# Public Opinion Analysis of Airlines

**Team 1:**

Fan Wu

Dayu Jia

Bowen Jiang

**Github link:**
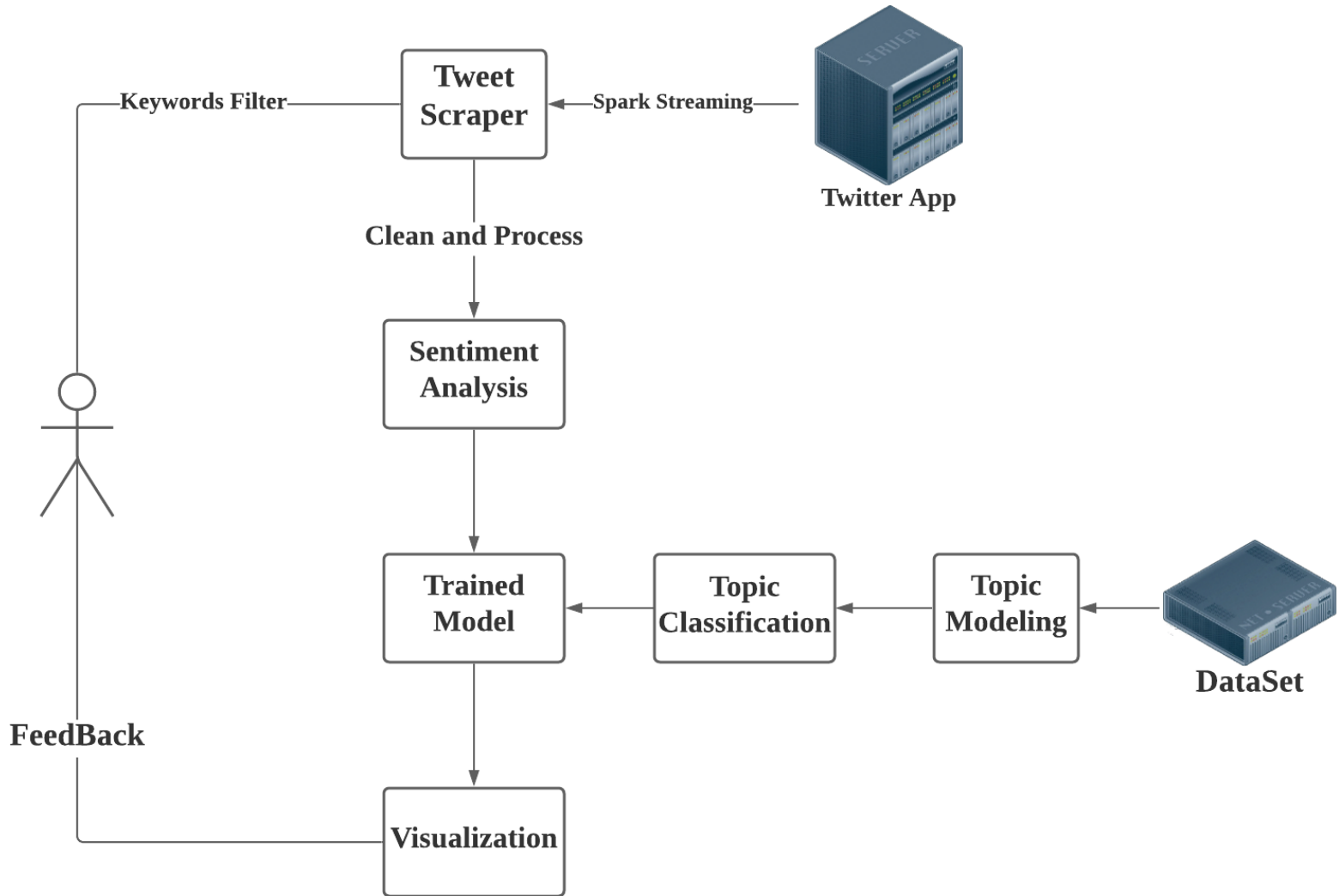
https://github.com/FanWu6/Spark-stream-twitter-analysis

# Use Cases

- User inputs the name of an airplane company

- and receives a comprehensive analysis on public opinion

# Methodology

- Tweets pulling-> Spark Streaming, twitter4j
- Data process -> StopWordsRemover, tf-idf, SparkSQL,etc...
- Topic modeling ->  Latent Dirichlet Allocation(LDA)
- Topic classification ->  Logistic Regression
- Sentiment analysis separately -> CoreNLP
- Visualization -> Elasticsearch, Kibana

# Data Sources



**Twitter US Airline Sentiment**

This dataset has 14485 rows and 30 columns

Source from:

https://www.kaggle.com/crowdflower/twitter-airline-sentiment

**Real-time data through Twitter API**

Approximately 150 tweets in 1 minute.

(Depend on keywords put in)

# Milestones

1st week:

Implement two ways to extract real-time data from Twitter; Choose LDA to do tweets topic modeling.

