

Building Real-time Analytics Dashboard using Apache Spark

Team #4:
Akshay Jain
Vinay Gor

Class: CSYE-7200 Big-Data Sys Engr Using Scala
Professor: Robin Hillyard

Github: https://github.com/akshaysjk/CSYE7200_Scala_Project_Team4

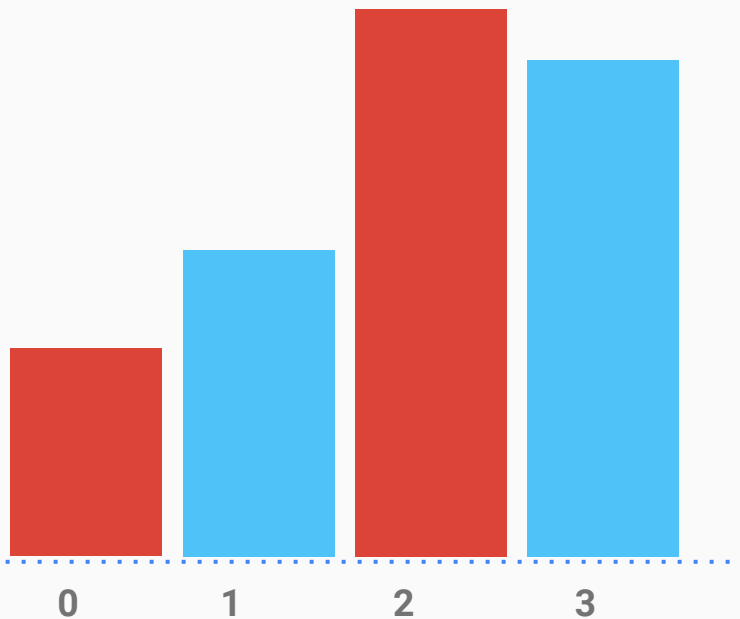
Goals of the Project

- To build a real-time analytics model using Stream Analytics.(Apache Spark)
- Reading of Data will be done in batches, to simulate real-time scenario.(Apache Kafka)
- Build a Real-time dashboard to display how the sales go on a particular day across different locations.
- Warehouse and inventory management at peak locations can be handled gracefully based on real-time analysis.

The problem

Batch Processing

The time delay between the collection of data and getting the result after the batch process.





The solution

Real-time Processing

Get the analytics in
real-time on Dashboard

Data Source

Data is taken from an ongoing competition on Analytics Vidhya website :
Practice Problem: Black Friday

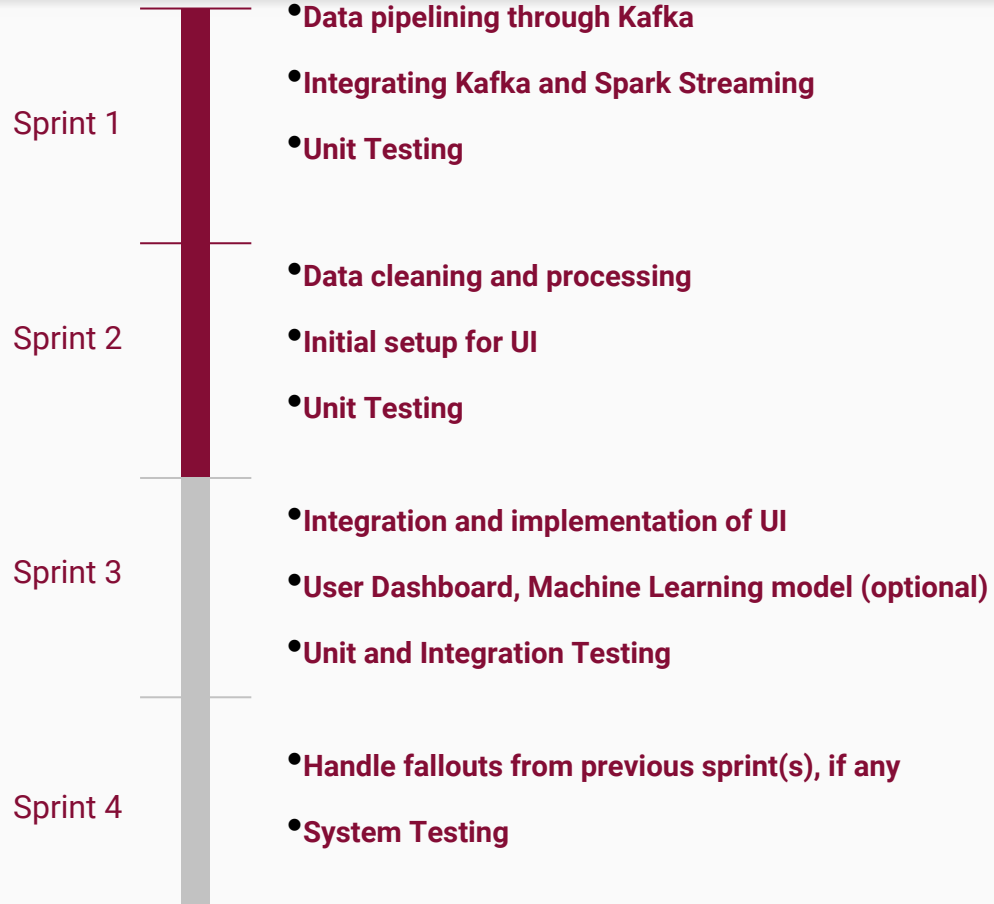
URL:

https://datahack.analyticsvidhya.com/contest/black-friday/#data_dictionary

Data :

train.csv : - consists of products and user details 0.5 million rows and 12 columns

Milestones / Sprints



Blockers faced

*Version compatibility with
Scala/Play/Spark*

Play Framework

Web Sockets

Methodology

Step 1

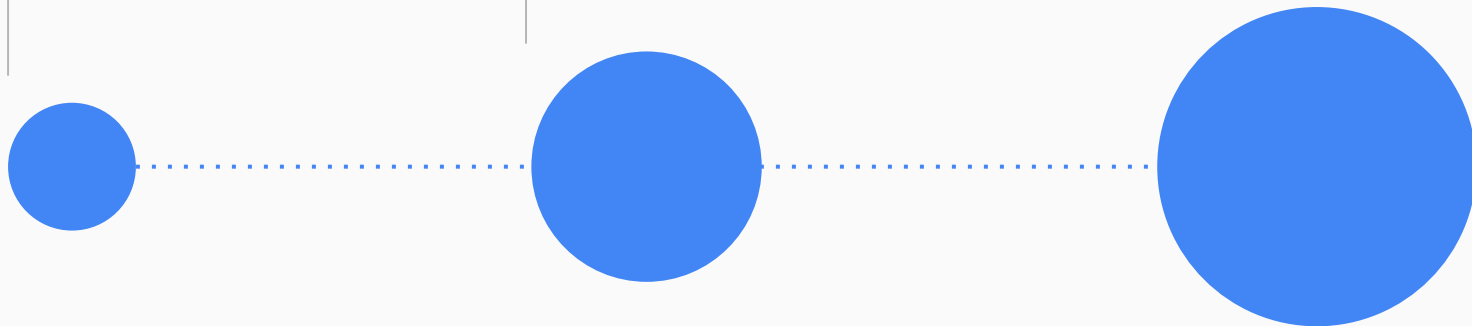
Read Data from CSV
through Apache Kafka

Step 2

Using Spark Streaming
to read data

Step 3

Clean the Data, run
analysis and pass data
to Dashboard through
Websockets



Technology Stack

Technologies:

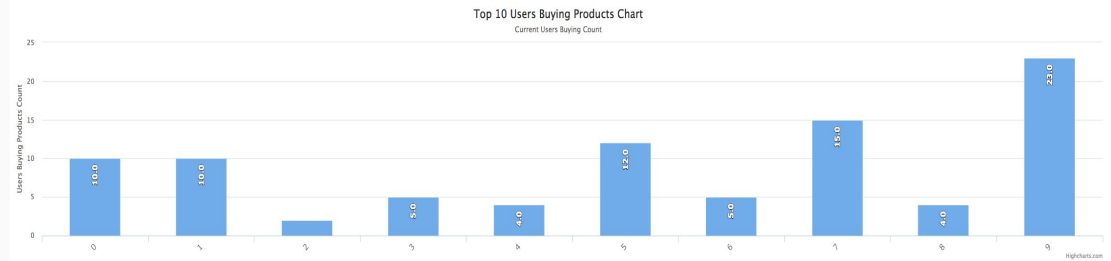
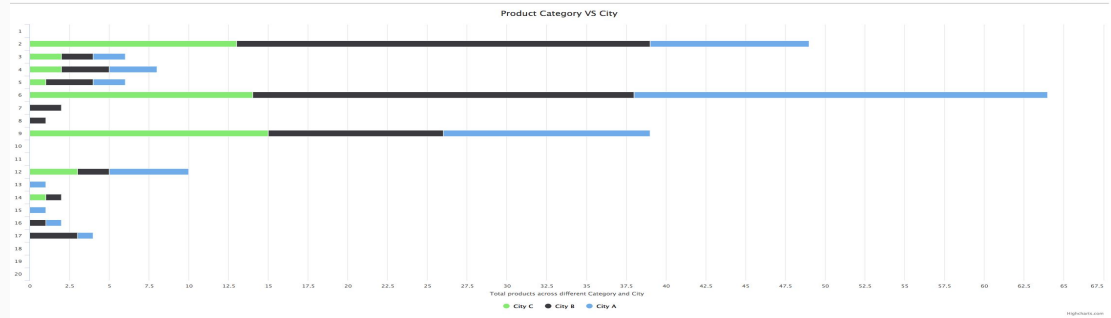
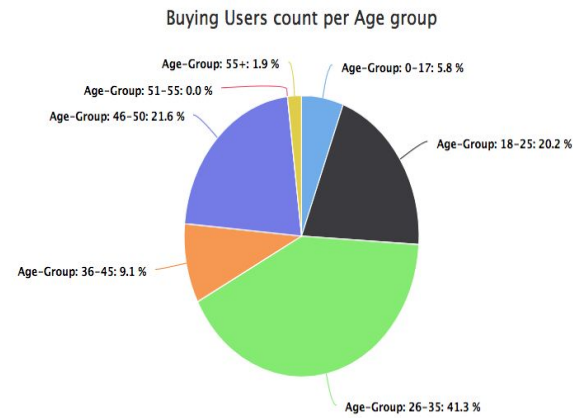
- Apache Kafka
- Spark Streaming
- Play Framework
- Web Socket Communication to pass data from Controller to Dashboard
- Highcharts.js for displaying charts
- Akka Actors for communication with Web Socket

Languages Used:

- Scala (66%)
- JavaScript (26%)
- HTML (6%)

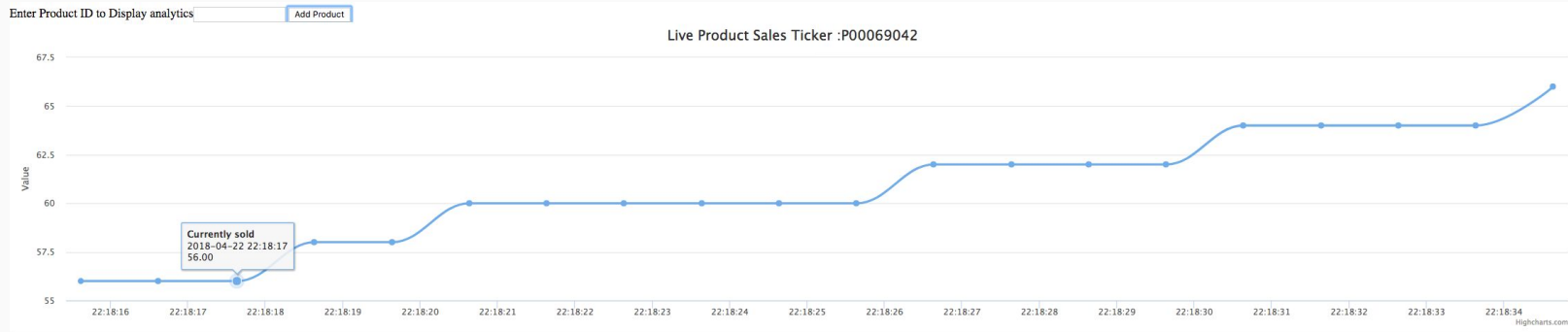
Use Case 1

Employee can see overview of the sales of products(eg: highest selling Product) on the dashboard homepage



