Real Time Twitter Streaming

Bhargav Pabbisetti

Introduction

High Level Project Encapsulation

- Distributed Stream Processing » Apache
 Spark
- Machine Learning » Naive Bayes Classifier
 [Apache Spark MLlib implementation]
- Publish data to Redis
- Visualization » Sentiment visualization on a World map using Datamaps

What is Spark and why we should use it?

- Spark is a lightning fast cluster computing platform
 - Spark lets you run programs upto 100x faster in memory and 10x faster on disk than Hadoop.
 - Spark is ideal for iterative, interactive and event stream processing.
 - We can take advantage of many additional libraries that are built on top of Spark Core.

Spark Ecosystem

Spark SQL & Shark Spark Streaming real-time processing

MLlib machine learning

GraphX graph processing

Spark Core

Standalone Scheduler

YARN

Mesos

Resilient Distributed Dataset(RDD)

- Spark's primary abstraction
- It is a collection of elements that is partitioned across a cluster and it can be operated on in parallel.
- Immutable
- Two types of RDD operations
 - Transformations
 - Creates Direct Acyclic Graph(DAG)
 - Lazy evaluation
 - No return value
 - Actions
 - Performs the transformations and the action that follows
 - Returns a value
- Fault tolerance
- Caching

Shared Variables

- Two types of variables
 - Broadcast variables
 - Read only copy on each machine
 - Distribute broadcast variable using efficient broadcast algorithm
 - Accumulators
 - Variables added through an associative operation
 - Implement counters and sums
 - Only the driver can read the accumulator

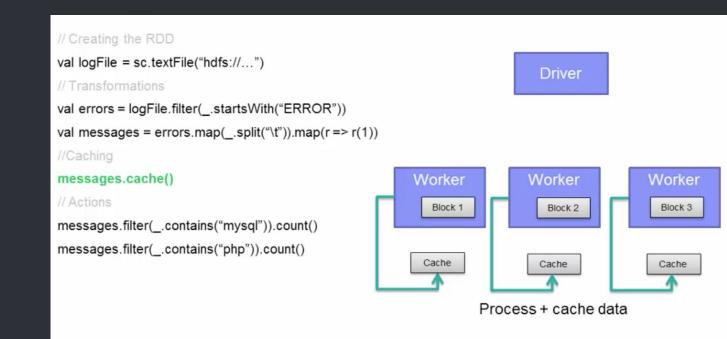
Basic Code Execution

```
// Creating the RDD
val logFile = sc.textFile("hdfs://...")
                                                                                  Driver
// Transformations
val errors = logFile.filter(_.startsWith("ERROR"))
val messages = errors.map(\_.split("\t")).map(r \Rightarrow r(1))
//Caching
messages.cache()
                                                             Worker
                                                                                Worker
                                                                                                     Worker
// Actions
                                                                                                       Block 3
                                                                Block 1
                                                                                    Block 2
messages.filter(_.contains("mysql")).count()
messages.filter(_.contains("php")).count()
                                                                       The data is partitioned into
```

