

Practical Machine Learning with R

Lecture 1

Objectives:

- ❖ Introduction to R programming Language for Machine Learning.
- ❖ Installation of R and RStudio. R Studio basics.
- ❖ Installation of packages from CRAN, Bioconductor and Github.
- ❖ R Introduction
- ❖ Fundamental concepts of Machine Learning.

Introduction To R

Origins:

- The R programming language was created by Ross Ihaka and Robert Gentleman at University of Auckland, New Zealand in 1996. R programming is one of the most powerful languages used in the scientific world as well as domain areas related to research and exploratory analysis.
- R programming language was remodeled utilizing the programming paradigms and principles of S, an object oriented language developed by John Chambers at the Bell Labs in 1970. Historically in the 1970s S was implemented for research and scientific endeavors but was not open source application.

Characteristics:

- **Open Source Framework:** R works under the GNU (General Public License) license which is not restrictive in its scope allowing any individual to download and implement it. It has its own community based features and forums.
- **Flexible:** Unlike many proprietary scientific and software systems R is not composed as a rigid architecture that mandates the analysis to be conducted only via a given schema. The analysis framework is not didactic in its implementation wherein the only way to conduct research based tests is through a menu driven graphical interface. Packages can be created by developers for different tasks as well as newer functions can be customized to achieve the objective of the research.
- **Extensible:** R is tailored to allow research not only in the domain areas of statistics but also provides packages and utilities for Bio informatics, Language processing, humanities, econometric, geo spatial arena to name a few. R implements the Object Oriented approach following in the footsteps of its parent programming language S.
- **Comprehensive:** R incorporates the full spectrum of computational and visualization frameworks that allow for holistic and novel research and analysis. It provides plethora of diverse tools and techniques to implement the full data science or machine learning project lifecycle from inception to completion. This exhaustive lifecycle entails importing

Cross Platform Compatible: The R programming language can be implemented seamlessly across platforms ranging from Windows OS, Mac OS, UNIX, Solaris to name a few. R programming has been used to create many industry based applications and companies like Google, Oracle, SAP are leveraging the power of this language. It can also be implemented in conjunction with other languages like Python.

Steps:

- # RStudio

The screenshot displays the RStudio environment with the following components:

- Source Editor:** Contains R code for data processing and visualization. Key lines include:


```

      34 library(readr)
      35 library(dplyr)
      36 library(plyr)
      37 library(ggplot2)
      38 library(stringr)
      39 library(cluster)
      40 #library(fpc)
      41 library(200)
      42 library(pca3d)
      43 library(factoextra)
      44 #library(FactoMineR)
      45 library(stats)
      46 library(tidyverse)
      47 #library(MASS)
      48 #library(vegan)
      49
      50 cytokinesPlasma <- read_excel("CytokinesClusterAnalysis.xlsx", sheet = "plasma")
      51 cytokinesHDL <- read_excel("CytokinesClusterAnalysis.xlsx", sheet = "hdl")
      52 cytokinesHDLPlasma <- bind_rows(cytokinesPlasma, cytokinesHDL)
      53 # cytokinesHDLPlasma <- read_excel("CytokinesHDLPlasma.xlsx", sheet = "hdl")
      54 cytokines_PCA <- prcomp(cytokinesHDLPlasma[, 4:60])
      55 plot(cytokines_PCA, type = "t")
      56 autoplot(cytokines_PCA, data = cytokinesHDLPlasma, colour =
      57 "Sample_Type_Condition", frame = TRUE, frame.colour = "Sample_Type_Condition", frame.type = "norm")
      58 autoplot(cytokines_PCA, data = cytokinesHDLPlasma, colour =
      59 "Sex", frame = TRUE, frame.colour = "Sex", frame.type = "norm")
      60 autoplot(cytokines_PCA, data = cytokinesHDLPlasma, colour =
      61 "Sample_Type", frame = TRUE, frame.colour = "Sample_Type", frame.type = "norm")
      62
      63 # Classical MDS
      64 # N rows (objects) x p columns (variables)
      65 # each row identified by a unique row name
      66
      67 # d <- dist(cytokinesHDLPlasma[,40:]) # euclidean distances between the rows
      68
      69 # plot(d, main = "MDS plot of cytokinesHDLPlasma data",
      70 #      xlab = "MDS1", ylab = "MDS2",
      71 #      col = "black", pch = 1, las = 1,
      72 #      xlim = c(-10, 10), ylim = c(-10, 10),
      73 #      main = "MDS plot of cytokinesHDLPlasma data")
      74
      75 # chunk 2
      76
      77 # reached getOption("max.print") -- omitted 187 rows
      78
      79 attr(,"scaled:center")
      80
      81 # MDS1
      82 # MDS2
      83 # MDS3
      84 # MDS4
      85 # MDS5
      86 # MDS6
      87 # MDS7
      88 # MDS8
      89 # MDS9
      90 # MDS10
      91 # MDS11
      92 # MDS12
      93 # MDS13
      94 # MDS14
      95 # MDS15
      96 # MDS16
      97 # MDS17
      98 # MDS18
      99 # MDS19
      100 # MDS20
      101 # MDS21
      102 # MDS22
      103 # MDS23
      104 # MDS24
      105 # MDS25
      106 # MDS26
      107 # MDS27
      108 # MDS28
      109 # MDS29
      110 # MDS30
      111 # MDS31
      112 # MDS32
      113 # MDS33
      114 # MDS34
      115 # MDS35
      116 # MDS36
      117 # MDS37
      118 # MDS38
      119 # MDS39
      120 # MDS40
      121 # MDS41
      122 # MDS42
      123 # MDS43
      124 # MDS44
      125 # MDS45
      126 # MDS46
      127 # MDS47
      128 # MDS48
      129 # MDS49
      130 # MDS50
      131 # MDS51
      132 # MDS52
      133 # MDS53
      134 # MDS54
      135 # MDS55
      136 # MDS56
      137 # MDS57
      138 # MDS58
      139 # MDS59
      140 # MDS60
      141 # MDS61
      142 # MDS62
      143 # MDS63
      144 # MDS64
      145 # MDS65
      146 # MDS66
      147 # MDS67
      148 # MDS68
      149 # MDS69
      150 # MDS70
      151 # MDS71
      152 # MDS72
      153 # MDS73
      154 # MDS74
      155 # MDS75
      156 # MDS76
      157 # MDS77
      158 # MDS78
      159 # MDS79
      160 # MDS80
      161 # MDS81
      162 # MDS82
      163 # MDS83
      164 # MDS84
      165 # MDS85
      166 # MDS86
      167 # MDS87
      168 # MDS88
      169 # MDS89
      170 # MDS90
