```
In [13]: from pyspark.streaming import StreamingContext
            from pyspark.sql import Row
           import time
In [14]: def save(time, rdd):
                try:
                     df = spark.createDataFrame(rdd.map(\
                      \textbf{lambda} \text{ row: } \texttt{Row}(\texttt{time=time, package=row[0], count=row[1])))}
                      df.show(5)
                     df.write.format("org.apache.spark.sql.cassandra")\
.options(table="cran_counts", keyspace="streamingbdm")\
.save(mode="append")
                 except:
                     pass
In [15]: ssc = StreamingContext(sc, 3)
 In [*]: cran_data = ssc.textFileStream("file:///home/victorbdm/assignStreaming/")
           cran_data = cran_data.map(lambda x: x.split(','))
           def remove_quotation(x):
    return([xx.replace('"', '') for xx in x])
            cran_data = cran_data.map(remove_quotation)
           package_download_count = cran_data.map(lambda package:(package[6], 1))
package_download_count = package_download_count.reduceByKey(lambda a,b: a+b)
            package_download_count.pprint()
            package_download_count.foreachRDD(save)
            ssc.start()
            time.sleep(120)
            ssc.stop(stopSparkContext=False)
```

victorbdm@bdm: ~

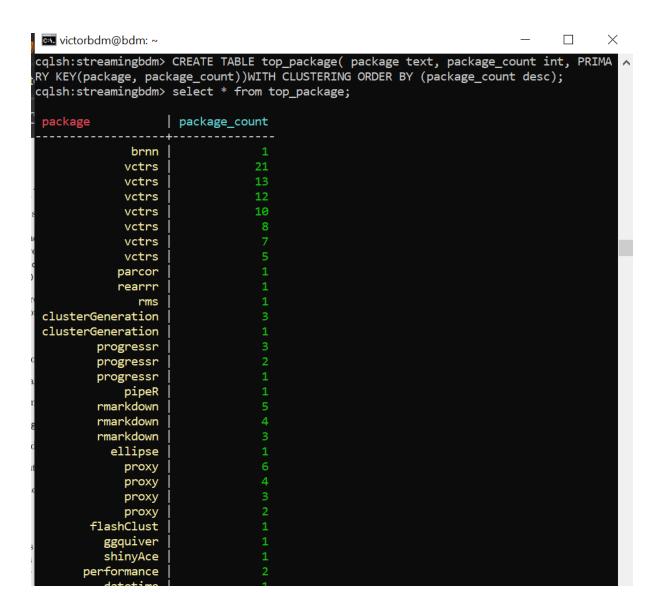
```
cqlsh:streamingbdm> drop table word_counts;
cqlsh:streamingbdm> cls
cqlsh:streamingbdm> CREATE TABLE cran_counts( time text, package
text, count int, PRIMARY KEY(time, package));
cqlsh:streamingbdm> select * from cran_counts;
time | package | count
(0 rows)
cqlsh:streamingbdm> select * from cran_counts;
 time
                          package
                                                   count
 2021-12-21 23:02:00+0000
                                        AutoDeskR
                                                        1
 2021-12-21 23:02:00+0000
                                               BH
 2021-12-21 23:02:00+0000
                                              C50
 2021-12-21 23:02:00+0000
                                              DBI
 2021-12-21 23:02:00+0000
                                                        1
                                               DΤ
 2021-12-21 23:02:00+0000
                                        DescTools
                                                        1
 2021-12-21 23:02:00+0000
                                         EnvStats
                                                        1
 2021-12-21 23:02:00+0000
                                                        1
                                            Exact
 2021-12-21 23:02:00+0000
                                              FNN
                                                        1
 2021-12-21 23:02:00+0000
                                          Formula
                                                        2
 2021-12-21 23:02:00+0000
                                        GADMTools
 2021-12-21 23:02:00+0000
                                                        1
                                           GGallv
 2021-12-21 23:02:00+0000
                                      GPArotation
 2021-12-21 23:02:00+0000
                                            GPfit
 2021-12-21 23:02:00+0000
                                            Hmisc
                                                        8
 2021-12-21 23:02:00+0000
                                        JSconsole
                                                        1
 2021-12-21 23:02:00+0000
                                           KMsurv
                                                        1
 2021-12-21 23:02:00+0000
                                          MLEcens
                                                        1
 2021-12-21 23:02:00+0000
                                     MatrixModels
 2021-12-21 23:02:00+0000
                                     ModelMetrics
                                                        4
 2021-12-21 23:02:00+0000
                                           NlinTS
                                                        1
 2021-12-21 23:02:00+0000
                                                       10
                                               R6
 2021-12-21 23:02:00+0000
                                     RColorBrewer
 2021-12-21 23:02:00+0000
                                            RCurl
                                                        1
2021-12-21 23:02:00+0000
                                           RMySQL
                                                        1
2021-12-21 23:02:00+0000
                                          RSQLite
                                                        1
2021-12-21 23:02:00+0000
                                           Rcgmin
```

```
2021-12-21 23:02:00+0000 | zip | 3
2021-12-21 23:02:00+0000 | zro | 1
2021-12-21 23:01:54+0000 | ATAforecasting | 1
2021-12-21 23:01:54+0000 | ATAforecasting | 1
2021-12-21 23:01:54+0000 | ATR | 1
2021-12-21 23:01:54+0000 | ATR | 1
2021-12-21 23:01:54+0000 | ATmet | 1
2021-12-21 23:01:54+0000 | Atment | 1
2021-12-21 23:01:54+0000 | BH | 1
2021-12-21 23:01:54+0000 | BH | 1
2021-12-21 23:01:54+0000 | BI | 2
2021-12-21 23:01:54+0000 | DistributionUtils | 1
2021-12-21 23:01:54+0000 | GeneralizedHyperbolic | 1
2021-12-21 23:01:54+0000 | Historal | 1
2021-12-21 23:01:54+0000 | MatrixModels | 4
2021-12-21 23:01:54+0000 | ModelNetrics | 3
```

