

G4.P-1

David Emanuel Craciunescu Laura Perez Medeiro

November 3, 2019

1 Data analysis using Hunt decision algorithm and linear regression

1.1

The data use in this exercise will be nine students marks and califications, as it has been done in the teorical classes. For this analysis, it wil be used the gain information meause, using Gini as the impurity measure.

First step in the analysis is read the data, which is contains in a txt file called qualifications.txt:

```
> library(utils)
> qualifications <- read.table("qualification.txt")
> sample = data.frame(qualifications)
```

In order to make the analysis, the package rpart will be used. this means, that it should be install before working with the dataset. In order to manage R packages, it will be used Packrat `library(utils) library(rpart)` clasification = `rpart(C.G .,data=sample, method="class", minsplit = 1)` clasification

Another package that can be used to do this analysis is tree:

```
library(tree) (clasificationTree = tree(C.G ., data = sample, mincut = 1,
minsize = 2)) clasificationTree
```

1.2

In this second part, tha dataset use is planets.txt. To this dataset, linear regression will be applied.

As it has been done before, the first step consist on reading data from a *txt* file:

```
> data <- read.table("planets.txt")
> data = data.frame(data)
> names(data)
```

```
[1] "Radius" "Density"
```

In order to quantify the correlation between the variables, it will be calculated the coefficient's matrix correlation:

```
> cor(data)
```

```

      Radius Density
Radius 1.000000 0.371063
Density 0.371063 1.000000

```

Then, it will be calculated and represented the minimum square error line:

```

> regression <- lm( Density~Radius, data)
> summary(regression)

```

Call:

```
lm(formula = Density ~ Radius, data = data)
```

Residuals:

```

Mercurio   Venus   Tierra   Marte
 0.70312 -0.01253  0.24566 -0.93624

```

Coefficients:

```

              Estimate Std. Error t value Pr(>|t|)
(Intercept)   4.3624      1.2050   3.620  0.0685 .
Radius         0.1394      0.2466   0.565  0.6289
---

```

```

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Residual standard error: 0.846 on 2 degrees of freedom

Multiple R-squared: 0.1377, Adjusted R-squared: -0.2935

F-statistic: 0.3193 on 1 and 2 DF, p-value: 0.6289

The equation's line is $y = 4.3624 + 0.1394x$

```

> library(gplots)
> par(mar = rep(2,4))
> plot(data$Density, data$Radius)
> abline(regression)

```

Finally, it is necessary to calculate ANOVA in order to analyze correctly the relation between variables.

```

> anova <- aov(Density~Radius, data)
> summary(anova)

```

```

      Df Sum Sq Mean Sq F value Pr(>F)
Radius    1  0.2286   0.2286   0.319  0.629
Residuals  2  1.4314   0.7157

```

2

The following part consist on doing the same analysis as it has been done before, but now with new datasets. `data <- read.table("vehiculos.txt")` `sample = data.frame(data)` `library(rpart)` `clasification = rpart(TipoVehiculo .,data=sample, method="class", minsplit = 1)` `clasification`