      171 # MDS91
      172 # MDS92
      173 # MDS93
      174 # MDS94
      175 # MDS95
      176 # MDS96
      177 # MDS97
      178 # MDS98
      179 # MDS99
      180 # MDS100
      181 # MDS101
      182 # MDS102
      183 # MDS103
      184 # MDS104
      185 # MDS105
      186 # MDS106
      187 # MDS107
      188 # MDS108
      189 # MDS109
      190 # MDS110
      191 # MDS111
      192 # MDS112
      193 # MDS113
      194 # MDS114
      195 # MDS115
      196 # MDS116
      197 # MDS117
      198 # MDS118
      199 # MDS119
      200 # MDS120
      201 # MDS121
      202 # MDS122
      203 # MDS123
      204 # MDS124
      205 # MDS125
      206 # MDS126
      207 # MDS127
      208 # MDS128
      209 # MDS129
      210 # MDS130
      211 # MDS131
      212 # MDS132
      213 # MDS133
      214 # MDS134
      215 # MDS135
      216 # MDS136
      217 # MDS137
      218 # MDS138
      219 # MDS139
      220 # MDS140
      221 # MDS141
      222 # MDS142
      223 # MDS143
      224 # MDS144
      225 # MDS145
      226 # MDS146
      227 # MDS147
      228 # MDS148
      229 # MDS149
      230 # MDS150
      231 # MDS151
      232 # MDS152
      233 # MDS153
      234 # MDS154
      235 # MDS155
      236 # MDS156
      237 # MDS157
      238 # MDS158
      239 # MDS159
      240 # MDS160
      241 # MDS161
      242 # MDS162
      243 # MDS163
      244 # MDS164
      245 # MDS165
      246 # MDS166
      247 # MDS167
      248 # MDS168
      249 # MDS169
      250 # MDS170
      251 # MDS171
      252 # MDS172
      253 # MDS173
      254 # MDS174
      255 # MDS175
      256 # MDS176
      257 # MDS177
      258 # MDS178
      259 # MDS179
      260 # MDS180
      261 # MDS181
      262 # MDS182
      263 # MDS183
      264 # MDS184
      265 # MDS185
      266 # MDS186
      267 # MDS187
      268 # MDS188
      269 # MDS189
      270 # MDS190
      271 # MDS191
      272 # MDS192
      273 # MDS193
      274 # MDS194
      275 # MDS195
      276 # MDS196
      277 # MDS197
      278 # MDS198
      279 # MDS199
      280 # MDS200
      281 # MDS201
      282 # MDS202
      283 # MDS203
      284 # MDS204
      285 # MDS205
      286 # MDS206
      287 # MDS207
      288 # MDS208
      289 # MDS209
      290 # MDS210
      291 # MDS211
      292 # MDS212
      293 # MDS213
      294 # MDS214
      295 # MDS215
      296 # MDS216
      297 # MDS217
      298 # MDS218
      299 # MDS219
      300 # MDS220
      301 # MDS221
      302 # MDS222
      303 # MDS223
      304 # MDS224
      305 # MDS225
      306 # MDS226
      307 # MDS227
      308 # MDS228
      309 # MDS229
      310 # MDS230
      311 # MDS231
      312 # MDS232
      313 # MDS233
      314 # MDS234
      315 # MDS235
      316 # MDS236
      317 # MDS237
      318 # MDS238
      319 # MDS239
      320 # MDS240
      321 # MDS241
      322 # MDS242
      323 # MDS243
      324 # MDS244
      325 # MDS245
      326 # MDS246
      327 # MDS247
      328 # MDS248
      329 # MDS249
      330 # MDS250
      331 # MDS251
      332 # MDS252
      333 # MDS253
      334 # MDS254
      335 # MDS255
      336 # MDS256
      337 # MDS257
      338 # MDS258
      339 # MDS259
      340 # MDS260
      341 # MDS261
      342 # MDS262
      343 # MDS263
      344 # MDS264
      345 # MDS265
      346 # MDS266
      347 # MDS267
      348 # MDS268
      349 # MDS269
      350 # MDS270
      351 # MDS271
      352 # MDS272
      353 # MDS273
      354 # MDS274
      355 # MDS275
      356 # MDS276
      357 # MDS277
      358 # MDS278
      359 # MDS279
      360 # MDS280
      361 # MDS281
      362 # MDS2
```

