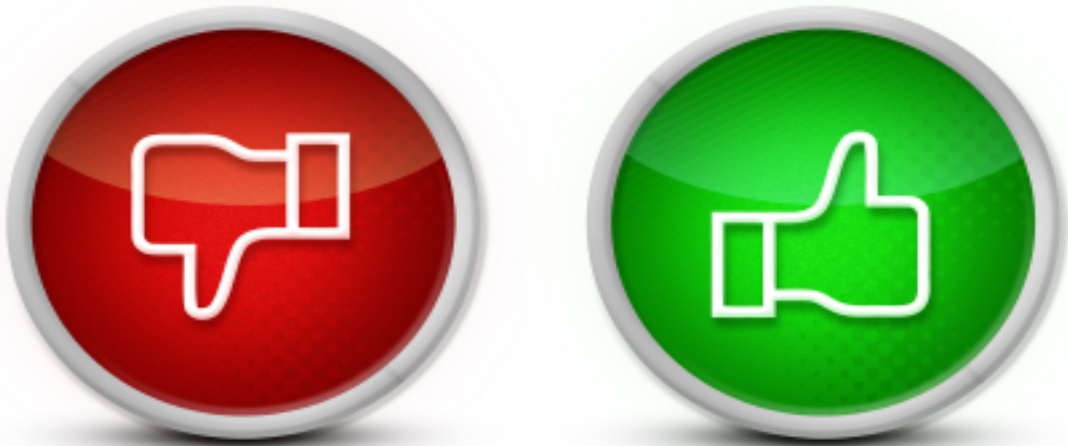


**BIG DATA AND INTELLIGENT ANALYTICS**  
**SENTIMENT ANALYSIS – IMDB MOVIE REVIEWS -**  
**REPORT**



**REPORT PREPARED BY**

MUBEEN

RASHMI KARUNANITHI

SANDEEP KUMAR RAMADOSS MAHENDRAN

## Table of Contents

<b>1. INTRODUCTION .....</b>	<b>3</b>
1.1 SENTIMENT ANALYSIS.....	3
<b>2. EXECUTIVE SUMMARY.....</b>	<b>3</b>
2.1 PROBLEM STATEMENT.....	3
2.2 APPROACH.....	5
2.2.1 Random Split .....	5
2.2.2 Exploratory Data Analysis.....	5
2.2.2.1 Case 1 .....	6
2.2.2.2 Case 2 .....	6
2.2.2.3 Case 3 .....	7
2.2.2.4 Case 4 .....	7
2.3 DATA STORAGE:.....	8
<b>3. FEATURE ENGINEERING .....</b>	<b>8</b>
3.1 Term frequency–inverse document frequency (TF-IDF).....	8
3.2 Word2Vec .....	10
<b>4. MODELS FOR TF - IDF .....</b>	<b>11</b>
4.1 RANDOM FOREST CLASSIFICATION .....	11
4.2 NAIVE BAYES .....	12
4.3 LOGISTIC REGRESSION WITH ML .....	13
<b>5. MODELS FOR WORD2VEC.....</b>	<b>15</b>
5.1 LOGISTIC REGRESSION WITH LBFGS.....	15
5.2 K MEANS CLUSTERING.....	16
5.3 RANDOM FOREST CLASSIFICATION .....	17
<b>6. COMPARISON.....</b>	<b>18</b>
6.1 TF - IDF.....	18
6.2 WORD2VEC .....	19
<b>7. FUTURE REVIEW PREDICTION.....</b>	<b>19</b>
<b>8. CONCLUSION .....</b>	<b>20</b>
<b>9. REFERENCES .....</b>	<b>20</b>

# 1. INTRODUCTION

In this project we will study the basics of sentiment analysis with Natural Language Processing techniques using both spark and gensim. We will try to build a predictive model for a real-world dataset and provide results and recommendations. We begin by looking at the basics of

## 1.1 SENTIMENT ANALYSIS

Sentiment Analysis aims to determine the attitude of the speaker or a writer with respect to some topic using NLP. A basic task in this analysis is classifying the polarity of a given text at the document, sentence or feature/aspect level - whether the expressed opinion in a document, sentence or an entity feature/aspect is positive or negative.

The rise of social media such as blogs and social networks have fueled interest in sentiment analysis. With the proliferation of reviews, ratings, recommendations etc, there is a great opportunity for businesses to identify new opportunities and manage their reputations.

NLP comprises of making the machine understand how humans speak, write and communicate. Humans communicate or convey messages in structured pattern, unstructured pattern, contain regional slangs and idioms. NLP tries to bridge the gap in communication. As our project deals with sentimental analysis of movie reviews we will be using NLP to build different models like TF-IDF and Word2Vec which will predict the polarity of the review.

# 2. EXECUTIVE SUMMARY

## 2.1 PROBLEM STATEMENT

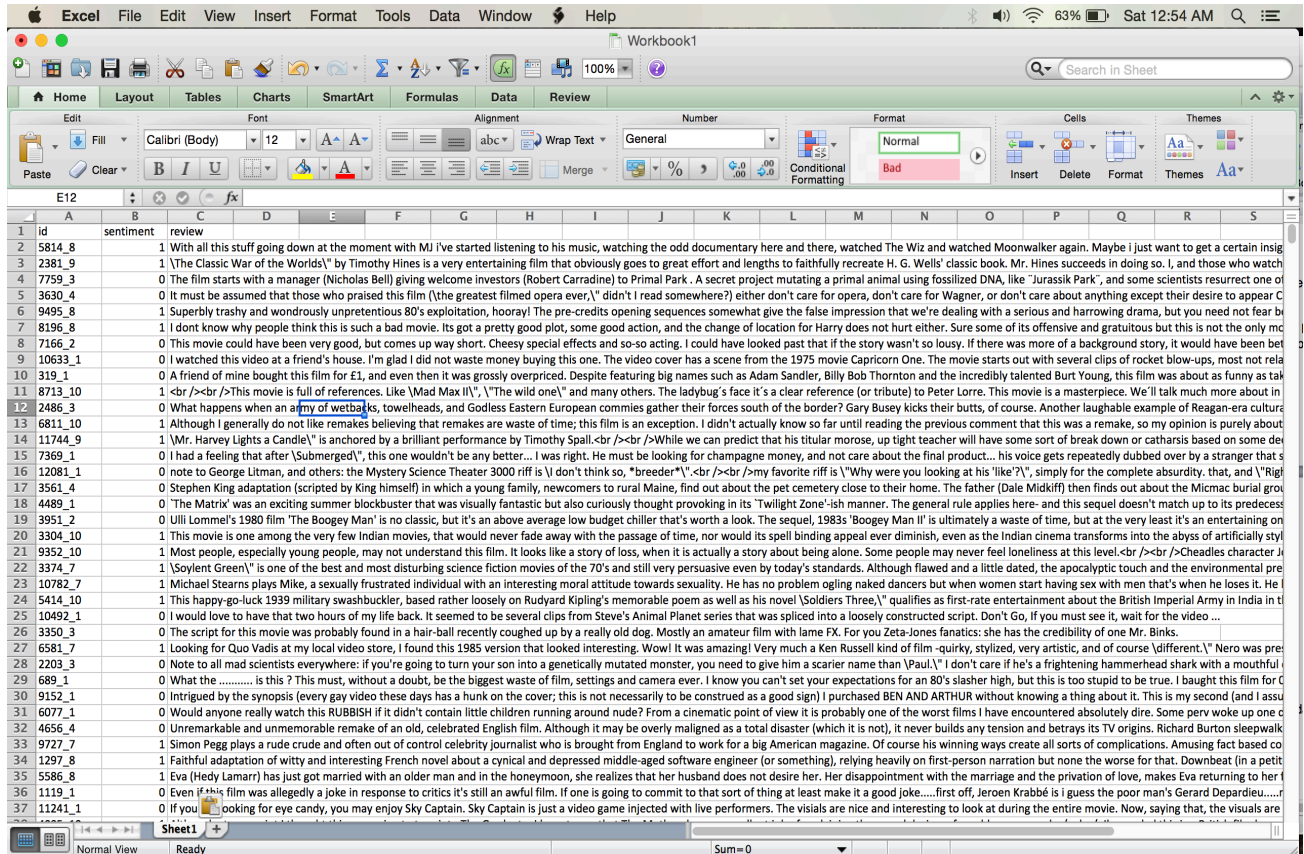
The given dataset provides the 25000 reviews of movies on IMDB. The sentiment is associated with each review as 0(thumbs down) or 1(thumbs up).

