

Nama : Syava Aprilia Puspitasari

NIM : 2241760129 / 18

Interaksi dengan Spark di Lingkungan Windows Menggunakan Docker

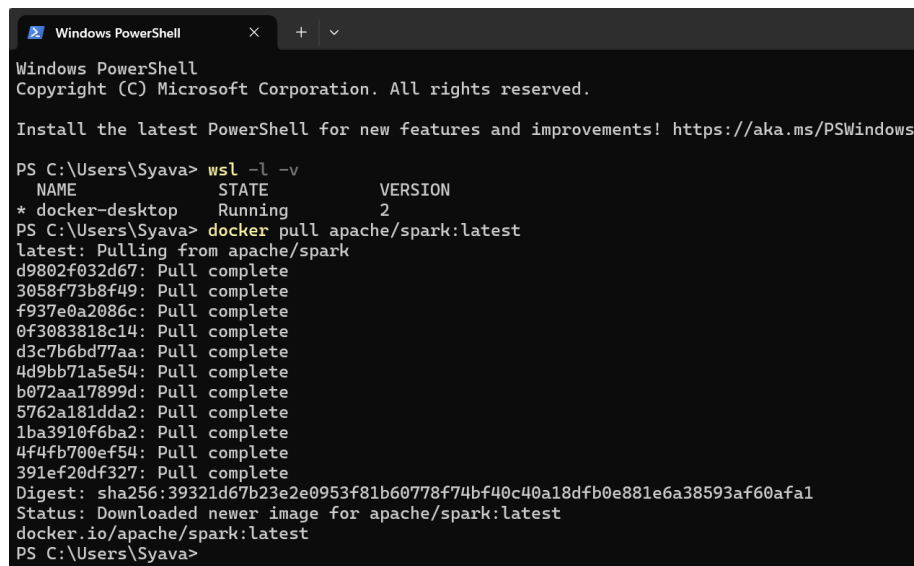
Dalam praktikum ini kita akan menjalankan Apache Spark di Windows menggunakan Docker dan mencoba membuat job sederhana dengan berbagai macam alternatif cara.

Prasyarat

1. Windows 10/11 (64-bit) dengan versi Pro, Enterprise, atau Education
2. Docker Desktop untuk Windows diinstal dan berjalan
3. WSL 2 (Windows Subsystem for Linux versi 2) diaktifkan

Langkah-langkah

1. Pull Image Spark Resmi



```
Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Install the latest PowerShell for new features and improvements! https://aka.ms/PSWindows

PS C:\Users\Syava> wsl -l -v
  NAME                STATE      VERSION
* docker-desktop      Running    2
PS C:\Users\Syava> docker pull apache/spark:latest
latest: Pulling from apache/spark
d9802f032d67: Pull complete
3058f73b8f49: Pull complete
f937e0a2086c: Pull complete
0f3083818c14: Pull complete
d3c7b6bd77aa: Pull complete
4d9bb71a5e54: Pull complete
b072aa17899d: Pull complete
5762a181dda2: Pull complete
1ba3910f6ba2: Pull complete
4f4fb700ef54: Pull complete
391ef20df327: Pull complete
Digest: sha256:39321d67b23e2e0953f81b60778f74bf40c40a18dfb0e881e6a38593af60afa1
Status: Downloaded newer image for apache/spark:latest
docker.io/apache/spark:latest
PS C:\Users\Syava>
```

