

R Notebook

Parametros:

Measure = Matthews correlation coefficient
Columns = learner
Performance = holdout_measure
Filter keys = sampling, weight_space, underbagging
Filter values = FALSE, FALSE, FALSE

```
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)){
  dots = paste0(params$filter_keys," == '",params$filter_values,"'")
  ds = filter_(ds, .dots = dots)
}
```

```
summary(ds)
```

```
##          learner      weight_space
## classif.ksvm      :684  Mode :logical
## classif.randomForest:684 FALSE:2052
## classif.rusboost   : 0  NA's :0
## classif.xgboost    :684
##
##
##
##          measure      sampling      underbagging
## Accuracy           : 0  ADASYN: 0  Mode :logical
## Area under the curve : 0  FALSE :2052 FALSE:2052
## F1 measure          : 0  SMOTE : 0  NA's :0
## G-mean              : 0
## Matthews correlation coefficient:2052
##
##
## tuning_measure      holdout_measure      holdout_measure_residual
## Min. : -0.02717      Min. : -0.05263      Min. : -0.4571
## 1st Qu.: 0.05889      1st Qu.: 0.00000      1st Qu.: 0.0000
## Median : 0.49619      Median : 0.56796      Median : 0.1669
## Mean : 0.47012      Mean : 0.47690      Mean : 0.2775
## 3rd Qu.: 0.81018      3rd Qu.: 0.86046      3rd Qu.: 0.4871
## Max. : 1.00000      Max. : 1.00000      Max. : 1.0000
## NA's :18            NA's :18            NA's :18
## iteration_count      dataset      imba.rate
## Min. :1      abalone      : 36      Min. :0.0010
## 1st Qu.:1      adult      : 36      1st Qu.:0.0100
## Median :2      bank      : 36      Median :0.0300
## Mean :2      car      : 36      Mean :0.0286
## 3rd Qu.:3      cardiocography-10clases: 36      3rd Qu.:0.0500
## Max. :3      cardiocography-3clases : 36      Max. :0.0500
## NA's :18      (Other)      :1836
```

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] 228 3
```

```
# Renomeando a variavel
df = df_tec_wide_residual

head(df)
```

```
##   classif.ksvm classif.randomForest classif.xgboost
## 1 -0.002681154      0.0000000      0.000000000
## 2 -0.002681154      0.0000000      0.000000000
## 3 -0.007549758      0.0000000     -0.003371752
## 4  0.103707362      0.0000000      0.055778009
## 5  0.011520780      0.3603999      0.449276758
## 6  0.011520780           NA      0.449276758
```

```
summary(df)
```

```
##   classif.ksvm      classif.randomForest classif.xgboost
## Min.      : -0.04044   Min.      : -0.02449   Min.      : -0.03192
## 1st Qu.:  0.00000   1st Qu.:  0.08717   1st Qu.:  0.18970
## Median :  0.19993   Median :  0.64741   Median :  0.66518
## Mean    :  0.32962   Mean    :  0.53999   Mean    :  0.56276
## 3rd Qu.:  0.63626   3rd Qu.:  0.87917   3rd Qu.:  0.88876
## Max.    :  1.00000   Max.    :  1.00000   Max.    :  1.00000
## NA's     : 6
```

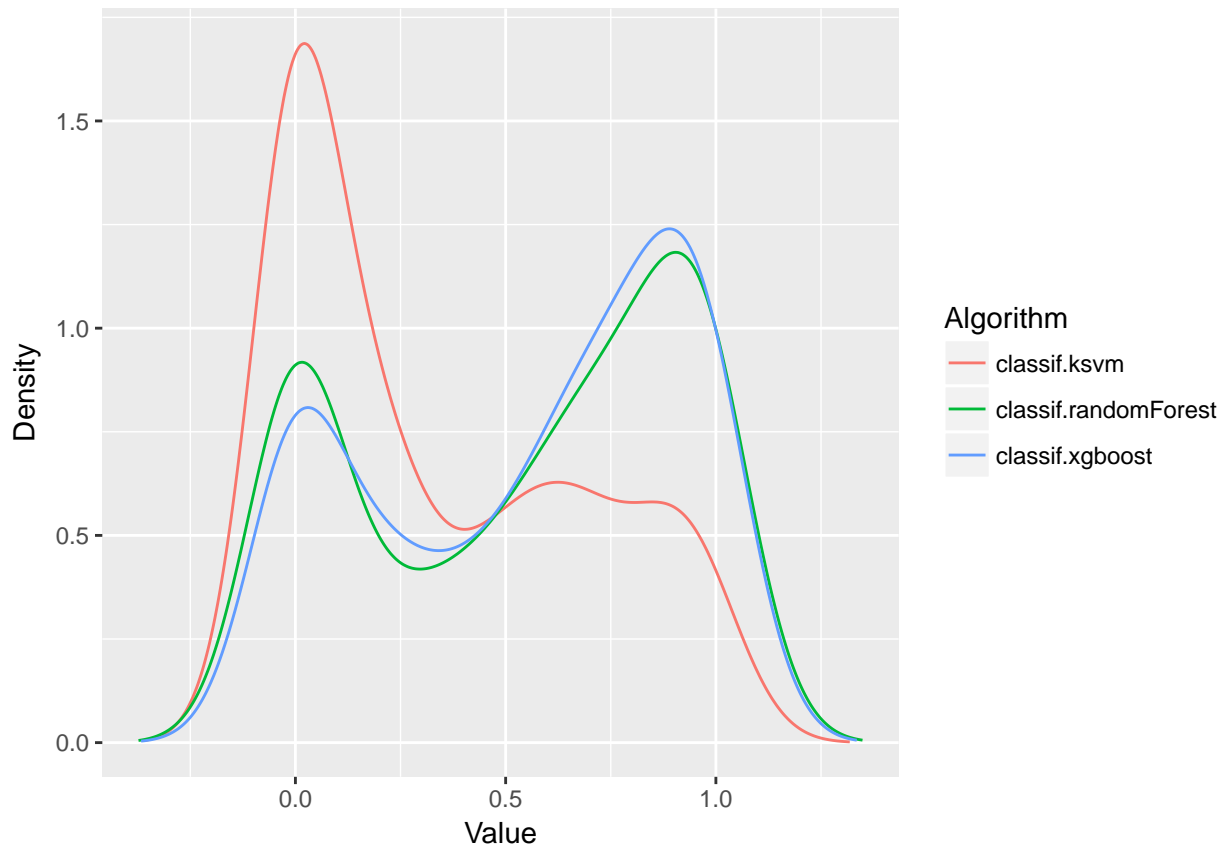
Verificando a média de cada coluna selecionada

```
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 classif.ksvm = 0.329621774060411"
## [1] "Media da coluna classif.randomForest = 0.539987503010754"
## [1] "Media da coluna classif.xgboost = 0.562761901495153"
```

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 = 79.042, df = 2, p-value < 2.2e-16
```

Testando as diferenças par a par

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

```
##      classif.ksvm classif.randomForest classif.xgboost
## [1,]          FALSE                TRUE             TRUE
## [2,]           TRUE                FALSE             TRUE
## [3,]           TRUE                TRUE              FALSE
```

Plotando os ranks

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

```
##      classif.ksvm classif.randomForest      classif.xgboost
##      2.462719      1.881579      1.655702
```

Plotando grafico de Critical Difference

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