Our approach aims to apply simple TF-IDF as well as Word2Vec models to try and build a system, which can predict the polarity of the review as thumbs down (0) or thumbs up (1).

# Sentiment Analysis - IMDB Movie Reviews

## Dataset Screenshot:

## Labeled TrainData



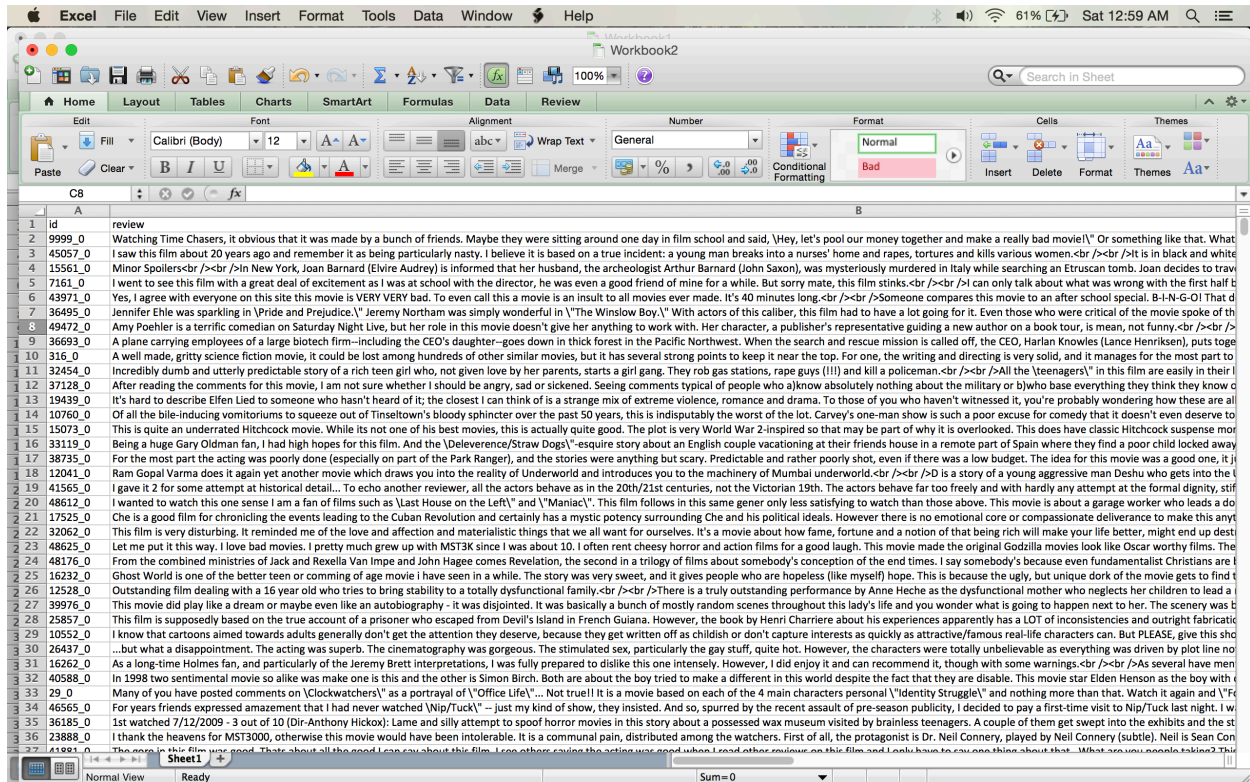
id	sentiment	review
5814_8	1	With all this stuff going down at the moment with MJ I've started listening to his music, watching the odd documentary here and there, watched The Wiz and watched Moonwalker again. Maybe i just want to get a certain insig
2381_9	1	The Classic War of the Worlds by Timothy Hines is a very entertaining film that obviously goes to great effort and lengths to faithfully recreate H. G. Wells' classic book. Mr. Hines succeeds in doing so. I, and those who watch
7759_3	0	The film starts with a manager (Nicholas Bell) giving welcome investors (Robert Carradine) to Primal Park. A secret project mutating a primal animal using fossilized DNA, like 'Jurassic Park', and some scientists resurrect one of
3630_4	0	It must be assumed that those who praised this film (the greatest filmed opera ever, I didn't read somewhere?) either don't care for opera, don't care for Wagner, or don't care about anything except their desire to appear C
9495_8	1	Superbly trashy and wondrously unpretentious 80's exploitation, hooray! The pre-credits opening sequences somewhat give the false impression that we're dealing with a serious and harrowing drama, but you need not fear b
8196_8	1	I dont know why people think this is such a bad movie. Its got a pretty good plot, some good action, and the change of location for Harry does not hurt either. Sure some of its offensive and gratuitous but this is not the only mc
7166_2	0	This movie could have been very good, but comes up way short. Cheesy special effects and so-so acting. I could have looked past that if the story wasn't so lousy. If there was more of a background story, it would have been bet
10633_1	0	I watched this video at a friend's house. I'm glad I did not waste money buying this one. The video cover has a scene from the 1975 movie Capricorn One. The movie starts out with several clips of rocket blow-ups, most not rela
319_1	0	A friend of mine bought this film for £1, and even then it was grossly overpriced. Despite featuring big names such as Adam Sandler, Billy Bob Thornton and the incredibly talented Burt Young, this film was about as funny as tak
8713_10	1	  This movie is full of references. Like 'Mad Max II', 'The wild one' and many others. The ladybug's face it's a clear reference (or tribute) to Peter Lorre. This movie is a masterpiece. We'll talk much more about in
2486_3	0	What happens when an army of wetbacks, towelheads, and Godless Eastern European commies gather their forces south of the border? Gary Busey kicks their butts, of course. Another laughable example of Reagan-era cultura
6811_10	1	Although I generally do not like remakes believing that remakes are waste of time; this film is an exception. I didn't actually know so far until reading the previous comment that this was a remake, so my opinion is purely about
11744_9	1	'Mr. Harvey Lights a Candle' is anchored by a brilliant performance by Timothy Spall.  While we can predict that his titular morose, up tight teacher will have some sort of break down or catharsis based on some de
7369_1	0	I had a feeling that after 'Submerged', this one wouldn't be any better... I was right. He must be looking for champagne money, and not care about the final product... his voice gets repeatedly dubbed over by a stranger that s
12081_1	0	note to George Litman, and others: the Mystery Science Theater 3000 riff is \I don't think so, "breeder"!\  my favorite riff is \"Why were you looking at his 'like'?\", simply for the complete absurdity, that, and \"Righ
3561_4	0	Stephen King adaptation (scripted by King himself) in which a young family, newcomers to rural Maine, find out about the pet cemetery close to their home. The father (Dale Midkiff) then finds out about the Micmac burial gro
14889_1	0	The Matrix' was an exciting summer blockbuster that was visually fantastic but also curiously thought provoking in its 'Twilight Zone'-ish manner. The general rule applies here- and this sequel doesn't match up to its predecess