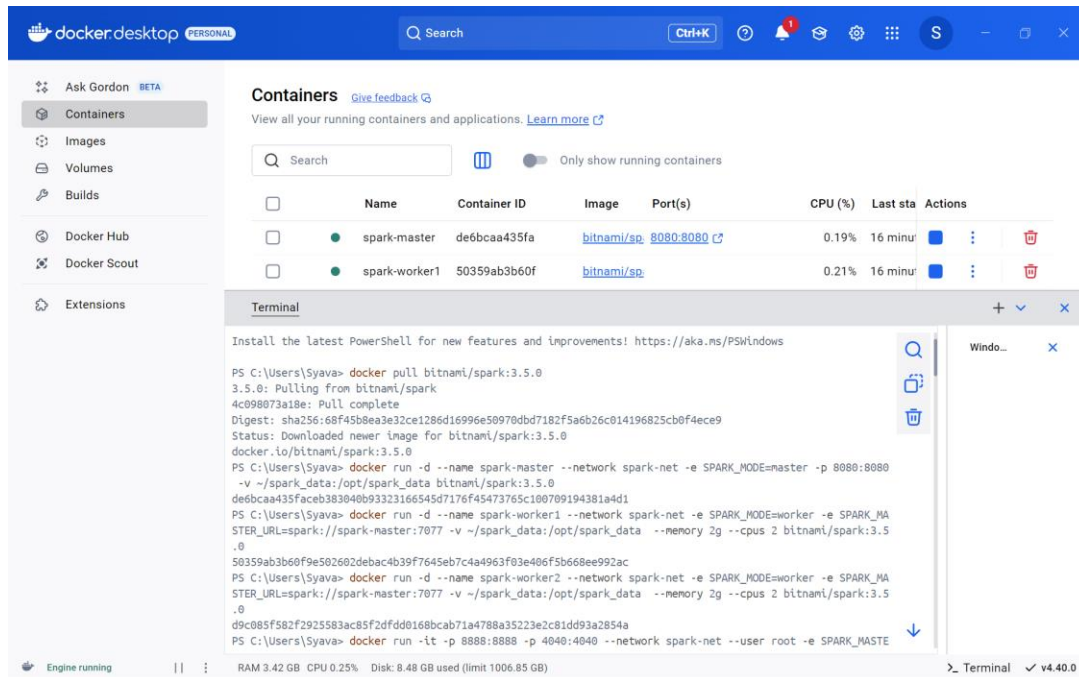
Menjalankan Master

```
docker run -d --name spark-master --network spark-net -e SPARK_MODE=master -p 8080:8080 -v ~/spark_data:/opt/spark_data bitnami/spark:3.5.0
```

Menjalankan Worker

```
docker run -d --name spark-worker1 --network spark-net -e SPARK_MODE=worker -e SPARK_MASTER_URL=spark://spark-master:7077 -v ~/spark_data:/opt/spark_data --memory 2g --cpus 2 bitnami/spark:3.5.0
```

```
docker run -d --name spark-worker2 --network spark-net -e SPARK_MODE=worker -e SPARK_MASTER_URL=spark://spark-master:7077 -v ~/spark_data:/opt/spark_data --memory 2g --cpus 2 bitnami/spark:3.5.0
```



Menjalankan Spark Shell lalu Menggunakan Jupyter Notebook dengan Spark

```
docker run -it -p 8888:8888 -p 4040:4040 --network spark-net --user root -e SPARK_MASTER=spark://spark-master:7077 -e GRANT_SUDO=yes -v ~/spark_data:/opt/spark_data jupyter/all-spark-notebook
```

docker desktop

PERSONAL

Ctrl+K

Ask Gordon BETA

Containers

Images

Volumes

Builds

Docker Hub

Docker Scout

Extensions

Containers [Give feedback](#)

View all your running containers and applications. [Learn more](#)

Only show running containers

<input type="checkbox"/>	Name	Container ID	Image	Port(s)	CPU (%)	Last sta	Actions
<input type="checkbox"/>	spark-master	de6bcaa435fa	bitnami/spark	8080:8080	0.19%	17 min	
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark		0.18%	16 min	

Terminal

```

STER_URL=spark://spark-master:7077 -v ~/spark_data:/opt/spark_data --memory 2g --cpus 2 bitnami/spark:3.5
.d
d9c88f592f2925583ac85f2dfdd0168bcab71a4788a35223e2c81dd93a2854a
PS C:\Users\Syava> docker run -it -p 8888:8888 -p 4040:4040 --network spark-net --user root -e SPARK_MASTER=spark://spark-master:7077 -e GRANT_SUDO=yes -v ~/spark_data:/opt/spark_data jupyter/all-spark-notebook
Unable to find image 'jupyter/all-spark-notebook:latest' locally
latest: Pulling from jupyter/all-spark-notebook
aeece8493d397: Pulling fs layer
fd92c719666c: Pulling fs layer
088f11eb1e74: Pulling fs layer
4f4fb780ef54: Waiting
ef8373d60b0b: Waiting
77e45ee945dc: Waiting
a30f89a0af6c: Waiting
dc42adc7eb73: Waiting
abaa8376a50: Waiting
aa099bb9e49a: Waiting
822c4bcfc6a6: Waiting
d25166dc7b: Waiting
964fc3e4ff9f: Waiting

```

Windo...

