# LETTERKENNY INSTITUTE OF TECHNOLOGY

# ASSIGNMENT COVER SHEET

Lecturer's Name: **James Connolly**

Assessment Title: **Data scraping**

Work to be submitted to: **08/01/19**

Date for submission of work: **James Connolly**

Place and time for submitting work:

---

## To be completed by the Student

Student's Name: **PRATEEK PARASHER**

Class: **MSc Big Data Analytics**

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: **08/01/19**
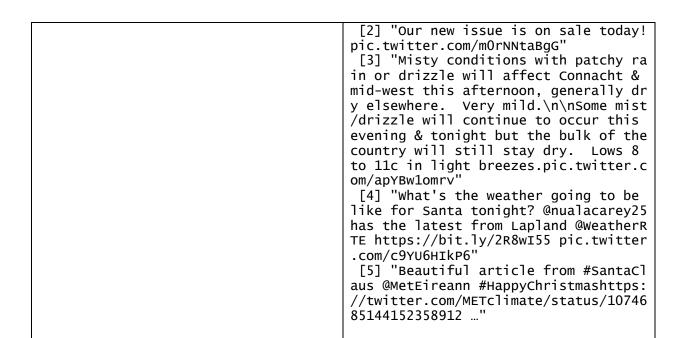
---

**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. In this lab I scrapped the data from twitter regarding climate warning like heatwaves warning, snow warning etc.

| Step Description: - | install packages<br>read link<br>head & str data |
|---|---|
| R code used to perform the scrape: - | install.packages("rvest")<br>library(rvest)<br>url <- 'https://twitter.com/weatherrte?lang=en'<br>web_page <- read_html(url)<br>head(web_page) |
| Structure of scraped data: - | str(web_page)<br>List of 2<br> $ node:<externalptr><br> $ doc :<externalptr><br>- attr(*, "class")= chr [1:2] "xml_document" "xml_node" |
| Result: - | $\`node\`<br><pointer: 0x000000000d9dd100><br><br>$doc<br><pointer: 0x0000000008ef38a0> |
| | |

| | ranking tweet data |
|---|---|
| Step Description: - | |
| R code used to perform the scrape: - | warning <- html_nodes(web_page,'.tweet-text') |
| Structure of scraped data: - | str(warning)<br><br>List of 24<br> $ :List of 2<br>  ..$ node:<externalptr><br>  ..$ doc :<externalptr><br>  ..- attr(*, "class")= chr "xml_node"<br> $ :List of 2<br>  ..$ node:<externalptr><br>  ..$ doc :<externalptr> |

| | |
|---|---|
| | ```..- attr(*, "class")= chr "xml_node"``` <br><br> ```  $ :List of 2``` <br> ```  ..$ node:<externalptr>``` <br> ```  ..$ doc :<externalptr>``` <br> ```..- attr(*, "class")= chr "xml_node"``` |
| **Result: -** | ```[1] <p class="TweetTextSize TweetTextSize--normal js-tweet-text tweet-text" lang="en" data-aria-label-part="4">With details of how to enter the Winter quarter of the RTE Wea ...``` <br> ```[2] <div class="QuoteTweet-text tweet-text u-dir js-ellipsis" lang="en" data-aria-label-part="2" dir="ltr">Our new issue is on sale today! <span class="twitter-timeline-link ...``` <br> ```[3] <p class="TweetTextSize TweetTextSize--normal js-tweet-text tweet-text" lang="en" data-aria-label-part="0">Misty conditions with patchy rain or drizzle will affect Conna ...``` <br> ```[4] <p class="TweetTextSize TweetTextSize--normal js-tweet-text tweet-text" lang="en" data-aria-label-part="0">What's the weather going to be like for Santa tonight? <a href ...``` <br> ```[5] <p class="TweetTextSize TweetTextSize--normal js-tweet-text tweet-text" lang="en" data-aria-label-part="4">Beautiful article from <a href="/hashtag/SantaClaus?src=hash"  ...``` |

| | |
|---|---|
| **Step Description: -** | length of tweet data |
| **R code used to perform the scrape: -** | length(warning) <br> data <- html_text(warning) <br> head(data, 30) <br> str(data) |
| **Structure of scraped data: -** | ```chr [1:24] "With details of how to enter the Winter quarter of the RTE Weather Photo competition...so snap to it! @RTEOne @"| __truncated__ ...``` |
| **Results: -** | ```[1] "With details of how to enter the Winter quarter of the RTE Weather Photo competition...so snap to it! @RTEOne @WeatherRTE #Weather #winter #photography #competitionhttps://twitter.com/RTE_GUIDE/status/107859904361237299 94 …"``` |

| | |
|---|---|
| | [2] "Our new issue is on sale today! pic.twitter.com/m0rNNtaBgG"<br><br>[3] "Misty conditions with patchy rain or drizzle will affect Connacht & mid-west this afternoon, generally dry elsewhere.  Very mild.\n\nSome mist/drizzle will continue to occur this evening & tonight but the bulk of the country will still stay dry.  Lows 8 to 11c in light breezes.pic.twitter.com/apYBw1omrv"<br><br>[4] "What's the weather going to be like for Santa tonight? @nualacarey25 has the latest from Lapland @WeatherRTE https://bit.ly/2R8wI55 pic.twitter.com/c9YU6HIkP6"<br><br>[5] "Beautiful article from #SantaClaus @MetEireann #HappyChristmashttps://twitter.com/METclimate/status/1074685144152358912 …" |

**GitHub Link :-  https://github.com/prateekparasher/web_scrap**