

In [ ]:

```
'''

(10 points) Part 3a: MapRedce
In this task, design a MapReduce program in python that reads all the original tweets (no retweets) in the sample tweets (congress-sample-10k.json.gz) and if a tweet is a reply to another tweet then output a record of the form <src_id, src_user, dst_id, dst_user>.

Create a small cluster (2 or 3 nodes) as per the AWS Guide and then ssh to your cluster and use Hadoop streaming to execute your mapreduce program.

Note: the Hadoop streaming jar file can be found at /usr/lib/hadoop-mapreduce/hadoop-p-streaming.jar

'''
```

In [21]:

```
congress_sample = spark.read.csv("s3://us-congress-tweets/congress-sample-10k.json.gz", header=True)
```

In [31]:

```
# your mapper function
# congress_sample_reply = congress_sample.select("in_reply_to_status_id_str")
# congress_sample_reply.show()

ak47 = congress_sample.limit(1).toPandas
```

In [35]:

```
ak47.head
```

```
'function' object has no attribute 'head'
Traceback (most recent call last):
AttributeError: 'function' object has no attribute 'head'
```

In [33]:

```
congress_sample.printSchema()
```



```

root
|-- {"extended_tweet":{"entities":{"urls":[]: string (nullable = true)
|-- "hashtags":[]1: string (nullable = true)
|-- "user_mentions":[{"indices":[02: string (nullable = true)
|-- 8]3: string (nullable = true)
|-- "screen_name":"jillwow"4: string (nullable = true)
|-- "id_str":"464222379"5: string (nullable = true)
|-- "name":"jillie"6: string (nullable = true)
|-- "id":464222379}7: string (nullable = true)
|-- {"indices":[98: string (nullable = true)
|-- 23]9: string (nullable = true)
|-- "screen_name":"RepAdamSchiff"10: string (nullable = true)
|-- "id_str":"29501253"11: string (nullable = true)
|-- "name":"Adam Schiff"12: string (nullable = true)
|-- "id":29501253}]13: string (nullable = true)
|-- "symbols":[]14: string (nullable = true)
|-- "full_text":"@jillwow @RepAdamSchiff That is what is so disturbing
& ominous 🤖 Also: string (nullable = true)
|-- we must remember our world history - didn\u2019t Hiler only have
30% or so support from the German ppl before he began his persecution o
f the Jewish ppl? 🤔 ": string (nullable = true)
|-- "display_text_range":[2417: string (nullable = true)
|-- 226}}: string (nullable = true)
|-- "in_reply_to_status_id_str":"1047522510839975936": string (nullabl
e = true)
|-- "in_reply_to_status_id":1047522510839975936: string (nullable = tr
ue)
|-- "created_at":"Wed Oct 03 17:30:03 +0000 2018": string (nullable =
true)
|-- "in_reply_to_user_id_str":"464222379": string (nullable = true)
|-- "source":"<a href=\"http://twitter.com/download/iphone\" rel=\"nof
ollow\">Twitter for iPhone</a>": string (nullable = true)
|-- "retweet_count":0: string (nullable = true)
|-- "retweeted":false: string (nullable = true)
|-- "geo":null: string (nullable = true)
|-- "filter_level":"low": string (nullable = true)
|-- "in_reply_to_screen_name":"jillwow": string (nullable = true)
|-- "is_quote_status":false: string (nullable = true)
|-- "id_str":"1047539288705880064": string (nullable = true)
|-- "in_reply_to_user_id":464222379: string (nullable = true)
|-- "favorite_count":0: string (nullable = true)
|-- "id":1047539288705880064: string (nullable = true)
|-- "text":"@jillwow @RepAdamSchiff That is what is so disturbing &am
p; ominous 🤖 Also: string (nullable = true)
|-- we must remember our world history - didn\u2019t\u2026 https://t.
co/h5SURkil2V": string (nullable = true)
|-- "place":null: string (nullable = true)
|-- "lang":"en"37: string (nullable = true)
|-- "quote_count":0: string (nullable = true)
|-- "favorited":false: string (nullable = true)
|-- "coordinates":null: string (nullable = true)
|-- "truncated":true: string (nullable = true)
|-- "timestamp_ms":"1538587803423": string (nullable = true)
|-- "reply_count":0: string (nullable = true)
|-- "entities":{"urls":[{"display_url":"twitter.com/i/web/status/1\u20
26": string (nullable = true)
|-- "indices":[121: string (nullable = true)