Question 2

ssc.start() time.sleep(120)

ssc.stop(stopSparkContext=False)



```
\times
victorbdm@bdm: ~
--MORE---
package
                 package_count
       pbkrtest
       pbkrtest
       ellipsis
       ellipsis
                              13
       ellipsis
                              12
       ellipsis
                              10
       ellipsis
       ellipsis
     rapidjsonr
          bit64
          bit64
          bit64
          bit64
          bit64
            DBI
            DBI
            DBI
            DBI
            DBI
        VALERIE
            fma
           ade4
         praise
         praise
            TSP
            bst
       tidytext
       tidytext
        stringr
        stringr
        stringr
        stringr
        whisker
        mvnfast
```

Question 3

question 3

```
cqlsh:streamingbdm> select * from country_counts;
 country | download_count
      IN
                         1
                         2
      IN
      IN
      IN
                        11
      PΕ
                         2
      PΕ
                         4
      PΕ
                         8
      PΕ
                        13
      PH
      ΑT
      ΑT
      JΡ
                         1
      JΡ
      JΡ
                        10
      JΡ
      JΡ
                        13
      HK
      HK
      HK
                        36
      HK
      HK
                        47
      FR
      FR
      FR
      FR
      NA
                        99
      NA
                       115
                       125
      NA
```

MORE		
country		
US	415	
US	436	
US	457	
US	505	
US	955	
MZ	3	
MZ	6	
MZ	11	
MZ	12	
MZ	15	
UA	1	
SE	2	
SE	4	
PK	1	
PK	2	
OM	1	
PL	4	
PL	9	
PL	13	
PL	14	
PL	25	
PL	27	
country	1	
MY	1	
MY	2	
BF	1	
BF	2	
BF	3	
RU	2	
RU	5	
RU	6	
KR	1	