Engine running

||

:

RAM 3.41 GB

CPU 0.25%

Disk: 8.48 GB used (limit 1006.85 GB)

Terminal

v4.40

Ask Gordon BETA

Containers

Images

Volumes

Builds

Docker Hub

Docker Scout

Extensions

Containers Give feedback

View all your running containers and applications. [Learn more](#)

Search

Only show running containers

	Name	Container ID	Image	Port(s)	CPU (%)	Last sta	Actions
<input type="checkbox"/>	spark-master	de6bcaa435fa	bitnami/spark	8080:8080	0.2%	17 min	
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark		0.2%	17 min	

Terminal

[I 2025-05-06 07:58:14.394 ServerApp] jupyter_lsp | extension was successfully loaded.
[I 2025-05-06 07:58:14.403 ServerApp] jupyter_server_mathjax | extension was successfully loaded.
[I 2025-05-06 07:58:14.406 ServerApp] jupyter_server_terminals | extension was successfully loaded.
[I 2025-05-06 07:58:14.423 LabApp] JupyterLab extension loaded from /opt/conda/lib/python3.11/site-packages/jupyterlab
[I 2025-05-06 07:58:14.423 LabApp] JupyterLab application directory is /opt/conda/share/jupyter/lab
[I 2025-05-06 07:58:14.424 LabApp] Extension Manager is 'pypi'.
[I 2025-05-06 07:58:14.435 ServerApp] JupyterLab | extension was successfully loaded.
[I 2025-05-06 07:58:14.448 ServerApp] jupyterlab_git | extension was successfully loaded.

Read the migration plan to Notebook 7 to learn about the new features and the actions to take if you are using extensions.

https://jupyter-notebook.readthedocs.io/en/latest/migrate_to_notebook7.html

Engine running RAM 3.42 GB CPU 0.67% Disk: 8.48 GB used (limit 1006.85 GB) Terminal v4.40.0

Ask Gordon BETA

Containers

Images

Volumes

Builds

Docker Hub

Docker Scout

Extensions

Containers Give feedback

View all your running containers and applications. [Learn more](#)

Search

Only show running containers

	Name	Container ID	Image	Port(s)	CPU (%)	Last sta	Actions
<input type="checkbox"/>	spark-master	de6bcaa435fa	bitnami/spark	8080:8080	0.2%	18 min	
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark		0.24%	17 min	

Terminal

[I 2025-05-06 07:58:16.411 ServerApp] Skipped non-installed server(s): bash-language-server, dockerfile-language-server-nodejs, javascript-typescript-languageserver, jedi-language-server, julia-language-server, pyright, python-language-server, python-lsp-server, r-languageserver, sql-language-server, texlab, typescript-language-server, unified-language-server, vscode-css-languageserver-bin, vscode-html-languageserver-bin, vscode-json-languageserver-bin, yamll-language-server
0.00s - Debugger warning: It seems that frozen modules are being used, which may make the debugger miss breakpoints. Please pass --x-frozen-modules=off to python to disable frozen modules.
0.00s - Note: Debugging will proceed. Set PYDEVD_DISABLE_FILE_VALIDATION=1 to disable this validation.
HTTPServerRequest(protocol='http', host='127.0.0.1:8888', method='GET', url='/lab/api/settings?1746518353140', version='HTTP/1.1', remote_ip='172.18.0.1')
Traceback (most recent call last):
 File "/opt/conda/lib/python3.11/site-packages/tornado/web.py", line 1786, in _execute
 result = await result
 ^^^^^^^^^^^^^
tornado.iostream.StreamClosedError: Stream is closed
[W 2025-05-06 07:59:17.395 LabApp] The extension "nbline-jupyterlab" is outdated.
[I 2025-05-06 07:59:17.395 LabApp] Build is up to date

Read the migration plan to Notebook 7 to learn about the new features and the actions to take if you are using extensions.

https://jupyter-notebook.readthedocs.io/en/latest/migrate_to_notebook7.html

Engine running RAM 3.43 GB CPU 0.00% Disk: 8.48 GB used (limit 1006.85 GB) Terminal v4.40.0

Ask Gordon BETA

Containers

Images

Volumes

Builds

Docker Hub

Docker Scout

Extensions

Images Give feedback

View and manage your local and Docker Hub images. [Learn more](#)