2nd week:

Implement data cleaning and LDA; Learn about topic classification algorithm and sentiment analysis.

3rd week:

Implement Logistic Regression and train the model; Sentiment analysis real-time tweets; Learn about Elasticsearch.
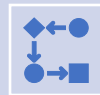
4th week:

Uploading data to Kibana and visualization; Make our code better.

# Acceptance criteria

- 60% of data can be classified correctly

- The accuracy of sentimental analysis should reach 80%

# Unit tests



```
val evaluator = new BinaryClassificationEvaluator().setMetricName("areaUnderROC")
println("The topic prediction accuracy of test data = " + (evaluator.evaluate(predictions)))
//
LogisticRegressionModel
```

LogisticRegressionModel ×

```
21/04/22 21:18:15 INFO SparkContext: Invoking stop() from shutdown hook
21/04/22 21:18:15 INFO SparkUI: Stopped Spark web UI at http://wufan-PC.lan:4040
The topic prediction accuracy of test data = 0.6761062958812633
```

```
val accuracy: Double = (accp * np + accn * nn + acco * no)/ (nn + no + np)
println("Accuracy of Sentiment: " + accuracy)
```

SentimentAccuracy

SentimentAccuracy ×

```
21/04/22 21:22:38 INFO SparkContext: Created broadcast 3 from rdd at Sentimen
21/04/22 21:22:38 INFO FileSourceScanExec: Planning scan with bin packing, ma
Accuracy of Sentiment: 0.9167457333333332
```

```
Preprocess success.
[info] PreprocessSpec:
[info] Preprocess
[info] - should Preprocess work
[info] Run completed in 37 seconds, 62 milliseconds.
[info] Total number of tests run: 15
[info] Suites: completed 6, aborted 0
[info] Tests: succeeded 15, failed 0, canceled 0, ignored 0, pending 0
[info] All tests passed.
[success] Total time: 57 s, completed 2021骞?鏈?1鏃?涓婂崍7:39:54:54
21/04/21 19:39:54 INFO ContextCleaner: Cleaned accumulator 5213
21/04/21 19:39:54 INFO ContextCleaner: Cleaned accumulator 5205
```
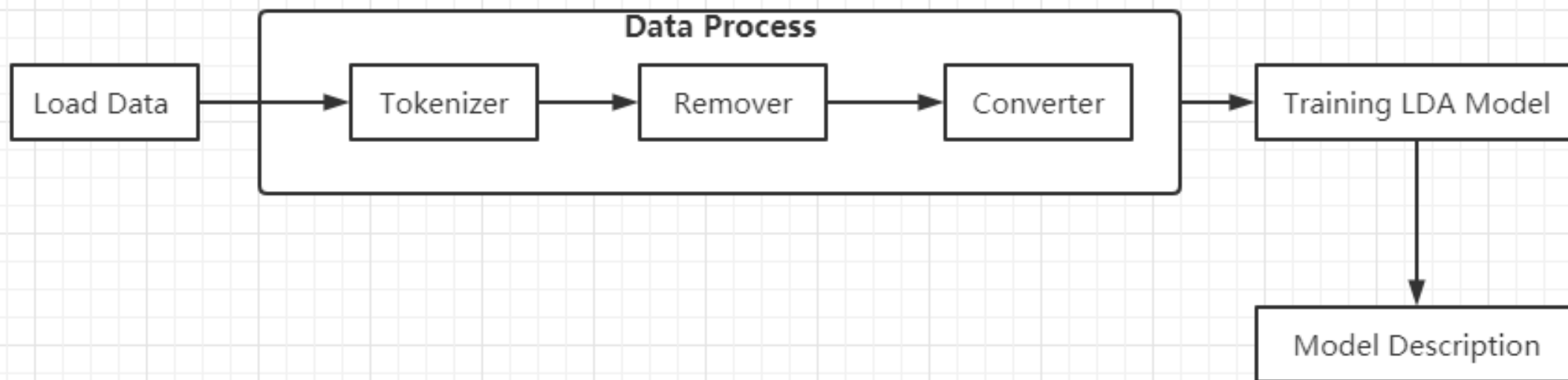
## Goals

- For any airline company, we can offer a real–time analysis on public opinion based on our model.

- For us, we want to learn the using of Scala and machine learning, and how to co-work on Github.