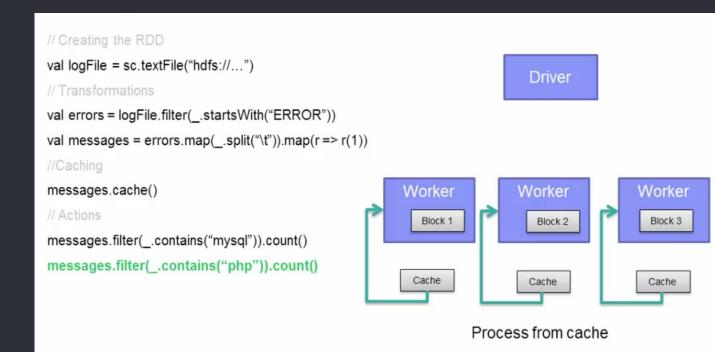
different blocks

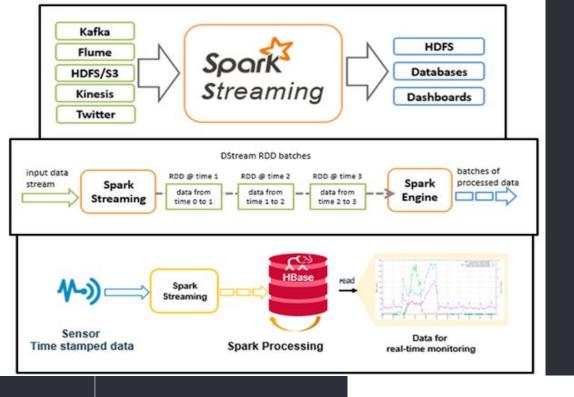
```
// Creating the RDD
val logFile = sc.textFile("hdfs://...")
                                                                                  Driver
// Transformations
val errors = logFile.filter(_.startsWith("ERROR"))
val messages = errors.map( .split("\t")).map(r \Rightarrow r(1))
//Cache
messages.cache()
                                                             Worker
                                                                                Worker
                                                                                                     Worker
// Actions
                                                                                                       Block 3
                                                                Block 1
                                                                                    Block 2
messages.filter(_.contains("mysql")).count()
messages.filter(_.contains("php")).count()
```

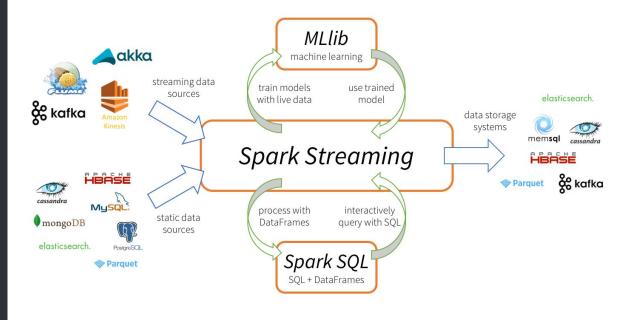
Driver sends the code to be executed on each block



```
// Creating the RDD
val logFile = sc.textFile("hdfs://...")
                                                                                Driver
// Transformations
val errors = logFile.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
//Caching
messages.cache()
                                                                              Worker
                                                                                                  Worker
                                                           Worker
// Actions
                                                                                                     Block 3
                                                              Block 1
                                                                                  Block 2
messages.filter(_.contains("mysql")).count()
messages.filter(_.contains("php")).count()
                                                             Cache
                                                                                Cache
                                                                                                     Cache
                                                                       Send the data back
                                                                       to the driver
```







Sample code using twitter4j

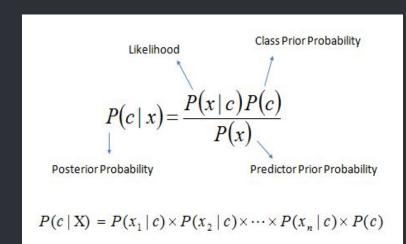
```
val oAuth: Some[OAuthAuthorization] = Authorization.TwitterAuth()
val rawTweets = TwitterUtils.createStream(ssc, oAuth) //twitter input dstreams
```

val classifiedTweets = rawTweets.filter(hasGeoLocation)
 .map(Sentimentprediction)

```
/**
* Predicts the sentiment of the tweet passed.
* Invokes Stanford Core NLP and MLlib methods for identifying the tweet sentiment.
 * # @param status -- twitter4j.Status object.
 * @return tuple with Tweet ID, Tweet Text, Core NLP Polarity, MLlib Polarity, Latitude, Longitude, Profile Image URL, Tweet Date.
*/
def Sentimentprediction(status: Status): (Long, String, String, Int, Int, Double, Double, String, String) = {
  val tweetText = replaceNewLines(status.getText)
  val (corenlpSentiment, mllibSentiment) = {
   // If tweet is in English, compute the sentiment by Mulib and also with Stanford CoreNLP.
   if (isTweetInEnglish(status)) {
      (CoreNLPSentimentAnalyzer.computeWeightedSentiment(tweetText),
       MLlibHelper.computeSentiment(tweetText, stopWordsList, naiveBayesModel))
    } else {
     // TODO: all non-English tweets are defaulted to neutral.
     // TODO: this is a workaround :: as we cant compute the sentiment of non-English tweets with our current model.
  (status.getId,
   status.getUser.getScreenName,
    tweetText,
    corenlpSentiment,
   mllibSentiment.
    status.getGeoLocation.getLatitude.
    status.getGeoLocation.getLongitude,
    status.getUser.getOriginalProfileImageURL,
    simpleDateFormat.format(status.getCreatedAt))
```