The top left window is the Script editor that contains the script of the code. It also incorporates features that help save the files using different formats, syntax highlighting , commenting as well as other enhanced features.

The top right window is the Workspace environment window that displays and saves detailed information related to the contextual data frames, vectors as well as the user defined objects.

Console window is at the lower left window. It serves as the space where the output of the code is displayed. A developer can also run a line code scripts.

The Misc displays is the lower right window The tabs provides with the information related to the files in the working directory, plots generated by the code, packages downloaded and installed as well as the help option .For a function help we can use the command `>help(function)`.

R Notebook framework is very useful since it has several chunks that can be combined to create a journal, book etc where diverse elements like images, code outputs, mark up as well as latex for formula can be incorporated.

Packages

Packages are bundles which encapsulate code, documentation, tests as well as data. Packages are created by R developers and uploaded on the CRAN, Bioconductor as well as Github. These packages can be downloaded by any individual from these website . As of now there are more than 2000 packages.

Installation of packages from CRAN, Bioconductor and Github.

CRAN:

➤ `install.packages ("packagename")`

BICONDUCTOR:

The command to get to the Bi conductor website is as follows:

➤ `source("http://bioconductor.org/biocLite.R")`

The command to install standard packages:

➤ `biocLite()`

The command to install specific packages:

➤ `biocLite("packagename")`

GitHub:

To install Github devtools package is used. The command is as follows:

➤ `install.packages("devtools")`

➤

The command to install a specific package is as follows:

- `devtools::install_github("username/packageName")`

Package Loading

The command to load a package is as follows:

- `library(packagename)`

To obtain help for a package

- `help(package="packagename")`

R INTRODUCTION

Every programming language has different Syntax (command grammar) and Semantics (Interpretation of the command).

- R is an object-oriented language. Everything in R are objects.
- The main data structure in R is a vector. `a<-4` is a 1*1 vector.
`v<-c(2,4,5,6)` is a 1*4 vector
This is assigning the values to the vector v

← is an assigning operator
- All the commands are functions for example `help()` is a function
- Some commands can have multiple outputs depending on the input arguments.

Main Data Structures

1) Vectors

Stores numerous elements. They can comprise of numbers, text , logical variables etc.

Vector data structure cannot simultaneously incorporate more than one data type . For instance , a vector cannot contain elements that are integers as well as text , both at the same time.

- `v<-c("Hello","Hi","Welcome")`
- `v`

`"Hello","Hi","Welcome"`

2) List

List stores elements that might be either of the same type or of different types.

`Record<-list(name="mary",gpa=4,major="History")`

To access an individuals record we have to provide the index.

3) Matrices and Arrays

Matrix is a two-dimensional framework with specific elements at the given indices. These elements in a matrix can be of one type only. Usually matrices are used for numerical computation therefore the stored data type is typically numerical.

```
Student_grades<-matrix(100,200,150,180,ncol=2)
```

Student_grades

```
      [,1] [,2]
[1,]  100  150
[2,]  200  180
```

4) Arrays

Are structures like matrices where the structure can be multidimensional.

5) Data Frames

Data Frames are the most paramount data structure within the R programming environment as it combines the functionality of lists and vectors. A data frame is analogous to a spreadsheet where it can consist of elements that are of different data types.

➤ `Student_frame<-data.frame(student_name,
student_grade,studentgpa,stringsAsFactors=FALSE)`

If we do not set the `stringsAsFactors=FALSE` then R converts every character vector to a factor.

➤ `Student_frame`

Student Name	Student grade	studentgpa

Data frames are very advantageous since they are displayed as a 2 dimensional matrix which can showcase disparate data types at the same time.

The columns of the data frames are defined as attributes or features and the rows are defined as examples.

Factors

Data can be either numerical or categorical. Categorical data can be further differentiated into nominal and ordinal type. Features of a data set are Nominal if there is no natural ordering whereas features of a data set are ordinal if there is a natural ordering in the dataset.

If data inherently is nominal or ordinal data then R uses a special data type called factors that saves the data with categories. Hypothetically, if we have a dataframe that has a column which defines the gender type Male and Female then R will create two categories. Since there are recurring records pertaining to males or females, the category label for male(male:label=1) and female(female:label=2) is stored once and subsequently this is stored as numbers 1,2,1,2 to identify every example/row. The R framework will automatically map 1 to Male and 2 to Female. This strategy not only saves the memory but also informs the reader that the data type has categories.

```
Gender<-factor (c("Male"."Female","Male"))
```

```
Gender
```

```
Male Female Male
```

```
Levels: Male Female
```

Useful Inbuilt Functions in R for exploring data features

- View Data Structures in Memory
 - ls()
- Remove data structure from memory
 - rm(Student_frame)
 - rm(list=ls())
- Reading Data from a file into a data frame
 - Student_record ← read.csv("student.csv", stringsAsFactors = FALSE)
- Visualizing seeing the structure of a data frame
 - str(Student_record)
- Summarizing the data in the data frame
 - summary(Student_record)
 - summary(Student_record\$gpa)
 - summary(Student_record[c("gpa","major")])
- mean
 - mean(Student_record\$gpa)
- median

