# LETTERKENNY INSTITUTE OF TECHNOLOGY

# **ASSIGNMENT COVER SHEET**

Lecturer's Name: James Connolly				
Assessment Title: Data scraping				
Work to be submitted to:James Connolly				
Date for submission of work: 29-08-2018				
Place and time for submitting work: 1:00 pm				
To be completed by the Student				
Student's Name: PRATEEK PARASHER	_			
Class: Msc in Big Data				
Subject/Module: Data Science				
Word Count (where applicable):				
I confirm that the work submitted has been produced solely through my own efforts.				
Student's signature: PRATEEK PARASHER Date: 28/8/2018				

# Notes

**Penalties:** The total marks available for an assessment is reduced by 15% for work submitted up to one week late. The total marks available are reduced by 30% for work up to two weeks late. Assessment work received more than two weeks late will receive a mark of zero. [Incidents of alleged plagiarism and cheating are dealt with in accordance with the Institute's Assessment Regulations.]

Plagiarism: Presenting the ideas etc. of someone else without proper acknowledgement (see section L1 paragraph 8).

**Cheating:** The use of unauthorised material in a test, exam etc., unauthorised access to test matter, unauthorised collusion, dishonest behaviour in respect of assessments, and deliberate plagiarism (see section L1 paragraph 8).

**Continuous Assessment:** For students repeating an examination, marks awarded for continuous assessment, shall normally be carried forward from the original examination to the repeat examination.

Description: - Twitter is a great source for sentiment data and social media mining furthermore it is quite easy to get significant amounts of data to be able to scrape data from Twitter. Last year we all experienced the different climate conditions and multiple climate warning announced by govt. this lab I scrapped the data from twitter regarding climate warning like heatwaves warning, snow warning etc. after that I visualize the wordcloud of scraped data from twitter.

GitHub Link :- https://github.com/prateekparasher/web scrap

# STEP DESCRIPTION AND CODE

```
#install package
install.packages("rvest")
library(rvest)
#read link
url <- 'https://twitter.com/weatherrte?lang=en'</pre>
web_page <- read_html(url)</pre>
#head & str data
head(web_page)
str(web_page)
#ranking tweet data
warning <- html_nodes(web_page,'.tweet-text')</pre>
head(warning, 30)
#leghth of tweet data
length(warning)
data <- html_text(warning)</pre>
#show first 30 tweets
head(data, 30)
#Install and load the required packages
# Install
```

install.packages("tm") # for text mining install.packages("SnowballC") # for text stemming

```
install.packages("wordcloud") # word-cloud generator
install.packages("RColorBrewer") # color palettes
# Load
library("tm")
library("SnowballC")
library("wordcloud")
library("RColorBrewer")
# Load the data as a corpus
docs <- Corpus(VectorSource(data))</pre>
#Inspect the content of the document
inspect(docs)
#Text transformation
toSpace <- content_transformer(function (x , pattern ) gsub(pattern, " ", x))
docs <- tm_map(docs, toSpace, "/")</pre>
docs <- tm_map(docs, toSpace, "@")</pre>
docs <- tm_map(docs, toSpace, "\\|")</pre>
# Convert the text to lower case
docs <- tm_map(docs, content_transformer(tolower))</pre>
# Remove numbers
docs <- tm_map(docs, removeNumbers)</pre>
# Remove english common stopwords
docs <- tm_map(docs, removeWords, stopwords("english"))</pre>
# Remove your own stop word
# stopwords as a character vector
docs <- tm_map(docs, removeWords, c("blabla1", "blabla2"))</pre>
# Remove punctuations
docs <- tm_map(docs, removePunctuation)</pre>
```

# # Eliminate extra white spaces docs <- tm\_map(docs, stripWhitespace) # Text stemming docs <- tm\_map(docs, stemDocument) #Build a term-document matrix dtm <- TermDocumentMatrix(docs) m <- as.matrix(dtm) v <- sort(rowSums(m),decreasing=TRUE) d <- data.frame(word = names(v),freq=v) head(d, 20) # Generate the Word cloud

# **STR (Structure of scraped data)**

```
str(d)
'data.frame': 222 obs. of 2 variables:
    $ word: Factor w/ 222 levels "across","advisory",..: 153 216 145 208 64 45 5
2 71 206 154 ...
    $ freq: num 12 12 9 9 7 6 6 5 5 4 ...
```

