

Monica Santana – Week 4 Exercise

SparkPi Output

```
santanamonica@dscbigdata: ~/dsc650-info/bellevue-bigdata/hadoop-hive-spark-hbase
36402 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Default ResourceProfile created, executor resources: Map(cores -> name: cores, amount: 1, script: /, vendors: memory -> name: memory, amount: 1024, script: /, vendors: ), task resources: Map(cpu -> name: cpu, amount: 1.0)
38161 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Registered executor NettyRpcEndpointRef(spark-client://executor) (172.28.1.2:55700) with ID 2
38520 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 0.0 in stage 0.0 (TID 0, worker1, executor 2, partition 0, PROCESS_LOCAL, 7002, None)
38677 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 1.0 in stage 0.0 (TID 1, worker1, executor 2, partition 1, PROCESS_LOCAL, 7004, bytes)
39199 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 2.0 in stage 0.0 (TID 2, worker1, executor 2, partition 2, PROCESS_LOCAL, 7006, bytes)
41692 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 3.0 in stage 0.0 (TID 3, worker1, executor 2, partition 3, PROCESS_LOCAL, 7008, bytes)
41747 [task-result-getter-0] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 0.0 in stage 0.0 (TID 0) in 3118 ms on worker1 (executor 2) (1/10)
41894 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 4.0 in stage 0.0 (TID 4, worker1, executor 2, partition 4, PROCESS_LOCAL, 7010, bytes)
41910 [task-result-getter-1] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 1.0 in stage 0.0 (TID 1) in 218 ms on worker1 (executor 2) (2/10)
42056 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 5.0 in stage 0.0 (TID 5, worker1, executor 2, partition 5, PROCESS_LOCAL, 7012, bytes)
42060 [task-result-getter-2] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 2.0 in stage 0.0 (TID 2) in 169 ms on worker1 (executor 2) (3/10)
42141 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 6.0 in stage 0.0 (TID 6, worker1, executor 2, partition 6, PROCESS_LOCAL, 7014, bytes)
42143 [task-result-getter-3] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 3.0 in stage 0.0 (TID 3) in 88 ms on worker1 (executor 2) (4/10)
42261 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 7.0 in stage 0.0 (TID 7, worker1, executor 2, partition 7, PROCESS_LOCAL, 7016, bytes)
42262 [task-result-getter-0] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 4.0 in stage 0.0 (TID 4) in 122 ms on worker1 (executor 2) (5/10)
42385 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 8.0 in stage 0.0 (TID 8, worker1, executor 2, partition 8, PROCESS_LOCAL, 7018, bytes)
42390 [task-result-getter-1] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 5.0 in stage 0.0 (TID 5) in 129 ms on worker1 (executor 2) (6/10)
42480 [task-result-getter-2] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 6.0 in stage 0.0 (TID 6) in 103 ms on worker1 (executor 2) (7/10)
42492 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 9.0 in stage 0.0 (TID 9, worker1, executor 2, partition 9, PROCESS_LOCAL, 7020, bytes)
42566 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 8.0 in stage 0.0 (TID 8, worker1, executor 2, partition 8, PROCESS_LOCAL, 7018, bytes)
42570 [task-result-getter-3] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 7.0 in stage 0.0 (TID 7) in 77 ms on worker1 (executor 2) (8/10)
42683 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Starting task 9.0 in stage 0.0 (TID 9, worker1, executor 2, partition 9, PROCESS_LOCAL, 7020, bytes)
42686 [task-result-getter-0] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 8.0 in stage 0.0 (TID 8) in 120 ms on worker1 (executor 2) (9/10)
42765 [task-result-getter-1] INFO org.apache.spark.scheduler.TaskSetManager - Finished task 9.0 in stage 0.0 (TID 9) in 93 ms on worker1 (executor 2) (10/10)
42769 [dag-scheduler-event-loop] INFO org.apache.spark.scheduler.DAGScheduler - ResultStage 0 (reduce at SparkPi.scala:38) finished in 7.314 s
42774 [task-result-getter-1] INFO org.apache.spark.scheduler.TaskSetManager - Removed TaskSet 0.0, whose tasks have all completed, from pool
42788 [dag-scheduler-event-loop] INFO org.apache.spark.scheduler.DAGScheduler - Job 0 in finished. Cancelling potential speculative or zombie tasks for this job
42800 [dag-scheduler-event-loop] INFO org.apache.spark.scheduler.TaskSetManager - Killing all running tasks in stage 0: Stage finished
42810 [main] INFO org.apache.spark.scheduler.DAGScheduler - Job 0 finished: reduce at SparkPi.scala:38, took 7.637181 s
42816 [main] INFO org.apache.spark.scheduler.DAGScheduler - Job 0 finished: reduce at SparkPi.scala:38, took 7.637181 s
42832 [main] INFO org.apache.spark.scheduler.DAGScheduler - Job 0 finished: reduce at SparkPi.scala:38, took 7.637181 s
42888 [main] INFO org.apache.spark.scheduler.DAGScheduler - Job 0 finished: reduce at SparkPi.scala:38, took 7.637181 s
42896 [YARN application state monitor] INFO org.apache.spark.scheduler.TaskSetManager - Interrupting monitor thread
43039 [main] INFO org.apache.spark.scheduler.TaskSetManager - Shutting down all executors
43040 [dispatcher-coarseGrainedScheduler] INFO org.apache.spark.scheduler.TaskSetManager - Asking each executor to shut down
43069 [main] INFO org.apache.spark.scheduler.TaskSetManager - YARN client scheduler backend stopped
43127 [dispatcher-event-loop-1] INFO org.apache.spark.scheduler.TaskSetManager - MapOutputTrackerMasterEndpoint - MapOutputTrackerMasterEndpoint stopped!
43267 [main] INFO org.apache.spark.scheduler.TaskSetManager - MemoryStore cleared
43268 [main] INFO org.apache.spark.scheduler.TaskSetManager - BlockManager stopped
43379 [main] INFO org.apache.spark.scheduler.TaskSetManager - BlockManagerMaster stopped
43384 [dispatcher-event-loop-0] INFO org.apache.spark.scheduler.TaskSetManager - OutputCommitCoordinatorEndpoint - OutputCommitCoordinator stopped!
43413 [main] INFO org.apache.spark.scheduler.TaskSetManager - Successfully stopped SparkContext
43431 [shutdown-hook-0] INFO org.apache.spark.scheduler.TaskSetManager - Shutdown hook called
43441 [shutdown-hook-0] INFO org.apache.spark.scheduler.TaskSetManager - Deleting directory /tmp/spark-a40f37fd-ff27-4c49-a2ad-f7cc0260110a
43457 [shutdown-hook-0] INFO org.apache.spark.scheduler.TaskSetManager - Deleting directory /tmp/spark-10f24bbc-c05c-4b07-b0f1-b00c6bac404c
bach-5.0#
```