- `median(Student_record$gpa)`
- range
 - `range(Student_record$gpa)`
- IQR
 - `IQR(Student_record$gpa)`
- Boxplot
 - `boxplot(Student_record$gpa, main="Boxplot of Used Student gpa", ylab = "GPA")`
- Histogram
 - `hist(Student_record$gpa, main="Boxplot of Used Student gpa", ylab = "GPA")`
- Variance, Sd
 - `var(Student_record$gpa)`
 - `sd(Student_record$gpa)`
- One way Table
 - `table(Student_record$major)`
- Proportions
 - `major_table ← table(Student_record$major)`
 - `prop.table(major_table)`
- Rounding
 - `round(major_table, digits = 1)`
- Scatterplot
 - `plot(x = Student_record$gpa, y = Student_record$grade,`

```
main = "Scatterplot of gpa and grade",  
xlab = "grade",  
ylab = "gpa")
```

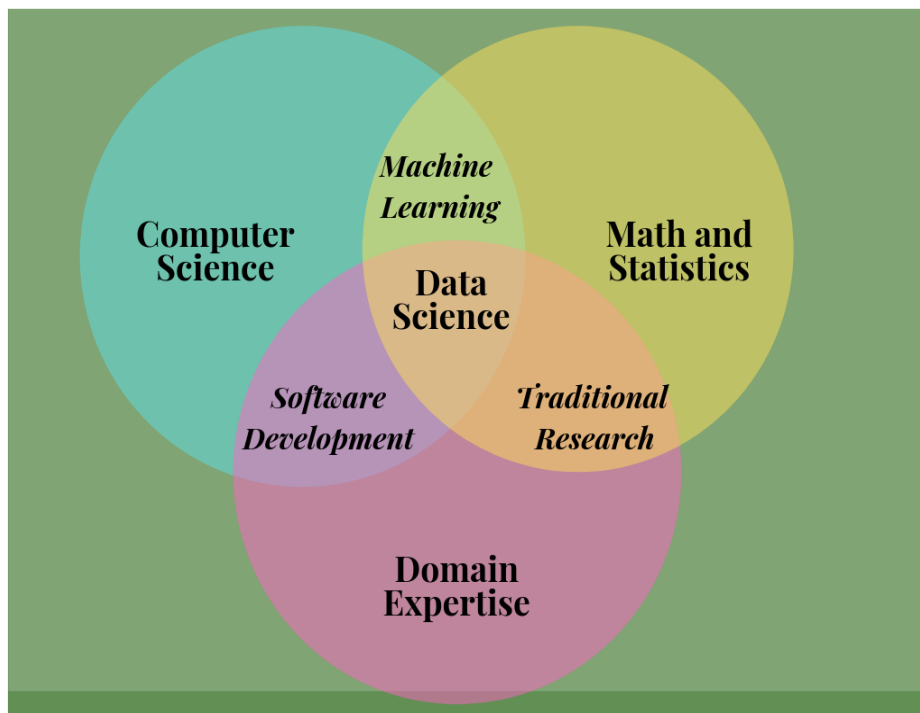
- Crosstabulation table

➤ `CrossTable(x = Student_record$major, y = Student_record $full_time)`

Foundation of Data Science --- An Introduction:

- In contemporary times the digital progression has caused a proliferation of data and data is electronically available across a whole spectrum media sources pertaining to diverse domain areas.
- Complex and large volume of data available must be translated into substantive inferences and actionable insights for it to be of any value. Data Science serves the explicit purpose of translating data into knowledge.
- Data Scientist utilizes numerous techniques from the field of Statistics, Mathematics, Computer Science, Machine Learning and Artificial Intelligence to convert information into knowledge.
- Data Science is a cross disciplinary set of skills that are applied to solve complex problems across diverse disciplines.