```

```

|-- 144]: string (nullable = true)
|-- "expanded_url": "https://twitter.com/i/web/status/10475392887058800
64": string (nullable = true)
|-- "url": "https://t.co/h5SURki12V"}]: string (nullable = true)
|-- "hashtags": []49: string (nullable = true)
|-- "user_mentions": [{"indices": [050: string (nullable = true)
|-- 8]51: string (nullable = true)
|-- "screen_name": "jillwow"52: string (nullable = true)
|-- "id_str": "464222379"53: string (nullable = true)
|-- "name": "jillie"54: string (nullable = true)
|-- "id": 464222379}55: string (nullable = true)
|-- {"indices": [956: string (nullable = true)
|-- 23]57: string (nullable = true)
|-- "screen_name": "RepAdamSchiff"58: string (nullable = true)
|-- "id_str": "29501253"59: string (nullable = true)
|-- "name": "Adam Schiff"60: string (nullable = true)
|-- "id": 29501253}61: string (nullable = true)
|-- "symbols": []62: string (nullable = true)
|-- "display_text_range": [2463: string (nullable = true)
|-- 140]: string (nullable = true)
|-- "contributors": null: string (nullable = true)
|-- "user": {"utc_offset": null: string (nullable = true)
|-- "friends_count": 1034: string (nullable = true)
|-- "profile_image_url_https": "https://pbs.twimg.com/profile_images/82
9148920223707137/L0JxcBmq_normal.jpg": string (nullable = true)
|-- "listed_count": 3: string (nullable = true)
|-- "profile_background_image_url": "http://abs.twimg.com/images/theme
s/themel/bg.png": string (nullable = true)
|-- "default_profile_image": false: string (nullable = true)
|-- "favourites_count": 47592: string (nullable = true)
|-- "description": "Proud citizen; ❤️ Obama: string (nullable = true)
|-- Michelle & HRC; #Resistance: string (nullable = true)
|-- #BoycottNRA: string (nullable = true)
|-- #BanAssaultWeapons #ActiveMeasures": string (nullable = true)
|-- "created_at": "Wed Jun 17 17:50:07 +0000 2015": string (nullable =
true)
|-- "is_translator": false: string (nullable = true)
|-- "profile_background_image_url_https": "https://abs.twimg.com/image
s/themes/themel/bg.png": string (nullable = true)
|-- "protected": false: string (nullable = true)
|-- "screen_name": "rposey42": string (nullable = true)
|-- "id_str": "3248076530": string (nullable = true)
|-- "profile_link_color": "1DA1F2": string (nullable = true)
|-- "translator_type": "none": string (nullable = true)
|-- "id": 3248076530: string (nullable = true)
|-- "geo_enabled": false: string (nullable = true)
|-- "profile_background_color": "C0DEED": string (nullable = true)
|-- "lang": "en"88: string (nullable = true)
|-- "profile_sidebar_border_color": "C0DEED": string (nullable = true)
|-- "profile_text_color": "333333": string (nullable = true)
|-- "verified": false: string (nullable = true)
|-- "profile_image_url": "http://pbs.twimg.com/profile_images/829148920
223707137/L0JxcBmq_normal.jpg": string (nullable = true)
|-- "time_zone": null: string (nullable = true)
|-- "url": null: string (nullable = true)
|-- "contributors_enabled": false: string (nullable = true)
|-- "profile_background_tile": false: string (nullable = true)

```

```
|-- "profile_banner_url":"https://pbs.twimg.com/profile_banners/324807
6530/1520646319": string (nullable = true)
|-- "statuses_count":70007: string (nullable = true)
|-- "follow_request_sent":null: string (nullable = true)
|-- "followers_count":660: string (nullable = true)
|-- "profile_use_background_image":true: string (nullable = true)
|-- "default_profile":true: string (nullable = true)
|-- "following":null: string (nullable = true)
|-- "name":"Rhoda🦋 ": string (nullable = true)
|-- "location":"California: string (nullable = true)
    USA": string (nullable = true)
|-- "profile_sidebar_fill_color":"DDEEF6": string (nullable = true)
|-- "notifications":null}}: string (nullable = true)
```