# Latent Dirichlet Allocation

# LDA Process Description

# Data Preprocess Code Description

```scala
//Tokenizing using the RegexTokenizer
val tokenizer = new RegexTokenizer()
  .setPattern("[\\W_]+")
  .setMinTokenLength(4)
  .setInputCol("corpus")
  .setOutputCol("tokens")

val tokenized_df: DataFrame = tokenizer.transform(corpus_df)

tokenized_df.select("tokens").show(10,false)

//Removing the Stop-words using the Stop Words remover
val add_stopwords = Array("http","jetblue","southwestair","americanair","flight",
  "usairways","thanks", "virginamerica","thank","today","flightled","united","please")
val stopwords = sparkSession.read.text("data/actualdata/stopwords.txt")
  .collect().map(row => row.getString(0)).union(add_stopwords)

val remover = new StopWordsRemover()
  .setStopWords(stopwords) // This parameter is optional
  .setInputCol("tokens")
  .setOutputCol("filtered")

val filtered_df: DataFrame = remover.transform(tokenized_df)

//Converting the Tokens into the CountVector
val vectorizer: CountVectorizerModel = new CountVectorizer()
  .setInputCol("filtered")
  .setOutputCol("features")
  .setVocabSize(10000)
  .setMinDF(5)
  .fit(filtered_df)

val countVectors: DataFrame= vectorizer.transform(filtered_df).select("id", "features")
countVectors.show(5,false)
import sparkSession.implicits._
val lda_countVector: RDD[(Long, linalg.Vector)] = countVectors.rdd.map {
  case Row(id: Long, countVector: Vector) => (id, Vectors.fromML(countVector)) }
(lda_countVector,vectorizer)
```

Using tokenizer filter away tokens with length <4

```
+------------------------------------------------------------------------------------+
|tokens                                                                              |
+------------------------------------------------------------------------------------+
|[virginamerica, what, dhepburn, said]                                              |
|[virginamerica, plus, added, commercials, experience, tacky]                       |
|[virginamerica, didn, today, must, mean, need, take, another, trip]                |
|[virginamerica, really, aggressive, blast, obnoxious, entertainment, your, guests, faces, they, have, little, recourse]|
|[virginamerica, really, thing, about]                                              |
|[virginamerica, seriously, would, flight, seats, that, didn, have, this, playing]  |
|[]                                                                                 |
|[virginamerica, nearly, every, time, this, worm, away]                             |
|[virginamerica, really, missed, prime, opportunity, without, hats, parody, there, https, mwpg7grezp]|
|[virginamerica, well, didn]                                                        |
+------------------------------------------------------------------------------------+
only showing top 10 rows
```

Remove stopwords

Set vocabulary size to 10000, and let the word show 5 times in each vocabulary

```
+---+-----------------------------------------------------------------------+
|id |features                                                               |
+---+-----------------------------------------------------------------------+
|0  |(2431,[92],[1.0])                                                      |
|1  |(2431,[74,331,839,1998],[1.0,1.0,1.0,1.0])                            |
|2  |(2431,[10,45,48,69,76,342,528],[1.0,1.0,1.0,1.0,1.0,1.0,1.0])         |
|3  |(2431,[37,288,698,2155],[1.0,1.0,1.0,1.0])                            |
|4  |(2431,[37,282],[1.0,1.0])                                             |
+---+-----------------------------------------------------------------------+
only showing top 5 rows
```

Completes the process of converting the documents into a vector of word counts

# LDA Code Description

# Some Optimizations To The Model

```
val lda = new LDA()
  .setOptimizer(new OnlineLDAOptimizer().setMiniBatchFraction(0.8))
  .setOptimizer("em")
  .setK(4)
  .setMaxIterations(100)
  .setDocConcentration(-1) // use default values
  .setTopicConcentration(-1)// use default values
```

- More iterations make the model more accurate

- Add new Stopwords to help filter

```
val add_stopwords  = Array("http","jetblue","southwestair","americanair","flight",
  "usairways","thanks", "virginamerica","thank","today","flightled","united","please")
val stopwords = sparkSession.read.text("data/actualdata/stopwords.txt")
  .collect().map(row => row.getString(0)).union(add_stopwords)
```

# Topic
# Classification

# Topic classification

For instance,

"We have the gold level plan and use it for everything, love the features! It is one of the best bang for buck possible."

A topic classification model that's been trained to understand these expressions (gold level plan, love the features, and best bang for buck) would be able to tag this review as topic of *Features* and *Price*.

In our case,

@VirginAmerica Is it me, or is your website down?  BTW, your new website isn't a great user experience.  Time for another redesign.

Customer Service Issue

# Data for training

```
flights 0.0167903690142048
cancelled    0.013848024722761337
hours    0.010398252437989924
help     0.009549276062319553
still    0.008408700272805191
late     0.00823532559798192
delayed  0.007862819032489875
hour     0.007674428383153015
like     0.007281068289706297
email    0.006798501598634885
==========
TOPIC 1
gate     0.014362314033549274
service  0.011501779434624122
time     0.01147332879552762
plane    0.01051470294529229
call     0.009846823449676802
cancelled    0.00818193145594124
hours    0.00812018480842092
hold     0.00796607052723924
help     0.00781181750360937
customer     0.007377330252671869
==========
TOPIC 2
```

late issue

customer service issue

- Late issue   ->   0
- Flight experience  issue  -> 2
- Customer  service  issue  -> 4
- Other type -> 6

# Logistic Regression

We have a bunch of data which has already been classified and it's used to train model.

