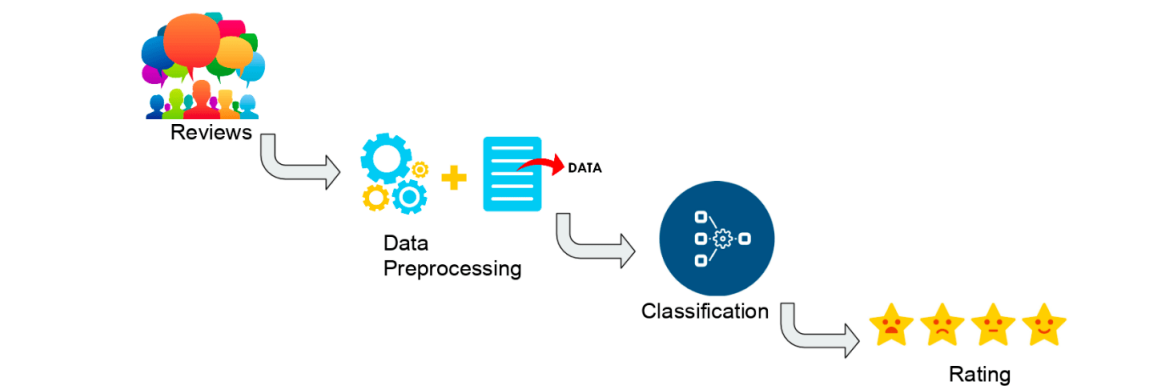


# CS 4830: Final Project

Team **Lockdown**: Kavish Shah, Akshara B, Vishal H, Vishnu Harshith  
*Indian Institute of Technology Madras*

## Introduction

The purpose of this project is to predict ratings based on massive data set containing Yelp reviews



**Fig 1:** Workflow for finding correlation b/w correlation between reviews and ratings

## Steps Involved

1. Batch Computation:  
Used a DataProc Cluster and submit a Spark job for data pre-processing and model training. Store a model in GCS bucket, and submit a spark job for evaluation on validation dataset stored in the bucket
2. Real Time Computation:  
Streamed data stored on the GCS bucket into Kafka and used spark streaming to read the data and make real-time predictions.

## Data-Set description

The data set is a Jason file, roughly of 9GB size and has the following columns:

Field name	Type	Mode	Policy tags ⓘ	Description
text	STRING	NULLABLE		
cool	INTEGER	NULLABLE		
funny	INTEGER	NULLABLE		
stars	FLOAT	NULLABLE		
date	TIMESTAMP	NULLABLE		bq-datetime
business_id	STRING	NULLABLE		
user_id	STRING	NULLABLE		
useful	INTEGER	NULLABLE		
review_id	STRING	NULLABLE		

**Fig: 2** Different columns in the yelp data set obtained using BigQuery

## Context

This dataset is a subset of Yelp’s businesses, reviews, and user data. Various platforms are present for people express their views on a matter. To get a general public opinion one may have to read all the views on the matter. This leads to a subclass of data mining problems known as sentiment analysis. We present below below the various pre-processing techniques adopted in this problem which may typically be used to work on any other massive data set as well.

Pre-processing

Loading the data

Loaded the JSON formatted data into the google cloud bucket using gs util command. Imported the data into BigQuery table and defined the schema.

Cleaning the Data

We perform the pre-processing of the data in a DataProc Cluster using Spark.The major hurdle dealing with text based models is that it consumes a lot of time pre-processing it. We have used many pre-processing techniques for this project.