3951_2	0	Ulli Lommel's 1980 film 'The Boogey Man' is no classic, but it's an above average low budget chiller that's worth a look. The sequel, 1983's 'Boogey Man II' is ultimately a waste of time, but at the very least it's an entertaining o
3304_10	1	This movie is one among the very few Indian movies, that would never fade away with the passage of time, nor would its spell binding appeal ever diminish, even as the Indian cinema transforms into the abyss of artificially styl
9352_10	1	Most people, especially young people, may not understand this film. It looks like a story of loss, when it is actually a story about being alone. Some people may never feel loneliness at this level.  Cheadles character J
3374_7	1	'Soylent Green' is one of the best and most disturbing science fiction movies of the 70's and still very persuasive even by today's standards. Although flawed and a little dated, the apocalyptic touch and the environmental pre
10782_7	1	Michael Stearns plays Mike, a sexually frustrated individual with an interesting moral attitude towards sexuality. He has no problem ogling naked dancers but when women start having sex with men that's when he loses it. He l
5414_10	1	This happy-go-lucky 1939 military swashbuckler, based rather loosely on Rudyard Kipling's memorable poem as well as his novel 'Soldiers Three,' qualifies as first-rate entertainment about the British Imperial Army in India in t
10492_1	0	I would love to have that two hours of my life back. It seemed to be several clips from Steve's Animal Planet series that was spliced into a loosely constructed script. Don't Go, If you must see it, wait for the video ...
3350_3	0	The script for this movie was probably found in a hair-ball recently coughed up by a really old dog. Mostly an amateur film with lame FX. For you Zeta-Jones fanatics: she has the credibility of one Mr. Binks.
5681_7	1	Looking for Quo Vadis at my local video store, I found this 1985 version that looked interesting. Wow! It was amazing! Very much a Ken Russell kind of film -quirky, stylized, very artistic, and of course 'different.' Nero was pre
2203_3	0	Note to all mad scientists everywhere: if you're going to turn your son into a genetically mutated monster, you need to give him a scarier name than 'Paul.' I don't care if he's a frightening hammerhead shark with a mouthful
689_1	0	What the ..... is this? This must, without a doubt, be the biggest waste of film, settings and camera ever. I know you can't set your expectations for an 80's slasher high, but this is too stupid to be true. I bought this film for C
9152_1	0	Intrigued by the synopsis (every gay video these days has a hunk on the cover; this is not necessarily to be construed as a good sign) I purchased BEN AND ARTHUR without knowing a thing about it. This is my second (and I assu
6077_1	0	Would anyone really watch this RUBBISH if it didn't contain little children running around nude? From a cinematic point of view it is probably one of the worst films I have encountered absolutely dire. Some perv woke up one c
4656_4	0	Unremarkable and unmemorable remake of an old, celebrated English film. Although it may be overly maligned as a total disaster (which it is not), it never builds any tension and betrays its TV origins. Richard Burton sleepwalk
9727_7	1	Simon Pegg plays a rude crude and often out of control celebrity journalist who is brought from England to work for a big American magazine. Of course his winning ways create all sorts of complications. Amusing fact based co
1297_8	1	Faithful adaptation of witty and interesting French novel about a cynical and depressed middle-aged software engineer (or something), relying heavily on first-person narration but none the worse for that. Downbeat (in a petit
5586_8	1	Eva (Hedy Lamarr) has just got married with an older man and in the honeymoon, she realizes that her husband does not desire her. Her disappointment with the marriage and the privation of love, makes Eva returning to her f
1119_1	0	Even if this film was allegedly a joke in response to critics it's still an awful film. If one is going to commit to that sort of thing at least make it a good joke.....first off, Jeroen Krabbé is I guess the poor man's Gerard Depardieu.....r
11241_1	0	If you looking for eye candy, you may enjoy Sky Captain. Sky Captain is just a video game injected with live performers. The visuals are nice and interesting to look at during the entire movie. Now, saying that, the visuals are