First 100 generated random numbers – First half

```
santanamonica@dscbigdata: ~/dsc650-info/bellevue-bigdata/hadoop-hive-spark-hbase
numbers: scala.collection.immutable.IndexedSeq[Int] = Vector(962, 78, 126, 159, 562, 848, 678, 91, 848, 197, 87, 386, 675, 913, 811, 563, 249, 110, 547, 291, 79, 3, 775, 798, 576, 285, 313, 99, 106, 831, 1, 628, 6, 680, 310, 490, 664, 252, 60, 4, 62, 405, 229, 610, 356, 667, 742, 393, 511, 413, 509, 729, 379, 480, 227, 479, 362, 609, 15, 987, 267, 355, 2, 114, 297, 902, 667, 200, 309, 90, 138, 113, 12, 6, 698, 587, 310, 330, 97, 26, 134, 784, 471, 600, 71, 947, 67, 879, 581, 21, 62, 3, 667, 573, 375, 779, 889, 9, 414, 991, 154, 880, 525, 816, 205, 729, 939, 662, 379, 112, 645, 116, 892, 37, 830, 332, 115, 122, 10, 711, 830, 384, 930, 53, 17, 6, 141, 554, 134, 584, 661, 287, 680, 934, 402, 308, 86, 718, 642, 707, 372, 861, 930, 0, 996, 560, 866, 873, 561, 873, 592, 994, 811, ...)
scala> val numbersRDD = sc.parallelize(numbers)
numbersRDD: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:26
scala> numbersRDD.take(100).foreach(println)
962
78
126
159
562
848
678
91
848
197
87
386
675
913
811
563
249
110
547
291
79
3
775
798
576
285
313
99
106
831
1
628
6
680
310
490
664
252
60
4
62
405
229
610
356
667
742
393
511
413
509
729
379
480
227
479
362
609
15
987
267
355
2
114
297
902
667
200
309
90
138
113
12
6
698
587
310
330
97
26
134
784
471
600
71
947
67
879
581
21
62
3
667
573
375
779
889
9
414
991
154
880
525
816
205
729
939
662
379
112
645
116
892
37
830
332
115
122
10
711
830
384
930
53
17
6
141
554
134
584
661
287
680
934
402
308
86
718
642
707
372
861
930
0
996
560
866
873
561
873
592
994
811
...
```

