# Tarefa Dois

Objetivo: Construir um *job* Spark por meio de um container Docker.

- Nesta atividade faremos uso da imagem **jupyter/all-spark-notebook** (https://registry.hub.docker.com/r/jupyter/all-spark-notebook) para criar um container e utilizar o recurso de shell oferecido pelo Spark. Os passos a executar são:

**1 - Realizar o *pull* da imagem jupyter/all-spark-notebook**

Comando utilizado:

### *sudo docker pull jupyter/all-spark-notebook*



Atenção: O tamanho total da imagem é 5.8 GB. Se você não tiver esse espaço disponível, recomendamos utilizar o **Google Colab** para codificar o exercício.

**2 - Criar um container a partir da imagem**

Comando utilizado:

### *docker run -it -p 8888:8888 jupyter/all-spark-notebook*

A url disponibilizada foi:

[http://127.0.0.1:8888/lab?token=2be3d8627f8fd5752dad5be97d9c3c718f5e7e1913786472](http://127.0.0.1:8888/lab?token=2be3d8627f8fd5752dad5be97d9c3c718f5e7e1913786472)



**3 - Em outro terminal, execute o comando `pyspark` no seu container. Pesquise sobre o comando *docker exec* para realizar esta ação. Utilize as flags *-i* e *-t* no comando.**

Comando utilizado para buscar o id do container:

**docker ps**



Comando utilizado para executar o pyspark no container:

**sudo docker exec -it 574075e1dec3 pyspark**

Dica: Você pode obter arquivos da Internet por meio do comando *wget* no seu container.

- Usando o *Spark Shell,* apresente a sequência de comandos Spark necessários para contar a quantidade de ocorrências de cada palavra contida no arquivo *README.md* de seu repositório *git.*

Passos:

Baixei o arquivo README pelo caminho temporário do raw usando o console do      JupyterLab:



Converti o arquivos para README.md, em seguida voltei para o terminal para verificar o caminho:

Utilizei os seguintes comandos para carregar o arquivo no pyspark e printar na tela:

**>>> from pyspark.sql import SparkSession**

**>>> spark = SparkSession.builder.appName("Cont").getOrCreate()**

**>>> readme_df = spark.read.text("/home/jovyan/work/README.md")**

**>>> readme_df.show(10, False)**



Os seguintes comandos foram utilizados para realizar a contagem e mostrar as 10 primeiras linhas:

**>>> from pyspark.sql.functions import explode, split, lower**

**>>> words_df = readme_df.select(explode(split(lower(readme_df.value), "\s+")).alias("word"))**

>>> *from pyspark.sql.functions import count*

*u* >>> *word_count=words_df.groupBy("word").agg(count("word").alias("count"))*

>>> *word_count = word_count.orderBy("count", ascending=False)*

>>> *word_count.show(10)*

```
                                                      lins@lins-Lenovo-G460: ~                    ×              lins@lins-Lenovo-G460: ~                    ×  ⊞  ▾
+------------------------------------------------------+
only showing top 10 rows

>>> from pyspark.sql.functions import explode, split, lower
>>> words_df = readme_df.select(explode(split(lower(readme_df.value), "\s+")).alias("word"))
>>> from pyspark.sql.functions import count
>>> word_count = words_df.groupBy("word").agg(count("word").alias("count"))
>>> word_count = word_count.orderBy("count", ascending=False)
>>> word_count.show(10)
+-------+-----+
|   word|count|
+-------+-----+
|      -|   67|
|       |   30|
|     de|   19|
|      e|   17|
|sprints|   11|
|     ##|    8|
|    aws|    7|
|     em|    6|
|     do|    6|
|    meu|    5|
+-------+-----+
only showing top 10 rows

>>>
```

Por último ultilisei o ".show()" semparametro para listar todas as palavras e contagens:

```
>>> word_count.show()
+--------+-----+
|    word|count|
+--------+-----+
|       -|   67|
|        |   30|
|      de|   19|
|       e|   17|
| sprints|   11|
|      ##|    8|
|     aws|    7|
|      em|    6|
|      do|    6|
|     meu|    5|
|      da|    5|
|    para|    4|
|       o|    4|
|    java|    4|
|       a|    4|
|     💻📚|    3|
|     com|    3|
|formação|    3|
|    data|    3|
|  trilha|    3|
+--------+-----+
only showing top 20 rows

>>>
```