## COLUMN DETAILS:

Column Name	Data Type	Description
id	String	Unique Identifiers
Sentiment	Number	Sentiment of rview
Review	String	Review Statements

# Sentiment Analysis - IMDB Movie Reviews

## UnLabeled TrainData



id	review
9999_0	Watching Time Chasers, it obvious that it was made by a bunch of friends. Maybe they were sitting around one day in film school and said, 'Hey, let's pool our money together and make a really bad movie!'
45057_0	I saw this film about 20 years ago and remember it as being particularly nasty. I believe it is based on a true incident: a young man breaks into a nurses' home and rapes, tortures and kills various women.
15561_0	Minor Spoilers->In New York, Joan Barnard (Elvire Audrey) is informed that her husband, the archeologist Arthur Barnard (John Saxon), was mysteriously murdered in Italy while searching an Etruscan tomb. Joan decides to travel to Italy to find out what happened to him.
7161_0	I went to see this film with a great deal of excitement as I was at school with the director, he was even a good friend of mine for a while. But sorry mate, this film stinks. I can only talk about what was wrong with the first half of the film.
43971_0	Yes, I agree with everyone on this site this movie is VERY VERY bad. To even call this a movie is an insult to all movies ever made. It's 40 minutes long. Someone compares this movie to an after school special. B-I-N-G-O! That's all I have to say about this movie.
36495_0	Jennifer Ehle was sparkling in 'Pride and Prejudice.' Jeremy Northam was simply wonderful in 'The Winslow Boy.' With actors of this caliber, this film had to have a lot going for it. Even those who were critical of the movie spoke of the quality of the acting.
49472_0	Amy Poehler is a terrific comedian on Saturday Night Live, but her role in this movie doesn't give her anything to work with. Her character, a publisher's representative guiding a new author on a book tour, is mean, not funny.
36693_0	A plane carrying employees of a large biotech firm—including the CEO's daughter—goes down in thick forest in the Pacific Northwest. When the search and rescue mission is called off, the CEO, Harlan Knowles (Lance Henriksen), puts together a team to find the plane.
316_0	A well made, gritty science fiction movie, it could be lost among hundreds of other similar movies, but it has several strong points to keep it near the top. For one, the writing and directing is very solid, and it manages for the most part to be a good movie.
32454_0	Incredibly dumb and utterly predictable story of a rich teen girl who, not given love by her parents, starts a girl gang. They rob gas stations, rape guys (!!!) and kill a policeman. All the 'teenagers' in this film are easily in their late teens or early 20s.
37128_0	After reading the comments for this movie, I am not sure whether I should be angry, sad or sickened. Seeing comments typical of people who know absolutely nothing about the military or base everything they think they know on the fact that the movie is not a war movie.
19439_0	It's hard to describe Elfen Lied to someone who hasn't heard of it; the closest I can think of is a strange mix of extreme violence, romance and drama. To those of you who haven't witnessed it, you're probably wondering how these are all in one movie.
10760_0	Of all the bile-inducing vomitoriums to squeeze out of Tinseltown's bloody splinter over the past 50 years, this is indisputably the worst of the lot. Carvey's one-man show is such a poor excuse for comedy that it doesn't even deserve to be mentioned.
15073_0	This is quite an underrated Hitchcock movie. While it's not one of his best movies, this is actually quite good. The plot is very World War 2-inspired so that may be part of why it is overlooked. This does have classic Hitchcock suspense moments.
33119_0	Being a huge Gary Oldman fan, I had high hopes for this film. And the 'Deleverage/Straw Dogs' -esque story about an English couple vacationing at their friends house in a remote part of Spain where they find a poor child locked away for the most part the acting was poorly done (especially on part of the Park Ranger), and the stories were anything but scary.
38735_0	For the most part the acting was poorly done (especially on part of the Park Ranger), and the stories were anything but scary. Predictable and rather poorly shot, even if there was a low budget. The idea for this movie was a good one, it just wasn't executed well.
12041_0	Ram Gopal Varma does it again yet another movie which draws you into the reality of Underworld and introduces you to the machinery of Mumbai underworld. It is a story of a young aggressive man Deshu who gets into the underworld and becomes a powerful figure.
41565_0	I gave it 2 for some attempt at historical detail... To echo another reviewer, all the actors behave as if in the 20th/21st centuries, not the Victorian 19th. The actors behave far too freely and with hardly any attempt at the formal dignity, stiff and unnatural.
48612_0	I wanted to watch this one sense I am a fan of films such as 'Vast House on the Left' and 'Maniac'. This film follows in this same genre only less satisfying to watch than those above. This movie is about a garage worker who leads a double life.
17525_0	Che is a good film for chronicling the events leading to the Cuban Revolution and certainly has a mystic potency surrounding Che and his political ideals. However there is no emotional core or compassionate deliverance to make this any more than a historical document.
32062_0	This film is very disturbing. It reminded me of the love and affection and materialistic things that we all want for ourselves. It's a movie about how fame, fortune and a notion of that being rich will make your life better, might end up destroying it.
48625_0	Let me put it this way. I love bad movies. I pretty much grew up with MST3K since I was about 10. I often rent cheesy horror and action films for a good laugh. This movie made the original Godzilla movies look like Oscar worthy films. The acting is terrible.
48176_0	From the combined ministries of Jack and Rexella Van Impe and John Hagee comes Revelation, the second in a trilogy of films about somebody's conception of the end times. I say somebody's because even fundamentalist Christians are I think a bit out of touch.
16232_0	Ghost World is one of the better teen or coming of age movie I have seen in a while. The story was very sweet, and it gives people who are hopeless (like myself) hope. This is because the ugly, but unique dork of the movie gets to find out that he is not alone.
12528_0	Outstanding film dealing with a 16 year old who tries to bring stability to a totally dysfunctional family. There is a truly outstanding performance by Anne Heche as the dysfunctional mother who neglects her children to lead a double life.
39976_0	This movie did play like a dream or maybe even like an autobiography - it was disjointed. It was basically a bunch of mostly random scenes throughout this lady's life and you wonder what is going to happen next to her. The scenery was beautiful.
25857_0	This film is supposedly based on the true account of a prisoner who escaped from Devil's Island in French Guiana. However, the book by Henri Charriere about his experiences apparently has a LOT of inconsistencies and outright fabrications. I know that cartoons aimed towards adults generally don't get the attention they deserve.
10552_0	I know that cartoons aimed towards adults generally don't get the attention they deserve, because they get written off as childish or don't capture interests as quickly as attractive/famous real-life characters can. But PLEASE, give this show a chance. It's not just a cartoon.
26437_0	...but what a disappointment. The acting was superb. The cinematography was gorgeous. The stimulated sex, particularly the gay stuff, quite hot. However, the characters were totally unbelievable as everything was driven by plot line no matter how stupid.
16262_0	As a long-time Holmes fan, and particularly of the Jeremy Brett interpretations, I was fully prepared to dislike this one intensely. However, I did enjoy it and can recommend it, though with some warnings. As several have mentioned, the acting is superb.
40588_0	In 1998 two sentimental movie so alike was make one is this and the other is Simon Birch. Both are about the boy tried to make a difference in this world despite the fact that they are disabled. This movie star Eiden Henson as the boy with cerebral palsy.
29_0	Many of you have posted comments on 'Clockwatchers' as a portrayal of 'Office Life'... Not true!! It is a movie based on each of the 4 main characters personal 'Identity Struggle' and nothing more than that. Watch it again and you'll see.
46565_0	For years friends expressed amazement that I had never watched 'Nip/Tuck' - just my kind of show, they insisted. And so, spurred by the recent assault of pre-season publicity, I decided to pay a first-time visit to Nip/Tuck last night. I was not disappointed.
36185_0	1st watched 7/12/2009 - 3 out of 10 (Dir-Anthony Hickox): Lame and silly attempt to spoof horror movies in this story about a possessed wax museum visited by brainless teenagers. A couple of them get swept into the exhibits and the story is a mess.
23888_0	I thank the heavens for MST3000, otherwise this movie would have been intolerable. It is a communal pain, distributed among the watchers. First of all, the protagonist is Dr. Neil Connery, played by Neil Connery (subtle). Neil is Sean Connery's son.
41881_0	The more I think about this movie, the more I love it. I can say about this film. I love others, giving the actor was good when I read other reviews on this film and I only knew to see one thing about that. What are you people taking? This is a masterpiece.

## 2.2 APPROACH

### 2.2.1 Random Split

We used random split to split the data into train data and test data so that the split data could be used for training the model and for testing the model.

### 2.2.2 Exploratory Data Analysis

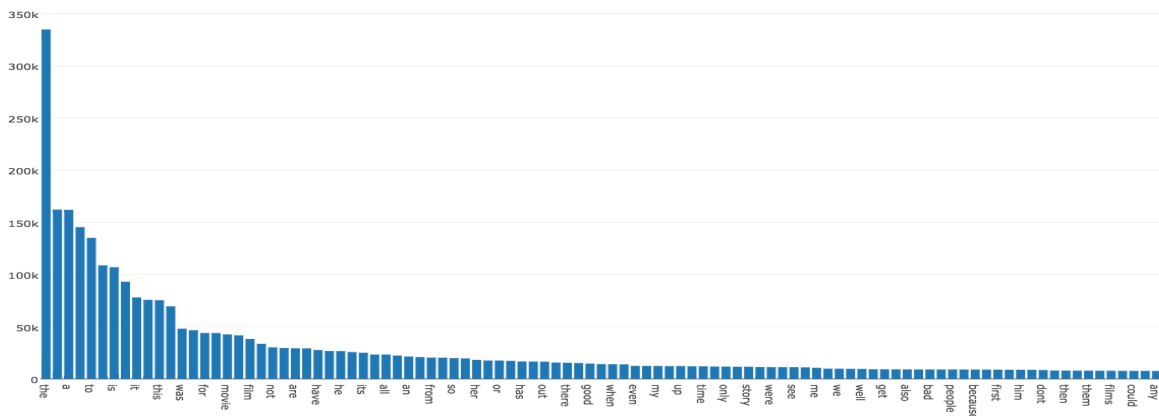
First the dataset is looked for irregularities. In our dataset the reviews will have special characters and html tags that needs to be removed as a part of data preprocessing step. We remove those html tags using The Beautiful Soup library. To remove the punctuation, Numbers we used NLTK.