Second half

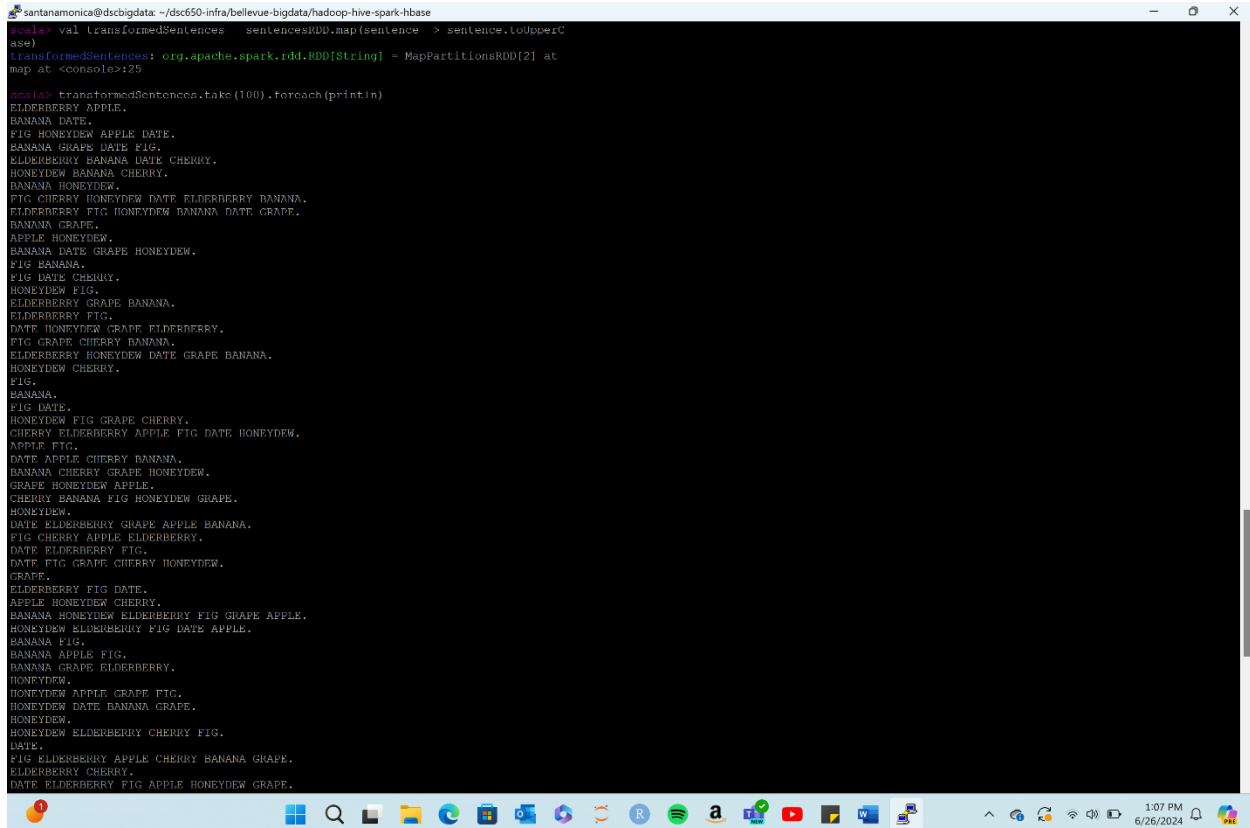
```
santanamonica@dscbigdata: ~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase
667
742
393
571
413
559
729
379
430
727
473
362
609
15
987
267
355
2
114
227
902
667
200
309
90
158
113
128
698
587
310
330
97
26
134
784
471
600
71
947
67
879
581
21
623
667
573
375
779
889
9
414
991
154
880
525
816

scala>
```

