

dplyr: powerful collection of R packages that are data tools for transforming and visualizing data Pipe Operator: %>%: together functions consectively

Functions (use logical sense to determine sequential): iris %> group by (Species) %>%

summarize (medianSL=median(Sepal.Length), maxSL=max(Sepal.Length))

filter() allows you to select a subset of rows in a data frame

arrange() sorts dataset in ascending or descending order mutate() allows you to update or create new columns summarize() turn many observations into single data point group by() allows you to summarize within groups distinct(),n distinct() Subset distinct/unique rows na if() Convert values to NA

### ggplot: implements grammar of graphics, can use it to visualize data

top\_n() If n is positive, chose top rows, otherwise if negative



# aes() Construct aesthetic mappings

geom\_abline(), geom\_hline() Ref lines: horizontal, vertical! labs(), xlab(), ylab() Modify axis, legend, and plot labels lims(), xlim(), ylim() Set scale limits theme() Modify theme such as grey, light, dark

by year <- gapminder %>% group by(year) %>% summarize (medianGdpPerCap=median(gdpPercap) ggplot(by\_year, aes(x=year, y=medianGdpPerCap))+ geom line()+ expand limits (y=0)

#### Date and Time in R: Abbreviations

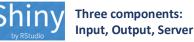
%d	Day of month (decimal)	%a	Weekday (abbv.)
%m	Month (decimal)	%A	Full weekday
%b	Month (abbv.)		Decimal
%B	Month (full name)	%w	Weekday
%у	Year (2 digit)		(0=Sunday)
%Y	Year (4 digit)	%j	Day of year
	Decimal week of the		(decimal)
%W	year (starting on	%W	Week of year
	Monday)		(starting Sun)

#### Transform date: Example code

> format(as.Date(x), "%Y/%m/%d")

**End-result format** 

- > strptime(as.character(X\$date), "%d/%m/%Y")<sup>Time-zone</sup>
- Difference of times: > difftime(time1, time2, tz, units
- = c("auto", "secs", "mins", "hours", "days", "weeks")) Interval: > interval(start = NULL, end = NULL, tzone
- = tz(start))
- Print Current Date/Time: Sys.Date(), Sys.time() Lubridate function: library(lubridate)
- > ymd hms("2012-12-31 23:59:59")
- ## [1] "2012-12-31 23:59:59 UTC"
- > ldate <- mdy\_hms("12/31/2012 23:59:59")
- > Idate + seconds(1)
- ## [1] "2013-01-01 UTC"



> library(shiny) ## Template

> ui <- fluidPage(

numericInput(inputId = "n", "Sample size", value = 25),

plotOutput(outputId = "hist") )

> server <- function(input, output) { output\$hist <- renderPlot({

hist(rnorm(input\$n)) }) }

> shinyApp(ui = ui, server = server)

submitButtion()	dataTableOutput()
checkboxGroupInput()	htmlOutput()
dateRangeinput()	imageOutput()
fileInput()	tableOutput()
numericInput()	plotOutput()
radioButtons()	uiOutput()
sliderInput()	verbatimTextOutput()

textInput()

Association Pattern Mining Rules in R: arules, eclat # Read the data with read transactions command

> txn <- read.transactions(file="baskets rdb.csv", rm.duplicates=F, format="single",sep=",")

# Look at the transactions data > inspect(txn) # Numerical and Graphical Analyze

> itemFrequency(txn); itemFrequencyPlot(txn)

# Running the Algorithm; finding patterns

> pattern <- apriori(txn, parameter = list(sup = 0.5))

> rules <- eclat(data = dataset, parameter = list (support = 0.003, minlen = 2))

> inspect(sort(rules, by = 'support')[1:10])

Support: fraction of transactions that contain both X and Y  $supp(X \rightarrow Y) = (\sigma(X \cup Y))/(|T|)$ Confidence: fraction of Y in transaction that contain X

textOutput()

