

LAPORAN JOBSHEET 7

BIG DATA



Oleh:

ALBANI RAJATA MALIK

2241760080/05

SIB 3E

PROGAM STUDI D-IV SISTEM INFORMASI BISNIS

JURUSAN TEKNOLOGI INFORMASI

POLITEKNIK NEGERI MALANG

```

PS C:\Users\USER> wsl -l -v
    NAME                STATE          VERSION
* docker-desktop      Running         2
PS C:\Users\USER> docker pull apache/spark:latest
latest: Pulling from apache/spark
b100ab72b3e9: Pull complete
c74c7c34259b: Pull complete
053e658edcf9: Pull complete
2fd40b3fb58e: Pull complete
215ed5a63843: Pull complete
ead058ffaa09: Pull complete
4f4fb700ef54: Pull complete
d4ebd433e6b9: Pull complete
3eaf8f3f646b: Pull complete
357990f05276: Pull complete
8f0df7ed423b: Pull complete
Digest: sha256:2ecf9cab4a1d0df1052731636f6029f6e80489fbd9ad0d24aaa31061f8aa8eea
Status: Downloaded newer image for apache/spark:latest
docker.io/apache/spark:latest
PS C:\Users\USER>

```

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

	Name	Container ID	Image	Port(s)	CPU (%)	Last sta	Actions
<input type="checkbox"/>	spark-master	de6bcaa435fa	bitnami/spark:3.5.0	8080:8080	0.19%	16 min	<div></div> <div></div> <div></div>
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark:3.5.0		0.21%	16 min	<div></div> <div></div> <div></div>

Terminal

+

-

x

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

PS C:\Users\USER> docker pull bitnami/spark:3.5.0
3.5.0: Pulling from bitnami/spark
4c980073a18e: Pull complete
Digest: sha256:68f45b8ea3e32ce1286d1699e50970dbd7182f5a6b26c814196825cb0f4ece9
Status: Downloaded newer image for bitnami/spark:3.5.0
docker.io/bitnami/spark:3.5.0
PS C:\Users\USER> 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
de6bcaa435faeb383040933216654507176f45473765c100709194381a4d1
PS C:\Users\USER> 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
50359ab3b60f9e502602debac4b39f7645eb7c4a4963f03e406f5b66ee992ac
PS C:\Users\USER> 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
d9c085f582f2925583ac85f2df0d0168bcab71a4788a35223e2c81d093a2854a
PS C:\Users\USER> docker run -it -p 8888:8888 -p 4040:4040 --network spark-net --user root -e SPARK_MASTER_URL=spark://spark-master:7077 -v ~/spark_data:/opt/spark_data jupyter/all-spark-notebook:latest
Unable to find image 'jupyter/all-spark-notebook:latest' locally
Latest: Pulling from jupyter/all-spark-notebook
aece4893e397: Pulling fs layer
fd92c71966dc: Pulling fs layer
088f1e1e74: Pulling fs layer
4f4fb700ef54: Waiting
ef83736000b0: Waiting
77e45ee945dc: Waiting
a30f89a9afcc: Waiting
dc42adc7eb73: Waiting
abaa3376a050: Waiting
aa099090e49a: Waiting
822c4c0cf6a6: Waiting
d25166dc0c7b: Waiting
964fc3e4ff9f: Waiting

Engine running

RAM 3.42 GB CPU 0.25% 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](#)

Only show running containers

	Name	Container ID	Image	Port(s)	CPU (%)	Last sta	Actions
<input type="checkbox"/>	spark-master	de6bcaa435fa	bitnami/spark:3.5.0	8080:8080	0.19%	17 min	<div></div> <div></div> <div></div>
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark:3.5.0		0.18%	16 min	<div></div> <div></div> <div></div>

Terminal

+

-

x

STER_URL=spark://spark-master:7077 -v ~/spark_data:/opt/spark_data --memory 2g --cpus 2 bitnami/spark:3.5.0

d9c085f582f2925583ac85f2df0d0168bcab71a4788a35223e2c81d093a2854a
PS C:\Users\USER> docker run -it -p 8888:8888 -p 4040:4040 --network spark-net --user root -e SPARK_MASTER_URL=spark://spark-master:7077 -e GRANT_SUDO=yes -v ~/spark_data:/opt/spark_data jupyter/all-spark-notebook:latest
Unable to find image 'jupyter/all-spark-notebook:latest' locally
Latest: Pulling from jupyter/all-spark-notebook
aece4893e397: Pulling fs layer
fd92c71966dc: Pulling fs layer
088f1e1e74: Pulling fs layer
4f4fb700ef54: Waiting
ef83736000b0: Waiting
77e45ee945dc: Waiting
a30f89a9afcc: Waiting
dc42adc7eb73: Waiting
abaa3376a050: Waiting
aa099090e49a: Waiting
822c4c0cf6a6: Waiting
d25166dc0c7b: Waiting
964fc3e4ff9f: Waiting

