

8/4/19

DATA ANALYTICS

→ Analytics : The science of using data to build models that lead to better decisions that add value to individuals, to companies, to institutions.

→ R \div A software environment for data analysis, statistical computing & graphics.
→ A programming language.

→ Why use R?

1. Nice graphics and visualizations
2. Free — open source project
3. Easy to re-run previous work

→ R console:

1. To get info about a particular function, we can type `?[fn-name]` eg: `?sqrt`.
It results document.
2. To save our value to a variable, we can type `var_name = sqrt(2)` (or) `var_name <- sqrt(2)`
result is stored in `var_name`.
3. Variable naming rules:— should not use spaces, shouldn't start with number.
4. `ls()` → to see list of variable names used in the R session.

→ R only allows 1 data type in 1 vector.

eg:- `country = c("India", "USA", "Brazil")`
`number = c("75", "60")`
 ↑ Combine

If we type `country[1]`, we get India.

→ `seq(0, 100, 2)`
 ↓ ↓ ↓
 starting last increment
 num num value.

→ All of the data in a single object → data.frame.

Eg: `CountryData = data.frame(Country, numberLifeExpectancy)`
 ↓
 To combine two ~~for~~ vectors.

→ To add another field/vector use '\$'

Eg: `CountryData$population = c(100, 200, 500)`

→ To combine two data frames, we use "rbind" function.

`Newfunction_name = rbind(dataframe1_name, dataframe2_name)`

→ `getwd()` - To find the path to the folder in console.

→ `WHO = read.csv("WHO.csv")` - To read csv file

→ `str(WHO)` - To get the structure of file

→ `summary(WHO)` - To get the numerical summary of file

→ `WHO-Europe = subset(WHO, Region == "Europe")`

↑
↑
 Dataset variable to be matched

→ write.csv(who-Europe , "who-Europe.csv")

↑ ↑
dataset file name

- To save the data into the csv file

→ $\text{sum}(\text{WFB} - \text{Europe})$ - To remove the data set

→ WHO & Under 15 - will print the under 15 data set
from WHO file

→ mean() - To find mean

→ sd () - To get the standard deviation

→ which mis () - will return the mis values index

→ which . max () - will return the max values index
or observation

→ plot (WHO d GNI , WHO d Fertility Rate) - B plot graph

\downarrow \downarrow
 x-axis y-axis

→ Outliers = subset (WHO, GNI > 1000 & Fertility Rate > 2.5)

- It will return the set with those specifications

→ now (Outliers) - will give the number of countries in the subset.

→ Outliers [c ("Country", "GDP", "FertilityRate")]

- it will give the table ^{branch} of the subset with those combinations.

→ plot (WHO's GNI, WHO's Fertility Rate) - scatter plot

→ hist (WHO & cellular subscribers) - histogram plot

→ boxplot (WHO's Life Expectancy ~ WHO's Region) - box plot

↑
we can also label the graph using

$xlab = "$	$"$	$xlim = c($
$ylab = "$	$"$	$breaks =$
$main = "$	$"$	

→ table (WHO & Region) - prints the number of observations of each region

→ Japply (who & over go, who & Region, mean)

↑ splits with this ↑ calculates this with 1st parameter

→ to apply (with Literacy Rate, with Region, min, max, min = True)

↑
to overcome the
non entry values
i.e. N/A

→ USDA = read.csv("USDA.csv")

→ `names (USDA)` - will give the names of the variables present in the file

→ match ("CAVIAR", USDA + Descriptions) - will give that particular index from Descriptions

→ HighSodium = USDA \$ Sodium > mean(USDA \$ Sodium ,

na.rm = TRUE)

↑ gives the output as TRUE FALSE

HighSodium = as.numeric()

↑
To get the output in numeric i.e 1 - TRUE
0 - FALSE

9/4/18

* Linear Regression :

Ashenfelter - tested the quality of wine

using a set of independent variables

↑ Age, weather, Harvest Rain
temperature, Winter Rain

⇒ One variable regression model

$$y_i = \beta_0 + \beta_1 x_i + e_i \leftarrow \text{error}$$

↑ dependent variable ↑ intercept co-efficient ↑ regression co-eff ↑ independent variable

↓ (for linear regression model)

