# Automated Tagging of Blogs Using LDA and Sentiment Analaysis

Anurag Mishra

Computer Science, NYU Tandon

aa9279@nyu.edu

Advised by: Professor Gustavo Sandoval

## ABSTRACT

Blog tagging is a concept where we add extra information usually in form of keywords or tags to the blogs to summarise it quickly and using just a cluster of tags.

In the generation of Machine Learning and AI and the growing need to handle massive dataset has become the need of the hour. Every day almost 5 million blogs are posted and it has shown an incremental pattern year-on-year.

It is understandable how much businesses can achieve if they leave no data untamed. But with such immense data comes the need of technology to handle the kind of data hand.

## KEYWORDS

Automated Tagging,Machine Learning, LDA,Sentiment Analysis, Kafka Queue, Zookeeper, Elastic Search, Word2Vec,Query

## 1 INTRODUCTION

The goal of this project was to generate the tags for the different kinds of blogs or articles. Users can specify the data and the query terms to search for the related blogs. All the related blogs are then be dumped into a Kafka queue which will process the streaming data and generate the associated tags with each blog Users can go to the specific blogs according to their choice. Its also include a pre-processing i.e, word2vec and TF-IDF. After the preprocessing the blogs are tagged using topic modeling- Latent Dirichlet Allocation(LDA) in pySpark. To minimize the storage requirement for the whole system we purge the resulting data after publishing the results to the user. We used pySpark, Apache Kafka and python which will be used to preprocess the massive data.

**Figure 1: Typical Json Tuple from Our Dataset**



## 2 HYPOTHESIS

"Well-defined tagging of articles can help us retrieve cluster of similar articles pertaining to a keyword which can be done by constructing an article database with tags that capture the essence of the article."

## 3 SYSTEM DESIGN AND OVERVIEW

### 3.1 Data Description and Source

The data consists of JSON files which have articles based on the following criteria:

News publishers: Bloomberg.com, CNBC.com, reuters.com, wsj.com, fortune.com Language: Only English Country: United States News Category: Financial news only The source for

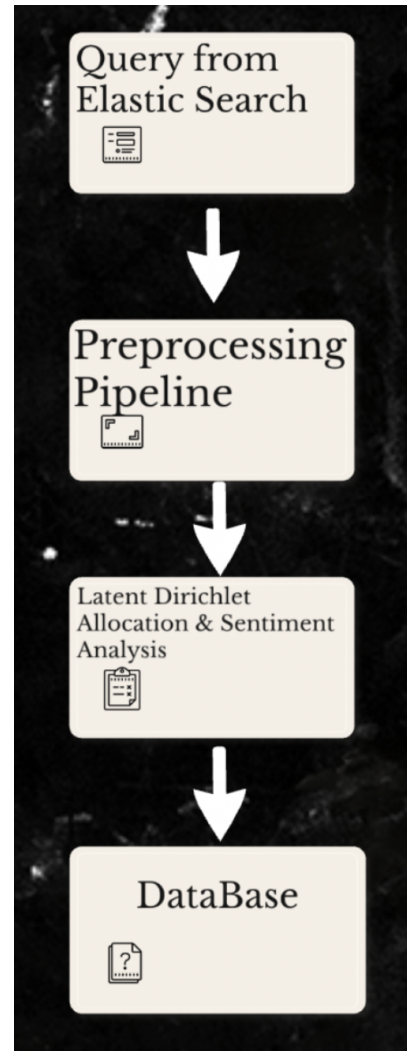**Figure 2: Shows the System Overview**



**Figure 3: System Architecture Diagram**



the articles (Archive source: httsps://Webhose.io/archive Total Data size is 1GB and the number of json files are 306K in 5 folders each folder for each news channel

## 3.2 System Flow

The UI allows the user to specify an umbrella query term. This query term is used to gather the relevant documents by querying an elastic search index. All the data collected using elastic search are then dumped into a Kafka queue, which helps in processing the streaming data and managing huge workloads due to the high amount of the relevant articles. The ML pipeline process consumes data from the Kafka queue. Incoming documents are subject to the cleaning of data for further analysis. It then goes through the core of the system for automated tagging. This is achieved by using Latent Dirichlet Allocation (LDA) to generate the tags relevant to each article. We also perform sentiment analysis in each article using NLTK. The generated tags, along with the sentiment, are then dumped into another ElasticSearch index. The UI will fetch the results from the ElasticSearch index and publish them. Now the user can perform searches from the search screen present in the UI and will be presented with the tags relevant to each article. The system also provides a link to every article in case the user wants to read the article.

## 4 TECHNOLOGIES USED FOR IMPLEMENTATION

(1) Elastic Search Elasticsearch is the distributed search and analytics engine at the heart of the Elastic Stack. Logstash and Beats facilitate collecting, aggregating, and enriching your data and storing it in Elasticsearch.
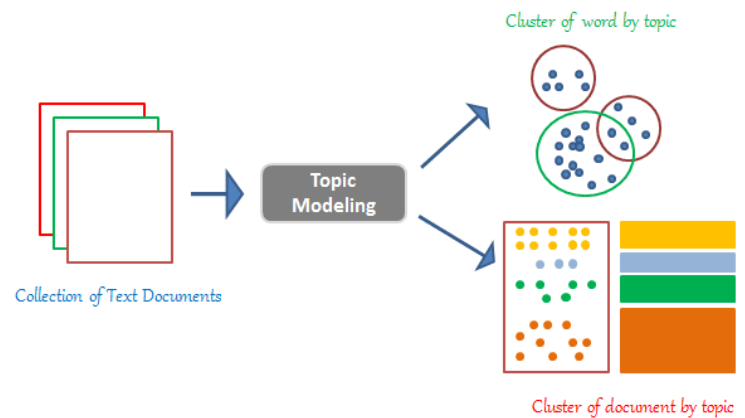
Kibana enables you to interactively explore, visualize, and share insights into your data and manage and monitor the stack. Elasticsearch is where the indexing, search, and analysis magic happens. Elasticsearch provides near real-time search and analytics for all types of data. Whether you have structured or unstructured text, numerical data, or geospatial data, Elasticsearch can efficiently store and index it in a way that supports fast searches.

(2) Kafka Apache Kafka is a community distributed event streaming platform capable of handling trillions of events a day. Initially conceived as a messaging queue, Kafka is based on an abstraction of a distributed commit log. Since being created and open sourced by LinkedIn

**Figure 4: Sequence Diagram**



**Figure 5: A high Level Overview of Latent Dirichlet Allocation**



in 2011, Kafka has quickly evolved from messaging queue to a full-fledged event streaming platform.

(3) PySpark PySpark is an interface for Apache Spark in Python. It not only allows you to write Spark applications using Python APIs, but also provides the PySpark shell for interactively analyzing your data in a distributed environment. PySpark supports most of Spark's features such as Spark SQL, DataFrame, Streaming, MLlib (Machine Learning) and Spark Core.

(4) React We wanted our system to be platform-independent and enable the user to access the system from a PC, Mac, or even a smartphone. Our frontend UI is, therefore, built using React. React helps in efficiently updating and rendering the right components of our UI when the data changes. It helps the system deliver beautiful visualizations of the data being processed, enabling the user to make more insightful decisions and ease his process of document discovery and understanding.

(5) Flask Flask is a lightweight WSGI web application framework. It is designed to make getting started quick and easy, with the ability to scale up to complex applications. It began as a simple wrapper around Werkzeug and Jinja and has become one of the most popular Python web application frameworks. Flask offers suggestions, but doesn't enforce any dependencies or project layout. It is up to the developer to choose the tools and libraries they want to use. There are many extensions provided by the community that make adding new functionality easy.

# 5 MACHINE LEARNING METHODS APPLIED

(1) Latent Dirichlet Allocation(LDA): A statistical model for discovering the abstract topics also known as topic modeling. Topic modeling is a method for unsupervised classification of documents, similar to clustering on numeric data, which finds some natural groups of items (topics) even when we're not sure what we're looking for.

(a) Assumptions:

 (i) Each document is just a collection of words or a "bag of words". Thus, the order of the words and the grammatical role of the words (subject, object, verbs, ...) are not considered in the model.

 (ii) Words like am/is/are/of/a/the/but/... don't carry any information about the "topics" and therefore can be eliminated from the documents as a pre-processing step. In fact, we can eliminate words that occur in at least For example, if our corpus contains only medical documents, words like human, body, health, etc might be present in most of the documents and hence can be removed as they don't add any specific information which would make the document stand out.

 (iii) We know beforehand how many topics we want. 'k' is pre-decided. All topic assignments except for the current word in question are correct, and then updating the assignment of the current word using our model of how documents are generated

## 5.1 How Does LDA work?

There are 2 parts in LDA:

(a) The words that belong to a document, that we already know.

**Figure 6: Code snippet of LDA Pipeline Handler**

```python
def create_lda_pipeline():
    # hyper-parameters
    num_topics = 10
    max_iterations = 10

    # pipeline
    tokenizer = Tokenizer(inputCol="content", outputCol="tokenized")
    cleaned = StopWordsRemover(inputCol="tokenized", outputCol="cleaned")
    tf = CountVectorizer(inputCol="cleaned", outputCol="raw_features", vocabSize=10, minDF=1.0)
    idf = IDF(inputCol="raw_features", outputCol="features")
    lda_model = LDA(featuresCol="features", k=num_topics, maxIter=max_iterations)
    return Pipeline(stages=[tokenizer, cleaned, tf, idf, lda_model])

spark = create_spark_session()

pipeline = create_lda_pipeline()
```

(b) The words that belong to a topic or the probability of words belonging into a topic, that we need to calculate.

*5.1.1* ***The Algorithm*** *.* Go through each document and randomly assign each word in the document to one of k topics (k is chosen beforehand). For each document d, go through each word w and compute :

(a) **p(topic t | document d): the proportion of words in document d that are assigned to topic t** tries to capture how many words belong to the topic t for a given document d. Excluding the current word. If a lot of words from d belongs to t, it is more probable that word w belongs to t.

**words in d with t +alpha/ words in d with any topic+ k*alpha**

(b) **(word w| topic t)**: the proportion of assignments to topic t over all documents that come from this word w. Tries to capture how many documents are in topic t because of word w. LDA represents documents as a mixture of topics. Similarly, a topic is a mixture of words. If a word has high probability of being in a topic, all the documents having w will be more strongly associated with t as well. Similarly, if w is not very probable to be in t, the documents which contain the w will be having very low probability of being in t, because rest of the words in d will belong to some other topic and hence d will have a higher probability for those topic. So even if w gets added to t, it won't be bringing many such documents to t.

(c) Update the probability for the word w belonging to topic t, as

```
p(word w with topic t) = p(topic t | document d) * p(word w | topic t)
```

## 5.2 Sentiment Analysis

Text sentiment classification is a fundamental sub-area in natural language processing. The sentiment classification algorithm is highly domain-dependent. For example, the
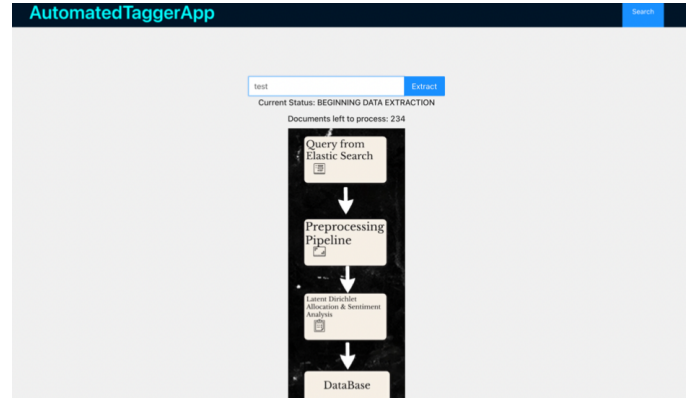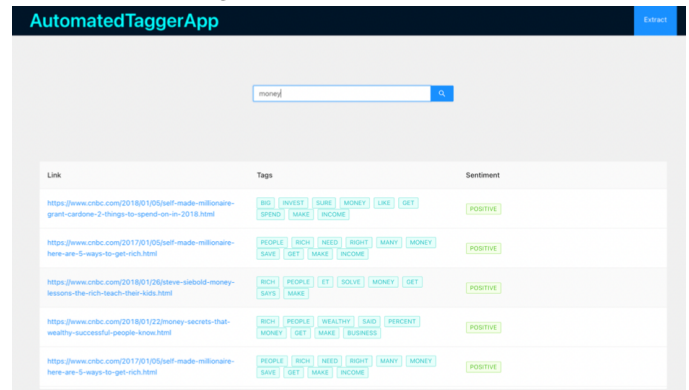
**Figure 7: Data Extraction Screen**



**Figure 8: Search Screen**



phrase "traffic jam" expresses negative sentiment in the sentence "I was stuck in a traffic jam on the elevated for 2 h." But in the domain of transportation, the phrase "traffic jam" in the sentence "Bread and water are essential terms in traffic jams" is without any sentiment. The most common method is to use the domain-specific data samples to classify the text in this domain. However, text sentiment analysis based on machine learning relies on sufficient labeled training data.

*5.2.1 Algorithm for sentiment Analysis.* Deciding to not make things complicated I just created a sentiment function which took the probabilty of each word and classifies it based on the if the probabilty was > 0.5 (Positive) or less than 0.5 (negative) and other wise (neutral)

A simple classfier function can be repersented as :

$$C(x) = \begin{cases} positive; x > 0 \cdot 5 \\ negative; x < 0 \cdot 5. \\ neutral; otherwise \end{cases}$$

**where C(x) is the sentiment classifier function**

**Figure 9: Code for C(x)**

```python
# sentiment analyzer function
def sentiment_analysis(data):
    analyzer = SentimentIntensityAnalyzer()
    score = analyzer.polarity_scores(data)
    if score['compound'] >= 0.05:
        return 'positive'
    elif score['compound'] <= -0.05:
        return 'negative'
    else:
        return 'neutral'
```

## 6   DEPLOYMENT

Install all dependencies (pip install -r requirements.txt )

(1) Python >= 3.6
(2) Spark
(3) Kafka
(4) ElasticSearch
(5) NPM or React for UI deployment

subsectionSetup

(1) Start the zookeeper, kafka and elastic search services using the following commands:
  (a) zookeeper-server-start
  (b) /usr/local/etc/kafka/zookeeper.properties
  (c) kafka-server-start /usr/local/etc/kafka/server.properties
  (d) elasticsearch
(2) Extract the project zip
(3) cd back into the project directory and run the following python scripts as shown below
(4) python3 backend.controller.ApiController.py
(5) python3 backend.ldapipeline.backendEngine.py
(6) To load the dataset into elasticsearch, run 'python -m backend.dataHandler.initES'
(7) Finally, cd into the UI directory and Start the react app by running the command " npm start". Then access the UI on the browser at http://localhost:3000/

## 7   CONCLUSION AND APPLICATIONS

Undocumented and poorly documented data causes problems for most data teams. Data users would love to have all of their data documented, but the effort to do this is far too large. Building an automated data documentation system to instantly document summarise the data is the need of the hour.

Documenting data is extremely important for any successful data initiative, but most people DESPISE documenting it. Documentation takes time and effort and is often created on the fly resulting in it not being a priority.This tool can this

approach by building an automated metadata tagging system to automatically classify data tags based on requirements.

There are many applications already in the run like instagram, twitter etc uses hashtags to store information about a post or tweet. This tools help to do the same by first extracting the keywords from the data and then user searching to find relevant documents.

Not only this, but it has gigantic application in optimising the Search engine where people are wasting time on catalogs and manually finding the keywords to link.

I also think that the app created during this project is not self sufficient for the application stated above but its a proof of concept from where amazing applications can be built.

## 8   REFERENCES

(1) What is Elasticsearch? | Elasticsearch Guide [7.14]
(2) https://medium.com/technofunnel/apache-kafka-in-5-minutes-c92c43ba3f39
(3) https://www.kaggle.com/jeet2016/us-financial-news-articles/version/1
(4) PySpark Documentation — PySpark 3.1.2 documentation
(5) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8509032/
(6) https://towardsdatascience.com/latent-dirichlet-allocation-lda-9d1cd064ffa2
(7) https://towardsdatascience.com/using-word2vec-to-analyze-news-headlines-and-predict-article-success-cdeda5f14751
(8) https://thinkinfi.com/latent-dirichlet-allocation-for-beginners-a-high-level-overview/