
News Data Classification

Student :

Afli RAMZI

Mhamdi OUMAYMA

Teacher :

Maher HENI

January 12, 2023

Contents

1	Introduction	3
2	Producer	3
3	Kafka Broker	4
4	Consumer	5
5	ML Algorithm	6
6	Elastic search & Kibana	9
7	Conclusion	11

List of Figures

1	Producer python scripts	3
2	Python kafka connector	3
3	Starting Zookeeper	4
4	Starting kafka	4
5	Kafka Ui	4
6	Calling the consumer	5
7	Receiving Data by the consumer	5
8	Libraries import	6
9	Creating pipeline	7
10	Fitting the model	7
11	Performance of the model	8
12	Docker compose elasticseach part	9
13	Docker compose kibana part	9
14	Starting docker containers	10
15	Indexing data	10
16	Kibana visualisations	11
17	Kibana visualisations	11

List of Tables

1 Introduction

Today, one of the most crucial aspects of our daily lives is news. People continue to concentrate on the news in order to develop a sense of the national perspective. However, because of the increasing number and variety of topics available, it can be difficult for listeners to focus on specific ones. Because of this, our main objective is to develop a real-time streamer that will be used to categorize news topics using news data. To complete this task, we attempted to implement a kind of seamless chain from the collection of the data to its visualization in the Kibana Dashboard.

2 Producer

In this section we will try to explain how we made a kafka producer that produces(streams) data from an api and store these data inside kafka broker.

The Guardian API, an API that provides real-time news, is our intended data source. The first session of our project involves writing Python scripts to stream the data and store it inside the Kafka broker. We must specify the api-url and api-key in order to retrieve the data from this API (this key is provided after creating account in the Guardian's web page) . After that, we connected our producer to the cloud-based kafka broker using a Python connector.

```
if __name__ == "__main__":

    key = '53baf88d-42eb-4df9-bf3e-f30d773b1726'
    fromDate = '2013-01-1'
    toDate = '2022-08-8'

    url = 'http://content.guardianapis.com/search?from-date='+ fromDate + '&to-date='+ toDate
    all_news = getData(url)

    if len(all_news)>0:
        prod = connect_kafka_producer()
        for story in all_news:
            publish_message(prod, 'test', story)
            time.sleep(1)
        if prod is not None:
            prod.close()
```

Figure 1: Producer python scripts

```
def connect_kafka_producer():
    _producer = None
    try:
        _producer = KafkaProducer(bootstrap_servers=['20.216.153.173:9092'], api_version=(0, 10
```

Figure 2: Python kafka connector

3 Kafka Broker

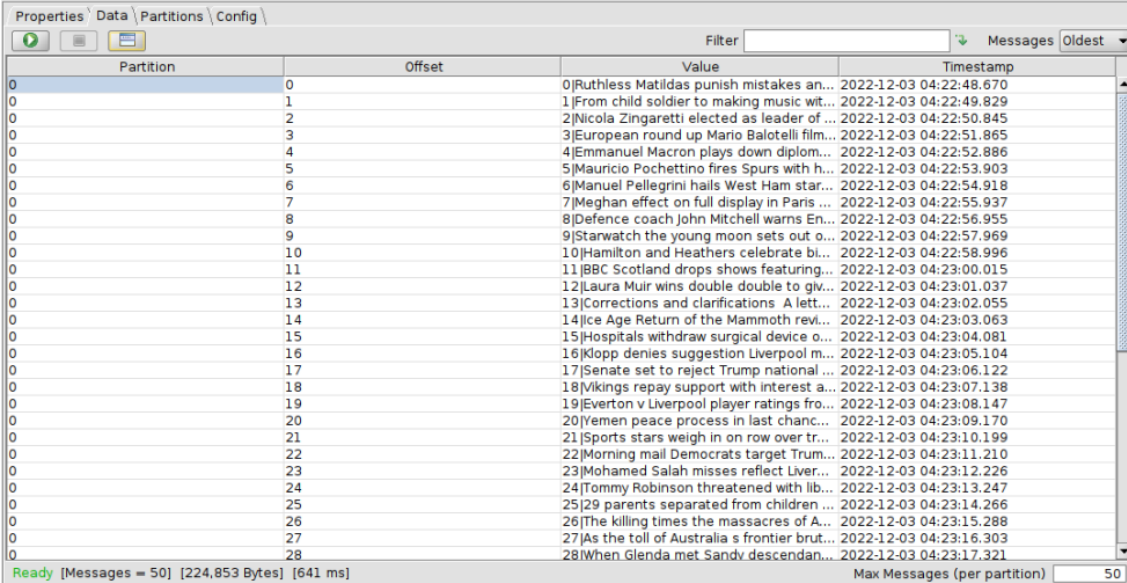
In order to store these news. Kafka broker will serve as good choice. Our Kafka broker is deployed on Azure cloud in a separate virtual machine.

```
ramzi-afl@kafka-maher-vm:~$ ~/kafka_2.13-3.0.0/bin/zookeeper-server-start.sh ~/kafka_2.13-3.0.0/config/zookeeper.properties
[2023-01-07 16:51:15,917] INFO Reading configuration from: /home/ramzi-afl/kafka_2.13-3.0.0/config/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2023-01-07 16:51:15,927] INFO clientPortAddress is 0.0.0.0:2181 (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2023-01-07 16:51:15,927] INFO secureClientPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2023-01-07 16:51:15,927] INFO observerMasterPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2023-01-07 16:51:15,927] INFO metricsProvider.className is org.apache.zookeeper.metrics.impl.DefaultMetricsProvider (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
```