Local

Docker Hub repositories

4.49 GB / 0 Bytes in use 2 images Last refresh: 57 minutes ago

Search

Only show running containers

	Name	Tag	Image ID	Created	Size	Actions
<input type="checkbox"/>	bitnami/spark	3.5.0	07dc860e5425	1 year ago	1.74 GB	
<input type="checkbox"/>	jupyter/all-spark-notebook	latest	0add020fa68c	2 years ago	5.67 GB	

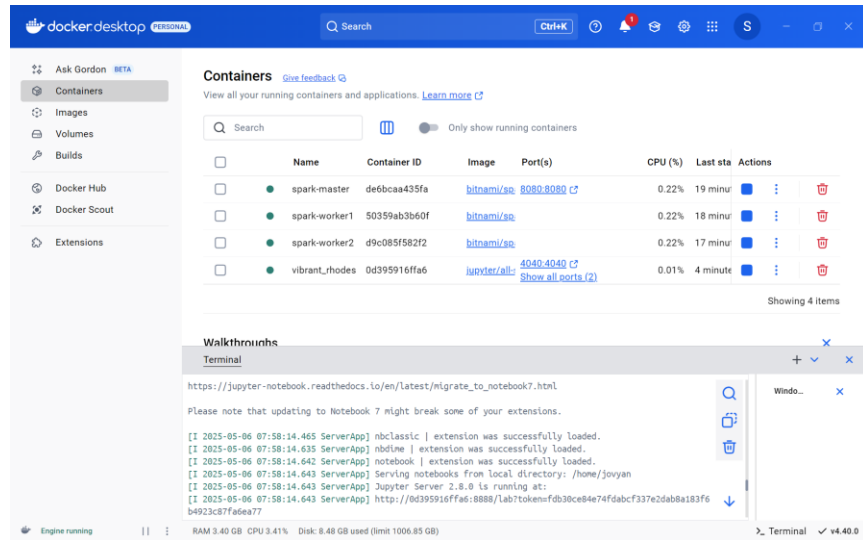
Terminal

HTTPServerRequest(protocol='http', host='127.0.0.1:8888', method='GET', url='/lab/api/settings?1746518353140', version='HTTP/1.1', remote_ip='172.18.0.1')
Traceback (most recent call last):
 File "/opt/conda/lib/python3.11/site-packages/tornado/web.py", line 1786, in _execute
 result = await result
 ^^^^^^^^^^^^^
tornado.iostream.StreamClosedError: Stream is closed
[W 2025-05-06 07:59:17.395 LabApp] The extension "nbline-jupyterlab" is outdated.
[I 2025-05-06 07:59:17.395 LabApp] Build is up to date

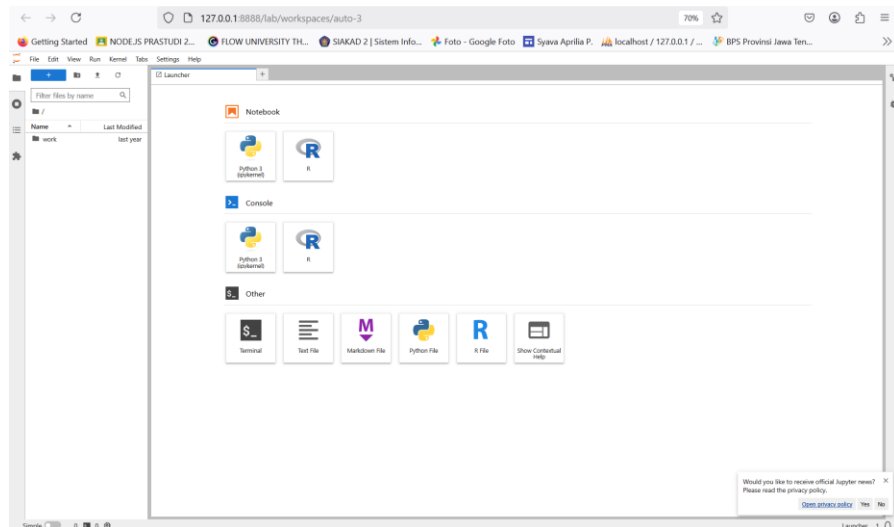
Read the migration plan to Notebook 7 to learn about the new features and the actions to take if you are using extensions.

https://jupyter-notebook.readthedocs.io/en/latest/migrate_to_notebook7.html

Engine running RAM 3.41 GB CPU 0.08% Disk: 8.48 GB used (limit 1006.85 GB) Terminal v4.40.0