Engine running

RAM 3.42 GB CPU 0.25% Disk: 8.48 GB used (limit 1006.85 GB)

Terminal v4.40.0

Ask Gordon

Containers

Images

Volumes

Builds

Docker Hub

Docker Scout

Extensions

Containers

Give feedback

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

Q 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	<div></div> <div></div> <div></div>
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark		0.2%	17 min	<div></div> <div></div> <div></div>

Terminal

+ v x

[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_nathjax | 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 'jupyterlab'.
[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

Windo... x

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

Ask Gordon

Containers

Images

Volumes

Builds

Docker Hub

Docker Scout

Extensions

Containers

Give feedback

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

Q 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	<div></div> <div></div> <div></div>
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark		0.24%	17 min	<div></div> <div></div> <div></div>

Terminal

+ v x

[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, unfld-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 --frozen-modules=off to python to disable frozen modules.
0.00s - Note: Debugging will proceed. Set PYTHON_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.tostream.StreamClosedError: Stream is closed
[W 2025-05-06 07:59:17.395 LabApp] The extension "nbdime-jupyterlab" is outdated.
[I 2025-05-06 07:59:17.395 LabApp] Build is up to date

Windo... x

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

Ask Gordon

Containers

Images

Volumes

Builds

Docker Hub

Docker Scout

Extensions

Containers

Give feedback

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

Q 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		0.22%	19 min	<div></div> <div></div> <div></div>
<input type="checkbox"/>	spark-worker1	50359ab3b60f	bitnami/spark		0.22%	18 min	<div></div> <div></div> <div></div>
<input type="checkbox"/>	spark-worker2	d9c085f582f2	bitnami/spark		0.22%	17 min	<div></div> <div></div> <div></div>
<input type="checkbox"/>	vibrant_rhodes	0d395916ffa6	jupyter/all:4040-4040	4040-4040	0.01%	4 minute	<div></div> <div></div> <div></div>