In [37]:

```
congress_sample.rdd.map(lambda r: r.notifications).collect()
```





```
An error occurred while calling z:org.apache.spark.api.python.PythonRDD.collectAndServe.
: org.apache.spark.SparkException: Job aborted due to stage failure: Task 0 in stage 5.0 failed 4 times, most recent failure: Lost task 0.3 in stage 5.0 (TID 11, ip-172-31-12-174.ec2.internal, executor 19): org.apache.spark.api.python.PythonException: Traceback (most recent call last):
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.py", line 1527, in __getattr__
    idx = self.__fields__.index(item)
ValueError: 'notifications' is not in list
```

During handling of the above exception, another exception occurred:

```
Traceback (most recent call last):
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.py", line 377, in main
    process()
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.py", line 372, in process
    serializer.dump_stream(func(split_index, iterator), outfile)
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/serializer.s.py", line 393, in dump_stream
    vs = list(itertools.islice(iterator, batch))
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/util.py", line 113, in wrapper
    return f(*args, **kwargs)
  File "<stdin>", line 1, in <lambda>
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.py", line 1532, in __getattr__
    raise AttributeError(item)
AttributeError: notifications
```

```
at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.handlePythonException(PythonRunner.scala:456)
at org.apache.spark.api.python.PythonRunner$$anon$1.read(PythonRunner.scala:592)
at org.apache.spark.api.python.PythonRunner$$anon$1.read(PythonRunner.scala:575)
at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.hasNext(PythonRunner.scala:410)
at org.apache.spark.InterruptibleIterator.hasNext(InterruptibleIterator.scala:37)
at scala.collection.Iterator$class.foreach(Iterator.scala:891)
at org.apache.spark.InterruptibleIterator.foreach(InterruptibleIterator.scala:28)
at scala.collection.generic.Growable$class.$plus$plus$eq(Growable.scala:59)
at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuffer.scala:104)
at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuffer
```

```

er.scala:48)
    at scala.collection.TraversableOnce$class.to(TraversableOnce.sc
ala:310)
    at org.apache.spark.InterruptibleIterator.to(InterruptibleItera
tor.scala:28)
    at scala.collection.TraversableOnce$class.toBuffer(TraversableO
nce.scala:302)
    at org.apache.spark.InterruptibleIterator.toBuffer(Interruptibl
eIterator.scala:28)
    at scala.collection.TraversableOnce$class.toArray(TraversableOn
ce.scala:289)
    at org.apache.spark.InterruptibleIterator.toArray(Interruptible
Iterator.scala:28)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)
    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkC
ontext.scala:2101)
    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkC
ontext.scala:2101)
    at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.sca
la:90)
    at org.apache.spark.scheduler.Task.run(Task.scala:123)
    at org.apache.spark.executor.Executor$TaskRunner$$anonfun$10.ap
ply(Executor.scala:408)
    at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:
1360)
    at org.apache.spark.executor.Executor$TaskRunner.run(Executor.s
cala:414)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPool
Executor.java:1149)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoo
lExecutor.java:624)
    at java.lang.Thread.run(Thread.java:748)

```

#### Driver stacktrace:

```

    at org.apache.spark.scheduler.DAGScheduler.org$apache$spark$sch
eduler$DAGScheduler$$failJobAndIndependentStages(DAGScheduler.scala:204
1)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage
$1.apply(DAGScheduler.scala:2029)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage
$1.apply(DAGScheduler.scala:2028)
    at scala.collection.mutable.ResizableArray$class.foreach(Resiza
bleArray.scala:59)
    at scala.collection.mutable.ArrayBuffer.foreach(ArrayBuffer.sca
la:48)
    at org.apache.spark.scheduler.DAGScheduler.abortStage(DAGSchedu
ler.scala:2028)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTasks
etFailed$1.apply(DAGScheduler.scala:966)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTasks
etFailed$1.apply(DAGScheduler.scala:966)
    at scala.Option.foreach(Option.scala:257)
    at org.apache.spark.scheduler.DAGScheduler.handleTaskSetFailed
(DAGScheduler.scala:966)