 $P(A \cap B) = P(A) \times P(B) \Rightarrow$  statistical indep.  $P(A \cap B) > P(A) \times P(B) \Rightarrow +ve \text{ correlated}$ 

 $conf(X \rightarrow Y) = (\sigma(X \cup Y))/(\sigma(X))$ 

 $P(A \cap B) < P(A) \times P(B) \Rightarrow -ve correlated$ Lift = (P(Y|X))/(P(Y))

Interest = (P(X,Y)) / (P(X)P(Y))Apriori Principle: If an itemset is freq,

X <- data[,3:ncol(data)];X <- scale(X,center=F,scale=F)

lambdas2 <- SVD results\$d # singular values of x

pmatrix <- scale(iris[,-5]) #scale variables

table(cboot.hclust\$result\$partition,iris[,5])

cboot.hclust\$bootmean # clusterwise means

cboot.hclust <- clusterboot(pmatrix,

U <- SVD results\$u # columns contain the left singular

V <- SVD\_results\$v # columns contain the right singular

clustermethod=hclustCBI,method="ward.D", k=3)

then all of its subsets must also be freq.

For m data points, there are  $\frac{m(m-1)/2}{2}$  pairs.

- $f_{00}$  = number of pairs that are from <u>different</u> classes and clustered into <u>different</u> groups
- $f_{01}$  = number of pairs that are from <u>different</u> classes and clustered into the <u>same</u> group •  $f_{10}$  = number of pairs that are from the <u>same</u> class and clustered into <u>different</u> groups
- $f_{11}$  = number of pairs that are from the <u>same</u> class and clustered into the <u>same</u> group

$${\rm Rand\ Statistic}\ =\ \frac{f_{00}+f_{11}}{f_{00}+f_{01}+f_{10}+f_{11}}\ \ {\rm Jaccard\ Coefficient}\ =\ \frac{f_{11}}{f_{01}+f_{10}+f_{11}}$$

graph based view for cohesion and proximity, same formula for prototype-based view  $cohesion(C_i) = \sum_{\boldsymbol{x} \in C_i, \boldsymbol{y} \in C_j} proximity(\boldsymbol{x}, \boldsymbol{y}) \quad separation(C_i, C_j) = \sum_{\boldsymbol{x} \in C_i, \boldsymbol{y} \in C_j} proximity(\boldsymbol{x}, \boldsymbol{y})$ 

**SVD Example Code** 

SVD\_results <- svd(X)

names(SVD results)

vectors (n by 1) of x

vectors (p by 1) of x

Y2 = as.matrix(X) %\*% V

Clusterboost Example Code

summary(cboot.hclust\$result)

cboot.hclust\$result\$partition

#### **PCA Example Code**

X <- data[,3:ncol(data)];X <- scale(X,center=F,scale=F) PCA results <- princomp(X)

names(PCA results) ## Names that are available

V <- PCA results\$loadings Y <- PCA\_results\$scores #transformed-each col is PC

# Check the variances

s1 <- apply(X,2,var) # variances of each column in X s2 <- apply(Y,2,var) # variances of each column in Y

### **Hierarchical Clustering**

data(iris) # load internal sample data mydata <- iris[,1:4] # select variables of interest mydata <- scale(mydata) # standardize variables d <- dist(mydata, method = "euclidean") #matrix fit <- hclust(d, method="ward.D2") # display dendogram groups2 <- cutree(fit, k=3)

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$Recall = \frac{tp}{tp + fn}$$

 $F = 2 \cdot \frac{\operatorname{precision} \cdot \operatorname{recall}}{}$ precision + recall An ROC curve (receiver operating characteristic curve) is graph showing the performance of a classification model at all classification thresholds. This curve plots two parameters: TPR

 $Accuracy = \frac{1}{tp + tn + fp + fn}$ 

# Confusion Matrix

# transformed

\* library(fpc)

	Actually	Actually
	Positive (1)	Negative (0)
Predicted Positive (1)	True Positives (TPs)	False Positives (FPs)
Predicted Negative (0)	False Negatives (FNs)	True Negatives (TNs)