Showing 4 items

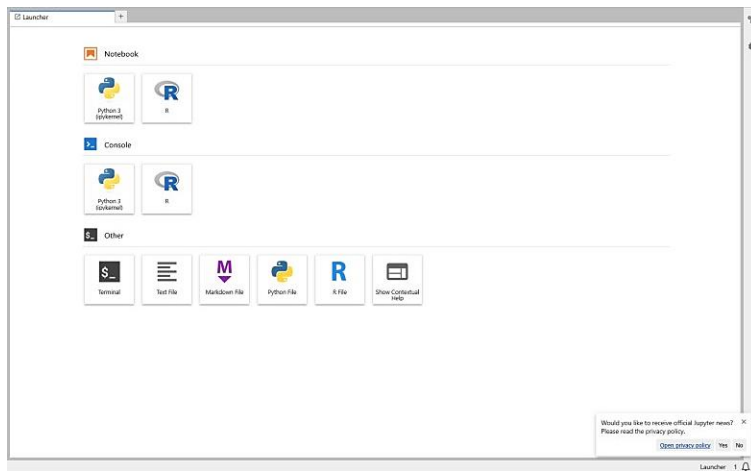
Walkthroughs

+ v x

Terminal

Windo... x

https://jupyter-notebook.readthedocs.io/en/latest/migrate_to_notebook7.html
Please note that updating to Notebook 7 might break some of your extensions.
[I 2025-05-06 07:58:14.465 ServerApp] nbclassic | extension was successfully loaded.
[I 2025-05-06 07:58:14.465 ServerApp] nbdime | extension was successfully loaded.
[I 2025-05-06 07:58:14.642 ServerApp] notebook | extension was successfully loaded.
[I 2025-05-06 07:58:14.643 ServerApp] Serving notebooks from local directory: /home/jovyan
[I 2025-05-06 07:58:14.643 ServerApp] Jupyter Server 2.8.0 is running at:
[I 2025-05-06 07:58:14.643 ServerApp] http://0d395916ffa6:8888/lab?token=fdb03ce64e7fdabc737e2dab8a183f6b4923c87f0dea77



```
CountWord.ipynb
[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)]
```

```
CountWord.ipynb
[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)]
```

```
CountWord.ipynb  X  CountWord.py  X  jovyan@7f403e1835da: ~  X  +

1 from pyspark.sql import SparkSession
2
3 if __name__ == "__main__":
4     spark = SparkSession.builder.appName("WordCount").getOrCreate()
5
6     data = ["Hello Spark", "Hello Docker", "Spark is awesome"]
7     lines = spark.sparkContext.parallelize(data)
8
9     counts = (
10         lines.flatMap(lambda x: x.split(' '))
11         .map(lambda x: (x, 1))
12         .reduceByKey(lambda a, b: a + b)
13     )
14
15     output = counts.collect()
16
17     for (word, count) in output:
18         print(f"{word}: {count}")
19
20     spark.stop()
```

```
CountWord.ipynb  X  CountWord.py  X  jovyan@7f403e1835da: ~  X  +

tor driver) (5/12)
25/06/14 16:14:56 INFO Executor: Finished task 11.0 in stage 1.0 (TID 23). 2100 bytes result sent to driver
25/06/14 16:14:56 INFO PythonRunner: Times: total = 289, boot = -718, init = 1006, finish = 1
25/06/14 16:14:56 INFO TaskSetManager: Finished task 11.0 in stage 1.0 (TID 23) in 479 ms on 7f403e1835da (exec
utor driver) (6/12)
25/06/14 16:14:56 INFO PythonRunner: Times: total = 291, boot = -645, init = 929, finish = 7
25/06/14 16:14:56 INFO TaskSetManager: Finished task 5.0 in stage 1.0 (TID 18) in 486 ms on 7f403e1835da (execu
tor driver) (7/12)
25/06/14 16:14:56 INFO PythonRunner: Times: total = 292, boot = -750, init = 1036, finish = 6
25/06/14 16:14:56 INFO TaskSetManager: Finished task 8.0 in stage 1.0 (TID 20) in 405 ms on 7f403e1835da (execu
tor driver) (8/12)
25/06/14 16:14:56 INFO Executor: Finished task 3.0 in stage 1.0 (TID 14). 2107 bytes result sent to driver
25/06/14 16:14:56 INFO TaskSetManager: Finished task 2.0 in stage 1.0 (TID 13) in 498 ms on 7f403e1835da (execu
tor driver) (9/12)
25/06/14 16:14:56 INFO Executor: Finished task 0.0 in stage 1.0 (TID 16). 2057 bytes result sent to driver
25/06/14 16:14:56 INFO Executor: Finished task 6.0 in stage 1.0 (TID 15). 2108 bytes result sent to driver
25/06/14 16:14:56 INFO TaskSetManager: Finished task 3.0 in stage 1.0 (TID 14) in 507 ms on 7f403e1835da (execu
tor driver) (10/12)
25/06/14 16:14:56 INFO TaskSetManager: Finished task 6.0 in stage 1.0 (TID 15) in 509 ms on 7f403e1835da (execu
tor driver) (11/12)
25/06/14 16:14:56 INFO TaskSetManager: Finished task 0.0 in stage 1.0 (TID 16) in 509 ms on 7f403e1835da (execu
tor driver) (12/12)
25/06/14 16:14:56 INFO TaskSchedulerImpl: Removed TaskSet 1.0, whose tasks have all completed, from pool
25/06/14 16:14:56 INFO DAGScheduler: ResultStage 1 (collect at /home/jovyan/CountWord.py:15) finished in 0.568
s
25/06/14 16:14:56 INFO DAGScheduler: Job 0 is finished. Cancelling potential speculative or zombie tasks for th
is job
25/06/14 16:14:56 INFO TaskSchedulerImpl: Killing all running tasks in stage 1: Stage finished
25/06/14 16:14:56 INFO DAGScheduler: Job 0 finished: collect at /home/jovyan/CountWord.py:15, took 3.142194 s
Hello: 2
is: 1
awesome: 1
Spark: 2
Docker: 1
25/06/14 16:14:56 INFO SparkContext: SparkContext is stopping with exitCode 0.
25/06/14 16:14:56 INFO SparkUI: Stopped Spark web UI at http://7f403e1835da:4041
25/06/14 16:14:56 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
25/06/14 16:14:57 INFO MemoryStore: MemoryStore cleared
25/06/14 16:14:57 INFO BlockManager: BlockManager stopped
25/06/14 16:14:57 INFO BlockManagerMaster: BlockManagerMaster stopped
25/06/14 16:14:57 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped
!
25/06/14 16:14:58 INFO SparkContext: Successfully stopped SparkContext
25/06/14 16:14:59 INFO ShutdownHookManager: Shutdown hook called
25/06/14 16:14:59 INFO ShutdownHookManager: Deleting directory /tmp/spark-b79a0a3a-18ec-44ef-9c3d-4dfb80f21ee2
25/06/14 16:14:59 INFO ShutdownHookManager: Deleting directory /tmp/spark-0a68e836-de0f-4aa6-b343-89fd70173cbd/
pyspark-f30240c1-174d-4fa0-80b3-7885754de392
25/06/14 16:14:59 INFO ShutdownHookManager: Deleting directory /tmp/spark-0a68e836-de0f-4aa6-b343-89fd70173cbd
(base) jovyan@7f403e1835da:~$
```