# Project 2 Part 3

## Kaushal Khatiwada

## 2024-03-18

## Word Cloud

```r
library(pdftools)
```

```
## Using poppler version 23.08.0
```

```r
library(tm)
```

```
## Loading required package: NLP
```

```r
library(wordcloud)
```

```
## Loading required package: RColorBrewer
```

```r
library(RColorBrewer)
```

lapply function is used to apply pdf_text function to all the element of list "pdfs"

```r
setwd("~/MDS503P2")
pdfs <- list.files(pattern="*.pdf")
alltext <- lapply(pdfs, pdf_text)
```

```
## PDF error: Expected the optional content group list, but wasn't able to find it, or it isn't an Array
```

```r
alltext <- unlist(alltext)
head(alltext,n=4)
```

```
## [1] "Big Data,\nMining, and\nAnalytics\n                    Components\n              of Strategic\n
## [2] " Big Data,\nMining, and\n Analytics\n        Components of\n Strategic Decision Making\n"
## [3] ""
## [4] " Big Data,\nMining, and\n Analytics\n        Components of\n Strategic Decision Making\n\n\n\n\n
```

create a single Corpus object containing all the text data, rather than separate Corpus objects for each PDF.

```
myCorpus <- Corpus(VectorSource(alltext))

myCorpus <- tm_map(myCorpus, content_transformer(tolower))
myCorpus <- tm_map(myCorpus,removePunctuation)
myCorpus <- tm_map(myCorpus,removeNumbers)
removeURL <- function(X) gsub("http[^[:space:]]*","",X)
myCorpus <- tm_map(myCorpus, content_transformer(function(x) gsub("\\\n", " ", x)))
myCorpus <- tm_map(myCorpus,removeURL)
myCorpus <- tm_map(myCorpus,removeWords,stopwords("english"))
myCorpus <- tm_map(myCorpus,removeWords,c("can","may","eg","ie","h","b","p","k","g","q","set","used"))
myCorpus <- tm_map(myCorpus, content_transformer(function(x) gsub("\n*","", x)))
myCorpus <- tm_map(myCorpus,stripWhitespace)

inspect(myCorpus[1:2])
```

```
## <<SimpleCorpus>>
## Metadata:  corpus specific: 1, document level (indexed): 0
## Content:   documents: 2
##
## [1] big data mining analytics components strategic decision making stephan kudyba foreword thomas dav
## [2]  big data mining analytics components strategic decision making
```

```
myTdm <- TermDocumentMatrix(myCorpus,control = list(wordLengths=c(2,Inf)))
m <- as.matrix(myTdm)
freq <- sort(rowSums(m),decreasing = T)
```

## Word Cloud

```
set.seed(1234)
wordcloud(words=names(freq),freq=freq, min.freq = 150,random.order = F,
          max.words = 120, colors = brewer.pal(8,"Dark2"),
          rot.per = 0.35,scale = c(5, 0.3),)
```

data mining

users general statistical threshold processing algorithms time discovery systems multiple system chapter pattern clusters concepts dimension han query itemsets applications well conf based two will han dimension networks outliers knowledge clustering given learning different training class methods attribute rule classification also cube classi large number measures edition section example analysis concept models multidimensional decision model minimum order use tuples object new cluster using table level databases network dimensions process rules value pp cation computer tree outlier objects algorithm however items frequent information one olap figure support patterns web sales include hierarchy user proc search database form method association many warehouse suppose customer space measure detection approach distance techniques values