Figure 3: Starting Zookeeper

```
ramzi-afl@kafka-maher-vm:~$ ~/kafka_2.13-3.0.0/bin/kafka-server-start.sh ~/kafka_2.13-3.0.0/config/server.properties
[2023-01-07 16:52:55,903] INFO Registered kafka:type=kafka.Log4jController MBean (kafka.utils.Log4jControllerRegistration$)
[2023-01-07 16:52:56,337] INFO Setting -D jdk.tls.rejectClientInitiatedRenegotiation=true to disable client-initiated TLS renegotiation (org.apache.zookeeper.common.X509Util)
[2023-01-07 16:52:56,506] INFO Registered signal handlers for TERM, INT, HUP (org.apache.kafka.common.utils.LoggingSignalHandler)
[2023-01-07 16:52:56,510] INFO starting (kafka.server.KafkaServer)
[2023-01-07 16:52:56,511] INFO Connecting to zookeeper on localhost:2181 (kafka.server.KafkaServer)
[2023-01-07 16:52:56,530] INFO [ZooKeeperClient Kafka server] Initializing a new session to localhost:2181. (kafka.zookeeper.ZooKeeperClient)
[2023-01-07 16:52:56,545] INFO Client environment:zookeeper.version=3.6.3--6401e4ad2087061bc6b9f80dec2d69f2e3c8660a, built on 04/08/2021 16:35 GMT (org.apache.zookeeper.ZooKeeper)
[2023-01-07 16:52:56,546] INFO Client environment:host.name=kafka-maher-vm.internal.cloudapp.net (org.apache.zookeeper.ZooKeeper)
```

Figure 4: Starting kafka



Partition	Offset	Value	Timestamp
0	0	0 Ruthless Matildas punish mistakes an...	2022-12-03 04:22:48.670
0	1	1 From child soldier to making music wit...	2022-12-03 04:22:49.829
0	2	2 Nicola Zingaretti elected as leader of ...	2022-12-03 04:22:50.845
0	3	3 European round up Mario Balotelli film...	2022-12-03 04:22:51.865
0	4	4 Emmanuel Macron plays down diplom...	2022-12-03 04:22:52.886
0	5	5 Mauricio Pochettino fires Spurs with h...	2022-12-03 04:22:53.903
0	6	6 Manuel Pellegrini hails West Ham star...	2022-12-03 04:22:54.918
0	7	7 Meghan effect on full display in Paris ...	2022-12-03 04:22:55.937
0	8	8 Defence coach John Mitchell warns En...	2022-12-03 04:22:56.955
0	9	9 Starwatch the young moon sets out o...	2022-12-03 04:22:57.969
0	10	10 Hamilton and Heathers celebrate bl...	2022-12-03 04:22:58.996
0	11	11 BBC Scotland drops shows featuring...	2022-12-03 04:23:00.015
0	12	12 Laura Muir wins double double to glw...	2022-12-03 04:23:01.037
0	13	13 Corrections and clarifications A lett...	2022-12-03 04:23:02.055
0	14	14 Ice Age Return of the Mammoth revi...	2022-12-03 04:23:03.063
0	15	15 Hospitals withdraw surgical device o...	2022-12-03 04:23:04.081
0	16	16 Klopp denies suggestion Liverpool m...	2022-12-03 04:23:05.104
0	17	17 Senate set to reject Trump national ...	2022-12-03 04:23:06.122
0	18	18 Vikings repay support with interest a...	2022-12-03 04:23:07.138
0	19	19 Everton v Liverpool player ratings fro...	2022-12-03 04:23:08.147
0	20	20 Yemen peace process in last chanc...	2022-12-03 04:23:09.170
0	21	21 Sports stars weigh in on row over tr...	2022-12-03 04:23:10.199
0	22	22 Morning mail Democrats target Trum...	2022-12-03 04:23:11.210
0	23	23 Mohamed Salah misses reflect Liver...	2022-12-03 04:23:12.226
0	24	24 Tommy Robinson threatened with lib...	2022-12-03 04:23:13.247
0	25	25 29 parents separated from children ...	2022-12-03 04:23:14.266
0	26	26 The killing times the massacres of A...	2022-12-03 04:23:15.288
0	27	27 As the toll of Australia s frontier brut...	2022-12-03 04:23:16.303
0	28	28 When Glenda met Sandy descendan...	2022-12-03 04:23:17.321

Figure 5: Kafka Ui

4 Consumer

We will now develop the consumer after first creating the producer that streams data to the Kafka broker. The principal function of the consumer is to retrieve data from the Kafka broker.

after importing all needed libraries , we call the consumer .

```
1 consumer = KafkaConsumer('test',
2                           bootstrap_servers=['20.216.153.173:9092'])
```

Command took 0.15 seconds -- by oumayma.mhamdi@studentambassadors.com at 07/01/2023 23:38:18 on Oumayma Mhamdi's Cluster

Figure 6: Calling the consumer

the consumer retrieve data from kafka broker . We can see data received by the consumer in the figure 7 .

```
1 for text in consumer:
2     print(text.value)
3
4
```

b'29|Plan to axe 91 000 civil servants only possible with cuts to services Doubts have been cast on claims by Boris Johnson that it will be possible to go ahead with plans to axe 91 000 civil servants without harming frontline services The prime minister wrote in May to civil servants justifying plans for a reduction in headcount of almost 20 saying the government must reduce its costs just as many families are doing But a review by Steve Barclay Johnson's former chief of staff is reported to have caused the Treasury to have second thoughts about the plans given the potential impact on wider services Johnson had tasked the cabinet with cutting staff by a fifth telling ministers during an away day in Stoke on Trent earlier this year that every bit of cash saved on government spending could be better used elsewhere Trade unions representing civil servants seized on reports about Treasury reticence while Labour said it was clear that the disastrous plan would cause historic staff shortages causing huge delays for people when trying to access appointments and renew passports and driving licences Now the government has not only belatedly come to the same conclusion but also found this shoddy proposal would have cost the taxpayer over 1bn in redundancy payments said Rachel Hopkins MP Labour's shadow minister in the Cabinet Office The plans have been backed by Liz Truss the Conservative leadership frontrunner who has been accused of making ludicrous claims as she vowed to cut civil service salaries and reduce expenditure to recoup 11bn a year in a war on Whitehall waste However a Whitehall insider who was said to have worked on the plans to axe tens of thousands of civil servants was quoted in the Financial Times as saying that the prime minister had announced the move without fully thinking through the implications You can only deliver 91 000 cuts by actual cuts to major frontline services they added There's no way you can get to that number through efficiency savings or reductions in HQ staff Another Whitehall source was quoted as suggesting that a figure of 91 000 is not realistic Canceled

Figure 7: Receiving Data by the consumer

5 ML Algorithm

To classify News into several categories , we use the Naive Bayes algorithm .
we choose Naive Bayes for many reasons :

- It doesn't require as much training data.
- It handles both continuous and discrete data.
- It is highly scalable with the number of predictors and data points.
- It is fast and can be used to make real-time predictions

To construct the model we firstly import required libraries

The screenshot shows a Databricks notebook interface. The top bar includes the Microsoft Azure logo, the Databricks logo, a search bar, and the keyboard shortcut CTRL + P. The notebook title is 'spark_pipeline_nb_databricks' and the language is set to Python. The code editor displays two code blocks. The first block, labeled 'Cmd 1', contains import statements for various PySpark and MLlib libraries. The second block, labeled 'Cmd 2', shows the creation of a SQLContext and a new DataFrame schema. A command execution bar below the first block indicates it took 0.93 seconds to run. A warning message is visible at the bottom of the notebook.

```

1 from pyspark.ml import Pipeline
2 from pyspark.ml.classification import LogisticRegression, NaiveBayes
3 from pyspark.ml.feature import HashingTF, Tokenizer, StopWordsRemover
4 from pyspark import SparkContext
5 from pyspark.sql import SQLContext
6 from pyspark.sql.types import StructType, StructField, IntegerType, StringType, DoubleType
7 from pyspark.ml.tuning import ParamGridBuilder, TrainValidationSplit
8 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
9 from pyspark.mllib.evaluation import MulticlassMetrics

Command took 0.93 seconds -- by oumayma.mhamdi@studentambassadors.com at 16/12/2022 21:01:11 on Oumayma Mhamdi's Cluster

1 sqlContext = SQLContext(sc)
2 newDF = [
3     StructField("id", IntegerType(), True),
4     StructField("text", StringType(), True),
5     StructField("label", DoubleType(), True)]
6 finalSchema = StructType(fields=newDF)

/databricks/spark/python/pyspark/sql/context.py:117: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
warnings.warn(

```

Figure 8: Libraries import

Then we Load the data and split it into training and testing datasets .
The training dataset presents 80 '%of the data .
The testing dataset presents 20 '%of the data . After that we applicate the tokenizer then
The stopwords remover and finnnally the hashing TF

Microsoft Azure | databricks | Search | CTRL + P | ML_BM | ? | oumayma.mhamdi@studentamba... |

spark_pipeline_nb_databricks Python | File Edit View Run Help | Last edit was 5 minutes ago | Give feedback | Run all | Oumayma Mhamdi's CL... | Schedule | Share |

```

1 dataset = sqlContext.read.format('csv').options(header='true', schema=finalSchema, delimiter='|').load('/FileStore/tables/dataset.csv')
2 #types = [f.dataType for f in dataset.schema.fields]
3 #print(types)
4 dataset = dataset.withColumn("label", dataset["label"].cast(DoubleType()))
5 dataset = dataset.withColumn("id", dataset["id"].cast(IntegerType()))
6 training, test = dataset.randomSplit([0.8, 0.2], seed=12345)

(1) Spark Jobs
dataset: pyspark.sql.dataframe.DataFrame = [id: integer, text: string ... 1 more field]
training: pyspark.sql.dataframe.DataFrame = [id: integer, text: string ... 1 more field]
test: pyspark.sql.dataframe.DataFrame = [id: integer, text: string ... 1 more field]
Command took 9.43 seconds -- by oumayma.mhamdi@studentambassadors.com at 16/12/2022 21:02:03 on Oumayma Mhamdi's Cluster

Cmd 4
1
2 tokenizer = Tokenizer(inputCol="text", outputCol="words")
3 remover = StopWordsRemover(inputCol="words", outputCol="filtered")
4 hashingTF = HashingTF(inputCol=remover.getOutputCol(), outputCol="features")
5 lr = LogisticRegression(maxIter=2, regParam=0.001)
6 nb = NaiveBayes(smoothing=1.0, modelType="multinomial")
7 pipeline = Pipeline(stages=[tokenizer, remover, hashingTF, nb])

Command took 0.44 seconds -- by oumayma.mhamdi@studentambassadors.com at 16/12/2022 21:02:18 on Oumayma Mhamdi's Cluster

```

Figure 9: Creating pipeline

The next step is to train the model so that it can be able to classify data correctly .

Microsoft Azure | databricks | Search | CTRL + P | ML_BM | ? | oumayma.mhamdi@studentamba... |

spark_pipeline_nb_databricks Python | File Edit View Run Help | Last edit was 6 minutes ago | Give feedback | Run all | Oumayma Mhamdi's CL... | Schedule | Share |

```

1 # Fit the pipeline to training documents.
2 model = pipeline.fit(training)
3 result = model.transform(test)\
4     .select("features", "label", "prediction")
5 correct = result.where(result["label"] == result["prediction"])
6 accuracy = correct.count()/test.count()
7 print("Accuracy of model = "+str(accuracy))
8 test_error = 1 - accuracy
9 print ("Test error = "+str(test_error))
10

(6) Spark Jobs
result: pyspark.sql.dataframe.DataFrame = [features: udt, label: double ... 1 more field]
correct: pyspark.sql.dataframe.DataFrame = [features: udt, label: double ... 1 more field]
Accuracy of model = 0.5104166666666666
Test error = 0.48958333333333337
Command took 14.92 seconds -- by oumayma.mhamdi@studentambassadors.com at 16/12/2022 21:02:53 on Oumayma Mhamdi's Cluster

```

Figure 10: Fitting the model

The final step is to evaluate the performance of the classifier . the metrics of performance used are :

- The F1 Score
- The Recall
- The Precision

We saved the pretrained model to applicate it on testing data .

The screenshot shows a Databricks notebook interface. The top bar includes the Microsoft Azure logo, the Databricks logo, a search bar, and the text 'CTRL + P'. The notebook title is 'spark_pipeline_nb_databricks' and the language is set to 'Python'. The code editor contains the following Python code:

```

1 evaluator = MulticlassClassificationEvaluator(labelCol="label", predictionCol="prediction", metricName="f1")
2 metric = evaluator.evaluate(result)
3 print("F1 metric = "+ str(metric))
4
5 evaluator = MulticlassClassificationEvaluator(labelCol="label", predictionCol="prediction", metricName="weightedRecall")
6 metric = evaluator.evaluate(result)
7 print("Recall = "+ str(metric))
8
9 evaluator = MulticlassClassificationEvaluator(labelCol="label", predictionCol="prediction", metricName="weightedPrecision")
10 metric = evaluator.evaluate(result)
11 print("Precision = "+ str(metric))
12 model.save("nbmodelNew")

```

Below the code editor, the output of the command is displayed:

```

(10) Spark Jobs
F1 metric = 0.46053494606126183
Recall = 0.5104166666666667
Precision = 0.5964778414606092
Command took 10.49 seconds -- by oumayma.mhamdi@studentambassadors.com at 16/12/2022 21:04:05 on Oumayma Mhamdi's Cluster

```

Figure 11: Performance of the model

6 Elastic search & Kibana

In order to index our Data and visualize them .ELK stack is mandatory that's why we were Deploying ELK stack as docker containers using docker compose in the cloud .

```
ramzi-afli@kafka-maher-vm:~/Elk-stack-docker$ vim docker-compose.yml
version: '3.7'

services:
  # Author Afli Ramzi
  # The 'setup' service runs a one-off script which initializes the
  # 'logstash_internal' and 'kibana_system' users inside Elasticsearch with the
  # values of the passwords defined in the '.env' file.
  #
  # This task is only performed during the *initial* startup of the stack. On all
  # subsequent runs, the service simply returns immediately, without performing
  # any modification to existing users.
  setup:
    build:
      context: setup/
      args:
        ELASTIC_VERSION: ${ELASTIC_VERSION}
    init: true
    volumes:
      - setup:/state:Z
    environment:
      ELASTIC_PASSWORD: ${ELASTIC_PASSWORD:-}
      LOGSTASH_INTERNAL_PASSWORD: ${LOGSTASH_INTERNAL_PASSWORD:-}
      KIBANA_SYSTEM_PASSWORD: ${KIBANA_SYSTEM_PASSWORD:-}
    networks:
      - elk
    depends_on:
      - elasticsearch

  elasticsearch:
    build:
```

Figure 12: Docker compose elasticsearch part

```
kibana:
  build:
    context: kibana/
    args:
      ELASTIC_VERSION: ${ELASTIC_VERSION}
  volumes:
    - ./kibana/config/kibana.yml:/usr/share/kibana/config/kibana.yml:ro,Z
  ports:
    - "5601:5601"
  environment:
    KIBANA_SYSTEM_PASSWORD: ${KIBANA_SYSTEM_PASSWORD:-}
  networks:
    - elk
  depends_on:
    - elasticsearch

networks:
  elk:
    driver: bridge

volumes:
  setup:
  elasticsearch:
```

Figure 13: Docker compose kibana part

```

ramzi-afli@kafka-maher-vm:~/Elk-stack-docker$ sudo docker compose up -d
WARN[0000] mount of type 'volume' should not define 'bind' option
[+] Running 1/2
[+] Running 5/5stack-docker_elk      Created           0.1s
# Network elk-stack-docker_elk      Created           0.1s
# Container elk-stack-docker-elasticsearch-1 Started           0.8s
# Container elk-stack-docker-setup-1 Started           2.9s
# Container elk-stack-docker-logstash-1 Started           2.2s
# Container elk-stack-docker-kibana-1 Started           1.9s
ramzi-afli@kafka-maher-vm:~/Elk-stack-docker$ sudo docker ps
CONTAINER ID   IMAGE                                COMMAND                  CREATED        STATUS        PORTS
b28d1cc29e64   elk-stack-docker_kibana             "/bin/tini -- /usr/l..." 17 seconds ago Up 14 seconds 0.0.0.0:5601→5601/tcp, :
::5601→5601/tcp
                elk-stack-docker-kibana-1
2a1c63859a6   elk-stack-docker_logstash           "/usr/local/bin/dock..." 17 seconds ago Up 14 seconds 0.0.0.0:5044→5044/tcp, :
::5044→5044/tcp, 0.0.0.0:9600→9600/tcp, ::9600→9600/tcp, 0.0.0.0:50000→50000/tcp, ::50000→50000/tcp, 0.0.0.0:50000→50000/ud
p, ::50000→50000/udp elk-stack-docker-logstash-1
d4511462d6be   elk-stack-docker_elasticsearch     "/bin/tini -- /usr/l..." 17 seconds ago Up 15 seconds 0.0.0.0:9200→9200/tcp, :
::9200→9200/tcp, 0.0.0.0:9300→9300/tcp, ::9300→9300/tcp
                elk-stack-docker-elasticsearch-1
ramzi-afli@kafka-maher-vm:~/Elk-stack-docker$

```

[illegible]

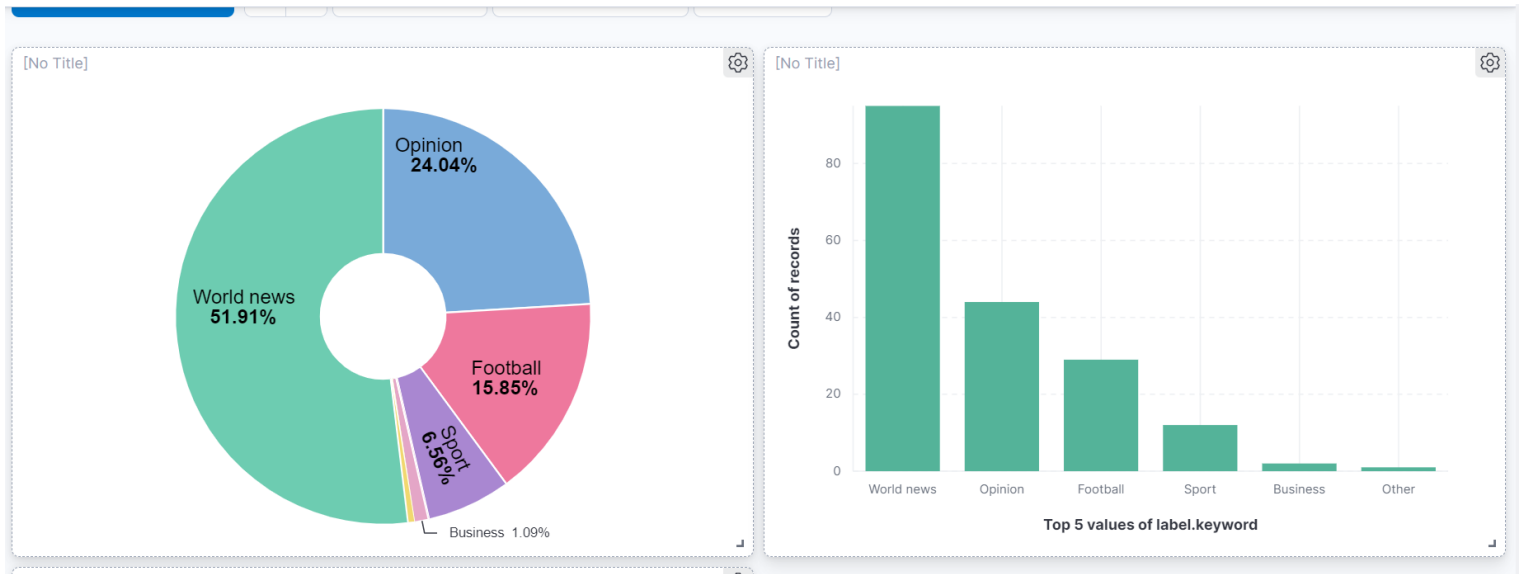


Figure 16: Kibana visualisations

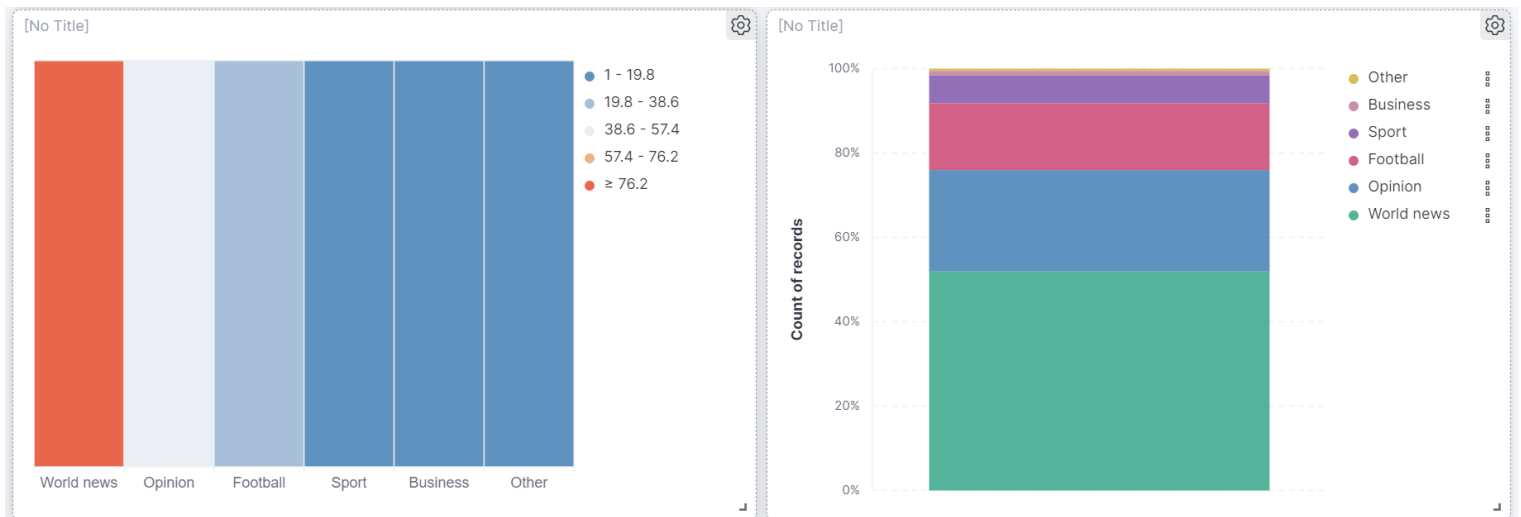


Figure 17: Kibana visualisations

==> These visualisations show the percentage of each news category based on the result given by the ML algorithm .

7 Conclusion

This project, named "News Data Classification", was conducted to create full intelligent chain of streaming , used to stream news in real time ,classify these data using machine learning algorithm , and indexing the news topics inside elastic search and visualize the most common using kibana dashboard .