	text cool funny stars			date	business_id		user_id useful		review_id
	I just had a pedi...	0	0	2.0 2014-08-03 21:38:45	MJ5xkXJ2UwJL60QH...	MXhsNk1DQw5qvwieZ...	3 912_cfGWNvf0nZudP...		
	We scheduled an a...	0	0	1.0 2015-05-18 19:22:54	r4cdabWEmTr2110Fu...	NXAX_SSLw2hxiAms...	3 V74eUostrSH5wpMgF...		
	[This old favorite...	1	1	4.0 2009-01-11 03:35:34	9DTypMuJ56GQZYdUW...	25N2f7A17Tgu9zQE-	2 6PyhrNqHroIvxZ6YD...		
	[This place was hu...	0	0	3.0 2017-07-03 03:15:45	7JWwexqil03Ua_wV9...	s72nhH1H1_lN_Ju4h...	0 Z_ShCWTEdEEvca9Fs...		
	[Went there to wat...	0	0	3.0 2013-09-20 04:24:55	Y7T9osP-hoop2nLnm...	xSt--5JG6QqTEOwXJ...	0 jsf1I9Sy1w3VvNkg1...		
	[Played cornhole f...	1	0	4.0 2017-10-07 23:35:14	3Jq5LfJ5fmJ5KmuA6...	dXP2z-Mqr1k2p2uKU...	1 yp5g42u10-bBEnUH...		
	[Invited a friend ...	0	0	1.0 2015-05-05 21:08:44	dQWE63u0iSfaf3Fup...	_IR48ok0ZkPMWJ2P1...	0 4kn_D8-XXseuNuGsv...		
	[What an incredibl...	0	0	5.0 2013-03-28 23:11:14	-4TMqnQJW1yd6NqGR...	cf-n5eNXJE7DG0iJb...	0 ytTNQ00foogMwofa1...		
	[Excited to give t...	0	0	2.0 2012-10-19 05:36:10	-4TMqnQJW1yd6NqGR...	7A87z0ycTC5D108bh...	0 Q2fHu4Zx0Ts115RaZ...		
	[This place is qui...	0	0	4.0 2017-10-08 23:37:14	z5KWNSDvgv-4I62P8...	AisJLgSqMlfkeae9y...	0 84jtzzhtl7Vvk8WBmi...		
only showing top 10 rows									

Fig 3: Original Dataframe

Most of them are imported from PySpark ML features.Tokenizing the words is where we split the entire sentences into words and the order of the words will not hold any importance any more. We performed following functions over the text reviews:

1. Convert to Lower case:
- Converted all the reviews to lower case to remove any discrepancies in the corpus

	text cool funny stars	date	business_id	user_id useful	review_id	words
	I just had a pedi...	0  0  2.0 2014-08-03 21:38:45	MJ5xkXJ2uWJL60QH...	MXhsNk1DQw5qvwieZ...	3 912_cfGWNvf0nZudP...	[i, just, had, a,...]
	We scheduled an a...	0  0  1.0 2015-05-18 19:22:54	r4cdabWEmTr2110Fu...	NXAX_SSLw2hxiAms...	3 V74eUostrSH5wpMgF...	[we, scheduled, a,...]
	[This old favorite...	1  1  4.0 2009-01-11 03:35:34	90TypmUJ56GQZYdUW...	25N2f7A17Tgu9zQE-	2 6PyhrNqHroIvxZ6YD...	[this, old, favor...
	[This place was hu...	0  0  3.0 2017-07-03 03:15:45	7JWwexqil03Ua_wV9...	s72nhH1Hl_lN_Ju4h...	0 Z_ShCWTEdEEvca9Fs...	[this, place, was...
	[Went there to wat...	0  0  3.0 2013-09-20 04:24:55	Y7T9osP-hoop2nLnm...	xSt--5JG6QqTEOwXJ...	0 jsf1I9Sy1w3VvNkg1...	[went, there, to,...]
	[Played cornhole f...	1  0  4.0 2017-10-07 23:35:14	3Jq5LfJ5fmJ5KmuA6...	dXP2z-Mqr1k2p2uKU...	1 yp5g42u10-bBEnUH...	[played, cornhole...
	[Invited a friend ...	0  0  1.0 2015-05-05 21:08:44	dQWE63u0iSfaf3Fup...	_IR48ok0ZkPMWJ2P1...	0 4kn_D8-XXseuNuGsv...	[invited, a, frie...
	[What an incredibl...	0  0  5.0 2013-03-28 23:11:14	-4TMqnQJW1yd6NqGR...	cf-n5eNXJE7DG0iJb...	0 ytTNQ00foogMwofa1...	[what, an, incred...
	[Excited to give t...	0  0  2.0 2012-10-19 05:36:10	-4TMqnQJW1yd6NqGR...	7A87z0ycTC5D108bh...	0 Q2fHu4Zx0Ts115RaZ...	[excited, to, giv...
	[This place is qui...	0  0  4.0 2017-10-08 23:37:14	zSKWNSDvgv-4I62P8...	AisJLgSqMlfkeae9y...	0 84jtzzhtl7Vvk8WBmi...	[this, place, is,...]