Visualization of Data Science domain specifications:



Fundamental Knowledge and Taxanomical Frameworks of Machine Learning

Machine learning is a repertoire of diverse algorithms that enable gaining substantive insights from the contextual data being analyzed which in turn can be leveraged for the decision making in manifold domain areas.

Machine learning is prevalently used in the field of Data Analytics which has proliferated the research sectors of society due to the advancement of information and technology which caused the generation of huge recorded data sets with disparate internal structures.

The big data which is characterized by attributes of volume, variety and velocity in conjunction with high computing power of the computer systems and the statistical models facilitated the advent and ubiquitous application of Machine learning algorithms.

Machine learning algorithm teaches the computer how to use the data to solve a problem.

The computational power of computers for implementation of Machine Learning can only be harnessed and utilized for Real life application by human intervention and direction. The process can only be initiated by humans as well as progressive incremental actions are at the discretion of humans. The algorithms might not be able to extrapolate and does not have the intuition of the human brain. In addition to the fact that it is not able to fully work like the brain of a human being, it also is not competent to always create and translate grammatically correct constructs.

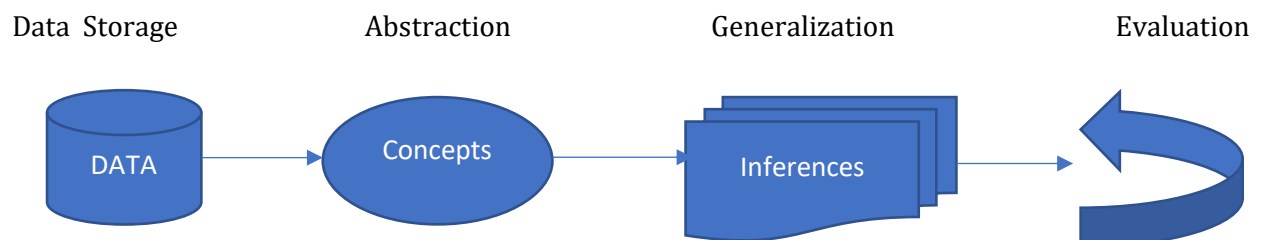
The Machine learning Algorithms is complimentary to the domain knowledge expert therefore ameliorated the decision making in a given scenario. A machine learning algorithm will help a medical doctor visually and numerically perceive the patterns in the biomarkers distributions consequent to which the doctor can make informed decisions related to his/her patient set.

Some important applications of machine learning include fraud detection, disease pathway exploration, weather prediction, auto pilot frameworks, crime pattern assessment etc

Machine Learning Process:

The machine learning process primarily comprises of four germane and interrelated components. These can be elucidated as follows:

- 1) Data Storage
- 2) Abstraction
- 3) Generalization
- 4) Evaluation



Data Storage:

The first component of the learning process is the contextual data in various storage devices and in various formats. Some of this data might be fully structured like a data base file from the Oracle company or it might be a unstructured data from a sequence of click stream.

Abstraction

The second component of the learning process is the conversion of the raw data into explicit knowledge representation and insightful patterns by mapping the data to a specific model . The model categories are as follows:

- Mathematical equations
- Relational diagrams
- Logical if/else
- Grouping into clusters

The conversion process entails fitting the model to the data set which is defined as Training. Subsequent to the training the data represents and summarizes the original data in a more abstract manner. An illustration to clarify this step is the generation of a linear regression equation from a data set composed of two correlated columns.

