

# News Classifier Project Report

## Objective:

To train machine learning model and classify upcoming news on the fly with good accuracy by building end to end machine learning pipeline.

To make containerized application which is scalable, robust, fault tolerant.

## Planning:

We are using agile methodology to build the project. Task level details are mentioned in below gantt chart.

#Sprint: 4

#People: 2

Project members:

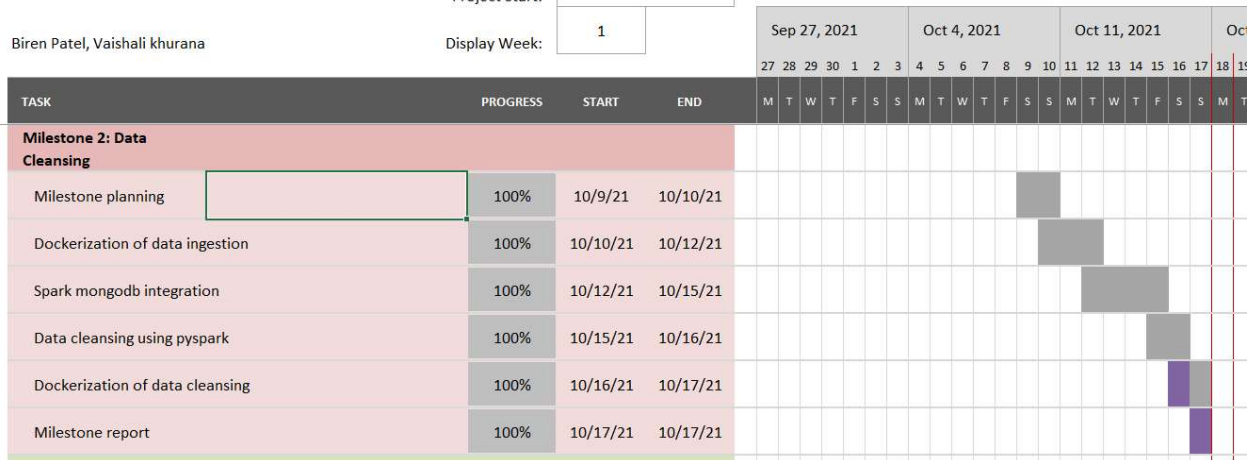
Biren Patel, Vaishali khurana

Project Start:

Tue, 9/28/2021

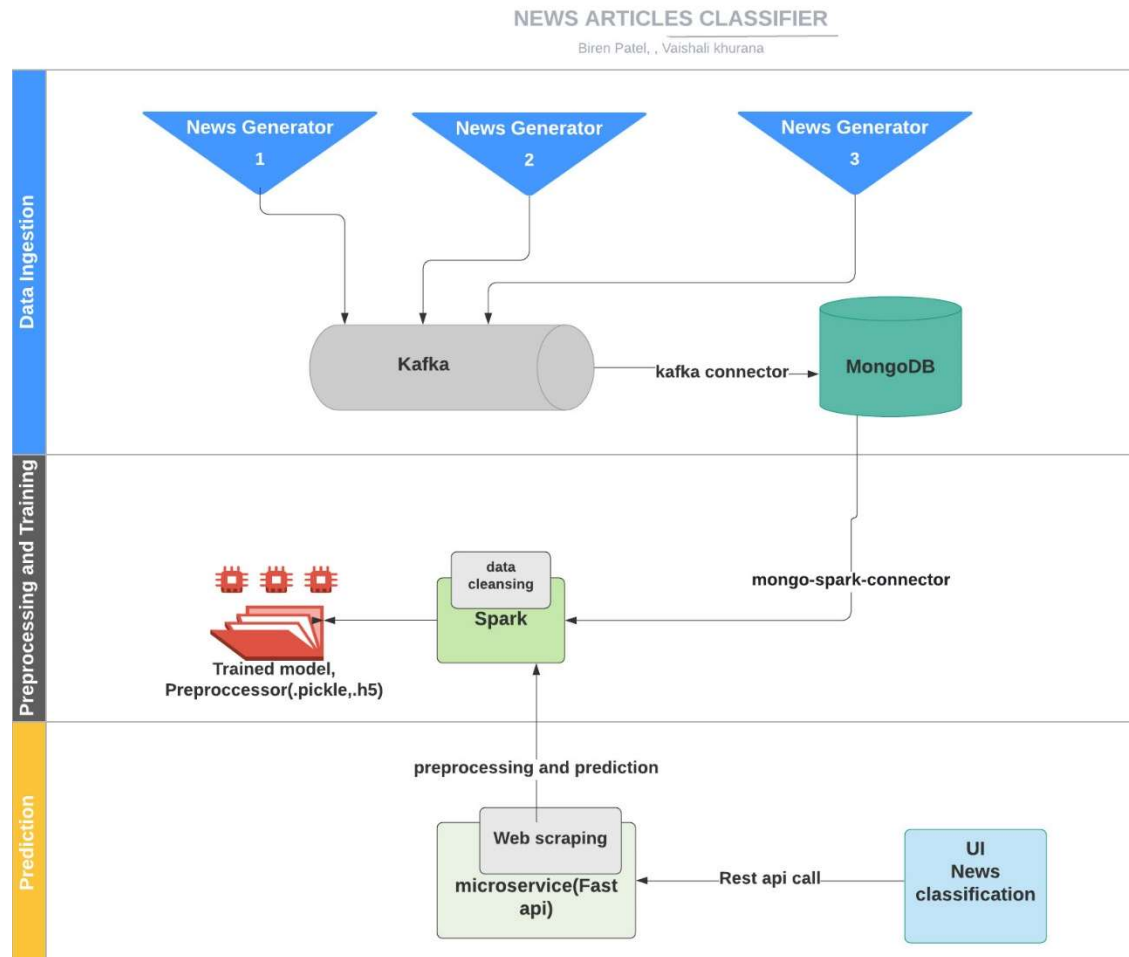
Display Week:

1



# News Classifier Project Report

## Architecture:



## Components and description:

Component	Description
Data Ingestion	Multithreaded service that is collecting the data from web using rapid api and custom news generator, passing it to kafka queue. Kafka consumer is used to sync the data between kafka and database. Finally, data is dumped into MongoDB database
Preprocessor and trainer	Reading data from mongo db to spark session using mongo-spark-connector. As a part of feature selection we are using category and summary columns.

## News Classifier Project Report

	Data cleansing including stop word removal, tokenization, tf-idf vectorization using pyspark. This cleaned data will be used for training.
Prediction	Scraping the weblink provided by user in UI, clean it and predict the news using rest api.

Since data cleansing is the second milestone we have explained it in detail here.

### Data Cleansing:

- **Environment details:**

Docker environment with kafka broker running on 9092 port, zookeeper running on 2181, mongo-db running on 27017, mongo-express running on 8082, spark-master running on 8080, spark-worker running on 8081, producer is connecting to kafka and producing news records to kafka queue, consumer is consuming news records from kafka and dumping to mongo db, preprocessor is cleansing the news data on spark node.

- **What goes in as an input:**

News records saved in news\_collection under news database MongoDB

### Viewing Collection: news\_collection

[New Document](#) [New Index](#)

Simple

Advanced

[Find](#)

Delete all 139 documents retrieved

[← First](#) [← Prev](#) [Next →](#) [Last →](#)

_id	title	date	summary	category	source
616d7fc390f3e115c28e3128	There Were 2 Mass Shootings In Texas Last Week, B...	2018-05-26	She left her husband. He killed their children. J...	CRIME	https://www.huffingtonpost.com/entry/texas-amanda...
616d7fc790f3e115c28e3129	Will Smith Joins Diplo And Nicky Jam For The 2018...	2018-05-26	Of course it has a song.	ENTERTAINMENT	https://www.huffingtonpost.com/entry/will-smith-j...
616d7fca90f3e115c28e312a	Hugh Grant Marries For The First Time At Age 57	2018-05-26	The actor and his longtime girlfriend Anna Eberst...	ENTERTAINMENT	https://www.huffingtonpost.com/entry/hugh-grant-m...

## News Classifier Project Report

- **How the input is being processed:**

Records from mongo db collection is processed using pyspark over spark node which is connected through mongo\_spark\_connector.

As a part of feature selection we are using category and summary columns.

Data cleansing including stop word removal, tokenization, tf-idf vectorization using pyspark

```
root
|-- _id: struct (nullable = true)
|   |-- oid: string (nullable = true)
|-- category: string (nullable = true)
|-- date: string (nullable = true)
|-- source: string (nullable = true)
|-- summary: string (nullable = true)
|-- title: string (nullable = true)

cleaning start
['_id', 'category', 'date', 'source', 'summary', 'title']
```

After feature selection and tf-idf vectorization:

```
root
|-- summary: string (nullable = true)
|-- category: string (nullable = true)
|-- tf: vector (nullable = true)
|-- idf: vector (nullable = true)
|-- label: double (nullable = false)
```

## News Classifier Project Report

- **What comes out as an output:** Processed and cleaned data, saved the preprocessed pipeline

```
2021-10-18 14:27:01 INFO DAGScheduler:54 - Job 3 finished: showString at NativeMethodAccessorImpl.java:0, took 0.050267 s
+-----+-----+-----+-----+
| category| summary| tf| idf|label|
+-----+-----+-----+-----+
| CRIME|She left her husb...(10000,[1662,3562...|(10000,[1662,3562...| 7.0|
| ENTERTAINMENT|Of course it has ...(10000,[1916,2460...|(10000,[1916,2460...| 1.0|
| ENTERTAINMENT|The actor and his...(10000,[512,2410...|(10000,[512,2410...| 1.0|
| ENTERTAINMENT|The actor gives D...(10000,[512,1384...|(10000,[512,1384...| 1.0|
| ENTERTAINMENT|The "Dietland" ac...(10000,[2678,3624...|(10000,[2678,3624...| 1.0|
+-----+-----+-----+-----+
only showing top 5 rows
```

- **Tools/libraries used:** Docker, Pycharm, kafka, MongoDB, zookeeper, Spark, kafka, pyspark, MongoExpress

## Challenges encountered and how we tackled them:

We faced below mentioned challenges. However, we have resolved them.

- Connecting MongoDB and spark
  - o Tried different ways of connecting spark to mongodb then implemented it with mongo-spark-connector
- Docker networking challenges
  - o Had to explore different docker-compose versions and related configuration for mac-os and windows.
- Authorization issue in MongoDB connection
  - o Configured appropriate parameters to resolve it
- We switched to Mongoddb to support auto scalability in future for large volume of text and unstructured data.
  - o Developed consumer service to connect kafka broker to mongoDb

## Future Scope:

- Further scale optimizations
- Implementing re-training mechanism using feedback feature
- Once we have large volume of labeled data, we will train our own model in place of transfer learning
- Adding multi language support

## Github link:

<https://github.com/biren162/Capstone/tree/master/Milestone-2>