- Sum of Squared Errors $SSE = e_1^2 + e_2^2 + \dots + e_n^2$

- Root mean square error $RMSE = \sqrt{\frac{SSE}{n}}$

- Total sum of squares (for base line model) =

↑ line will be || to x-axis

point is found by taking the mean of all data points

$$R^2 = 1 - \frac{SSE}{SST}$$

$$0 \leq SSE \leq SST, 0 \leq SST$$

if $R^2 = 0$ means no improvement over baseline
 $R^2 = 1$ means a perfect predictive model

⇒ Multiple Linear Regression model with K variables

$$y_i = \beta_0 + \beta_1 x_1^i + \beta_2 x_2^i + \dots + \beta_K x_K^i + e^i$$

where y_i = dependent variable

x_j^i = j^{th} independent variable for i^{th} observation

e^i = error terms

wine = read.csv("wine.csv")

→ model 1 = $\text{lm}(\text{Price} \sim \text{AGST}, \text{data} = \text{wine})$
↑
linear model

→ model 1 \$ residuals - will give the error terms

→ SSE = sum(model 1 \$ residuals ^ 2)

→ model 2 = $\text{lm}(\text{Price} \sim \text{AGST} + \text{HarvestRain}, \text{data} = \text{wine})$

* Correlation :

A measure of the linear relationship between variables

if correlation = +1 → perfect positive linear relationship
0 → no linear relationship
-1 → perfect negative linear relationship

→ cor(wine \$ WinterRain, wine \$ Price)

* Predictive ability :

{ training data - data that we already used
test data - new data

wineTest = read.csv("wine-test.csv")

→ ~~predictTest~~

→ model4 = lm(Price ~ AGST + HarvestRains + WinterRains
+ Age, data = wine)

→ predictTest = predict(model4, newdata = wineTest)

- which predicts the model with new data

→ $SSE = \sum \left(\overset{\text{actual value}}{\text{wineTest\$Price}} - \overset{\text{prediction value}}{\text{predictTest}} \right)^2$

$SST = \sum \left(\overset{\text{actual value}}{\text{wineTest\$Price}} - \overset{\text{mean of training set}}{\text{mean(wine\$Price)}} \right)^2$

$R^2 = 1 - SSE/SST$

- `qualityTrain = subset(quality, split == TRUE)`
- `qualityTest = subset(quality, split == FALSE)`

- `QualityLog = glm(PoorCare ~ OfficeVisits + Narcotics,`

↑
general Logistic
Regression model

`data = qualityTrain, family = binomial)`

- `predictTrain = predict(QualityLog, type = "response", newdata = qualityTest)`
- will give the prediction probabilities

- `apply(predictTrain, qualityTrain$PoorCare, mean)`

- mean of prediction of true outcomes

→ Threshold value - to get the binary prediction
• in place of probabilities.

→ Confusion Matrix / Classification Matrix

	Predicted = 0	Predicted = 1
Actual = 0	True Negatives (TN)	False Positives (FP)
Actual = 1	False Negatives (FN)	True Positives (TP)

0 - Goodcare 1 - Poorcare

TN - Actual goodcare predicted goodcare

TP - Actual poorcare predicted poorcare

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

True Positive Rate

$$\text{Specificity} = \frac{TN}{TN + FP}$$

True Negative Rate

- A model with high threshold $\uparrow \rightarrow$ Sensitivity \downarrow
Specificity \uparrow

table (quality T trains & Poor Care, predict T trains > 0.5)

\uparrow
Row = true outcome

\uparrow
Column = Predictions

$$\text{Overall Accuracy} = (TN + TP) / N$$

$$\text{Overall Error Rate} = (FP + FN) / N$$

$$\text{False Negative Error Rate} = FN / (FN + TP)$$

$$\text{False Positive Error Rate} = FP / (FP + TN)$$

$$\text{Base line model accuracy} = (TN + FP) / N$$

- install packages ("ROCR")

$$\text{ROCRPredict} = \text{prediction} \left(\underset{\substack{\downarrow \\ \text{predicted value}}}{\text{predictTest}}, \underset{\substack{\downarrow \\ \text{Actual value}}}{\text{test \& Ten Year CHO}} \right)$$

- as numeric (performance (ROCRpred, "auc") @ y.values)
- will give the auc value