

R Notebook

Parametros:

Measure = Matthews correlation coefficient
Columns = sampling, weight_space, underbagging
Performance = holdout_measure
Filter keys = imba.rate
Filter values = 0.03

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

Tratamento dos dados

Carregando data set compilado

```
ds = read.csv("/home/rodrigo/Dropbox/UNICAMP/IC/estudo_cost_learning/SummaryResults/summary_compilation.  
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      cardiocography-10clases: 900 3rd Qu.:0.0500
## Max. :3      cardiocography-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)){
  ds = filter_at(ds, .vars = params$filter_keys, .vars_predicate = any_vars(. == params$filter_values))
}
```

```
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              : 0  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:2970
##
##
## tuning_measure      holdout_measure      holdout_measure_residual
## Min.      :-0.05673  Min.      :-0.1757  Min.      :-0.4658
## 1st Qu.: 0.33347    1st Qu.: 0.0000  1st Qu.: 0.0391
## Median : 0.83196    Median : 0.5030  Median : 0.2116
## Mean      : 0.66187    Mean      : 0.4753  Mean      : 0.3111
## 3rd Qu.: 0.98596    3rd Qu.: 0.8126  3rd Qu.: 0.5286
## Max.      : 1.00000    Max.      : 1.0000  Max.      : 1.0000
## NA's      :48        NA's      :48        NA's      :48
## 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      :48        (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

```
# 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$performance)))

# 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] 198 5
```

```
# Renomeando a variavel
df = df_tec_wide_residual

summary(df)
```

```
## ADASYN, FALSE, FALSE FALSE, FALSE, FALSE FALSE, FALSE, TRUE
## Min. : -0.05271 Min. : -0.02749 Min. : -0.06218
## 1st Qu.: 0.13467 1st Qu.: 0.05858 1st Qu.: 0.20568
## Median : 0.55465 Median : 0.55050 Median : 0.39109
## Mean : 0.50021 Mean : 0.48068 Mean : 0.42513
## 3rd Qu.: 0.85705 3rd Qu.: 0.84007 3rd Qu.: 0.63304
## Max. : 1.00000 Max. : 1.00000 Max. : 1.00000
## NA's :6 NA's :2 NA's :2
## FALSE, TRUE, FALSE SMOTE, FALSE, FALSE
## Min. : -0.03234 Min. : -0.05605
## 1st Qu.: 0.05397 1st Qu.: 0.14262
## Median : 0.51095 Median : 0.55878
## Mean : 0.47394 Mean : 0.49714
## 3rd Qu.: 0.82366 3rd Qu.: 0.82510
## Max. : 1.00000 Max. : 1.00000
## NA's :2 NA's :4
```

Verificando a média de cada coluna selecionada

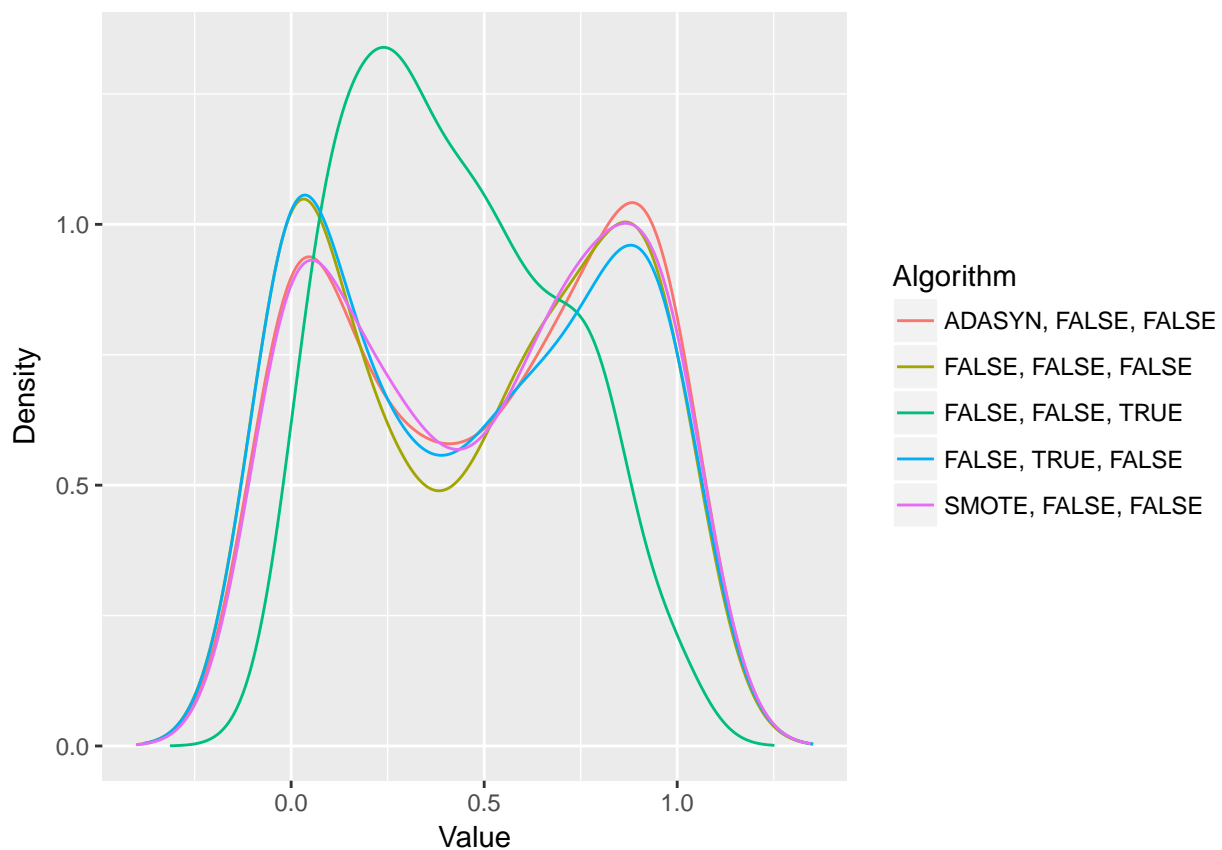
```
for(i in (1:dim(df)[2])){
  #print(df[,i])
}
```

```
print(paste("Media da coluna ", colnames(df)[i], " = ", mean(df[,i], na.rm = TRUE), sep=""))
}
```

```
## [1] "Media da coluna ADASYN, FALSE, FALSE = 0.500213486555043"
## [1] "Media da coluna FALSE, FALSE, FALSE = 0.480682919573954"
## [1] "Media da coluna FALSE, FALSE, TRUE = 0.425127343849434"
## [1] "Media da coluna FALSE, TRUE, FALSE = 0.473944555499608"
## [1] "Media da coluna SMOTE, FALSE, FALSE = 0.497144268022409"
```

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 = 15.155, df = 4, p-value = 0.004391
```

Testando as diferencas par a par

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

```
##      ADASYN, FALSE, FALSE FALSE, FALSE, FALSE FALSE, FALSE, TRUE
## [1,]          FALSE          FALSE          TRUE
## [2,]          FALSE          FALSE          FALSE
## [3,]           TRUE          FALSE          FALSE
## [4,]          FALSE          FALSE          FALSE
## [5,]          FALSE          FALSE          TRUE
##      FALSE, TRUE, FALSE SMOTE, FALSE, FALSE
## [1,]          FALSE          FALSE
## [2,]          FALSE          FALSE
## [3,]          FALSE          TRUE
## [4,]          FALSE          FALSE
## [5,]          FALSE          FALSE
```

Plotando grafico de Critical Difference

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

