**Shell scripting assignment**
A basic introduction to Shell scripting for Msc GeoData science students at UBS as part of the
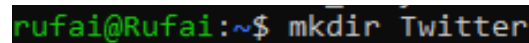Copernicus Master in Digital Earth.
Facilitated by: Dr. Charlotte Petellier
The following highlights below shows the steps used to achieve certain tasks using the UNIX-based
shell. The shell scripts and the rationale for using them were also mentioned.

1. Create a new Directory
   As a first step, it is important to create a singular working folder for organizing the files, data
   and other output of the task. Here, a twitter folder was created using the shell command:
   
   *mkdir Twitter*

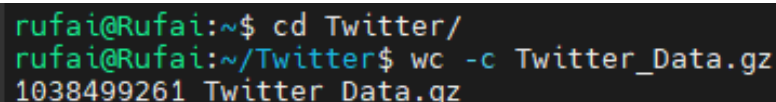   rufai@Rufai:~$ mkdir Twitter

   Fig 1: Shell command for creating a directory.

2. Download the file. How big is it?
   The twitter data was uploaded on the Unix environment using a third party software --
   MobaXterm and decompressed afterwards using the shell command:
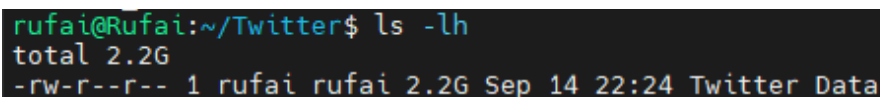   
   *gzip - d Twitter.gz*
   
   To get an idea of the size of the file, the *wc -c* command was used to check the file size prior
   to decompressing it.

   rufai@Rufai:~$ cd Twitter/
   rufai@Rufai:~/Twitter$ wc -c Twitter_Data.gz
   1038499261 Twitter_Data.gz

   Fig 2: Shell command for extracting the number of bytes occupied by a file.

   From the information displayed, we can tell that the twitter data in its compressed format is
   about **1038499261 bytes (~ 1 GB) in** size. After decompression, the file size was checked using
   the *ls -lh* command

   rufai@Rufai:~/Twitter$ ls -lh
   total 2.2G
   -rw-r--r-- 1 rufai rufai 2.2G Sep 14 22:24 Twitter_Data

   Fig 3: Shell command for showing human-readable information (including the size) about the files in a folder.

   This shows that the twitter file in its decompressed format i**s approximately 2.2GB.**

3. What delimiter is used to separate the columns in the file?
   A quick visual inspection of the twitter data, shows that the columns are tab-delimited;
   separated by tabs. A See a sample of the data in the image shown below.

4. How many columns are there? What do the columns describe?

   To extract the number of columns in the data, the awk shell command shown below was used:

   <div align="center">

   *awk -F "\t"'{print NF; exit}' Twitter_Data*

   </div>

   

   The output shows that there are **4 columns i**n the data. These columns describe the individual identity of the tweets (tweet id), the username, data of the tweets and the content of the tweets respectively.

5. How many Tweets are there in the file?

   The data contains at least 15,089,920 tweets as extracted by the shell command shown below.

   

6. What is the date range for Tweets in this file?

   To extract the date range from the twitter data, the awk command was also used but this time, the output was redirected to a date file.

   <div align="center">

   *awk -F'\t' '{print$3}' Twitter_Data > date*

   </div>

   

   The top and last 10 content of this date file was examined to get an idea of the tweet date range. A sample of these two sets of information are shown below:

```
rufai@Rufai:~/Twitter$ head date
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
Tue Feb 11 12:18:36 +0000 2014
rufai@Rufai:~/Twitter$ tail date
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
Tue Feb 18 23:15:00 +0000 2014
rufai@Rufai:~/Twitter$
```

Based on this, we can conclude that the tweet date ranges from **Tue, Feb 11 12:18:36 +0000 2014** to **Tue, Feb 18 23:15:00 +0000 2014.**

7. How many unique users are there?

The number of unique users in the file was computed using the shell command:

*awk '{print $2}' Twitter_Data | sort | uniq |wc -l*

```
rufai@Rufai:~/Twitter$ awk '{print $2}' Twitter_Data | sort | uniq |wc -l
8977904
```

From this output we can conclude that there are **8977904 unique users**.

8. When was the first mention in the file of Donald Trump and what was the tweet?

To extract the first mention of Donald Trump in the twitter file, the shell command shown below was used.

*grep -w Twitter_Data -e 'Donald Trump' | head -n1*

```
rufai@Rufai:~/Twitter$ grep -w Twitter_Data -e 'Donald Trump' | head -n1
433215995134476289      Maddog4U_1st    Tue Feb 11 12:28:36 +0000 2014  RT @
/qAcG…
```

The tweet that first mention Donald Trump was:

```
rufai@Rufai:~/Twitter$ awk {'print $4} Twitter_Data | grep Twitter_Data -e "Donald Trump" | head -n1 | cut -f4
RT @aedan_smith: Be interesting to see the detail on this one:  BBC News - Donald Trump loses offshore wind farm challenge http://t.co/qAcG…
rufai@Rufai:~/Twitter$
```

**RT @aedan_smith: Be interesting to see the details on this one:  BBC News - Donald Trump loses offshore wind farm challenge http://t.co/qAcG…**

The command used to extract this tweet is:

*awk {'print $4'} Twitter_Data | grep Twitter_Data -e "Donald Trump" | head -n1 | cut -f4*

9. How many times has Donald Trump been mentioned? What about Barack Obama? Hillary Clinton?

The respective shell commands shown below were used to extract the number of times each aforementioned string was found in the twitter file.

    I.   Donald Trump: A total of 130 references were found as shown in the graphical display below.

*grep -o -i Twitter_Data -o -i "Donald Trump" | wc -l*

```
RT @aedan_smith: Be interesting to see the detail on this one. BBC News - Donald
rufai@Rufai:~/Twitter$ grep -o -i "Donald Trump" Twitter_Data | wc -l
130
```

    II.   Barack Obama: A total of 482 references were found as shown in the graphical display below.

*grep -o -i Twitter_Data  -o -i "Barack Obama" | wc -l*

```
rufai@Rufai:~/Twitter$ grep -o -i "Barack Obama" Twitter_Data | wc -l
482
```

    III.   Hillary Clinton: A total of 127 references were found as shown in the graphical display below.

*grep -o -i Twitter_Data  -o -i "Hillary  Clinton" | wc -l*

```
rufai@Rufai:~/Twitter$ grep -o -i "Hillary Clinton" Twitter_Data | wc -l
127
```

10. Do you think we have captured all the references to Donald Trump, Barack Obama, and Hilary Clinton? What other strings might we need to try? What problems might we face?

No, we haven't captured all the references to the aforementioned names in the field. Largely because of different sentence cases, incomplete and incorrect spellings in the file. For instance, there is still a reference to Barack Obama using **Barack obam, Barack Obam,** for **Hillary Clinton** we still have **Hilary Clinton. Donald trump, and donald trump.** There are also possibilities for people to make tweets about these people using their nicknames, which might be specific to certain regions or groups of people.