Use Logistic Regression（**Spark MLlib**） to do this.

Use Logistic regression to predict the topic of text.

# Pre-process

```scala
//load training data
val newsgroupsRawData: RDD[String] = sc.textFile( path = "data/actualdata/train.csv")
//counts
println("The number of documents read in is " + newsgroupsRawData.count() + ".")

case class newsgroupsCaseClass(text: String, topic: String)

//remove all other character but words
val newsgroups: DataFrame = newsgroupsRawData.map{case (lines) =>
    val topic = lines.split( regex = ",").take(1)(0)
    val text = TrainingUtils.processText(lines)
    newsgroupsCaseClass(text,topic)}.toDF()
newsgroups.cache()

newsgroups.printSchema()

newsgroups.sample( withReplacement = false, fraction = 0.001, seed = 10L).show( numRows = 10, truncate = false)

newsgroups.groupBy( col1 = "topic").count().show()

//transform to another dataframe
val Labelednewsgroups = newsgroups.withColumn( colName = "label", newsgroups("topic").cast( to = "double"))

Labelednewsgroups.sample( withReplacement = false, fraction = 0.003, seed = 10L).show( numRows = 5, truncate = false)

//Split documents from a list of (id, text, label) tuples¶
val Array(training, test) = Labelednewsgroups.randomSplit(Array(0.9, 0.1), seed = 12345)
```
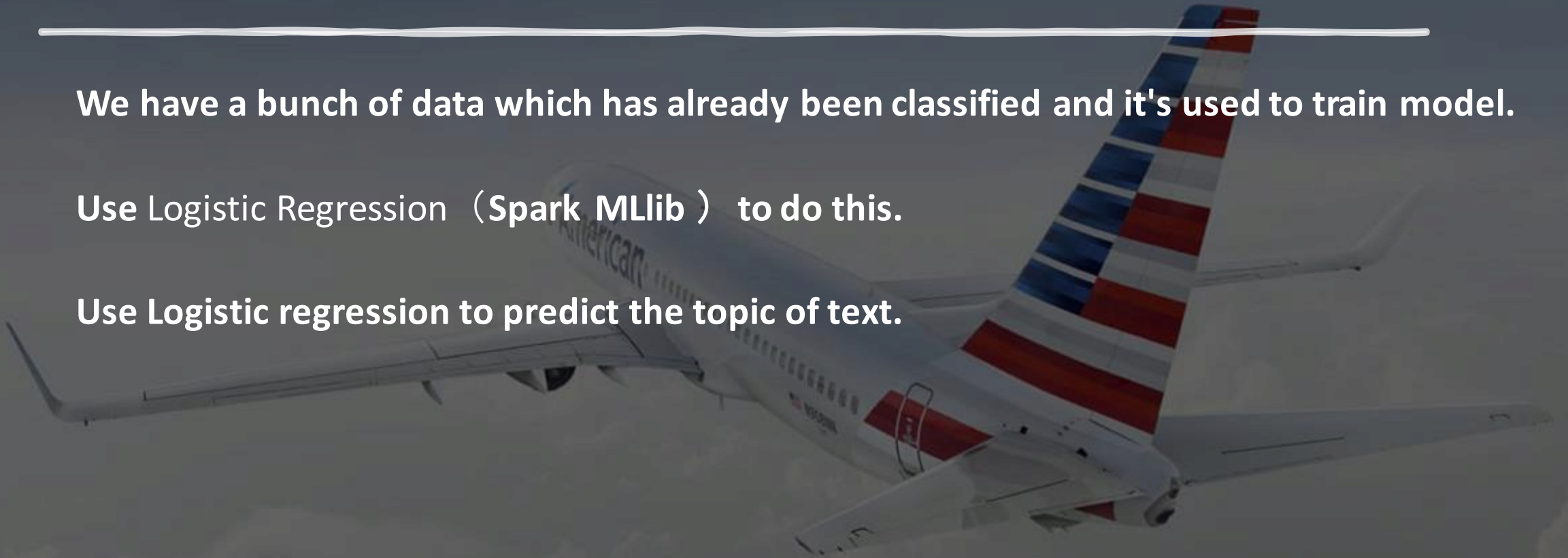
1. load data

2. extract column text, and remove all other charactor but only leave words

```
+--------------------------------------------------------------------------------------------------+-----+
|text                                                                                              |topic|
+--------------------------------------------------------------------------------------------------+-----+
| he was trying to take stuff from the under the seat in front of him and bugging out throughout the flight. feel safe. |2    |
| this means within one week i will have filed 2 compensation complaints to your website            |4    |
| my wife is trying to get a group of clients to their destination and just got disconnected after2 hours holding. Help.|4    |
```