Custom transformations: 1st one is changing the words to uppercase

CODE:

```
val transformedSentences = sentencesRDD.map(sentence => sentence.toUpperCase)
transformedSentences.take(100).foreach(println)
```



```
santanamonica@dscbigdata: ~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase
scala> val transformedSentences = sentencesRDD.map(sentence => sentence.toUpperCase)
transformedSentences: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[12] at map at <console>:25
scala> transformedSentences.take(100).foreach(println)
ELDERBERRY APPLE.
BANANA DATE.
FIG HONEYDEW APPLE DATE.
BANANA GRAPE DATE FIG.
ELDERBERRY BANANA DATE CHERRY.
HONEYDEW BANANA CHERRY.
BANANA HONEYDEW.
FIG CHERRY HONEYDEW DATE ELDERBERRY BANANA.
ELDERBERRY FIG HONEYDEW BANANA DATE GRAPE.
BANANA GRAPE.
APPLE HONEYDEW.
BANANA DATE GRAPE HONEYDEW.
FIG BANANA.
FIG DATE CHERRY.
HONEYDEW FIG.
ELDERBERRY GRAPE BANANA.
ELDERBERRY FIG.
DATE HONEYDEW GRAPE ELDERBERRY.
FIG GRAPE CHERRY BANANA.
ELDERBERRY HONEYDEW DATE GRAPE BANANA.
HONEYDEW CHERRY.
FIG.
BANANA.
FIG DATE.
HONEYDEW FIG GRAPE CHERRY.
CHERRY ELDERBERRY APPLE FIG DATE HONEYDEW.
APPLE FIG.
DATE APPLE CHERRY BANANA.
BANANA CHERRY GRAPE HONEYDEW.
GRAPE HONEYDEW APPLE.
CHERRY BANANA FIG HONEYDEW GRAPE.
HONEYDEW.
DATE ELDERBERRY GRAPE APPLE BANANA.
FIG CHERRY APPLE ELDERBERRY.
DATE ELDERBERRY FIG.
DATE FIG GRAPE CHERRY HONEYDEW.
GRAPE.
ELDERBERRY FIG DATE.
APPLE HONEYDEW CHERRY.
BANANA HONEYDEW ELDERBERRY FIG GRAPE APPLE.
HONEYDEW ELDERBERRY FIG DATE APPLE.
BANANA FIG.
BANANA APPLE FIG.
BANANA GRAPE ELDERBERRY.
HONEYDEW.
HONEYDEW APPLE GRAPE FIG.
HONEYDEW DATE BANANA GRAPE.
HONEYDEW.
HONEYDEW ELDERBERRY CHERRY FIG.
DATE.
FIG ELDERBERRY APPLE CHERRY BANANA GRAPE.
ELDERBERRY CHERRY.
DATE ELDERBERRY FIG APPLE HONEYDEW GRAPE.
```