```

```

    at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.doOn
Receive(DAGScheduler.scala:2262)
    at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.onRe
ceive(DAGScheduler.scala:2211)
    at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.onRe
ceive(DAGScheduler.scala:2200)
    at org.apache.spark.util.EventLoop$$anon$1.run(EventLoop.scala:
49)
    at org.apache.spark.scheduler.DAGScheduler.runJob(DAGScheduler.
scala:777)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:206
1)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:208
2)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:210
1)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:212
6)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1.apply(RDD.scala:
945)
    at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperati
onScope.scala:151)
    at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperati
onScope.scala:112)
    at org.apache.spark.rdd.RDD.withScope(RDD.scala:363)
    at org.apache.spark.rdd.RDD.collect(RDD.scala:944)
    at org.apache.spark.api.python.PythonRDD$.collectAndServe(Pytho
nRDD.scala:166)
    at org.apache.spark.api.python.PythonRDD.collectAndServe(Python
RDD.scala)
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAcce
ssorImpl.java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMe
thodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at py4j.reflection.MethodInvoker.invoke(MethodInvoker.java:244)
    at py4j.reflection.ReflectionEngine.invoke(ReflectionEngine.jav
a:357)
    at py4j.Gateway.invoke(Gateway.java:282)
    at py4j.commands.AbstractCommand.invokeMethod(AbstractCommand.j
ava:132)
    at py4j.commands.CallCommand.execute(CallCommand.java:79)
    at py4j.GatewayConnection.run(GatewayConnection.java:238)
    at java.lang.Thread.run(Thread.java:748)
Caused by: org.apache.spark.api.python.PythonException: Traceback (most
recent call last):
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.
py", line 1527, in __getattr__
    idx = self.__fields__.index(item)
ValueError: 'notifications' is not in list

```

During handling of the above exception, another exception occurred:

Traceback (most recent call last):

```
File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
```

```

1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.p
y", line 377, in main
    process()
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.p
y", line 372, in process
    serializer.dump_stream(func(split_index, iterator), outfile)
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/serializer
s.py", line 393, in dump_stream
    vs = list(itertools.islice(iterator, batch))
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/util.py",
line 113, in wrapper
    return f(*args, **kwargs)
  File "<stdin>", line 1, in <lambda>
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.
py", line 1532, in __getattr__
    raise AttributeError(item)
AttributeError: notifications

```

```

    at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.
handlePythonException(PythonRunner.scala:456)
    at org.apache.spark.api.python.PythonRunner$$anon$1.read(Python
Runner.scala:592)
    at org.apache.spark.api.python.PythonRunner$$anon$1.read(Python
Runner.scala:575)
    at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.
hasNext(PythonRunner.scala:410)
    at org.apache.spark.InterruptibleIterator.hasNext(Interruptible
Iterator.scala:37)
    at scala.collection.Iterator$class.foreach(Iterator.scala:891)
    at org.apache.spark.InterruptibleIterator.foreach(Interruptible
Iterator.scala:28)
    at scala.collection.generic.Growable$class.$plus$plus$eq(Growab
le.scala:59)
    at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuff
er.scala:104)
    at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuff
er.scala:48)
    at scala.collection.TraversableOnce$class.to(TraversableOnce.sc
ala:310)
    at org.apache.spark.InterruptibleIterator.to(InterruptibleItera
tor.scala:28)
    at scala.collection.TraversableOnce$class.toBuffer(TraversableOn
ce.scala:302)
    at org.apache.spark.InterruptibleIterator.toBuffer(Interruptibl
eIterator.scala:28)
    at scala.collection.TraversableOnce$class.toArray(TraversableOn
ce.scala:289)
    at org.apache.spark.InterruptibleIterator.toArray(Interruptible
Iterator.scala:28)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)

```

```

    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkContext.scala:2101)
    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkContext.scala:2101)
    at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.scala:90)
    at org.apache.spark.scheduler.Task.run(Task.scala:123)
    at org.apache.spark.executor.Executor$TaskRunner$$anonfun$10.apply(Executor.scala:408)
    at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:1360)
    at org.apache.spark.executor.Executor$TaskRunner.run(Executor.scala:414)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1149)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
    ... 1 more