3.cast topic to double

```
+--------------------------------------------------------------------------------------------------+-----+-----+
|text                                                                                              |topic|label|
+--------------------------------------------------------------------------------------------------+-----+-----+
| he was trying to take stuff from the under the seat in front of him and bugging out throughout the flight. feel safe. |2    |2.0  |
| this means within one week i will have filed 2 compensation complaints to your website            |4    |4.0  |
```

4.split data to training and test

# Main stages

- Tokenizer

- StopWordsRemover

- TF-IDF

- Logistic Regression

```
a
about
above
after
again
against
all
am
an
and
any
are
aren't
arent
as
at
be
because
```

```scala
val tokenizer = new Tokenizer().setInputCol("text").setOutputCol("words")
val remover = new StopWordsRemover().setInputCol("words").setOutputCol("filtered").setCaseSensitive(false)
val hashingTF = new HashingTF().setNumFeatures(1000).setInputCol("filtered").setOutputCol("rawFeatures")
val idf = new IDF().setInputCol("rawFeatures").setOutputCol("features").setMinDocFreq(0)
val lr = new LogisticRegression().setRegParam(0.01).setThreshold(0.5)
val pipeline = new Pipeline().setStages(Array(tokenizer, remover, hashingTF, idf, lr))
```

# Pipeline

# How to Optimize?

```
|                text|prediction|topic|
+--------------------+----------+-----+
| 2 and a half hou...|       0.0|    4|
| 2 hours on hold ...|       4.0|    4|
| I would love to ...|       6.0|    2|
| Im just praying ...|       0.0|    2|
| Instructions say...|       4.0|    4|
| Third flight in ...|       6.0|    2|
| Why offer automa...|       4.0|    4|
| Why you released...|       4.0|    4|
| Yo yo yo stuck o...|       0.0|    0|
| You respond to m...|       4.0|    4|
| airport and 2 ex...|       4.0|    4|
| anyone there to ...|       4.0|    4|
| but not sufficie...|       6.0|    4|
| flight 2031 wors...|       2.0|    2|
| has the worst cu...|       4.0|    4|
| have time and th...|       4.0|    4|
| if it is ever co  |       4.0|    4|
```