## Sentiment Analysis - IMDB Movie Reviews

Following are the plots of few data analysis

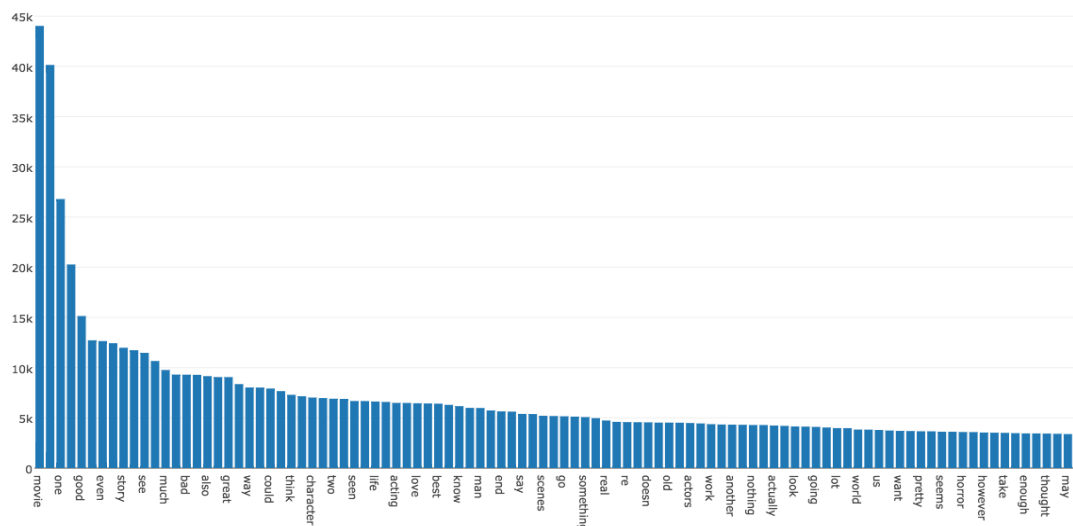
### 2.2.2.1 Case 1

The following plot tells about the word count in the dataset that is a very naive word count:



### 2.2.2.2 Case 2

In the previous plot the words like: the, of, is, I etc show up, which are not really useful. So we are using NLTK (Natural Language ToolKit) which will remove words like the, of, is and keep words that don't repeat often. After using NLTK the data analysis of the reviews is as follows:

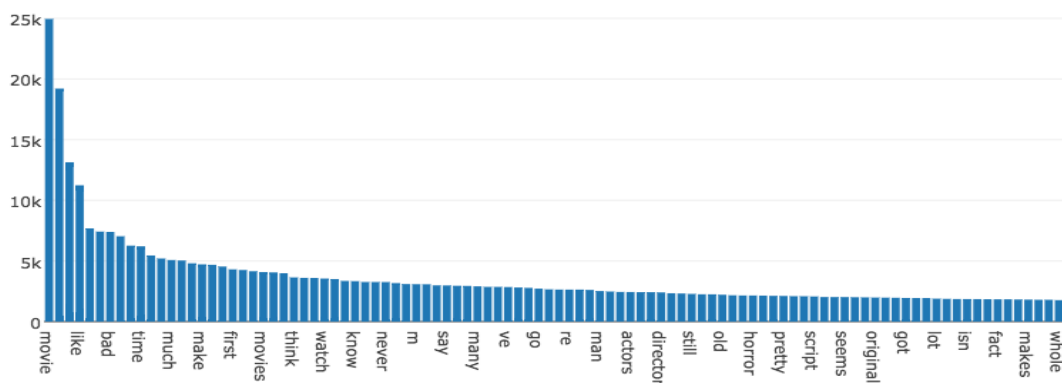


## Sentiment Analysis - IMDB Movie Reviews

The plot has the following words - movie, like, even, really, much, people, great, make, think, watch

### 2.2.2.3 Case 3

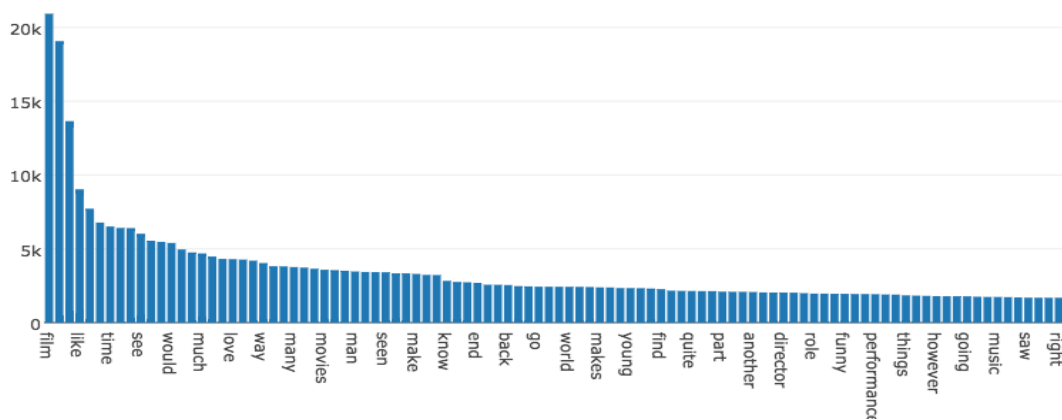
Below is the plot of negative reviews which tells about the negative word count



And the plot has words like bad, never, old, horror etc.

### 2.2.2.4 Case 4

Below is the plot of positive reviews which tells about the positive word count



And the plot has words like love, funny, right etc.



## 2.3 DATA STORAGE:

We store the test data and train data in Amazon S3. In future if a new data comes, it will also be stored in S3. Amazon Simple Storage Service (Amazon S3), provides developers and IT teams with secure, durable, highly-scalable object storage. Amazon S3 is easy to use, with a simple web services interface to store and retrieve any amount of data from anywhere on the web. Amazon S3 can be used alone or together with other AWS services such as Amazon Elastic Compute Cloud (Amazon EC2), Amazon Elastic Block Store (Amazon EBS), and Amazon Glacier, as well as third party storage repositories and gateways.

## 3. FEATURE ENGINEERING

Feature engineering is the process of creating 'features' that could be used in machine learning algorithms. For the purpose of sentiment analysis, the text that we use to train and predict has to be converted to usable and efficient numerical vectors to be used in the classification algorithms. Unlike other process, for text analysis, the assigning of vectors cannot be done manually and natural language processing (NLP) has to be used. For this specific problem, we used two different techniques for feature engineering, with varied results.

### 3.1 Term frequency–inverse document frequency (TF-IDF)

Term frequency–inverse document frequency (TF-IDF), is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus. TF-IDF is available as part of `pyspark.ml.features` package and helps in converting the text to vectors depending on the occurrence and repetition of the word. Before using TF-IDF, The data has to be cleaned. The first step in cleaning the text is the removal of HTML tags from the text. For this purpose, we use a package called Beautiful Soup available in python. The next step would be removed the punctuations and non-alphabetical characters, these characters will affect the efficiency of the created vectors. This is done using regular expressions by replacing anything that is not alphabetical with a space.

Every sentence in a spoken language has words, which repeat frequently. In the English language it would be the preposition, articles, pronouns etc. These words are called stopwords in natural language processing. It is important we remove these words before the features could be created. We use the package called `nltk` for this purpose. We download the stop words and remove those words from our text by looping through it.