Use case 2

Employee inputs query (such as product number) and gets the query specific real-time values.



```
[info] - should Summation of prodcategorySufferList
[info] FunctionalSpec:
[info] Routes
[info] - should send 404 on a bad request
[info] ScalaTest
[info] Run completed in 10 seconds, 407 milliseconds.
[info] Total number of tests run: 21
[info] Suites: completed 2, aborted 0
[info] Tests: succeeded 21, failed 0, canceled 0, ignored 0, pending 0
[info] All tests passed.
[info] Passed: Total 21, Failed 0, Errors 0, Passed 21
[success] Total time: 54 s, completed Apr 23, 2018 3:31:10 AM
[info] [04/23/2018 03:31:10.120] [Thread-3] [CoordinatedShutdown(akka://sbt-web)] Starting coordinated shutdown from JVM shutdown hook
```

```

The command "sbt ++$TRAVIS_SCALA_VERSION clean" exited with 0.
478 $ sbt ++$TRAVIS_SCALA_VERSION test
479 Picked up _JAVA_OPTIONS: -Xmx2048m -Xms512m
480 [info] Loading project definition from /home/travis/build/akshaysjk/CSYE7200_Scala_Project_Team4
481 [info] Loading settings from build.sbt ...
482 [info] Set current project to CSVKafka (in build file:/home/travis/build/akshaysjk/CSYE7200_Scala_Project_Team4)
483 [info] Setting Scala version to 2.10.4 on 0 projects.
484 [info] Excluded 1 projects, run ++ 2.10.4 -v for more details.
485 [info] Reapplying settings...
486 [info] Set current project to CSVKafka (in build file:/home/travis/build/akshaysjk/CSYE7200_Scala_Project_Team4)
487 [info] Updating ...
488 [info] Done updating.
489 [info] Compiling 2 Scala sources to /home/travis/build/akshaysjk/CSYE7200_Scala_Project_Team4/target/classes
490 [info] Done compiling.
491 [info] Compiling 1 Scala source to /home/travis/build/akshaysjk/CSYE7200_Scala_Project_Team4/target/classes
492 [info] Done compiling.
493 log4j:WARN No appenders could be found for logger (org.apache.kafka.clients.producer.ProducerConfig).
494 log4j:WARN Please initialize the log4j system properly.
495 log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
496 [info] KafkaSpec:
497 [info] - should match localhost
498 [info] - should match 9092
499 [info] - should match org.apache.kafka.common.serialization.StringSerializer for key
500 [info] - should match org.apache.kafka.common.serialization.StringSerializer for value
501 [info] - should match CSVKafka
502 [info] Run completed in 1 second, 471 milliseconds.
503 [info] Total number of tests run: 5
504 [info] Suites: completed 1, aborted 0
505 [info] Tests: succeeded 5, failed 0, canceled 0, ignored 0, pending 0
506 [info] All tests passed.
507 [success] Total time: 21 s, completed Apr 23, 2018 3:29:20 AM
508
509
510 The command "sbt ++$TRAVIS_SCALA_VERSION test" exited with 0.

```

Acceptance criteria

- 85% of the time, Spark Streaming will clean the data received, process it and generate/update the dashboard within 10 sec

$1524356551266 - 1524356550685 = 605 \text{ Milliseconds}$
~ 0.6 seconds

$1524356552174 - 1524356550728 = 1446$
~ 1.5 seconds

Time range ~ 0.6 to 5.3

Criteria Met!

```
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
Streaming Data Started!
[info] application - Received a message
[info] application - Logging input : ProductID P00102542 at time :1524356550661
[info] application - Logging input : ProductID P00273442 at time :1524356550685
[info] application - Logging input : ProductID P00281542 at time :1524356550686
[info] application - Logging input : ProductID P00367542 at time :1524356550687
[info] application - Logging input : ProductID P00253042 at time :1524356550688
```

Developer Tools - http://localhost:9000/startStreaming

Elements Console Sources Network Performance Memory Application

top Filter Default levels ▾ Gr

I am here

Logging result :=> P00273442 at time: 1524356551266

Logging result :=> P00367542 at time: 1524356551363

Logging result :=> P00273442 at time: 1524356551518

```
[info] application - Logging input : ProductID P00000142 at time :1524356550723
[info] application - Logging input : ProductID P00284642 at time :1524356550727
[info] application - Logging input : ProductID P00313342 at time :1524356550728
[info] application - Logging input : ProductID P00288342 at time :1524356550732
```

Logging result :=> P00273442 at time: 1524356552055

Logging result :=> P00313342 at time: 1524356552174

Logging result :=> P00313342 at time: 1524356552254

Future Scope

- Recommendation of product to Users with the help of Machine learning Algorithm
- User's Dashboard to display Analytics of the products purchased

Thank You!