Question 4

question 4

```
In [3]:
    ssc = StreamingContext(sc, 3)
        cran_data = ssc.textFileStream("file:///home/victorbdm/assignStreaming,
        cran_data = cran_data.map(lambda x: x.split(','))

    def remove_quotation(x):
        return([xx.replace('"', '') for xx in x])
        cran_data = cran_data.map(remove_quotation)
        ggplot2_package = cran_data.map(lambda package:(package[6], 1))
        ggplot2_package = ggplot2_package.reduceByKey(lambda a,b: a+b)
        ggplot2_package= ggplot2_package.filter(lambda a: 'ggplot2' in a)
        ggplot2_package.foreachRDD(save)

    ssc.start()
    time.sleep(120)
    ssc.stop(stopSparkContext=False)
```

```
victorbdm@bdm: ~
```

```
(0 rows)
cqlsh:streamingbdm> select * from gplot_counts;
time
                          package gp count
2021-12-22 02:17:18+0000
                          ggplot2
                                           22
                           ggplot2
2021-12-22 02:17:09+0000
                                           54
                           ggplot2
2021-12-22 02:17:15+0000
                                           34
(3 rows)
cqlsh:streamingbdm> select * from gplot_counts;
time
                          package gp_count
2021-12-22 02:17:18+0000
                           ggplot2
                                           22
                          ggplot2
2021-12-22 02:17:09+0000
                                           54
2021-12-22 02:17:24+0000
                           ggplot2
                                          15
2021-12-22 02:17:21+0000
                          ggplot2
                                           17
2021-12-22 02:17:15+0000 | ggplot2 |
                                           34
(5 rows)
cqlsh:streamingbdm> select * from gplot_counts;
time
                          package gp_count
2021-12-22 02:17:18+0000
                           ggplot2
                                           22
2021-12-22 02:17:36+0000
                           ggplot2
                                           17
                           ggplot2
2021-12-22 02:17:09+0000
                                           54
2021-12-22 02:17:39+0000
                           ggplot2
                                           17
2021-12-22 02:17:27+0000
                          ggplot2
                                           23
2021-12-22 02:17:24+0000
                           ggplot2
                                           15
2021-12-22 02:17:45+0000
                          ggplot2
                                           27
                           ggplot2
2021-12-22 02:17:33+0000
                                           28
                          ggplot2
2021-12-22 02:17:21+0000
                                           17
2021-12-22 02:17:30+0000
                           ggplot2
                                           19
                           ggplot2
2021-12-22 02:17:42+0000
                                           19
                          ggplot2
2021-12-22 02:17:15+0000
                                           34
```

victorbdm@bdm: ~

```
(26 rows)
cqlsh:streamingbdm> select * from gplot_counts;
time
                           package gp_count
 2021-12-22 02:18:30+0000
                            ggplot2
                                             28
 2021-12-22 02:17:54+0000
                             ggplot2
                                             20
                            ggplot2
 2021-12-22 02:18:21+0000
                                             27
 2021-12-22 02:17:18+0000
                             ggplot2
                                             22
                            ggplot2
 2021-12-22 02:17:57+0000
                                             23
                            ggplot2
 2021-12-22 02:18:24+0000
                                             19
                            ggplot2
 2021-12-22 02:17:36+0000
                                             17
 2021-12-22 02:17:09+0000
                            ggplot2
                                             54
                            ggplot2
 2021-12-22 02:17:39+0000
                                             17
                            ggplot2
 2021-12-22 02:17:27+0000
                                             23
                            ggplot2
 2021-12-22 02:17:24+0000
                                             15
 2021-12-22 02:18:03+0000
                            ggplot2
                                             14
 2021-12-22 02:17:45+0000
                             ggplot2
                                             27
                            ggplot2
 2021-12-22 02:18:15+0000
                                             17
                            ggplot2
 2021-12-22 02:17:51+0000
                                             13
 2021-12-22 02:18:27+0000
                            ggplot2
                                             28
                            ggplot2
 2021-12-22 02:18:48+0000
                                             30
 2021-12-22 02:18:42+0000
                             ggplot2
                                             22
 2021-12-22 02:18:39+0000
                            ggplot2
                                             27
                            ggplot2
 2021-12-22 02:19:03+0000
                                             14
                            ggplot2
 2021-12-22 02:18:36+0000
                                             33
 2021-12-22 02:17:33+0000
                            ggplot2
                                             28
 2021-12-22 02:17:21+0000
                            ggplot2
                                             17
                            ggplot2
 2021-12-22 02:18:00+0000
                                             18
                            ggplot2
 2021-12-22 02:17:30+0000
                                             19
                            ggplot2
 2021-12-22 02:18:54+0000
                                             27
                             ggplot2
 2021-12-22 02:18:51+0000
                                             27
 2021-12-22 02:17:48+0000
                            ggplot2
                                             20
                             ggplot2
 2021-12-22 02:17:42+0000
                                             19
 2021-12-22 02:18:09+0000
                             ggplot2
                                             16
                             ggplot2
 2021-12-22 02:18:57+0000
                                             16
```