2nd transformation: Using the filter() and contains() to filter out only the words that have honey in them

CODE:

```
val transformedSentences = sentencesRDD.filter(sentence => sentence.contains("honey"))
transformedSentences.take(100).foreach(println)
```

```
santanamonica@dscbigdata: ~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase

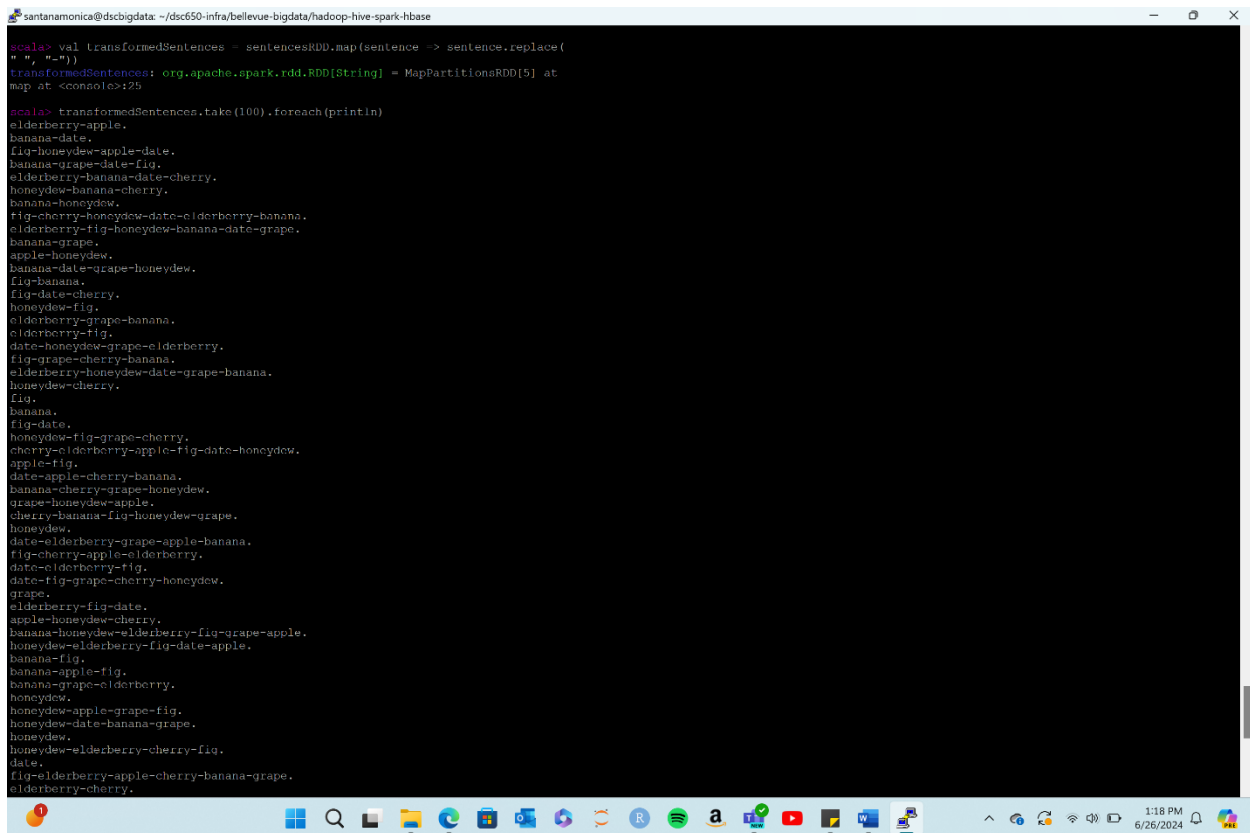
scala> val transformedSentences = sentencesRDD.filter(sentence => sentence.contains("honey"))
transformedSentences: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[4] at filter at <console>:125

scala> transformedSentences.take(100).foreach(println)
fig honeydew apple date.
honeydew banana cherry.
banana honeydew.
fig cherry honeydew date elderberry banana.
elderberry fig honeydew banana date grape.
apple honeydew.
banana date grape honeydew.
honeydew fig.
date honeydew grape elderberry.
elderberry honeydew date grape banana.
honeydew cherry.
honeydew fig grape cherry.
cherry elderberry apple fig date honeydew.
banana cherry grape honeydew.
grape honeydew apple.
cherry banana fig honeydew grape.
honeydew.
date fig grape cherry honeydew.
apple honeydew cherry.
banana honeydew elderberry fig grape apple.
honeydew elderberry fig date apple.
honeydew.
honeydew apple grape fig.
honeydew date banana grape.
honeydew.
honeydew elderberry cherry fig.
date elderberry fig apple honeydew grape.
banana apple honeydew date fig grape.
elderberry grape honeydew.
date elderberry apple honeydew.
honeydew.
grape banana apple honeydew fig.
apple date grape honeydew.
cherry fig honeydew apple.
date honeydew apple banana cherry elderberry.
grape honeydew elderberry cherry fig banana.
honeydew.
apple honeydew grape banana fig elderberry.
cherry elderberry grape banana honeydew apple.
elderberry honeydew apple.
date apple honeydew.
grape fig honeydew cherry elderberry banana.
apple honeydew grape date fig.
elderberry apple honeydew cherry banana grape.
date honeydew cherry.
grape banana elderberry honeydew.
cherry honeydew.
banana fig date honeydew apple.
elderberry date apple honeydew fig banana.
honeydew banana.
fig elderberry date honeydew banana apple.
```

