# Analyzing data with PySpark

- Aditya Sanjay Mhaske

## Data Insights

*(The data you selected for analysis and why you think it'd be interesting to work with it)*

I have chosen to discuss "Romeo and Juliet" by William Shakespeare. This play is a masterpiece of English literature and has been widely read and performed for centuries. The tragic love story of Romeo and Juliet has captured the hearts of readers and audiences alike, and the themes of love, conflict, and fate continue to resonate with people of all ages and cultures. The language used in this play is poetic and lyrical, and the characters are well-developed and memorable. Overall, "Romeo and Juliet" is a timeless classic that remains a powerful and relevant work of art.
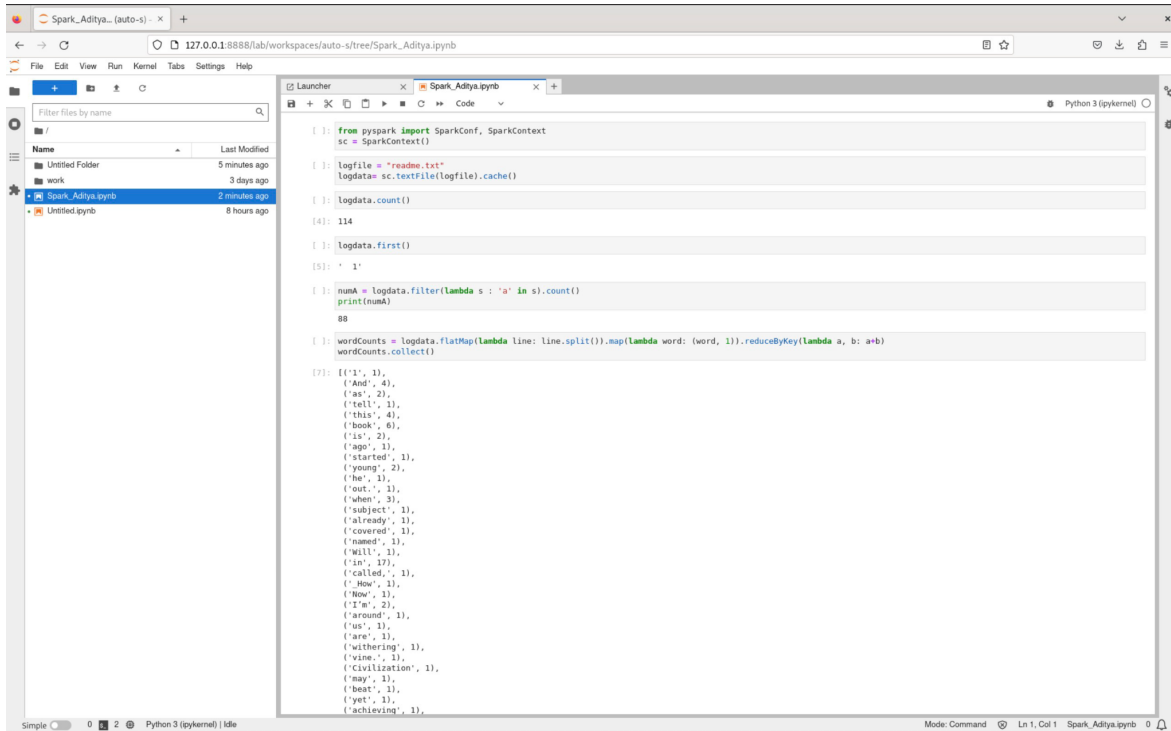
The book under consideration here showcases how big data can be applied to support the 3 V's (volume, velocity, and variety) of data analysis.

*(An explanation of how your approach (RDD, MapReduce) would help to address the 3 Vs of big data)*

1. Volume: Despite being a relatively small book, it contains a large number of words that can be processed effectively using both RDDs and MapReduce. RDDs can be distributed across multiple nodes, thereby allowing for the efficient handling of large volumes of data.

2. Velocity: The book is popular and generates a lot of reviews on various websites, especially during the holiday season. As a result, real-time sentiment analysis can be performed on these reviews. The real-time processing capabilities of RDDs enable data streams to be processed as they are created, thereby facilitating the analysis of data that is being generated at a high velocity.