Generalization

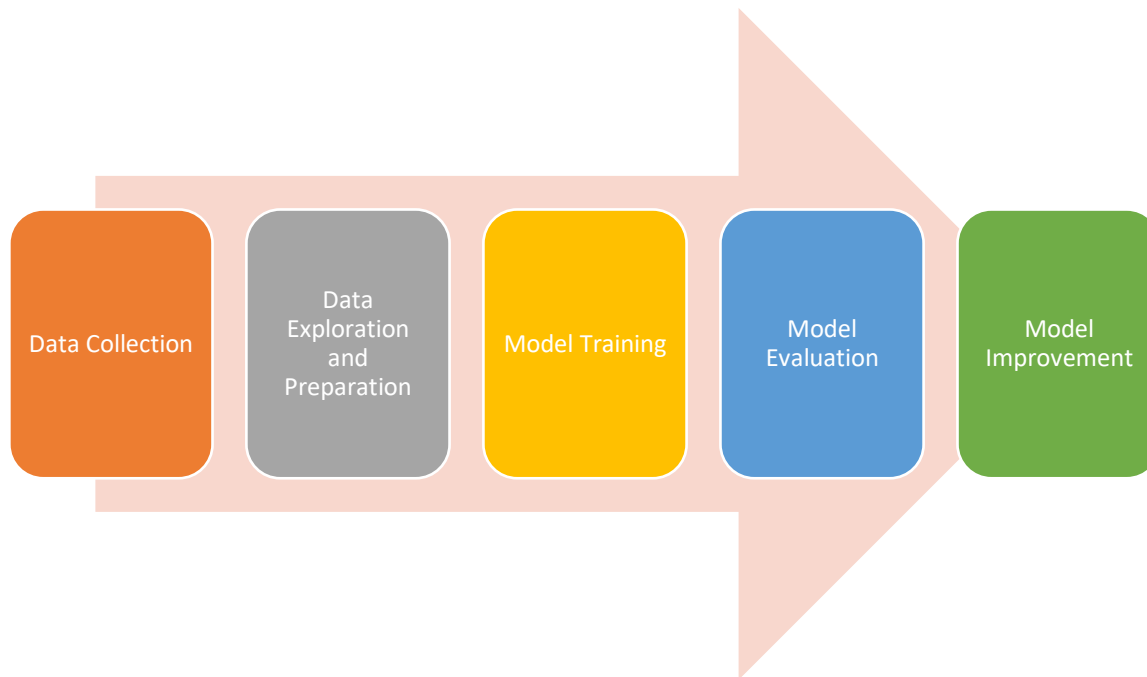
The third component of the learning process is Generalization which entails targeted identification of patterns that align with the objective of the study or are significant in comparison to all possible patterns to effectuate the decision process. After the algorithm completes its run it heuristically (short cut mechanisms) might conjecture related to the most useful inferences. If the heuristics applied are errorneous then it is defined as bias. For instance if a human is identified as standing on two feet then a toddler will not be identified as a human. Bias needs to be minimized in order to be able to make substantive and credible decisions.

Evaluation

Is the fourth component in the learning process which ascertains the success of the machine learning algorithm to inform decision making in the face of the varying levels of bias observed inherently in all of them. Subsequent to the training of the data ,t he model created by the generalization process is tested on a new data set defined as the testing data. This will determine how well the characterization of the training data generalizes to the testing data. The generalization mapping can never be a perfect mapping but data scientist endeavor to obtain the best fitted model.

The generalization is not perfect due to the explained and unexplainable variations in the data. The reasons for these variations might include errorneous measurements, missing data or complexity of the phenomenon. It is important to identify the underlying pattern without taking cognizance of the underlying noise because if we model the noise defined as overfitting , the generalization to the test data will inaccurate. If a model works well during training but categorizes imperfectly during testing is due to the issue of overfitting.

Lifecycle of Machine learning Framework



Taxonomy of Machine Learning

The diverse ML algorithms can be disaggregated at the topmost level into two categories namely Supervised Learning Algorithms and Unsupervised Learning Algorithms.

Supervised Learning Algorithms : This category of algorithms is used to either numerically predict a target feature by mapping its relationship with the other features or generating a categorical classification of the target feature in context of the other features. This classification generates classes with levels.

Supervised Algorithms are provided explicit instructions as to what it needs to learn and how the algorithm will learn. The algorithm optimizes the model to find the best combination of features to obtain the target feature.