Spark MLlib

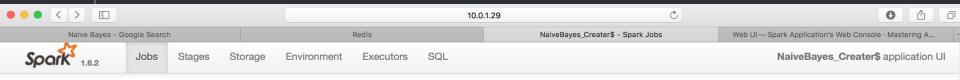
- Multinomial Naive Bayes is used to train 1.6 million tweets collected by Stanford University.
- Data is in csv format with 6 fields
 - 0 the polarity of the tweet(0=negative, 2=neutral, 4=positive)
 - 1 the id of the tweet
 - 2 the date of the tweet
 - 3 the user that tweeted
 - 4 the text of the tweet
- Term Frequency are the features.
- Achieved 78% accuracy.
- Pre trained model from stanford was also used to determine sentiment of a text
 - This model uses Recursive
 Neural Network that builds on top of grammatical structures



MLlib Code Snapshots

```
/**
* Creates a Naive Bayes Model of Tweet and its Sentiment from the Sentiment140 file.
* @param sc
                   — Spark Context.
* @param stopWordsList -- Broadcast variable for list of stop words to be removed from the tweets.
def createAndSaveNBModel(sc: SparkContext, stopWordsList: Broadcast[List[String]]): Unit = {
 val trainDF = loadTrainingFile(sc, PropertiesLoader.SENTIMENT140_TRAIN_DATA_PATH)
 val labeled_rdd = trainDF.select("POLARITY", "STATUS").rdd.map{
    case Row(polarity: Int, tweet: String) =>
     val tweetInWords: Seg[String] = MLlibHelper.cleanTweets(tweet, stopWordsList.value)
     LabeledPoint(polarity, MLlibHelper.transformFeatures(tweetInWords))
 labeled rdd.cache()
 val naiveBayesModel: NaiveBayesModel = NaiveBayes.train(labeled_rdd, lambda = 1.0, modelType = "multinomial")
 naiveBayesModel.save(sc, PropertiesLoader.NAIVEBAYES MODEL PATH)
```

Spark UI Mllib



Spark Jobs (?)

Total Uptime: 1.6 min Scheduling Mode: FIFO Active Jobs: 1 Completed Jobs: 11

▶ Event Timeline

Active Jobs (1)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
11	count at NaiveBayes_Creater.scala:107	2016/12/11 15:03:07	4 s	0/2	7/9

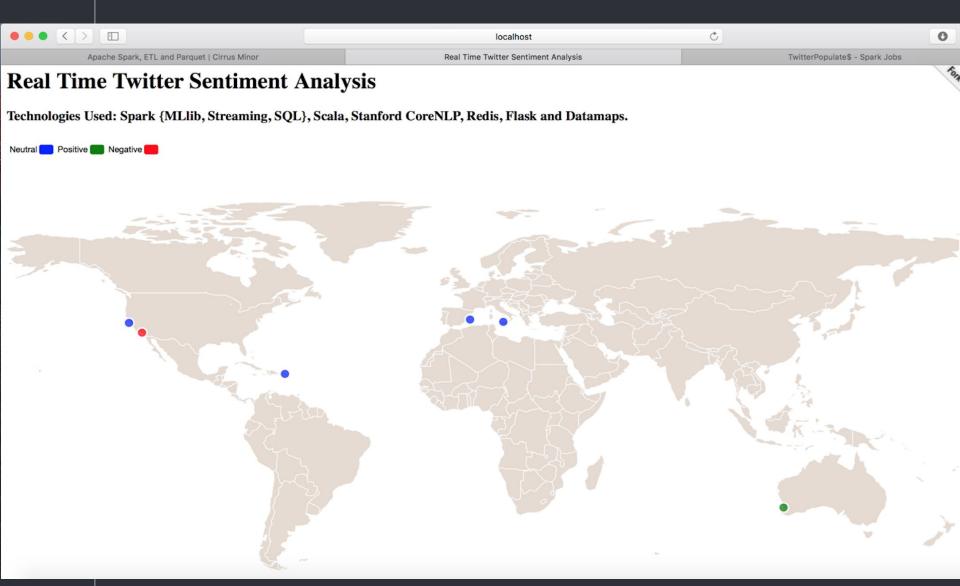
Completed Jobs (11)

