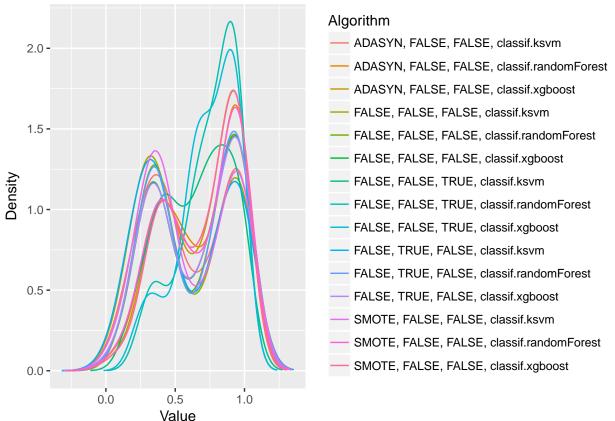
```
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.ksvm = 0.587209270677754"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.randomForest = 0.645394442377463"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.xgboost = 0.652228617182511"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.ksvm = 0.605755380600871"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.randomForest = 0.682624754358001"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.xgboost = 0.697862502217016"
```

# 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 = 107.35, df = 14, p-value = 2.22e-16
```

```
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.ksvm = 0.587209270677754"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.randomForest = 0.645394442377463"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.xgboost = 0.652228617182511"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.ksvm = 0.605755380600871"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.randomForest = 0.682624754358001"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.xgboost = 0.697862502217016"
```

# 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 = 107.35, df = 14, p-value = 2.22e-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,]                             TRUE
##  [4,]                            FALSE
##  [5,]                            FALSE
##  [6,]                            FALSE
##  [7,]                            FALSE
##  [8,]                             TRUE
##  [9,]                             TRUE
## [10,]                            FALSE
## [11,]                            FALSE
## [12,]                            FALSE
## [13,]                            FALSE
## [14,]                            FALSE
## [15,]                             TRUE
##         ADASYN, FALSE, FALSE, classif.randomForest
##  [1,]                                  FALSE
##  [2,]                                  FALSE
##  [3,]                                  FALSE
##  [4,]                                   TRUE
##  [5,]                                  FALSE
##  [6,]                                  FALSE
##  [7,]                                  FALSE
##  [8,]                                  FALSE
##  [9,]                                  FALSE
## [10,]                                   TRUE
## [11,]                                  FALSE
## [12,]                                  FALSE
## [13,]                                  FALSE
## [14,]                                  FALSE
## [15,]                                  FALSE
##         ADASYN, FALSE, FALSE, classif.xgboost
##  [1,]                              TRUE
##  [2,]                             FALSE
##  [3,]                             FALSE
##  [4,]                              TRUE
##  [5,]                             FALSE
##  [6,]                             FALSE
##  [7,]                             FALSE
##  [8,]                             FALSE
##  [9,]                             FALSE
## [10,]                              TRUE
## [11,]                             FALSE
## [12,]                             FALSE
## [13,]                              TRUE
## [14,]                             FALSE
## [15,]                             FALSE
```

8

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

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

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

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

# Plotando os ranks

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

```
##          ADASYN, FALSE, FALSE, classif.ksvm
##                                    9.323171
## ADASYN, FALSE, FALSE, classif.randomForest
##                                    7.274390
##       ADASYN, FALSE, FALSE, classif.xgboost
##                                    6.402439
##           FALSE, FALSE, FALSE, classif.ksvm
##                                    9.957317
##  FALSE, FALSE, FALSE, classif.randomForest
##                                    8.585366
##        FALSE, FALSE, FALSE, classif.xgboost
##                                    8.213415
##            FALSE, FALSE, TRUE, classif.ksvm
##                                    7.579268
##    FALSE, FALSE, TRUE, classif.randomForest
##                                    6.225610
##        FALSE, FALSE, TRUE, classif.xgboost
##                                    6.426829
##            FALSE, TRUE, FALSE, classif.ksvm
##                                   10.195122
##    FALSE, TRUE, FALSE, classif.randomForest
##                                    8.567073
##        FALSE, TRUE, FALSE, classif.xgboost
##                                    8.512195
##           SMOTE, FALSE, FALSE, classif.ksvm
##                                    9.353659
##   SMOTE, FALSE, FALSE, classif.randomForest
##                                    7.274390
##       SMOTE, FALSE, FALSE, classif.xgboost
##                                    6.109756
```

# Plotando grafico de Critical Diference

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