Once we have reached this state, the next step would be to tokenize the text into words. This is done using the `tokenizer` function available in `pyspark.ml.features`. Once the tokenization is done, we use a hashing function to convert the words into numerical

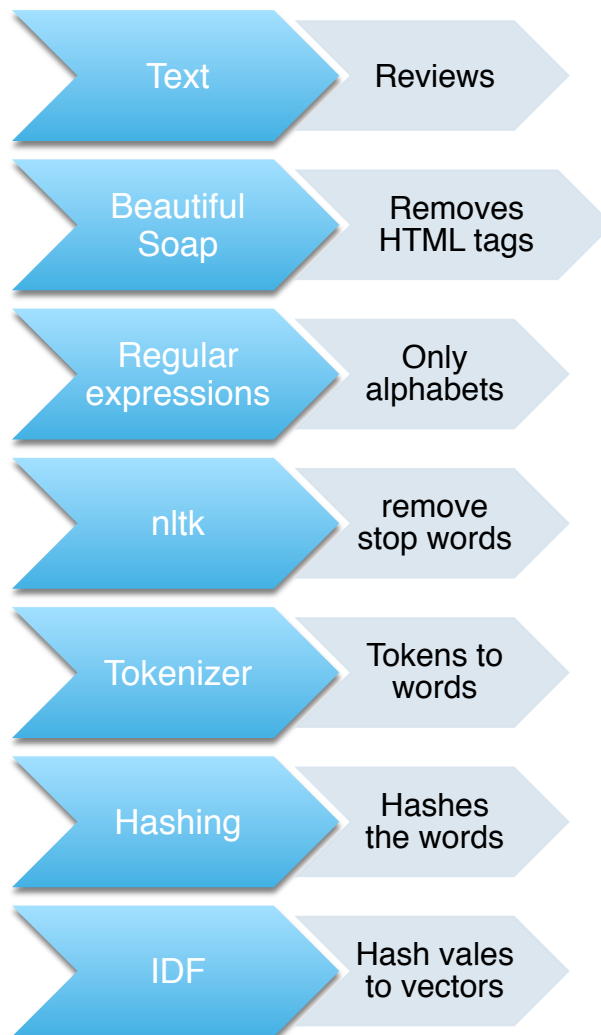


## Sentiment Analysis - IMDB Movie Reviews

values for the creation of TF-IDF model. The converted vectors are then fed to the IDF class for fitting. Once the model is fitted, it can be used for transforming our hashed vales to vectors. The IDF model on transformation gives a SparseVector of features. These features can be used in the ml algorithm. The IDF model also gives us options to choose the number of features we would want to use.

We tried the IDF model for 50 features and 300 features. On evaluation of the classification model, the model with 300 features gave an area under ROC value of about 80% whereas the model with 50 features performed under an area under ROC of 67%.

The steps from text to vectors using TF-IDF are given below.

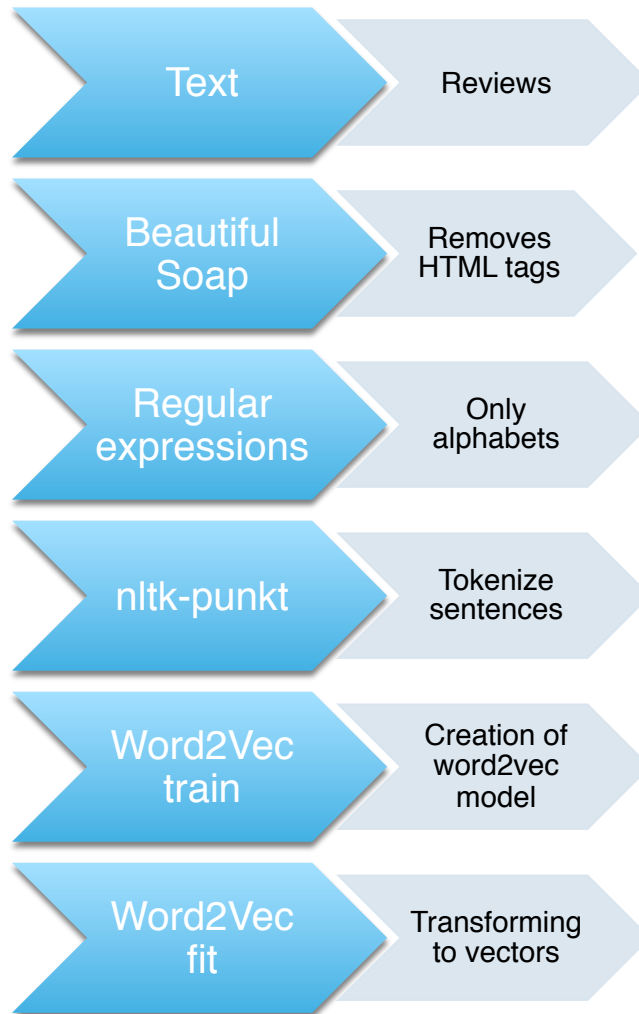


## 3.2 Word2Vec

The Word2Vec also converts the words to vectors, but taking into consideration the semantics and grammar of the language used. Word2Vec is also available in the `pyspark.ml.feature` package. The Word2Vec works by creating a model with the language training data to learn the semantics of the language used in processing. This requires a lot of text for training and it is to be noted that this training data need not be labeled for the sentiment it represents. So, we used the unlabeled reviews that were available for this process. The unlabeled data also goes through the process of cleaning. Beautiful Soup and regular expressions are used to remove html tags and non-alphabetic characters in the text.

For the word2vec model to be fitted it is important that this data has to be in the form of sentences instead of words. English language has various ways in which its sentences would be ended. To create sentences from paragraphs, we make use of yet another feature from the `nltk` package called `punkt`. This `punkt` is fed to a tokenizer and the reviews are tokenized into sentences. This is then fed to the Word2Vec model for fitting. Once our model is fitted, we can use it to train our actual training data for the machine learning algorithm. This data is also cleaned for html tags and non-alphabets, but in both the word2vec training data and the ml training data is not treated for stop words, as the stop words are integral part of language semantics. We were able to attain an area under ROC of about 89% with the vectors created using the word2Vec model (300 features)

The process diagram is as follows



## 4. MODELS FOR TF - IDF

### 4.1 RANDOM FOREST CLASSIFICATION

Random forests train a set of decision trees separately, so the training can be done in parallel. The algorithm injects randomness into the training process so that each decision tree is a bit different. Combining the predictions from each tree reduces the variance of the predictions, improving the performance on test data.

The randomness injected into the training process includes:

- Subsampling the original dataset on each iteration to get a different training set (a.k.a. bootstrapping).
- Considering different random subsets of features to split on at each tree node.

Apart from these randomizations, decision tree training is done in the same way as for individual decision trees.

### CODE

```
stringIndexer = StringIndexer(inputCol="label", outputCol="indexed")
si_model = stringIndexer.fit(selectData)
td = si_model.transform(selectData)
rfc = RandomForestClassifier(maxDepth=2, labelCol="indexed")

pipeline = Pipeline(stages=[tokenizer, hashingTF,idf,stringIndexer , rfc])

paramGrid = ParamGridBuilder().addGrid(hashingTF.numFeatures, [300,
400]).addGrid(rfc.maxDepth, [2, 5, 10]).build()

cv =
CrossValidator().setNumFolds(3).setEstimator(pipeline).setEstimatorParamMaps(param
Grid).setEvaluator(BinaryClassificationEvaluator())
cvModel = cv.fit(pipelineTrainingData)
testTransform = cvModel.transform(pipelineTestData)
predictions = testTransform.select('review', 'label', 'prediction')
predictionsAndLabels = predictions.map(lambda x : (x[1], x[2]))
trainErr = predictionsAndLabels.filter(lambda r : r[0] != r[1]).count() /
float(testData.count())
print("TrainErr: "+str(trainErr))
BinaryClassificationEvaluator().evaluate(testTransform)
```

### RESULT

Accuracy = 55.600362

## 4.2 NAIVE BAYES

Naive Bayes is a simple multiclass classification algorithm with the assumption of independence between every pair of features. Naive Bayes can be trained very efficiently. Within a single pass to the training data, it computes the conditional probability distribution of each feature given label, and then it applies Bayes' theorem to compute the conditional probability distribution of label given an observation and use it for prediction.