```

Traceback (most recent call last):

```

File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/rdd.py", line 816, in collect
    sock_info = self.ctx._jvm.PythonRDD.collectAndServe(self._jrdd.rdd())
File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/java_gateway.py", line 1257, in __call__
    answer, self.gateway_client, self.target_id, self.name)
File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/sql/utils.py", line 63, in deco
    return f(*a, **kw)
File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/protocol.py", line 328, in get_return_value
    format(target_id, ".", name), value)
py4j.protocol.Py4JJavaError: An error occurred while calling z:org.apache.spark.api.python.PythonRDD.collectAndServe.
: org.apache.spark.SparkException: Job aborted due to stage failure: Task 0 in stage 5.0 failed 4 times, most recent failure: Lost task 0.3 in stage 5.0 (TID 11, ip-172-31-12-174.ec2.internal, executor 19): org.apache.spark.api.python.PythonException: Traceback (most recent call last):
File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.py", line 1527, in __getattr__
    idx = self.__fields__.index(item)
ValueError: 'notifications' is not in list

```

During handling of the above exception, another exception occurred:

Traceback (most recent call last):

```

File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.py", line 377, in main
    process()
File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_0001/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.py", line 372, in process
    serializer.dump_stream(func(split_index, iterator), outfile)

```

```

File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/serializer
s.py", line 393, in dump_stream
    vs = list(itertools.islice(iterator, batch))
File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/util.py",
line 113, in wrapper
    return f(*args, **kwargs)
File "<stdin>", line 1, in <lambda>
File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.
py", line 1532, in __getattr__
    raise AttributeError(item)
AttributeError: notifications

```

```

    at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.
handlePythonException(PythonRunner.scala:456)
    at org.apache.spark.api.python.PythonRunner$$anon$1.read(Python
Runner.scala:592)
    at org.apache.spark.api.python.PythonRunner$$anon$1.read(Python
Runner.scala:575)
    at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.
hasNext(PythonRunner.scala:410)
    at org.apache.spark.InterruptibleIterator.hasNext(Interruptible
Iterator.scala:37)
    at scala.collection.Iterator$class.foreach(Iterator.scala:891)
    at org.apache.spark.InterruptibleIterator.foreach(Interruptible
Iterator.scala:28)
    at scala.collection.generic.Growable$class.$plus$plus$eq(Growab
le.scala:59)
    at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuff
er.scala:104)
    at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuff
er.scala:48)
    at scala.collection.TraversableOnce$class.to(TraversableOnce.sc
ala:310)
    at org.apache.spark.InterruptibleIterator.to(InterruptibleItera
tor.scala:28)
    at scala.collection.TraversableOnce$class.toBuffer(TraversableO
nce.scala:302)
    at org.apache.spark.InterruptibleIterator.toBuffer(Interruptibl
eIterator.scala:28)
    at scala.collection.TraversableOnce$class.toArray(TraversableOn
ce.scala:289)
    at org.apache.spark.InterruptibleIterator.toArray(Interruptible
Iterator.scala:28)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)
    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkC
ontext.scala:2101)
    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkC
ontext.scala:2101)
    at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.sca
la:90)
    at org.apache.spark.scheduler.Task.run(Task.scala:123)

```

```

    at org.apache.spark.executor.Executor$TaskRunner$$anonfun$10.ap
ply(Executor.scala:408)
    at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:
1360)
    at org.apache.spark.executor.Executor$TaskRunner.run(Executor.s
cala:414)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPool
Executor.java:1149)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoo
lExecutor.java:624)
    at java.lang.Thread.run(Thread.java:748)