3rd transformation: Using map() and the replace() to change the spaces to hyphens in the words

CODE:

```
val transformedSentences = sentencesRDD.map(sentence => sentence.replace(" ", "-"))
transformedSentences.take(100).foreach(println)
```



```
santanamonica@dscbigdata: ~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase

scala> val transformedSentences = sentencesRDD.map(sentence => sentence.replace(
" ", "-"))
transformedSentences: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[5] at
map at <console>:25

scala> transformedSentences.take(100).foreach(println)
elderberry-apple.
banana-date.
fig-honeydew-apple-date.
banana-grape-date-fig.
elderberry-banana-date-cherry.
honeydew-banana-cherry.
banana-honeydew.
fig-cherry-honeydew-date-elderberry-banana.
elderberry-fig-honeydew-banana-date-grape.
banana-grape.
apple-honeydew.
banana-date-grape-honeydew.
fig-banana.
fig-date-cherry.
honeydew-fig.
elderberry-grape-banana.
elderberry-fig.
date-honeydew-grape-elderberry.
fig-grape-cherry-banana.
elderberry-honeydew-date-grape-banana.
honeydew-cherry.
fig.
banana.
fig-date.
honeydew-fig-grape-cherry.
cherry-elderberry-apple-fig-date-honeydew.
apple-fig.
date-apple-cherry-banana.
banana-cherry-grape-honeydew.
grape-honeydew-apple.
cherry-banana-fig-honeydew-grape.
honeydew.
date-elderberry-grape-apple-banana.
fig-cherry-apple-elderberry.
date-elderberry-fig.
date-fig-grape-cherry-honeydew.
grape.
elderberry-fig-date.
apple-honeydew-cherry.
banana-honeydew-elderberry-fig-grape-apple.
honeydew-elderberry-fig-date-apple.
banana-fig.
banana-apple-fig.
banana-grape-elderberry.
honeydew.
honeydew-apple-grape-fig.
honeydew-date-banana-grape.
honeydew.
honeydew-elderberry-cherry-fig.
date.
fig-elderberry-apple-cherry-banana-grape.
elderberry-cherry.
```

4th transformation: Using filter() and startsWith() to filter out the words only starting with date and prints a true or false. True if the word starts with date and false if it doesn't

CODE:

```
val transformedSentences = sentencesRDD.filter(sentence => sentence.startsWith("date"))
transformedSentences.take(100).foreach(println)
```

[illegible]