only showing too 10 rows

2. Remove Stop Words
- Stop words such as prepositions or articles ("a", "the") are also removed so that we only keep the more meaningful words in the corpus.

	text cool funny stars	date	business_id	user_id useful	review_id	words	words_nsw
	I just had a pedi...	0	0	2.0 2014-08-03 21:38:45 MJ5xkXJ2uWJL60QH...	MXhsNk1DQw5qvwieZ...	3 912_cfGWNvf0nZudP...	[i, just, had, a,...] [pedi, & acrylic...
	We scheduled an a...	0	0	1.0 2015-05-18 19:22:54 r4cdabWEmTr2110Fu...	NXAX_SSLw2hxiAms...	3 V74eUostrSH5wpMgF...	[we, scheduled, a,...] [scheduled, appoi...
	[This old favorite...	1	1	4.0 2009-01-11 03:35:34 90TypmUJ56GQZYdUW...	25N2f7A17Tgu9zQE-	2 6PyhrNqHroIvxZ6YD...	[this, old, favor...] [old, favorite, m...
	[This place was hu...	0	0	3.0 2017-07-03 03:15:45 7JWwexqil03Ua_wV9...	s72nhH1Hl_lN_Ju4h...	0 Z_ShCWTEdEEvca9Fs...	[this, place, was...] [place, watch, u...
	[Went there to wat...	0	0	3.0 2013-09-20 04:24:55 Y7T9osP-hoop2nLnm...	xSt--5JG6QqTEOwXJ...	0 jsf1I9Sy1w3VvNkg1...	[went, there, to...] [went, watch, eag...
	[Played cornhole f...	1	0	4.0 2017-10-07 23:35:14 3Jq5LfJ5fmJ5KmuA6...	dXP2z-Mqr1k2p2uKU...	1 yp5g42u10-bBEnUH...	[played, cornhole...] [played, cornhole...
	[Invited a friend ...	0	0	1.0 2015-05-05 21:08:44 dQWE63u0iSfaf3Fup...	_IR48ok0ZkPMWJ2P1...	0 4kn_D8-XXseuNuGsv...	[invited, a, frie...] [invited, friend...
	[What an incredibl...	0	0	5.0 2013-03-28 23:11:14 -4TMqnQJW1yd6NqGR...	cf-n5eNXJE7DG0iJb...	0 ytTNQ00foogMwofa1...	[what, an, incred...] [incredibly, gorg...
	[Excited to give t...	0	0	2.0 2012-10-19 05:36:10 -4TMqnQJW1yd6NqGR...	7A87z0ycTC5D108bh...	0 Q2fHu4Zx0Ts115RaZ...	[excited, to, giv...] [excited, give, h...
	[This place is qui...	0	0	4.0 2017-10-08 23:37:14 zSKWNSDvgv-4I62P8...	AisJLgSqMlfkeae9y...	0 84jtzzhtl7Vvk8WBmi...	[this, place, is,...] [place, quite, ha...