Setelah itu, akses Jupyter Notebook di: <http://localhost:8888>



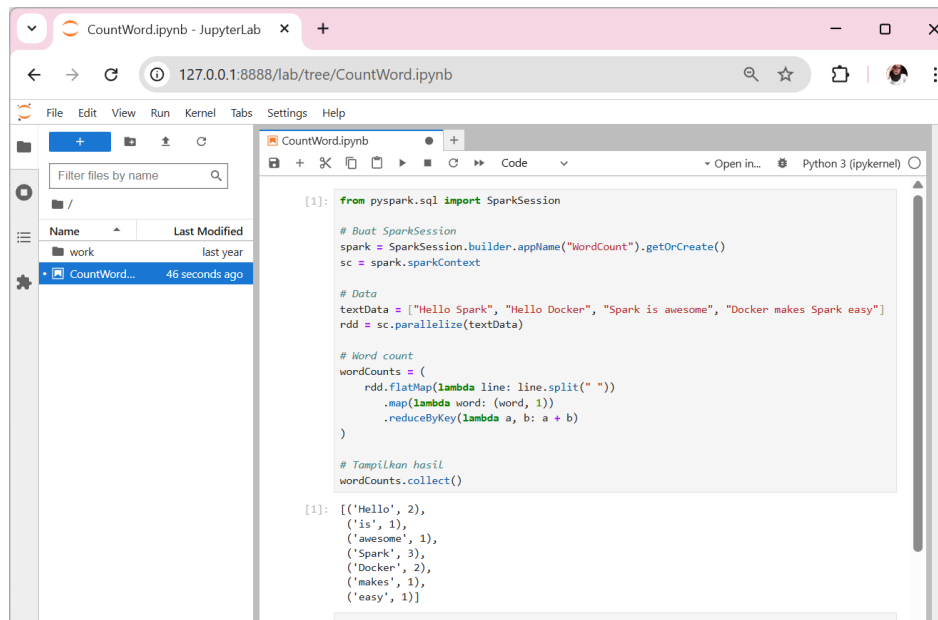
Contoh Program Word Count dengan Spark di Docker

Berikut adalah contoh program Word Count (menghitung kemunculan kata) menggunakan Apache Spark

yang bisa dijalankan di lingkungan Docker:

Cara 3: Menggunakan Jupyter Notebook

Jika Anda menggunakan Jupyter Notebook (seperti di container jupyter/all-spark-notebook):



```
[1]: from pyspark.sql import SparkSession

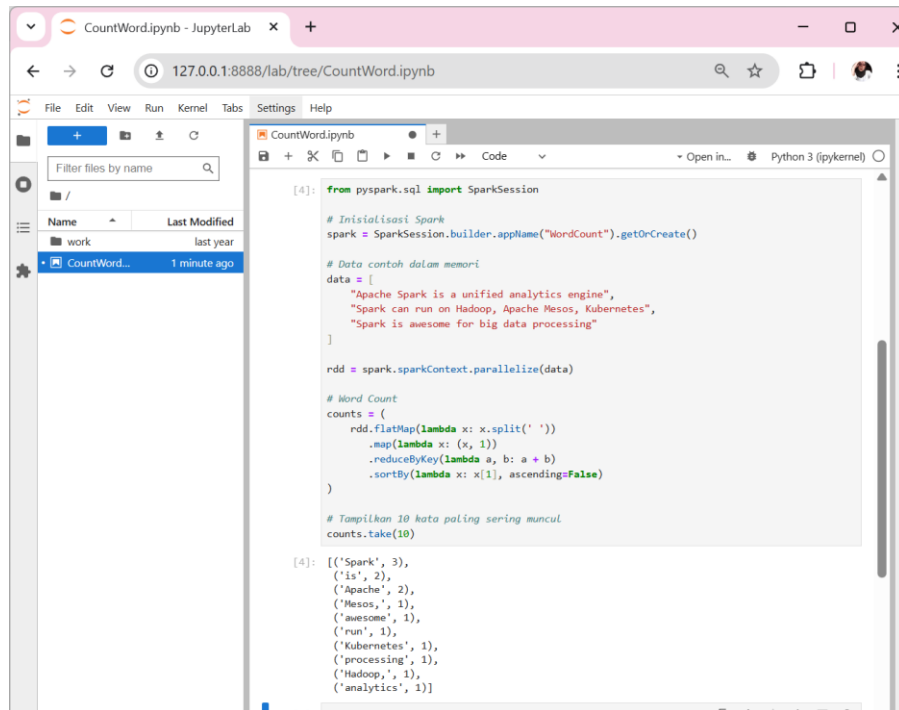
# Buat SparkSession
spark = SparkSession.builder.appName("WordCount").getOrCreate()
sc = spark.sparkContext

# Data
textData = ["Hello Spark", "Hello Docker", "Spark is awesome", "Docker makes Spark easy"]
rdd = sc.parallelize(textData)

# Word count
wordCounts = (
    rdd.flatMap(lambda line: line.split(" "))
    .map(lambda word: (word, 1))
    .reduceByKey(lambda a, b: a + b)
)

# Tampilkan hasil
wordCounts.collect()
```