## 01

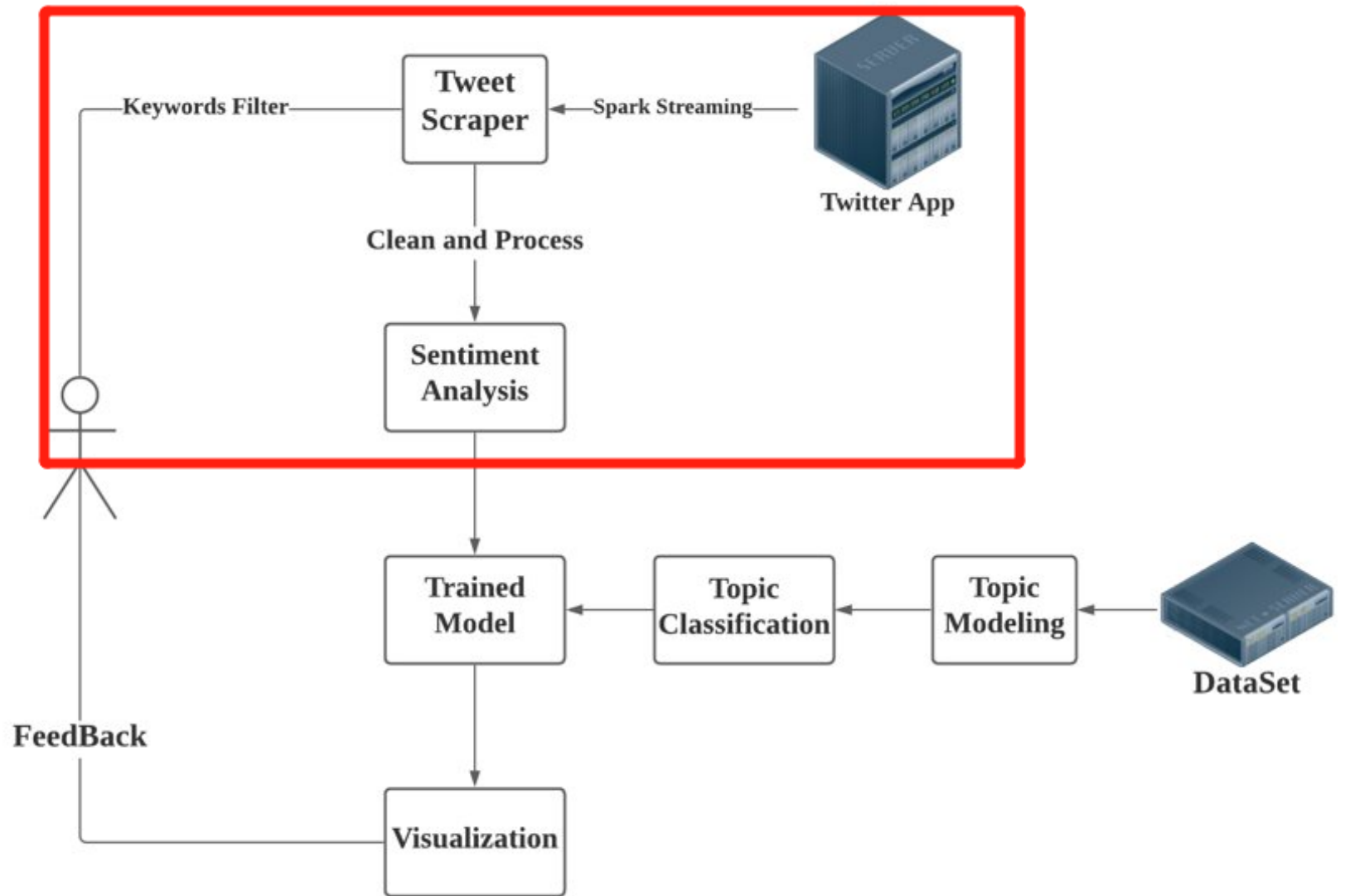**1.Tune the hyperparameters of model.**

## 02

**2.Make data more accurate.**
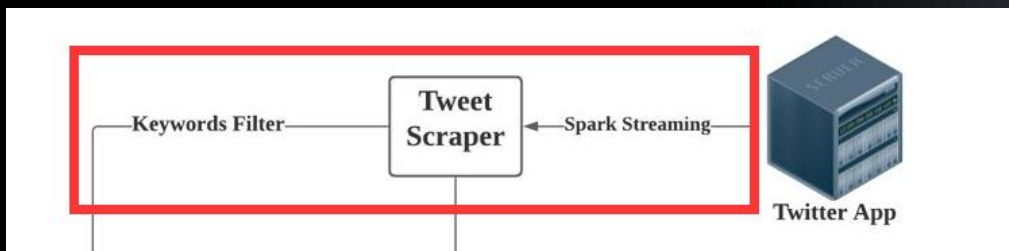
## 03

**3.Training with more data...**

Twitter Real-time Data

# Twitter Scraper

# Twitter Scraper

```scala
object TweetScrapper {
  def setupLogging(): Unit = {
    import org.apache.log4j.{Level, Logger}
    val rootLogger = Logger.getRootLogger
    rootLogger.setLevel(Level.ERROR)
  }

  /** Configures Twitter service credentials using twitter.txt
  def setupTwitter(): Unit = {
    import scala.io.Source

    val lines = Source.fromFile("data/actualdata/twitter.txt")
    for (line <- lines.getLines) {
      val fields = line.split( regex = " ")
      if (fields.length == 2) {
        System.setProperty("twitter4j.oauth." + fields(0), fields(1))
      }
    }
    lines.close()
  }
```

```scala
def main(args: Array[String]) {

  // Configure Twitter credentials using twitter.txt
  setupTwitter()

  // Set up a Spark streaming context named "PopularHashtags" that runs locally using
  // all CPU cores and one-second batches of data
  val ssc = new StreamingContext( master = "local[*]", appName = "PopularHashtags", Seconds(1))

  // Get rid of log spam (should be called after the context is set up)
  setupLogging()

  // Create a DStream from Twitter using our streaming context
  val keywords =  Configure.tweetfiltersc.getString( path = "KEYWORDS").split( regex = ",").toSeq
  println(keywords)
  val tweets = TwitterUtils.createStream(ssc, None,keywords)
  // Now extract the text of each status update into DStreams using map()
  val statuses: DStream[String] = tweets.filter(t=>t.getLang()=="en").map(status => status.getText)


  val spark = SparkSession.builder
    .master(Configure.sparkc.getString( path = "MASTER_URL"))
    .appName( name = "TweetStream")
    .getOrCreate()

  spark.sparkContext.setLogLevel("ERROR")
```