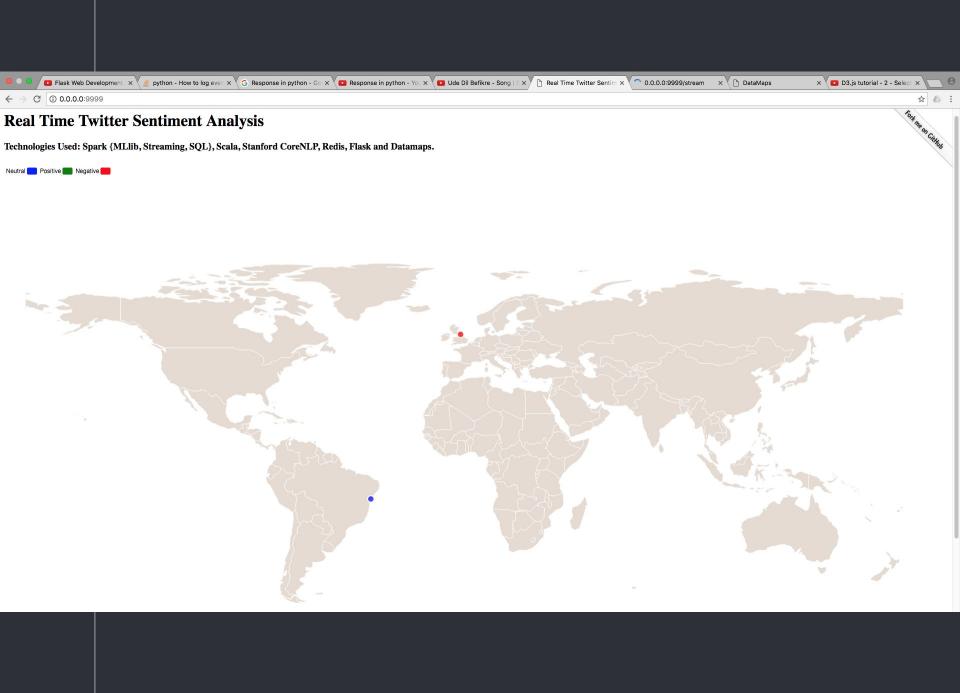
Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
10	count at NaiveBayes_Creater.scala:107	2016/12/11 15:02:28	38 s	1/1	8/8
9	aggregate at InferSchema.scala:41	2016/12/11 15:02:24	4 s	1/1	8/8
8	first at CsvRelation.scala:267	2016/12/11 15:02:24	8 ms	1/1	1/1
7	take at NaiveBayes.scala:216	2016/12/11 15:02:23	0.6 s	1/1	1/1
6	parquet at NaiveBayes.scala:213	2016/12/11 15:02:23	39 ms	1/1	2/2
5	first at modelSaveLoad.scala:129	2016/12/11 15:02:23	9 ms	1/1	1/1
4	parquet at NaiveBayes.scala:206	2016/12/11 15:02:22	1.0 s	1/1	1/1
3	saveAsTextFile at NaiveBayes.scala:202	2016/12/11 15:02:22	65 ms	1/1	1/1
2	collect at NaiveBayes.scala:401	2016/12/11 15:01:43	39 s	2/2	16/16
1	aggregate at InferSchema.scala:41	2016/12/11 15:01:37	5 s	1/1	8/8
0	first at CsvRelation.scala:267	2016/12/11 15:01:37	0.1 s	1/1	1/1

