

## Project 4 - Data Analytics with Hadoop

Project ini merupakan tugas individu dan ditujukan untuk melatih dan mengukur pemahaman terhadap konsep pada materi Data Analytics with Hadoop .

Soal:

- Analisis code project
- Jalankan project dilocal machine masing-masing
- Tambahkan code mapreduce untuk mendapatkan total pendapatan transaksi setiap bulannya
- Kumpulkan tugas dalam bentuk link Github.

Upload project yang sudah dibahas dikelas ke dalam Google Classroom paling lambat Hari Rabu, Jam 23.59 WIB.

### Jawaban

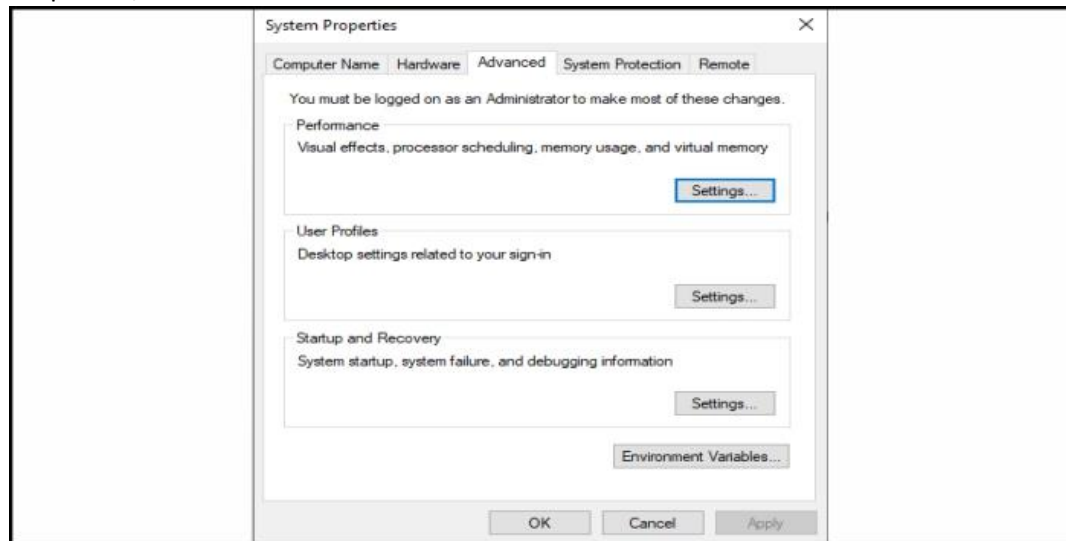
#### **Panduan Install Apache Hadoop 3.2.1 pada OS Windows**

Prasyarat sebelum melakukan instalasi Hadoop versi 3.2.1 pada windows 10 adalah melakukan instalasi Java. Semua versi hadoop hanya support pada Java versi 8. Berikut langkah-langkahnya:

1. Masuk ke website oracle berikut <https://www.oracle.com/java/technologies/javase/javase-jdk8-downloads.html>
2. Download Java Development Kit 8 (JDK 8) windows x64
3. Instalasi Java
  - a. Buka file jdk .exe yang telah download, kemudian ikuti proses instalasi sampai selesai

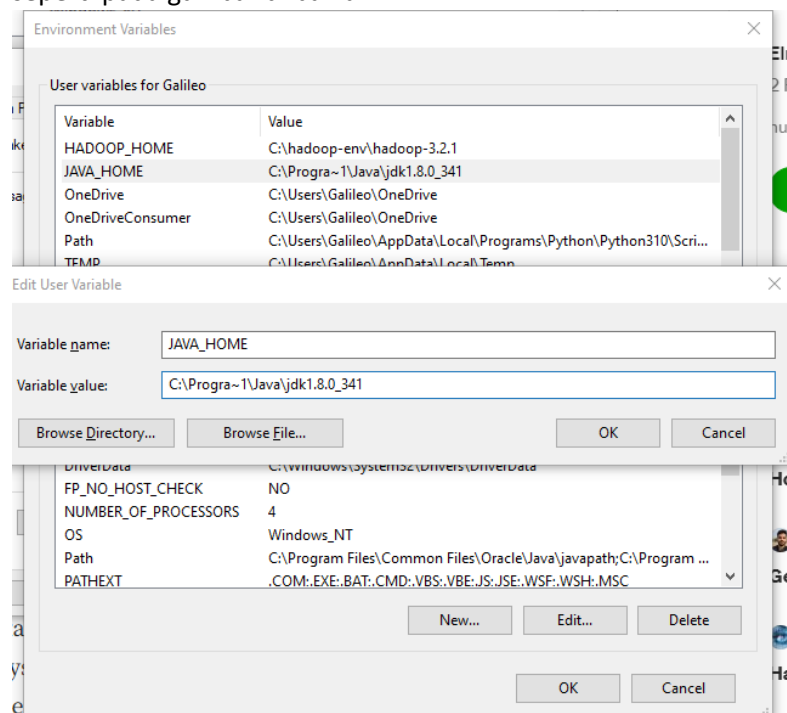


- b. Lakukan Settings Environment Variables Java. Pertama, buka Control Panel – System and Security – System – Advanced System Settings. Kemudian akan muncul dialog box System Properties, lalu klik Environment Variables.