```

#### Driver stacktrace:

```

    at org.apache.spark.scheduler.DAGScheduler.org$apache$spark$sch
eduler$DAGScheduler$$failJobAndIndependentStages(DAGScheduler.scala:204
1)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage
$1.apply(DAGScheduler.scala:2029)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage
$1.apply(DAGScheduler.scala:2028)
    at scala.collection.mutable.ResizableArray$class.foreach(Resiza
bleArray.scala:59)
    at scala.collection.mutable.ArrayBuffer.foreach(ArrayBuffer.sca
la:48)
    at org.apache.spark.scheduler.DAGScheduler.abortStage(DAGSchedu
ler.scala:2028)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskS
etFailed$1.apply(DAGScheduler.scala:966)
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskS
etFailed$1.apply(DAGScheduler.scala:966)
    at scala.Option.foreach(Option.scala:257)
    at org.apache.spark.scheduler.DAGScheduler.handleTaskSetFailed
(DAGScheduler.scala:966)
    at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.doOn
Receive(DAGScheduler.scala:2262)
    at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.onRe
ceive(DAGScheduler.scala:2211)
    at org.apache.spark.scheduler.DAGSchedulerEventProcessLoop.onRe
ceive(DAGScheduler.scala:2200)
    at org.apache.spark.util.EventLoop$$anon$1.run(EventLoop.scala:
49)
    at org.apache.spark.scheduler.DAGScheduler.runJob(DAGScheduler.
scala:777)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:206
1)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:208
2)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:210
1)
    at org.apache.spark.SparkContext.runJob(SparkContext.scala:212
6)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1.apply(RDD.scala:
945)
    at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperati
onScope.scala:151)
    at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperati
onScope.scala:112)

```

```

    at org.apache.spark.rdd.RDD.withScope(RDD.scala:363)
    at org.apache.spark.rdd.RDD.collect(RDD.scala:944)
    at org.apache.spark.api.python.PythonRDD$.collectAndServe(Python
nRDD.scala:166)
    at org.apache.spark.api.python.PythonRDD.collectAndServe(Python
RDD.scala)
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAcce
ssorImpl.java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMe
thodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at py4j.reflection.MethodInvoker.invoke(MethodInvoker.java:244)
    at py4j.reflection.ReflectionEngine.invoke(ReflectionEngine.jav
a:357)
    at py4j.Gateway.invoke(Gateway.java:282)
    at py4j.commands.AbstractCommand.invokeMethod(AbstractCommand.j
ava:132)
    at py4j.commands.CallCommand.execute(CallCommand.java:79)
    at py4j.GatewayConnection.run(GatewayConnection.java:238)
    at java.lang.Thread.run(Thread.java:748)
Caused by: org.apache.spark.api.python.PythonException: Traceback (most
recent call last):
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.
py", line 1527, in __getattr__
    idx = self.__fields__.index(item)
ValueError: 'notifications' is not in list

```

During handling of the above exception, another exception occurred:

```

Traceback (most recent call last):
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.p
y", line 377, in main
    process()
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/worker.p
y", line 372, in process
    serializer.dump_stream(func(split_index, iterator), outfile)
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/serializer
s.py", line 393, in dump_stream
    vs = list(itertools.islice(iterator, batch))
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/util.py",
line 113, in wrapper
    return f(*args, **kwargs)
  File "<stdin>", line 1, in <lambda>
  File "/mnt/yarn/usercache/livy/appcache/application_1572637375434_000
1/container_1572637375434_0001_01_000087/pyspark.zip/pyspark/sql/types.
py", line 1532, in __getattr__
    raise AttributeError(item)
AttributeError: notifications

```

```

    at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.
handlePythonException(PythonRunner.scala:456)

```



```

    at org.apache.spark.api.python.PythonRunner$$anon$1.read(Python
Runner.scala:592)
    at org.apache.spark.api.python.PythonRunner$$anon$1.read(Python
Runner.scala:575)
    at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.
hasNext(PythonRunner.scala:410)
    at org.apache.spark.InterruptibleIterator.hasNext(Interruptible
Iterator.scala:37)
    at scala.collection.Iterator$class.foreach(Iterator.scala:891)
    at org.apache.spark.InterruptibleIterator.foreach(Interruptible
Iterator.scala:28)
    at scala.collection.generic.Growable$class.$plus$plus$eq(Growab
le.scala:59)
    at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuff
er.scala:104)
    at scala.collection.mutable.ArrayBuffer.$plus$plus$eq(ArrayBuff
er.scala:48)
    at scala.collection.TraversableOnce$class.to(TraversableOnce.sc
ala:310)
    at org.apache.spark.InterruptibleIterator.to(InterruptibleItera
tor.scala:28)
    at scala.collection.TraversableOnce$class.toBuffer(TraversableO
nce.scala:302)
    at org.apache.spark.InterruptibleIterator.toBuffer(Interruptibl
eIterator.scala:28)
    at scala.collection.TraversableOnce$class.toArray(TraversableOn
ce.scala:289)
    at org.apache.spark.InterruptibleIterator.toArray(Interruptible
Iterator.scala:28)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)
    at org.apache.spark.rdd.RDD$$anonfun$collect$1$$anonfun$13.appl
y(RDD.scala:945)
    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkC
ontext.scala:2101)
    at org.apache.spark.SparkContext$$anonfun$runJob$5.apply(SparkC
ontext.scala:2101)
    at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.sca
la:90)
    at org.apache.spark.scheduler.Task.run(Task.scala:123)
    at org.apache.spark.executor.Executor$TaskRunner$$anonfun$10.ap
ply(Executor.scala:408)
    at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:
1360)
    at org.apache.spark.executor.Executor$TaskRunner.run(Executor.s
cala:414)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPool
Executor.java:1149)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoo
lExecutor.java:624)
    ... 1 more