# Twitter Scraper with Kafka

```scala
def main(args: Array[String]): Unit ={
  // set log level
  setupLogging()


  val cb = new ConfigurationBuilder
  cb.setDebugEnabled(true)
    .setOAuthConsumerKey(Configure.twitter.getString( path = "CONSUMER_KEY"))
    .setOAuthConsumerSecret(Configure.twitter.getString( path = "CONSUMER_KEY_SECRET"))
    .setOAuthAccessToken(Configure.twitter.getString( path = "ACCESS_TOKEN"))
    .setOAuthAccessTokenSecret(Configure.twitter.getString( path = "ACCESS_TOKEN_SECRET"))
    .setJSONStoreEnabled(true)


  //create kafka props
  val props = new Properties()
  props.put("bootstrap.servers", Configure.kafkac.getString( path = "BOOTSTRAP_SERVERS"));
  props.put("key.serializer", "org.apache.kafka.common.serialization.StringSerializer"); // Str
  props.put("value.serializer", "org.apache.kafka.common.serialization.StringSerializer"); // S

  val producer = new KafkaProducer[String, String](props)
  val kafkatopic = Configure.kafkac.getString( path = "TOPIC")
  val statusListener = new StatusListener {
    /*
    StatusListener defines what to do with the tweets as they stream
    */
```
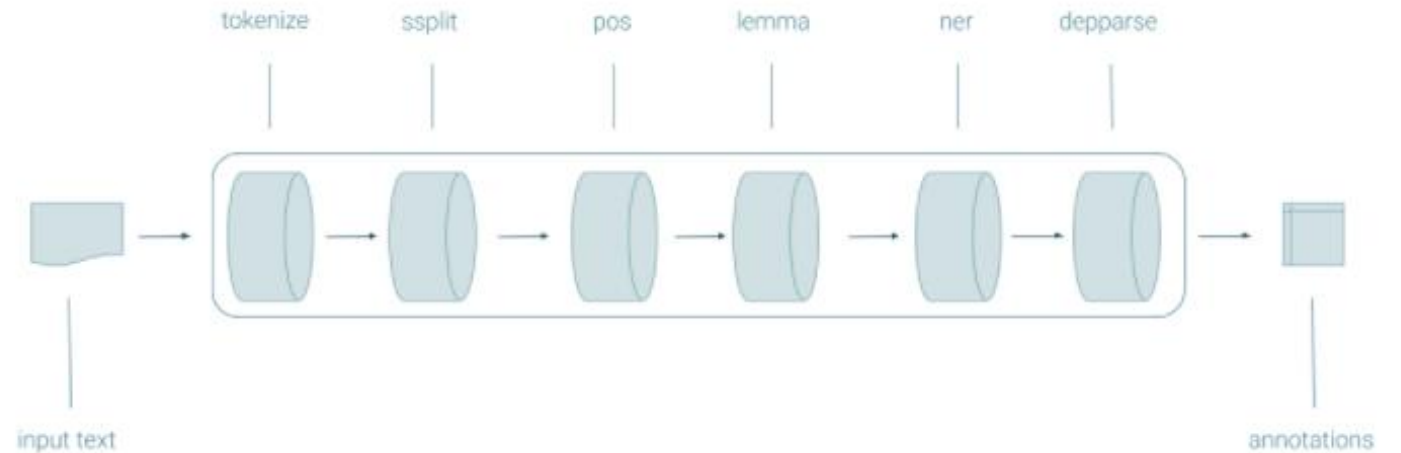
# Sentimental analysis with CoreNLP

```
object CleanTweets {

    // Data from twitter would be like "RT @User: xxx", clean it
    def clean(input: String): String = input.substring(input.indexOf(": ") + 2)

}
```

# Upload to Elastic



```scala
/**
 * Upload the data processed to Elastic Search
 * @param dataFrame Data whose type is dataframe
 * @param path Elastic Search path
 */
def dataFrameToElastic(dataFrame: DataFrame, path: String): Unit = {
  dataFrame.write
    .format( source = "org.elasticsearch.spark.sql")
    .option("es.port", 9200)
    .option("es.nodes", "localhost")
    .mode( saveMode = "append")
    .save(path)
}
```

# Upload to Elastic

```
Text,senti,classify
"This is a test0",positive,0.0
"This is a test2",positive,2.0
```

GET localhost:9200/test0419/doc/_search

Params  Authorization  Headers (7)  Body  Pre-request Script  Tests  Se

Body  Cookies  Headers (4)  Test Results

Pretty  Raw  Preview  Visualize  JSON

```
16          "hits": [
17              {
18                  "_index": "test0419",
19                  "_type": "doc",
20                  "_id": "GWWA6HgBAtphAarWfeOY",
21                  "_score": 1.0,
22                  "_source": {
23                      "Text": "This is a test0",
24                      "senti": "positive",
25                      "classify": "0.0"
26                  }
27              },
28              {
29                  "_index": "test0419",
30                  "_type": "doc",
31                  "_id": "GmWA6HgBAtphAarWfeOY",
32                  "_score": 1.0,
33                  "_source": {
34                      "Text": "This is a test2",
35                      "senti": "positive",
36                      "classify": "2.0"
37                  }
38              },
39
```

Download as CSV    Cancel    Save to library    Save and return

Search                                                                    KQL    📅 ⌄    Last 4 days                        Show dates    🔄 Refresh

⊜ —  + Add filter

**Quick select**                                                              ‹  ›

**tweetsairline**                        ⌄

▦ Stacked bar ⌄          #  ≣          Last ⌄        4          days ⌄        **Apply**

🔍 Search field names              ✕

**Commonly used**

Field filters  0                  ⌄
                                                        Today                    Last 24 hours

# Records                                               This week                Last 7 days

                                                        Last 15 minutes          Last 30 days

⌄ Available fields ⑦           3                        Last 30 minutes          Last 90 days

t  airline_sentiment.keyword                            Last 1 hour              Last 1 year

#  prediction

t  text.keyword                        **Kibana**        **Recently used date ranges**