only showing top 10 rows

3. Unigrams & trigrams - Frequency > 20
- We are going attempt two modeling approaches. The first one is to split the model into trigrams (e.g. every three words in a single review will be split together. We will then pick out the trigrams that appear more than 20 times in the corpus, so that we can eliminate any random phrases that only appear once or twice.
- These phrases, will be then joined together, using an underscore "", and be replaced in the original text. The pipeline of Tokenize → CountVectorizer (BagOfWords) → TF-IDF . Model will then be applied using the new texts. This step ensures that we have a well-mixed combinations of unigrams and trigrams in our training data.

	text cool funny stars		date	business_id	user_id useful	review_id	words	words_nsw	ngram
I just had a pedi...	0	0	2.0 2014-08-03 21:38:45	MJ5xkXJ2uWJL60QH...	MXhsNk1DQw5qvwieZ...	3 912_cfGWNvf0nZudP...	[i, just, had, a,...]	[pedi, & acrylic...	[I just had, just
We scheduled an a...	0	0	1.0 2015-05-18 19:22:54	r4cdabWEmTr2110Fu...	NXAX_SSLw2hxiAms...	3 V74eUostrSH5wpMgF...	[we, scheduled, a...	[scheduled, apol...	[We scheduled an...
[This old favorit...	1	1	4.0 2009-01-11 03:35:34	90TypmUJ56GQZYdUW...	25N2f7A17Tgu9zQE-	2 6PyhrNqHroIvxZ6YD...	[this, old, favor...	[old, favorite m...	[This old favor...
[This place was hu...	0	0	3.0 2017-07-03 03:15:45	7JWwexqil03Ua_wV9...	s72nhH1Hl_lN_Ju4h...	0 Z_ShCWTEdEEvca9Fs...	[this, place, was...	[place, hugely u...	[This place was,
[Went there to wat...	0	0	3.0 2013-09-20 04:24:55	Y7T9osP-hoop2nLnm...	xSt--5JG6QqTEOwXJ...	0 jsf1I9Sy1w3VvNkg1...	[went, there, to...	[went, watch, eag...	[Went there to, t...
[Played cornhole f...	1	0	4.0 2017-10-07 23:35:14	3Jq5LfJ5fmJ5KmuA6...	dXP2z-Mqr1k2p2uKU...	1 yp5g42u10-bBEnUH...	[played, cornhole...	[played, cornhole...	[Iplayed cornhole ...
[Invited a friend ...	0	0	1.0 2015-05-05 21:08:44	dQWE63u0iSfaf3Fup...	_IR48ok0ZkPMWJ2P1...	0 4kn_D8-XXseuNuGsv...	[invited, a, frie...	[invited, friend...	[Iinvited a friend...
[What an incredibl...	0	0	5.0 2013-03-28 23:11:14	-4TMqnQJW1yd6NqGR...	cf-n5eNXJE7DG0iJb...	0 ytTNQ00foogMwofa1...	[what, an, incre...	[incredibly, gorg...	[What an incredib...
[Excited to give t...	0	0	2.0 2012-10-19 05:36:10	-4TMqnQJW1yd6NqGR...	7A87z0ycTC5D108bh...	0 Q2fHu4Zx0Ts115RaZ...	[excited, to, giv...	[excited, give, h...	[Excited to give,
[This place is qui...	0	0	4.0 2017-10-08 23:37:14	zSKWNSDvgv-4I62P8...	AisJLgSqMlfkeae9y...	0 84jtzzhtl7Vvk8WBmi...	[this, place, is...	[place, quite, ha...	[This place is, p...
*****									
only showing top 10 rows									

# Model Architecture

We have used two basic methods to determine the rating of the reviews:

## Support Vector Machine

With a train-test split ratio of 80:20 we then run 50 iterations of Support Vector Machine (SVM) model on the training data. The parallel computing of the Spark attributes to speeding up the entire process. Given an imbalanced data-set, it is important know which classification metrics we are going to optimize. Accuracy in such a scenario will not be representitaive of how well the model is performing, while F1 score – as a weighted average of precision and recall could reveal how well the model performs in identifying both the prediction relevancy and what % of truly relevant results are correctly predicted. The parameters used here are:

- 1. Number of Iterations = 50
- 2. Regularization Parameter = 0.3
- 3. Learning Algorithm: Stochastic Gradient Descent

## Elastic Net Logistic Regression

Using the same training data, we instead applied a regularized logistic regression. SVM focuses on finding the separating plane that maximizes the distance of the closest points to the margin, whereas Logistic Regression maximizing the probability of the data (i.e. the further it lies away from the hyperplane the better). In addition to a normal logistic regression, we used a linear combination of L1 and L2 regularization, i.e. Elastic Net, to prevent the model from overfitting. Since LASSO (L1) tend to select only one significant coefficients the other, EN adds in the penalty from Ridge Regression (L2) that helps to overcome the disadvantages.

- 1. Reg Parameter = 0.3
- 2. ElasticNet Parameter = 0.8
- 3. Max Iterations = 100

## Model performance

The dataset is quite imbalanced and hence we used the F1 score to evaluate the model performance.

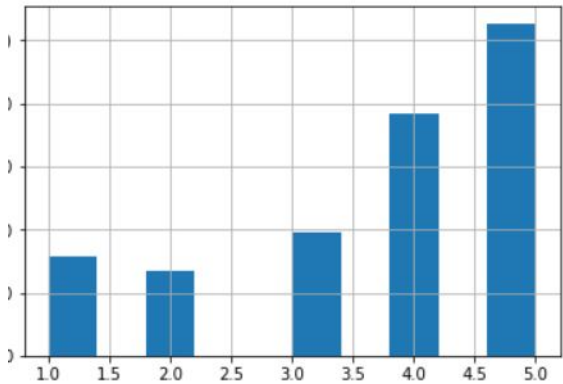


Fig 4: F1 score using logistic regression

## Logistic Regression

Below is the code used -

```
1 from __future__ import print_function
2 from pyspark.ml.classification import LogisticRegression
3 from pyspark.ml import Pipeline, PipelineModel
4 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
5 import os
6 from pyspark import SparkContext
7 from pyspark.sql.session import SparkSession
8 from pyspark.streaming import StreamingContext
9 import pyspark.sql.types as tp
10 from pyspark.ml.feature import *
11 from pyspark.ml import Pipeline
12 from pyspark.ml.feature import StringIndexer, OneHotEncoderEstimator, VectorAssembler
13 from pyspark.ml.feature import StopWordsRemover, Word2Vec, RegexTokenizer
14 from pyspark.ml.classification import LogisticRegression
15 #from pyspark.mllib.feature import HashingTF, IDF
16
17 sc = SparkContext(appName="Project ")
```

```

18 spark = SparkSession(sc)
19
20 spark_df = spark.read.format("bigquery").option("table", "bdl_project20.yelp").load().
    toDF("text", "cool", "funny", "stars", "date", "business_id", "user_id", "useful", "
        review_id")
21
22 spark_df.show(10)
23
24 tok = RegexTokenizer(inputCol= "text" , outputCol= "words", pattern= '\\W')
25 stopword_rm = StopWordsRemover(inputCol="words", outputCol="words_nsw")
26 hashingTF = HashingTF(inputCol="words_nsw", outputCol="rawFeatures")
27 idf = IDF(inputCol="rawFeatures", outputCol="features")
28 lr = LogisticRegression(featuresCol="features", labelCol="stars")
29 pipe = Pipeline(stages=[tok, stopword_rm, hashingTF, idf, lr])
30 model = pipe.fit(spark_df)
31 model.save("gs://bdl__project20/model/")
32 predictions = model.transform(spark_df)
33
34 predictions.show(10)
35
36 evaluator = MulticlassClassificationEvaluator(labelCol="stars", predictionCol="
    prediction", metricName="f1")
37
38 f1 = evaluator.evaluate(predictions)
39
40 print("F1 score =", f1)