```

In [19]:

6

In [ ]:

```
# map
rdd = congress_sample
sorted(rdd.map(lambda x: (x, 1)).collect())
```

In [ ]:

```
#Reduce
wordsRDD = sc.parallelize(congress_sample, 100)
wordCountsCollected = (wordsRDD
                        .map(lambda w: (w, 1))
                        .reduceByKey(lambda x,y: x+y)
                        .collect())
```

In [ ]:

In [39]:

```
congress_sample = spark.read.text("s3://us-congress-tweets/raw/files.txt")  
congress_sample.collect()
```



[illegible]

[illegible]

[illegible]

[illegible]



```
ress-tweets/raw/part-00269.snappy'), Row(value='s3://us-congress-tweet
s/raw/part-00270.snappy'), Row(value='s3://us-congress-tweets/raw/part-
00271.snappy'), Row(value='s3://us-congress-tweets/raw/part-00272.snapp
y'), Row(value='s3://us-congress-tweets/raw/part-00273.snappy'), Row(va
lue='s3://us-congress-tweets/raw/part-00274.snappy'), Row(value='s3://u
s-congress-tweets/raw/part-00275.snappy'), Row(value='s3://us-congress-
tweets/raw/part-00276.snappy'), Row(value='s3://us-congress-tweets/raw/
part-00277.snappy'), Row(value='s3://us-congress-tweets/raw/part-00278.
snappy'), Row(value='s3://us-congress-tweets/raw/part-00279.snappy'), R
ow(value='s3://us-congress-tweets/raw/part-00280.snappy'), Row(value='s
3://us-congress-tweets/raw/part-00281.snappy'), Row(value='s3://us-cong
ress-tweets/raw/part-00282.snappy'), Row(value='s3://us-congress-tweet
s/raw/part-00283.snappy'), Row(value='s3://us-congress-tweets/raw/part-
00284.snappy'), Row(value='s3://us-congress-tweets/raw/part-00285.snapp
y'), Row(value='s3://us-congress-tweets/raw/part-00286.snappy'), Row(va
lue='s3://us-congress-tweets/raw/part-00287.snappy'), Row(value='s3://u
s-congress-tweets/raw/part-00288.snappy'), Row(value='s3://us-congress-
tweets/raw/part-00289.snappy'), Row(value='s3://us-congress-tweets/raw/
part-00290.snappy'), Row(value='s3://us-congress-tweets/raw/part-00291.
snappy'), Row(value='s3://us-congress-tweets/raw/part-00292.snappy'), R
ow(value='s3://us-congress-tweets/raw/part-00293.snappy'), Row(value='s
3://us-congress-tweets/raw/part-00294.snappy'), Row(value='s3://us-cong
ress-tweets/raw/part-00295.snappy'), Row(value='s3://us-congress-tweet
s/raw/part-00296.snappy'), Row(value='s3://us-congress-tweets/raw/part-
00297.snappy'), Row(value='s3://us-congress-tweets/raw/part-00298.snapp
y'), Row(value='s3://us-congress-tweets/raw/part-00299.snappy'), Row(va
lue='s3://us-congress-tweets/raw/part-00300.snappy'), Row(value='s3://u
s-congress-tweets/raw/part-00301.snappy'), Row(value='s3://us-congress-
tweets/raw/part-00302.snappy'), Row(value='s3://us-congress-tweets/raw/
part-00303.snappy'), Row(value='s3://us-congress-tweets/raw/part-00304.
snappy'), Row(value='s3://us-congress-tweets/raw/part-00305.snappy'), R
ow(value='s3://us-congress-tweets/raw/part-00306.snappy'), Row(value='s