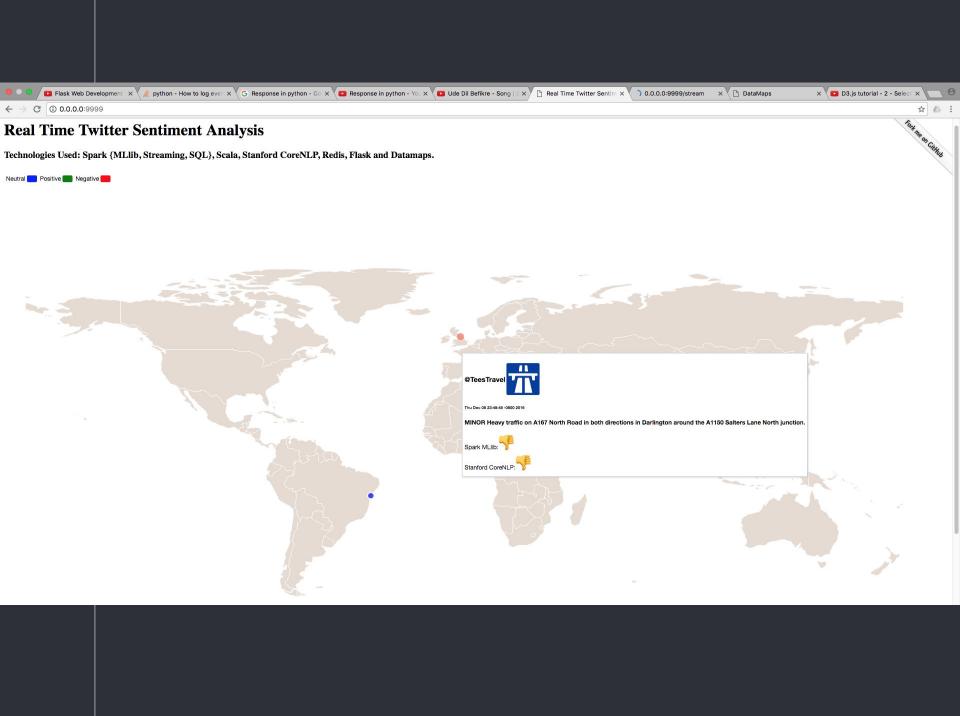
Web Server and Visualization

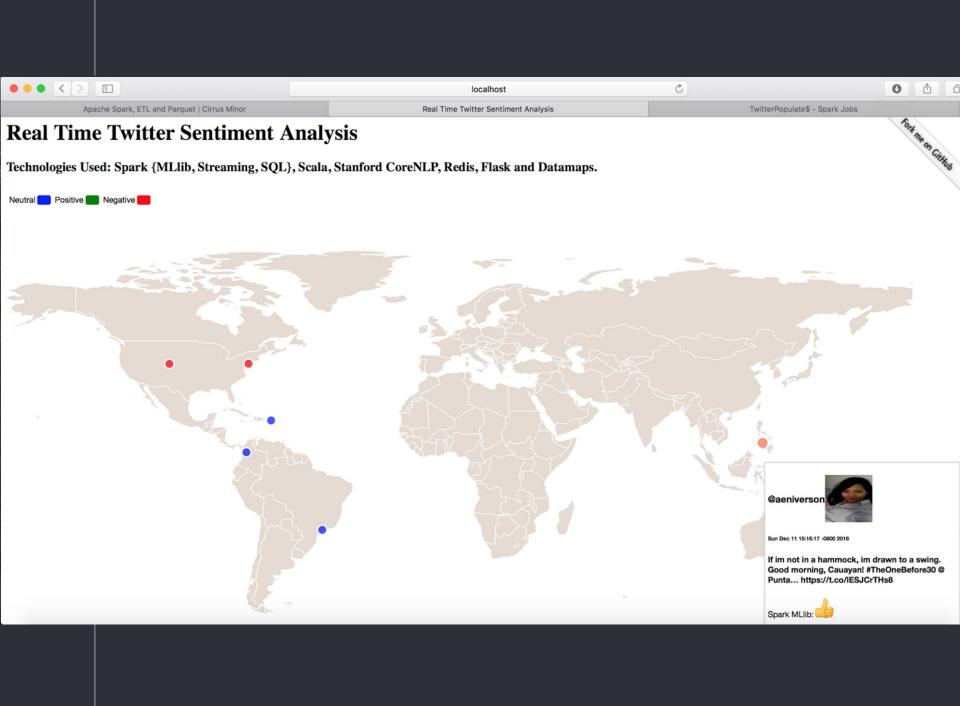
- Messaging System is built using Redis.
- Result of tweets are published to Redis which is subscribed by front end webapp.
- Web application uses server sent event.
- Python Flask web server
- HTML5, Javascript, D3 js are used to visualize real time tweets
- World map is created using datamaps.
- Bubbles represent live tweet.
- Hovering over a bubble reveals details of the tweet

Snapshots of Real Time Visualization









Further work and Improvement Areas

- Update project using Apache spark 2.0
 - Utilize Data Frames and Datasets
 - Use machine learning package org.apache.spark.ml
- Find methods to process non english tweets
- Visualization can be improved

Thanks!

ANY QUESTIONS?