```

```

20/09/20 06:47:35 INFO com.google.cloud.spark.bigquery.direct.DirectBigQueryRelation: Querying table codelearn-281105.bdl_project20.
20/09/20 06:47:35 INFO com.google.cloud.spark.bigquery.direct.DirectBigQueryRelation: Going to read from codelearn-281105.bdl_project
20/09/20 06:47:35 INFO com.google.cloud.spark.bigquery.direct.DirectBigQueryRelation: Created read session for table 'codelearn-281105.bdl_project20.yelp'
20/09/20 06:47:35 INFO com.google.cloud.spark.bigquery.direct.DirectBigQueryRelation: Requested 14 max partitions, but only received 10
F1 score = 0.7076878667
20/09/20 07:00:05 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@1bf2ab2b(HTTP/1.1,[http/1.1]){0.0.0.0:4040}

```

textbfFig 4: F1 score using logistic regression

## Real-time computation

We used the trained model to perform real-time predictions of test data streaming to Kafka. We have used a part of the data for streaming.

Below is the code used -

```

1 from __future__ import print_function
2 from pyspark import SparkContext
3 import pyspark.sql.types as tp
4 from pyspark.sql import Row
5 from pyspark.streaming import StreamingContext
6 from pyspark.ml.classification import LogisticRegression
7 from pyspark.streaming.kafka import KafkaUtils
8 import os
9 from pyspark.sql import *
10 from pyspark.sql.types import StructType, StructField, FloatType
11 from pyspark.ml.feature import VectorAssembler, StringIndexer, VectorIndexer,
    StandardScaler, IndexToString
12 from pyspark.ml.classification import LogisticRegression
13 from pyspark.sql.session import SparkSession
14 from pyspark.ml import Pipeline, PipelineModel
15 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
16 import json
17 from pyspark.sql import SQLContext
18
19 kafka_topic = 'from-pubsub'
20 zk = '10.182.0.2:2181' # Apache ZooKeeper quorum
21 app_name = 'from-pubsub' # Can be some other name
22 sc = SparkContext(appName="KafkaPubsub")
23 ssc = StreamingContext(sc, 20)
24 model = PipelineModel.load("gs://bdl__project20/model/")
25 schema = StructType([StructField("sl", FloatType(), True), StructField("sw", FloatType(), True), StructField("pl", FloatType(), True), StructField("pw", FloatType(), True)])
26 #spark = SparkSession.builder.appName(app_name).config("spark.master", "local").getOrCreate()
27 spark = SparkSession(sc)
28 kafkaStream = KafkaUtils.createStream(ssc, zk, app_name, {kafka_topic: 1})
29 #kafkaStream.pprint()
30 print("Starting")
31 def get_prediction(temp):
32     print("Inside IF")
33     temp.toDF().show()
34     dstream = temp.map(lambda x: json.loads(x[1]))
35     dstream.toDF().show()
36     dat = dstream.map(lambda x:x[0])
37     dat.toDF().show()

```

```

38 #dat1 = dat.map(lambda x: x.split(','))
39 #dat2 = dat1.map(lambda x: [float(i) for i in x])
40 #dat3 = dat2.map(lambda x: x[1:])
41 dat4 = dat.map(lambda x: Row(text=str(x)))
42 #df = dat.toDF().toPandas()
43 #df.columns = ['text']
44 #print(df)
45 df1 = spark.createDataFrame(dat4)
46 df1.show()
47 pred = model.transform(df1)
48 pred.show()
49 kafkaStream.foreachRDD(lambda x: get_prediction(x))
50 print("Ending")
51 ssc.start()
52 ssc.awaitTermination()

```

Publishing rate is 10 seconds.Kafka was ingesting at a rate of 20 seconds.

## Conclusion

The Logistic Regression model performs better than SVM. Nowadays, real-time computing generally needs to process large amounts of data, in addition to meeting some of the requirements of non-real-time computing (e.g., accurate results). The most important requirement of real-time computing is the response to computing results in real time. This project enabled us to understand the importance of real time computing.