[1]: [('Hello', 2), ('is', 1), ('awesome', 1), ('Spark', 3), ('Docker', 2), ('makes', 1), ('easy', 1)]



```
[4]: from pyspark.sql import SparkSession

# Inisialisasi Spark
spark = SparkSession.builder.appName("WordCount").getOrCreate()

# Data contoh dalam memori
data = [
    "Apache Spark is a unified analytics engine",
    "Spark can run on Hadoop, Apache Mesos, Kubernetes",
    "Spark is awesome for big data processing"
]

rdd = spark.sparkContext.parallelize(data)

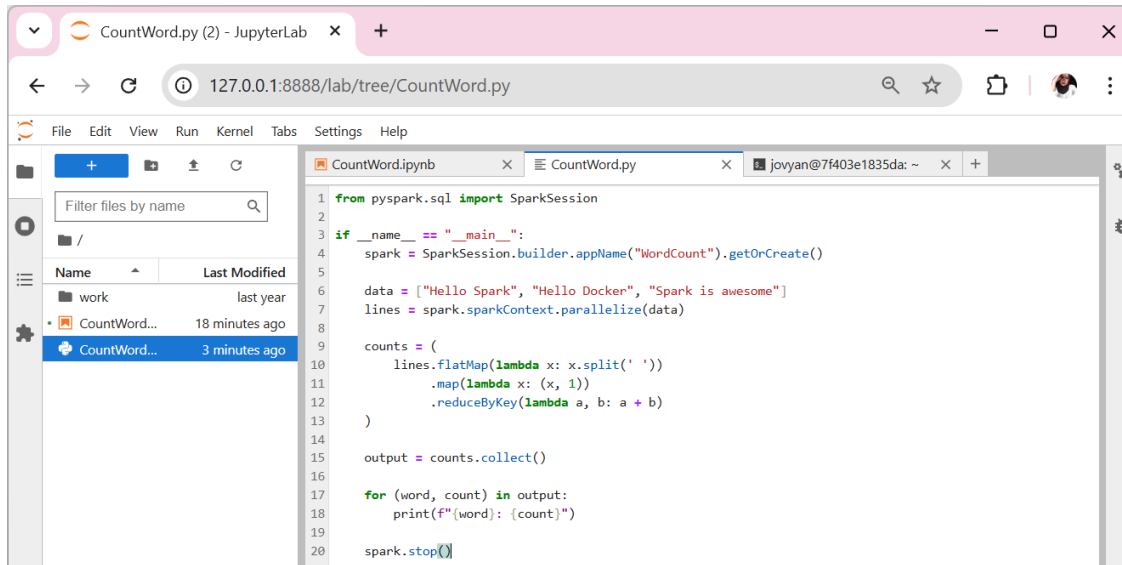
# Word Count
counts = (
    rdd.flatMap(lambda x: x.split(' '))
    .map(lambda x: (x, 1))
    .reduceByKey(lambda a, b: a + b)
    .sortBy(lambda x: x[1], ascending=False)
)

# Tampilkan 10 kata paling sering muncul
counts.take(10)
```

[4]: [('Spark', 3), ('is', 2), ('Apache', 2), ('Mesos', 1), ('awesome', 1), ('run', 1), ('Kubernetes', 1), ('processing', 1), ('Hadoop', 1), ('analytics', 1)]

Menjalankan Program sebagai Script

1. Buat file wordcount.py dengan isi berikut:



2. Jalankan jika sudah memastikan skrip di atas benar:

- Simpan file CountWord.py
- Jalankan:
spark-submit CountWord.py