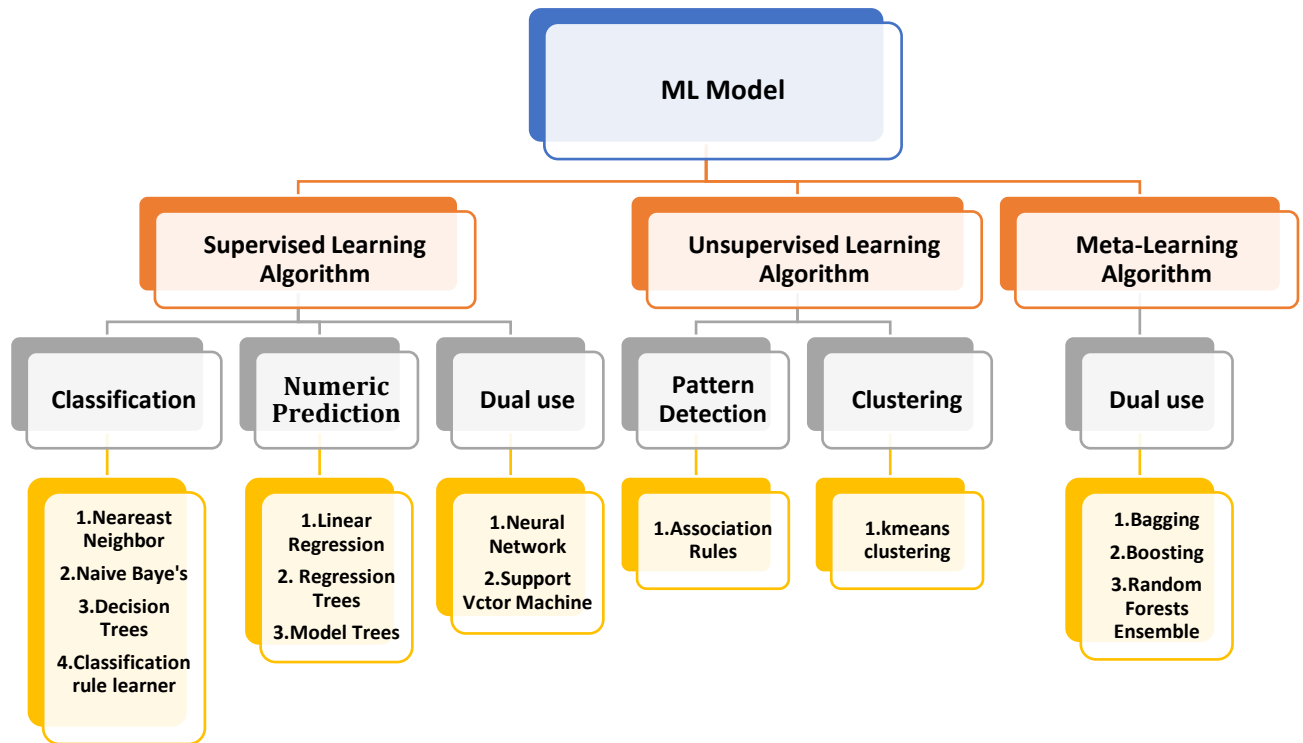
Examples prediction of risk of cancer patients, spam email prediction.

Unsupervised Learning Algorithms

Is the identification of the inherent pattern in a data set. The pattern discovery for the unsupervised learning does not entail having a target feature , on the contrary it is training a descriptive model with no feature more important than the other. For instances if we wish to ascertain the shopping patterns of individuals people who buy soda might also buy chips. Gamblers , thieves might have the same habits. The differentiation of descriptive model into homogeneous groups creates clusters which have to be interpreted and analyzed by data scientists. This segmentation procedure is called Clustering.

Meta Learners

These are algorithms that learning how to learn more effectively. This optimizes the final prediction.



The basic strategy of implementing a machine learning algorithm is identifying from amongst four learning tasks namely Classification, Numeric prediction, Pattern recognition and clustering. The more nuanced and detailed oriented review will be performed in subsequent lectures.

Class Work

- 1) Install R , Install R Studio
- 2) Download the data from the book
- 3) Create a github account
- 4) Choose data set from Kaggle, create a project group