These models are typically used for document classification. Within that context, each observation is a document and each feature represents a term whose value is the frequency of the term (in multinomial naive Bayes) or a zero or one indicating whether the term was found in the document (in Bernoulli naive Bayes). Feature values must be nonnegative. The model type is selected with an optional parameter "multinomial" or

“bernoulli” with “multinomial” as the default. Additive smoothing can be used by setting the parameter  $\lambda$  (default to 1.0). For document classification, the input feature vectors are usually sparse, and sparse vectors should be supplied as input to take advantage of sparsity. Since the training data is only used once, it is not necessary to cache it.

### CODE

```
lp = selectData.map(lambda x : LabeledPoint(x.label,x.features))
(trainingData, testData) = lp.randomSplit([0.6, 0.4])
nb = NaiveBayes.train(trainingData,1.0)
pipeline = Pipeline(stages=[tokenizer, hashingTF,idf, nb])
model = pipeline.fit(trainingData)
selected = model.transform(testData).select('review', 'label', 'prediction')

# Build a parameter grid.
paramGrid = ParamGridBuilder().addGrid(hashingTF.numFeatures, [300,
400]).addGrid(nb.regParam, [0.01, 0.1, 1.0]).build()

#Set up cross-validation.
cv =
CrossValidator().setNumFolds(3).setEstimator(pipeline).setEstimatorParamMaps(param
Grid).setEvaluator(BinaryClassificationEvaluator())

#Fit a model with cross-validation.
cvModel = cv.fit(trainingData)
testTransform = cvModel.transform(testData)
predictions = testTransform.select("review", "prediction", "label")
predictionsAndLabels = predictions.map(lambda x : (x[1], x[2]))
trainErr = predictionsAndLabels.filter(lambda r : r[0] != r[1]).count() /
float(testData.count())
print("TrainErr: "+str(trainErr))
BinaryClassificationEvaluator().evaluate(testTransform)
```

### RESULT

```
accuracy = 72
recall = 71.5
precision: 72.05
```

## 4.3 LOGISTIC REGRESSION WITH ML

Logistic Regression is a regression model where the nominal variable is categorical. Use Binary Logistic regression when you have one nominal variable with two values like

male/female or good/bad and one measurement variable, which is an independent variable. Use Multiple Logistic regression when you have one nominal variable and more than one measurement variable. The goal of logistic regression is to predict the probability of getting the desired nominal value given the measurement value and also to predict that getting a particular nominal value is dependent on the measurement value.

### CODE

```
(trainingData, testData) = selectData.randomSplit([0.6, 0.4])
lr = LogisticRegression(maxIter=5, regParam=0.01)
model = lr.fit(trainingData)
result = model.transform(testData)
evaluator = BinaryClassificationEvaluator()
evaluator.evaluate(result, {evaluator.metricName: "areaUnderPR"})
evaluator.evaluate(result, {evaluator.metricName: "areaUnderROC"})
```

Here we do a random split of the data that is available to us. We train the model with 60 percent of data and test the model with 40 percent of the data. We set the maximum no of iterations as 5 and the regularization parameter to 0.01. We fit the model according to the training data. Then the test data is transformed based on the model. We predict whether the sentiment of the review is 1 (good) or 0 (Bad). AreaUnderROC values plots false positive rate against true positive rate. So higher the auROC value better the model.

### RESULT

```
areaUnderPR = 0.8065552181670259
Area under ROC = 0.8275891274493837
training error = 0.244788029925
```

### CROSS VALIDATOR

An important task in ML is *model selection*, or using data to find the best model or parameters for a given task. This is also called *tuning*. Pipelines facilitate model selection by making it easy to tune an entire Pipeline at once, rather than tuning each element in the Pipeline separately.

Currently, spark.ml supports model selection using the CrossValidator class, which takes an Estimator, a set of ParamMaps, and an Evaluator. CrossValidator begins by splitting the dataset into a set of *folds* which are used as separate training and test datasets; e.g., with k=3 folds, CrossValidator will generate 3 (training, test) dataset pairs, each of which uses 2/3 of the data for training and 1/3 for testing. CrossValidator iterates through the set of ParamMaps. For each ParamMap, it trains the

given Estimator and evaluates it using the given Evaluator. The ParamMap that produces the best evaluation metric (averaged over the k folds) is selected as the best model. CrossValidator finally fits the Estimator using the best ParamMap and the entire dataset.

### CROSS VALIDATOR CODE

```
# Build a parameter grid.
paramGrid=ParamGridBuilder().addGrid(hashingTF.numFeatures,[300, 400])
.addGrid(lr.regParam, [0.01, 0.1, 1.0]).build()

#Set up cross-validation.
cv=CrossValidator().setNumFolds(3).setEstimator(pipeline).setEstimatorParamMaps(pa
ramGrid).setEvaluator(BinaryClassificationEvaluator())

#Fit a model with cross-validation.
cvModel = cv.fit(trainingData)
testTransform = cvModel.transform(testData)
predictions = testTransform.select("review", "prediction", "label")
predictionsAndLabels = predictions.map(lambda x : (x[1], x[2]))
trainErr = predictionsAndLabels.filter(lambda r : r[0] != r[1]).count() /
float(testData.count())
print("TrainErr: "+str(trainErr))
BinaryClassificationEvaluator().evaluate(testTransform)
```

For our model we have decided to go with 300 features to run the cross validator class.

### RESULT

```
areaUnderPR = 0.8065552181670259
Area under ROC = 0.8275891274493837
training error = 0.244788029925
```

## 5. MODELS FOR WORD2VEC

### 5.1 LOGISTIC REGRESSION WITH LBFGS

Logistic Regression is a regression model where the nominal variable is categorical. Use Binary Logistic regression when you have one nominal variable with two values like male/female or good/bad and one measurement variable, which is an independent variable. Use Multiple Logistic regression when you have one nominal variable and more than one measurement variable. The goal of logistic regression is to predict the



probability of getting the desired nominal value given the measurement value and also to predict that getting a particular nominal value is dependent on the measurement value.

### CODE

```
#Creating RDD of LabeledPoints
lpSelectData = selectData.map(lambda x : (x.id, LabeledPoint(x.label,x.features)))
#Splitting the data for training and test
(trainingData, testData) = lpSelectData.randomSplit([0.9, 0.1])
# training the Logistic regression with LBFGS model
lrm = LogisticRegressionWithLBFGS.train(trainingData.map(lambda x: x[1]),
iterations=10)
#fetching the labels and predictions for test data
labelsAndPreds = testData.map(lambda p: (p[0],p[1].label, lrm.predict(p[1].features)))
#calculating the accuracy and printing it.
accuracy = labelsAndPreds.filter(lambda (i, v, p): v == p).count() / float(testData.count())
print("Accuracy = " + str(accuracy))
```

### RESULT:

Accuracy = 0.843182696022

## 5.2 K MEANS CLUSTERING

k-means is one of the most commonly used clustering algorithms that clusters the data points into a predefined number of clusters. The MLlib implementation includes a parallelized variant of the k-means++ method called kmeans. The implementation in MLlib has the following parameters:

- k is the number of desired clusters.
- Max iterations are the maximum number of iterations to run.
- Initialization Mode specifies either random initialization or initialization via k-meansII.
- Runs are the number of times to run the k-means algorithm (k-means is not guaranteed to find a globally optimal solution, and when run multiple times on a given dataset, the algorithm returns the best clustering result).
- Initialization Steps determines the number of steps in the k-meansII algorithm.
- Epsilon determines the distance threshold within which we consider k-means to have converged.

### CODE

```
selectRDD = selectData.map(lambda s: s.features)
(trainingData, testData) = selectRDD.randomSplit([0.6, 0.4])
clusters = KMeans.train(trainingData, 2, maxIterations=10,
```

```
runs=10, initializationMode="random")
```

```
def error(point):  
    center = clusters.centers[clusters.predict(point)]  
    return sqrt(sum([x**2 for x in (point - center)]))
```

```
WSSSE = trainingData.map(lambda point: error(point)).reduce(lambda x, y: x + y)  
print("Within Set Sum of Squared Error = " + str(WSSSE))
```

### RESULT

Within Set Sum of Squared Error = 3347.68

## 5.3 RANDOM FOREST CLASSIFICATION

Random forests train a set of decision trees separately, so the training can be done in parallel. The algorithm injects randomness into the training process so that each decision tree is a bit different. Combining the predictions from each tree reduces the variance of the predictions, improving the performance on test data.

The randomness injected into the training process includes:

- Subsampling the original dataset on each iteration to get a different training set (a.k.a. bootstrapping).
- Considering different random subsets of features to split on at each tree node.

Apart from these randomizations, decision tree training is done in the same way as for individual decision trees.

### CODE

```
#Creating RDD of LabeledPoints  
lpSelectData = selectData.map(lambda x : (x.id, LabeledPoint(x.label,x.features)))  
#Instantiating string indexer for random forest  
stringIndexer = StringIndexer(inputCol="label", outputCol="indexed")  
#fitting the data in stringindexer  
si_model = stringIndexer.fit(selectData)  
#transforming the data  
transformData = si_model.transform(selectData)  
#Splitting the data for training and test  
(trainingData, testData) = transformData.randomSplit([0.6, 0.4])  
#instantiating Random forest model  
randomForest = RandomForestClassifier(numTrees=2, maxDepth=2,  
labelCol="indexed", seed=42)  
#training the model  
randomForestModel = randomForest.fit(trainingData)
```

```
#transforming test data
result = randomForestModel.transform(testData)
#calculating the accuracy and printing it.
accuracy = result.filter(result.label == result.prediction).count() / float(testData.count())
print("Accuracy = " + str(accuracy))
```

## RESULT

Accuracy = 61.215326

## 6. COMPARISON

### 6.1 TF - IDF

We did run the above said three models for TF – IDF and the results are listed below

Algorithm	Result
Random Forest Classification	Accuracy = 55.600362
Naive Bayes	accuracy = 72 recall = 71.5 precision: 72.05
TF_IDF LR	areaUnderPR = 79.20658874070523 Area under ROC = 80.61931789957208 training error = 26.133306765
ML TF_IDF LR cross validator	areaUnderPR = 80.65552181670259 Area under ROC = 82.75891274493837 training error = 24.4788029925

Logistic Regression with ML was better of all three because of the ML implementation with cross validator.

### FEATURE SELECTION:

Based on the results we decided to run the Logistic Regression With LBFGS of mllib because the results were better in Logistic Regression with ML implementation.

### CODE:

```
#Creating RDD of LabeledPoints
lpSelectData = selectData.map(lambda x : (x.id, LabeledPoint(x.label,x.features)))
#Splitting the data for training and test
(trainingData, testData) = lpSelectData.randomSplit([0.9, 0.1])
# training the Logistic regression with LBFGS model
```

```
lrm = LogisticRegressionWithLBFGS.train(trainingData.map(lambda x: x[1]),
iterations=10)
#fetching the labels and predictions for test data
labelsAndPreds = testData.map(lambda p: (p[0],p[1].label, lrm.predict(p[1].features)))
#calculating the accuracy and printing it.
accuracy = labelsAndPreds.filter(lambda (i, v, p): v == p).count() / float(testData.count())
print("Accuracy = " + str(accuracy))
```

### Result

Accuracy = 74.8927

## 6.2 WORD2VEC

We did run the above said three models for word2vec and the results are listed below

Algorithm	Result
Random Forest Classification	Accuracy = 61.215326
K means Clustering	Within Set Sum of Squared Error = 3347.68
Logistic Regression with LBFGS	Accuracy = 84.3182696022

Logistic Regression with LBFGS gave better result than the other two models.

## 7. FUTURE REVIEW PREDICTION

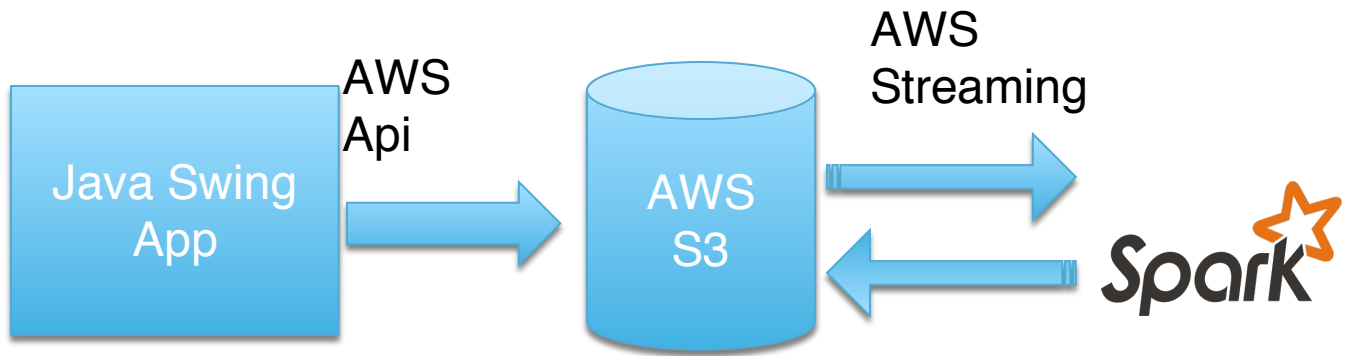
If we want to predict the sentiment of a new review we have a User Interface where the user has to enter his/her review and click on Submit button.

### WORKFLOW

The UI is developed is using JAVA Swing. From the Java application the review is read and stored in Amazon S3 using Amazon S3 API. From the storage using AWS streaming the review is given for prediction to the model.

### AWS STREAMING

Amazon CloudFront, the easy-to-use content delivery service, now supports the ability to stream audio and video files. Traditionally, world-class streaming has been out of reach of for many customers – running streaming servers was technically complex, and customers had to negotiate long- term contracts with minimum commitments in order to have access to the global streaming infrastructure needed to give high performance.



After prediction the result is written back to S3 using AWS Streaming.

## 8. CONCLUSION

Logistic Regression with ML gave better result in TF – IDF. So we decided to implement Logistic Regression with LBFGS in ml. Word2Vec is a deep learning model that is better than TF – IDF. Since we got better result for Logistic regression with LBFGS in TF – IDF we decided to implement Logistic Regression with LBFGS in Word2Vec, which gave better result as expected than every other model.

## 9. REFERENCES

<http://spark.apache.org/docs/latest/>

<https://databricks.com/blog/2015/07/29/new-features-in-machine-learning-pipelines-in-spark-1-4.html>

<http://alexminnaar.com/word2vec-tutorial-part-i-the-skip-gram-model.html>

<https://code.google.com/p/word2vec/>