3://us-congress-tweets/raw/part-00307.snappy'), Row(value='s3://us-cong
ress-tweets/raw/part-00308.snappy'), Row(value='s3://us-congress-tweet
s/raw/part-00309.snappy'), Row(value='s3://us-congress-tweets/raw/part-
00310.snappy'), Row(value='s3://us-congress-tweets/raw/part-00311.snapp
y'), Row(value='s3://us-congress-tweets/raw/part-00312.snappy'), Row(va
lue='s3://us-congress-tweets/raw/part-00313.snappy'), Row(value='s3://u
s-congress-tweets/raw/part-00314.snappy'), Row(value='s3://us-congress-
tweets/raw/part-00315.snappy'), Row(value='s3://us-congress-tweets/raw/
part-00316.snappy'), Row(value='s3://us-congress-tweets/raw/part-00317.
snappy'), Row(value='s3://us-congress-tweets/raw/part-00318.snappy'), R
ow(value='s3://us-congress-tweets/raw/part-00319.snappy'), Row(value='s
3://us-congress-tweets/raw/part-00320.snappy'), Row(value='s3://us-cong
ress-tweets/raw/part-00321.snappy'), Row(value='s3://us-congress-tweet
s/raw/part-00322.snappy'), Row(value='s3://us-congress-tweets/raw/part-
00323.snappy'), Row(value='s3://us-congress-tweets/raw/part-00324.snapp
y'), Row(value='s3://us-congress-tweets/raw/part-00325.snappy'), Row(va
lue='s3://us-congress-tweets/raw/part-00326.snappy'), Row(value='s3://u
s-congress-tweets/raw/part-00327.snappy'), Row(value='s3://us-congress-
tweets/raw/part-00328.snappy'), Row(value='s3://us-congress-tweets/raw/
part-00329.snappy'), Row(value='s3://us-congress-tweets/raw/part-00330.
snappy'), Row(value='s3://us-congress-tweets/raw/part-00331.snappy'), R
ow(value='s3://us-congress-tweets/raw/part-00332.snappy'), Row(value='s
3://us-congress-tweets/raw/part-00333.snappy'), Row(value='s3://us-cong
ress-tweets/raw/part-00334.snappy'), Row(value='s3://us-congress-tweet
s/raw/part-00335.snappy'), Row(value='s3://us-congress-tweets/raw/part-
```

```
00336.snappy'), Row(value='s3://us-congress-tweets/raw/part-00337.snappy'), Row(value='s3://us-congress-tweets/raw/part-00338.snappy'), Row(value='s3://us-congress-tweets/raw/part-00339.snappy'), Row(value='s3://us-congress-tweets/raw/part-00340.snappy'), Row(value='s3://us-congress-tweets/raw/part-00341.snappy'), Row(value='s3://us-congress-tweets/raw/part-00342.snappy'), Row(value='s3://us-congress-tweets/raw/part-00343.snappy'), Row(value='s3://us-congress-tweets/raw/part-00344.snappy'), Row(value='s3://us-congress-tweets/raw/part-00345.snappy'), Row(value='s3://us-congress-tweets/raw/part-00346.snappy'), Row(value='s3://us-congress-tweets/raw/part-00347.snappy'), Row(value='s3://us-congress-tweets/raw/part-00348.snappy'), Row(value='s3://us-congress-tweets/raw/part-00349.snappy'), Row(value='s3://us-congress-tweets/raw/part-00350.snappy'), Row(value='s3://us-congress-tweets/raw/part-00351.snappy'), Row(value='s3://us-congress-tweets/raw/part-00352.snappy'), Row(value='s3://us-congress-tweets/raw/part-00353.snappy'), Row(value='s3://us-congress-tweets/raw/part-00354.snappy')]
```

In [ ]: