

TEXT CLASSIFICATION USING APACHE SPARK

Overview

Apache spark

Text classification is critical task in processing of large data and documents into several categories for managing text classification process. It provides an outline related to message, email, tweets and text using the machine learning algorithm. Apache Spark is a multi-language engine that is possible to use in the text classification process. This big data analytics tools has been implemented using the source analysis for real time workload management. This tool is used for the big data analysis on very large dataset. In this text classification process, pyspark will be used by using multiple Apache Spark features including DataFrame, Spark SQL, MLlib and Spark Core [2]. In email, spam detection is important which will be solved through this machine learning model. It has developed a strong binding among different types of software language using machine learning application.

The Apache Spark has spark core which is playing the principle role in managing the task transmission, functionality and scheduling in challenging business environment. *Resilient Distributed Dataset (RDD)* concepts have been used in the Apache Spark for advanced computation systems for the analytical process. The RDD can be used in operational process to manage the processing of user information into *Directed Acyclic Graph (DAG)* to analyse the sequence of different nodes. Spark SQL has been utilised in the creation of RDD to enable mapping framework for filtration and aggregation of sample datasets. Core drivers of Spark has used for simplification of tasks according to the distribution channel for managing the sustainability. The Spark SQL has been applied for project interface management to deal with the standard structure to create better interface development [3]. Along with the SQL interface, the data frame is used for examine the queries and types of data sets across different locations. Using the Python language, transmission of data has been unified for creation different batches. Here, the Exploratory Data Analysis (EDA) has been used for creation of better data scaling process to manage the fundamental challenges.

Machine learning algorithm has been implemented to develop appropriate ecosystems for project data classification process. Jupyter notebook can develop the best environment for any platform to run effective computational. It consists of mathematical equations, statistical modelling, codes, narrative text and media visualisation framework to manage the easy classification process. Here, the Jupyter has used advanced interface which can support data

tools like Apache Spark to manage analytical process for improving classification process. Python is the best language in respect to the Apache Spark to manage the robust data classification process. Therefore, *interactive developing environment (IDE)* has been used for creation of integrated notebook creation to get better translation to manage scripts [1]. Different types of extensions have been used for selecting the better execution for functionality management.

Key features of Apache Spark

The Apache Spark has several key features for managing the functionality of data classification.

Data Streaming: Real time streaming is necessary for managing effective data processing for increasing data classification. Here, different types languages have been used for managing the batches.

SQL analytics: Different types of queries and distribution channel have been used for improving the reporting and dash boarding framework for managing efficiency in warehouse management.

Data science: Effective sampling has been used for performing better scaling of largedatasets for improving the fundamental challenges. Distributed framework has been applied in the processing of parallel for managing the accuracy of classification process.

Machine learning: Machine learning and python coding have been used for Apache Spark to create clusters regarding consistent machine operations [3]. Various set of machine learning tools have been used for effective filtering and regression process to increase the critical activities in competitive data modelling.

Apache spark ecosystem: The ecosystems has been used in the Apache Spark manage the simultaneous operations of different machines. Open sources and distributions have been used for effective documentation in better experience to assist better contribution.

Benefits of Apache spark

- **Speed:** The Apache Spark has several benefits in the context of data classification process. It can read large data sets 10 or 100 times faster than other data classification process such as disk, memory and card reader. In memory engine has been used in effective data processing to execute multiple operations.

- **Real time stream:** Integration of advanced framework and performance management is necessary for managing the fundamental challenges in different size of small batches. Apache has the highest quality in different language management and accessibility for natural community development.
- **Support workload:** It has seen Apache Spark can control several workloads using the real time data analytics and machine learning. Queries interaction strategies have been used in the Apache Spark [2]. Data engineers can opportunities in algorithm development for improving the functional effectiveness in complex file processing and management.
- **Usability:** Dynamics of multiple programming languages increased the usability in effective programming process. It has found that programming language can manage the application of dynamic culture and streaming of graphical processing. Spark machine learning manage all kind of tasks like clustering, regression, collaborative and classification to deliver active data management.

Design

Text classification model development needs dataset that has been collected from kaggle. It has been ensured that dataset has open source license so that it can be used in the development purpose. Spam text message classification data has been identified with CC0: public domain license [1]. The classification practices have been used to evaluate the toxic text from a content based on effective classification framework. In this process whole information has been imported from large data sets to run ML based observation process in the jupyter platform. Numpy and pandas algorithm have been used for processing of data and apply the pyspark to manage the dependencies. Here, unwanted expressions have been removed from data sets to address spam contents of any message and texts. String level has been converted into the integer format for managing the advanced featuring. Different kind of algorithm such as logistics classification, k mean clustering, native bayes and forest classification can manage the accuracy of spam detection. The below figure has shown that around 87% of text represents the normal text and 13% are spam in whole information.

Category		Message
ham spam	87% 13%	5157 unique values
ham		Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got a...
ham		Ok lar... Joking wif u oni...
spam		Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entr...
ham		U dun say so early hor... U c already then say...
ham		Nah I don't think he goes to usf he

Figure 1: text classification dataset

(Source: [1])

It is needed to integrated pyspark in the jupyter notebook for text classification and spam detection machine learning model development. Tokenizer process also has been used for converting sentences to list of words. The classification of text depends on the advanced machine learning for categorising the better prediction. After the importing of data sets, special characters and features have been utilised to develop multiple classification process. Therefore, pyspark algorithm is used for streamline of different types of data for transformation development to evaluate the numerical features. Matching different characters is necessary for defining the fundamental issues and challenges. Further, evaluation and seaborn for better development can provide consistency in creation of classified results.

```
In [45]: import pyspark
        from pyspark.sql import SparkSession
        from pyspark.sql.functions import *
        from pyspark.ml.feature import *
        from pyspark.ml.linalg import Vector
        from pyspark.ml import Pipeline
        from pyspark.ml.evaluation import MulticlassClassificationEvaluator

In [67]: import tensorflow as tf
        from tensorflow.keras.preprocessing.text import Tokenizer
        from tensorflow.keras.preprocessing.sequence import pad_sequences
        from sklearn.preprocessing import OneHotEncoder
        import numpy as np
```

Figure 2: Python library import for text classification

Architecture

In this spam detection process, python machine learning has been used. Jupyter notebook has been used as machine learning model development tool. Apache spark has been used for the text classification so that email spam can be detected and categorized. UI components and functions have classified the interface of Jupyter platform to manage the functional issues. Here, the notebook elements can manage by the interface design and selection of specific options in the classification process. The interface can change during the editing phase to manage the simplicity of Notebook and manage the dashboard. In the code section of the below figure, any developer can write python codes for Apache Spark to develop textclassification application. After successful writing, run can be used for testing the contentsfor creation of plot and images. Here, kemal option can be used for edited codes for removingthe issues in the coding to clear whole output and get re-explanation [6]. Markdown cell can be used for effective documentation of code and information. Therefore, several options are existed in the Jupyter interface for calculation, heading, equation, image, and integers for the file management.

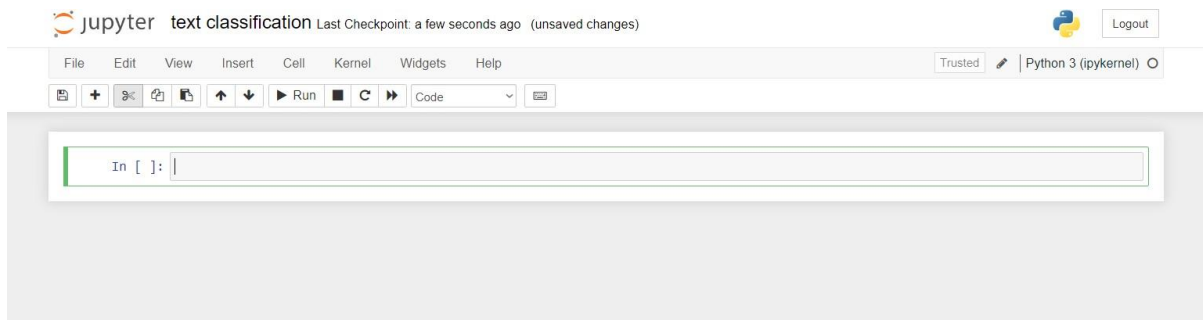


Figure 3: Jupyter notebook interface

In the jupyter notebook interface dataset has been imported using pandas framework. This dataset has been divided into train dataset and text dataset. Here, the Jupyter notebook also needed for installing the pyspark for text analyse. After that, it is needed identify proper machine learning model for the text classification. In this situation, logistic regression model has been identified as this is most suitable model for spam detection. Logistic regression model has been used through accuracy and confusion matrix process. Accuracy is effective for the identification accuracy of the developed model and confusion matrix helps in the visual representation of actual vs. predicted values [5]. Accuracy will help to understand classified text so that spam from the email can be detected. Performance evaluation has been completed using the splitting techniques structure creation. After splitting CSV file, classification function model has been used for the cell mode management to run the file classification.

```
In [2]: pip install pyspark

Collecting pyspark
  Downloading pyspark-3.5.0.tar.gz (316.9 MB)
Collecting py4j==0.10.9.7
  Downloading py4j-0.10.9.7-py2.py3-none-any.whl (200 kB)
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py): started
  Building wheel for pyspark (setup.py): still running...
  Building wheel for pyspark (setup.py): finished with status 'done'
  Created wheel for pyspark: filename=pyspark-3.5.0-py2.py3-none-any.whl size=317425366 sha256=de9e3f3dbd4de1ecaed7ce905e0bb621ad2c744757af69eec142264cef243115
  Stored in directory: c:\users\91882\appdata\local\pip\cache\wheels\57\bd\14\ce9e21f2649298678d011fb8f71ed38ee70b42b94fef0be142
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully installed py4j-0.10.9.7 pyspark-3.5.0
Note: you may need to restart the kernel to use updated packages.
```

Figure 4: Apache PySpark installation for text classification

The balancing is necessary for managing the sustainability of text categorisation process to manage the efficiency of data detection. It has found that text classification can manage the efficiency in real time textual data management. Further, nature of text can be fined using the

categorisation of data into ham and spam segments. Toxicity nature of contents can be detected using the intent detection and binary classification process [7]. Here, the below figure has shown that messages are processed into small tokens to extract features of messages. Decision making and learning mechanism have to be used for addressing category output management. Based on the below image, it can state that around 50% of messages are normal and 50% are spam types. It can reduce the complexity trends analysis for managing the classification process.

```
In [64]: textdata=textdata.filter(textdata['_c0']!='Category')
textdata=textdata.withColumnRenamed('_c0','Category').withColumnRenamed('_c1','Message')
textdata.show()
```

```
+-----+-----+
|Category|      Message|
+-----+-----+
|      ham|Go until jurong p...|
|      ham|Ok lar... Joking ...|
|     spam|Free entry in 2 a...|
|      ham|U dun say so earl...|
|      ham|Nah I don't think...|
|     spam|FreeMsg Hey there...|
|      ham|Even my brother i...|
|      ham|As per your reque...|
|     spam|WINNER!! As a val...|
|     spam|Had your mobile 1...|
|      ham|I'm gonna be home...|
|     spam|SIX chances to wi...|
|     spam|URGENT! You have ...|
|      ham|I've been searchi...|
|      ham|I HAVE A DATE ON ...|
|     spam|XXXMobileMovieClu...|
|      ham|Oh k...i'm watchi...|
|      ham|Eh u remember how...|
|      ham|Fine if thats th...|
|     spam|England v Macedon...|
+-----+-----+
only showing top 20 rows
```

Figure 5: Balance of texts in classification process

Keras model has been used in this classification process to manage the sequence of layers. There are two method of vectorisation layer which are used in the creation of model for evaluation of effective result. Here, the sentence encoders such as dense, dropout are used for addressing the optimised results of information detection. Variables should be analysed for development of optimised distribution for assembling the performance sequence in defining the Apache Spark in accuracy management.

[illegible]

```
In [76]: trainingSize=int(len(messages)*0.7)
trainingSeq=padded_sequences[:trainingSize]
trainingLabel=np.array(messageLabels[:trainingSize])
testingSeq=padded_sequences[trainingSize:]
testingLabel=np.array(messageLabels[trainingSize:])
```

The below image shows the model has used testing of several sentences for managing the validation. The training validation loss and accuracy have been analysed using the keras model which mentioned in the image [6]. Here, 5 data sets have been applied for managing the evaluation framework to address validation of data and splitting process using the keras model. The two validation testing has been organised during the testing process where the accuracy percentage are 92.86% and 97.07% that means acceptability and accuracy of the model is suitable for text classification process.

```

In [79]: tfmodel = tf.keras.Sequential([\
    tf.keras.layers.Embedding(input_dim=len(word_index) + 1,output_dim=16,input_length=tra\
    tf.keras.layers.Flatten()),\
    tf.keras.layers.Dense(32, activation='relu'),\
    tf.keras.layers.Dense(1, activation='sigmoid')
])

In [80]: tfmodel.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

In [81]: tfmodel.fit(trainingSeq,trainingLabel,epochs=5,validation_data=(testingSeq,testingLabel))

Epoch 1/5
122/122 [=====] - 2s 9ms/step - loss: 0.2786 - accuracy: 0.8977 - val_loss: 0.1303 - val_acc
uracy: 0.9713
Epoch 2/5
122/122 [=====] - 1s 5ms/step - loss: 0.0756 - accuracy: 0.9828 - val_loss: 0.0568 - val_acc
uracy: 0.9827
Epoch 3/5
122/122 [=====] - 1s 5ms/step - loss: 0.0331 - accuracy: 0.9915 - val_loss: 0.0434 - val_acc
uracy: 0.9868
Epoch 4/5
122/122 [=====] - 1s 5ms/step - loss: 0.0177 - accuracy: 0.9949 - val_loss: 0.0415 - val_acc
uracy: 0.9886
Epoch 5/5
122/122 [=====] - 1s 5ms/step - loss: 0.0109 - accuracy: 0.9974 - val_loss: 0.0366 - val_acc
uracy: 0.9898

Out[81]: <keras.src.callbacks.History at 0x7a4c6cb0dc90>

In [83]: loss,accuracy=tfmodel.evaluate(testingSeq,testingLabel)
print("accuracy : {}".format(accuracy*100))

53/53 [=====] - 0s 1ms/step - loss: 0.0366 - accuracy: 0.9898
accuracy : 98.98386001586914

```

Figure 7: Accuracy analysis using keras model

Bibliography

- [1] T. AI, "Spam text message classification," Kaggle, <https://www.kaggle.com/datasets/team-ai/spam-text-message-classification/data> (accessed Nov. 20, 2023).
- [2] Nazari, Elham, Mohammad Hasan Shahriari, and Hamed Tabesh. "BigData analysis in healthcare: apache hadoop, apache spark and apache flink." *Frontiers in Health Informatics* vol. 8. ed. 1, pp. 14, 2019.
- [3] Haggag, Mohamed, Mohsen M. Tantawy, and Magdy MS El-Soudani. "Implementing a deep learning model for intrusion detection on apache spark platform." *IEEE Access* pp. 163660-163672, Aug 2020.
- [4] Pintye, István, et al. "Big data and machine learning framework for clouds and its usage for text classification." *Concurrency and Computation: Practice and Experience* vol. 33. Ed.19, pp 6164, 2021.
- [5] Shetty, Sujala D. "Sentiment analysis, tweet analysis and visualization on big data using Apache Spark and Hadoop." IOP Conference Series: Materials Science and Engineering. Vol. 1099. No. 1. IOP Publishing, 2021.
- [6] Salihoun, Mohammed. "State of art of data mining and learning analytics tools in higher education." *International Journal of Emerging Technologies in Learning (iJET)* vol 15, ed. 21, pp. 58-76, 2020.
- [7] Wang, Shuang, Jian Luo, and Liangfu Luo. "Large-scale text multiclass classification using spark ML packages." *Journal of Physics: Conference Series*. Vol. 2171. No. 1. IOP Publishing, 2022.

My GitHub Repository

<https://github.com/YFA23SCM78K/big-data-final-project/issues/1>

Conclusion

This project demonstrates an end-to-end implementation of a spam detection model using Apache Spark's machine learning capabilities. By harnessing Spark's ability to rapidly process large amounts of text data in parallel, the developed logistic regression model can effectively identify spam messages with high accuracy. The project provides a valuable template for applying distributed machine learning techniques to text classification tasks for spam filtering and other real-world applications.