## **RESULT & VISUALIZATION**

```
.oncenc. documents. 21
[1] With live coverage on #RTEOne of the #PopeInIreland and Mass from the #PhoenixPark, there may be changes to the s
thedule today. If so you can see this week's #ISL forecast on #RTEPlayer. Please RT @RTEOne
[2] There's a Weather Advisory in place for heavy rain this evening and early tonight - get the forecast for that and
the next few days at 1.25 on #RTEOne
[3] Heat warnings remain in operation across Turkey will highs of 35 degrees. @nualacarey25 has more in your #europea
nforecast at 5.50pm this evening on #RTEOne
[4] Come hail, rain or shine #RTEWeather is a constant in the #RTENewSeason schedule. Our office is open 365 days a y
ear! @nualacarey25https://twitter.com/helcurran/status/1030049983716892672 ...
[5] Myself and @nualacarey25 at #RTENewSeason pic.twitter.com/yftEkgx9BD
[6] WATCH: After a long dry spell the rain pours down in Dublin #OnThisDay in 1968 https://bit.ly/2BOvbH7 pic.twitter
com/b6vCoKhcBE
[7] Mostly dry with some sunny spells in the east, but cloudier in the west and northwest with patchy rain and drizzl
  extending slowly eastwards but tending to die away. Top temperatures will be 17 to 22 degrees.pic.twitter.com/cfw4I
[8] Dry with some sunshine this morning but cloudier and humid in the southwest with rain gradually spreading northea
stwards, turning persistent and heavy later in the afternoon and evening. Top temperatures of 18 to 21 degrees.
[9] Cloudy with some outbreaks of rain today but generally dry, sunshine will develop in west & north but cooler air moving into west ...highs of 16C there- up to 22C for the east & southeast, join me for an update on @rteone at 11.50a
[10] Very hot weather in Europe now, especially in Iberia. High temperature warnings have been issued as the figures r
each over 40 degrees in Spain and Portugal. More details in the #europeanforecast on @RTEOne at 5.50pm approx.
[11] Sunscreen that is spread on the skin too thinly may provide less than half the expected degree of protection, a sawdy has shown https://www.rte.ie/news/newslens/2018/0725/980981-sun-lotion-application/ ... via @rtenews
[12] Solar UV Index for Wednesday: HIGH in all areas.
[13] Dry in most places today with long sunny spells. However there is the chance of a few showers along Atlantic coas all counties. Warm with highs of 20 to 25 degrees. Light to moderate southerly winds will increase strong and gusty al
ong the Atlantic seaboard later in the day.pic.twitter.com/OQy9Sd7ojf
[14] See your weekly #ISL weather forecast back on #RTEOne at 5.50pm approx. this evening. Available shortly thereafte
on @rteplayer #irishsignlanguage #weather #rteplayer @SenanDunne @electricginger @Carabellew @CarolineMcTweet @RTEOn
[15] Weather warnings have been issued for rain and thunder in eastern parts of Europe, while in southern Italy they a
e in place for high temperatures. More details this evening #RTEOne at 5.50pm in your European forecast. #Weather
[16] Solar UV index: HIGH in sunshine nationwide on Sunday.
[17] A series of weather warnings remain in place across Europe. They concern a mix of rain & thunderstorms for Croati
i, Bosnia, Serbia and Italy. In Sardinia and Italy there are also wind warnings in operation . \nFull details at 5.50p
n on #RTEOne with @nualacarey25
[18] Apart from well scattered showers, many places will remain dry this afternoon and evening. Rather cloudy general
```

Fig 1 – warning tweets scraped from twitter using R

> head(d, 20	)	
	word	freq
rteone	rteone	12
will	will	12
rain	rain	9
weather	weather	9
evening	evening	7
degrees	degrees	6
dry	dry	6
forecast	forecast	5
warnings	warnings	5
rteplayer	rteplayer	4
today	today	4
highs	highs	4
temperatures	temperatures	4
west	west	4
high	high	4
showers	showers	4
isl	isl	3
place	place	3
nualacarey	nualacarey	3
remain	remain	3
> cot cood/1	2247	

Fig  $\,2\,$  - Table containing the frequency of the words

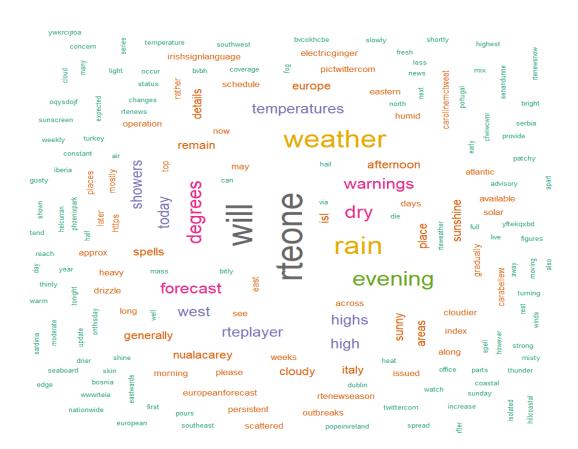


Fig 3 - Data cloud visualization