                                                        Last 4 days

› Empty fields ⑦              0                         This week

› Meta fields                3                          Last 2 days

                                                        Today

                                                        Last 24 hours

        **Drop some fields here to start**

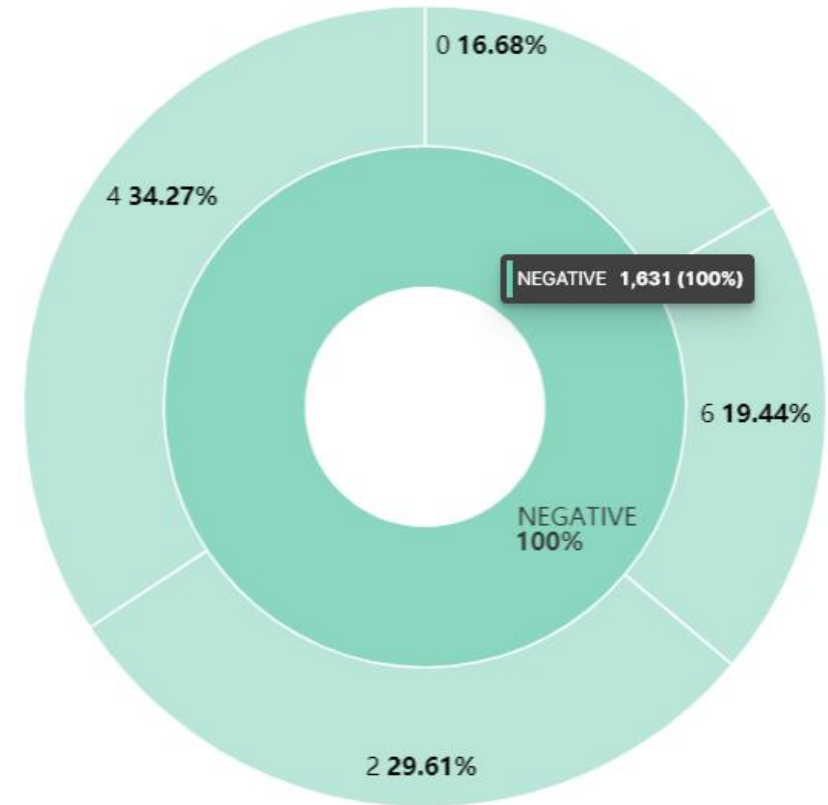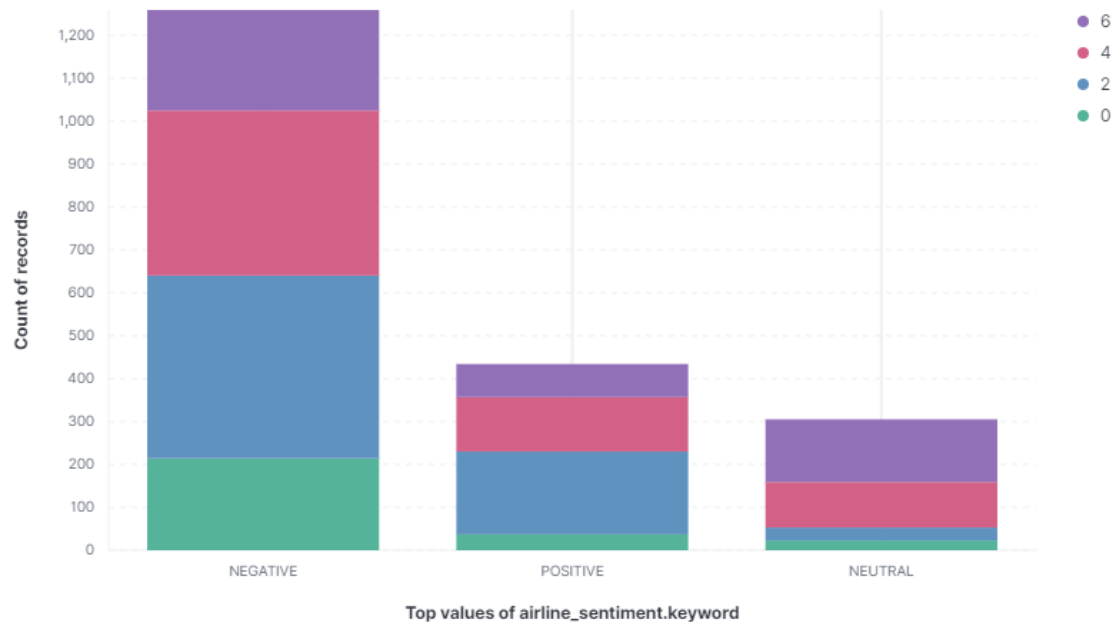                                                        **Refresh every**

                                                        2            seconds ⌄        ▷ **Start**

        Lens is a new tool for creating visualization

        Make requests and give feedback ⧉

# Example

Thank you for watching !