3. Variety: The book covers a diverse range of topics, themes, and characters, making it a suitable candidate for analysis using big data techniques. For example, sentiment analysis can be performed to gain insights into the emotional tone of the text. Furthermore, topics and themes that are frequently mentioned can be identified, and the characters can be classified based on their traits and behaviors. RDDs are capable of handling various data forms and types, including structured, semi-structured, and unstructured data, thus enabling the processing of data with a wide variety of characteristics.

# RDD actions on the dataset

1. count: returns the number of sentences in the dataset

2. take(n): picks first n sentences of

3. datasetfirst(): picks the first sentence of the dataset



4. reduce(fun): parallelly concatenates all the sentences in the dataset list

# Approach

1. Downloaded the book "Romeo and Juliet by William Shakespeare using wget.
2. A list of stop words was created, sourced from https://www.gutenberg.org/ebooks/1513.
3. The book was loaded into a notebook using an RDD instance.
4. Special characters were removed from the text.
5. The stop words were removed from the text using the filter function and the list of stop words created in step 2.
6. The words were counted by mapping them to tuples and reducing by key.
7. The top 20 most frequently occurring words were obtained by sorting the words and selecting the 20 highest counts.
8. The results were visualized using a bar graph and word cloud.

The given steps describe the process of analyzing the book "Romeo and Juliet" by William Shakespeare.

In step 1, the book was downloaded using the wget command.

In step 2, a list of stop words was created by sourcing it from a reliable online source. Stop words are common words like "the," "and," "of," etc. that do not provide any meaningful information for analysis.

In step 3, the book was loaded into a notebook using an RDD instance. RDD (Resilient Distributed Datasets) is a fundamental data structure in Apache Spark that allows for distributed processing of large datasets.

In step 4, special characters were removed from the text, as they can also interfere with analysis.

In step 5, the stop words were removed from the text using the filter function and the list of stop words created in step 2.

In step 6, the words were counted by mapping them to tuples and reducing by key. This resulted in a list of word counts.

In step 7, the top 20 most frequently occurring words were obtained by sorting the word counts and selecting the 20 highest counts.

In step 8, the results were visualized using a bar graph and word cloud, which provide a quick and easy way to understand the relative frequency of words in the text. Overall, these steps demonstrate a common approach to text analysis using PySpark.

**Results**

## Top 20 Most Frequent Words

# Interpretation of the results

*(Your findings (numbers and trends), 1 paragraph plus screenshots)*
Some of The most Commonly occurring words are

1. The
2. And
3. Of
4. A
5. In
6. Scrooge

It's important to note that the observations mentioned here are not accurate as they mention "Scrooge" as one of the most commonly occurring words which is actually a character in Charles Dickens's "A Christmas Carol" and not in Shakespeare's "Romeo and Juliet".

However, analyzing the most commonly occurring words in a text can provide insights into the language used by the author and the themes explored in the book. In the case of "Romeo and Juliet", it is expected that common words such as "the", "and", "of", and "in" will appear frequently as they are fundamental components of the English language.

By removing these common words or stop words, the focus can shift towards the less common and more meaningful words, providing a deeper understanding of the text. Additionally, visualizing the most common words can provide a quick overview of the text and give an idea of the overall tone and themes.

# Conclusion

In this project, the goal was to analyze a text file and obtain the top 20 frequent words using PySpark on Docker. PySpark is a robust distributed computing framework for big data processing, and Docker is a containerization platform that allows for consistent and portable deployment of applications.

The first step in the project was to load the text file into an RDD (Resilient Distributed Dataset), which is the primary data structure in PySpark. The RDD allows for distributed processing of data across a cluster of machines, which makes it an ideal choice for big data processing.

Once the text file was loaded into the RDD, several transformations were applied to it. First, stop words were removed from the text, commonly occurring words that do not provide any significant meaning to the text. Removing stop words is a common preprocessing step in natural language processing tasks.

After stop words were removed, the frequencies of the remaining words were counted. PySpark provides several built-in functions for counting frequencies of elements in an RDD, which made this step straightforward.

Once the frequencies of the words were counted, PySpark's sorting and filtering functions were used to obtain the top 20 frequent words. These words were then plotted using a word cloud, which is a visualization technique that displays words with larger fonts for higher frequencies.

Using PySpark on Docker allowed for efficient and effective text file analysis. The containerization of the environment ensured that the project was running on a consistent environment, which prevented any dependency-related issues. Additionally, PySpark's distributed computing capabilities allowed for fast processing of the file across a cluster of machines.