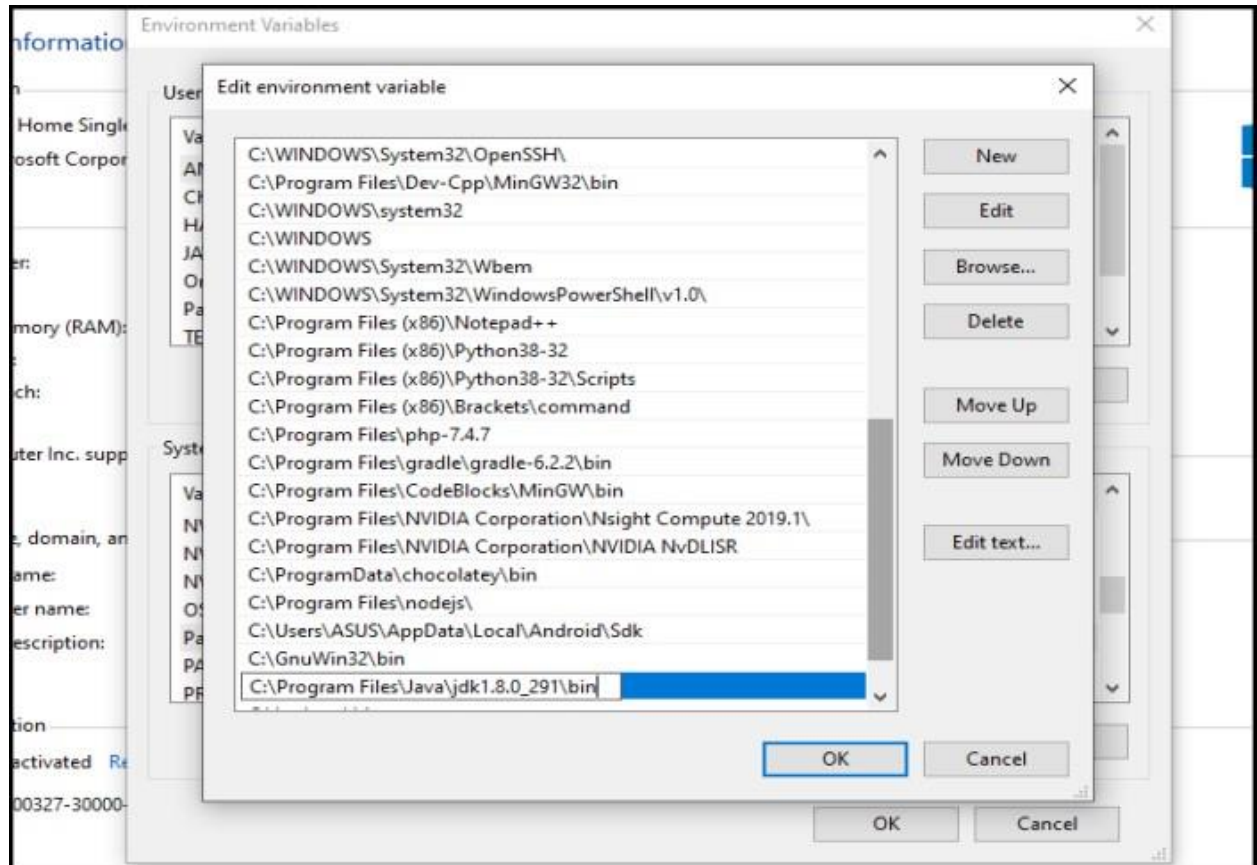


- c. Atur Home Java. Pada User variables, klik New kemudian isi Variable name dengan JAVA\_HOME dan Variable value dengan direktori C:\Program Files\Java\jdk1.8.0\_341. Kemudian klik OK.

Di tahap ini, perlu ada sedikit modifikasi dimana penulisan Variable value tidak diperkenankan ada spasi (Program Files) karena akan menyebabkan perintah tidak berjalan. Maka lakukan sedikit editing, ubah Program Files menjadi Progra~1. Atau jika file jdk kamu berada di Program Files(x86) maka ubah menjadi Progra~2. Sehingga tampilannya akan seperti pada gambar di bawah.



- d. Atur path Java. Pada System variables, klik Path. Kemudian klik New dan isi dengan direktori `jdk\bin`. Lalu klik OK pada Edit environment variable, Environment Variables dan System Properties.



- e. Buka cmd, cek versi java dengan perintah **java -version**. Apabila muncul versi java yang diinstall, maka proses instalasi berhasil.

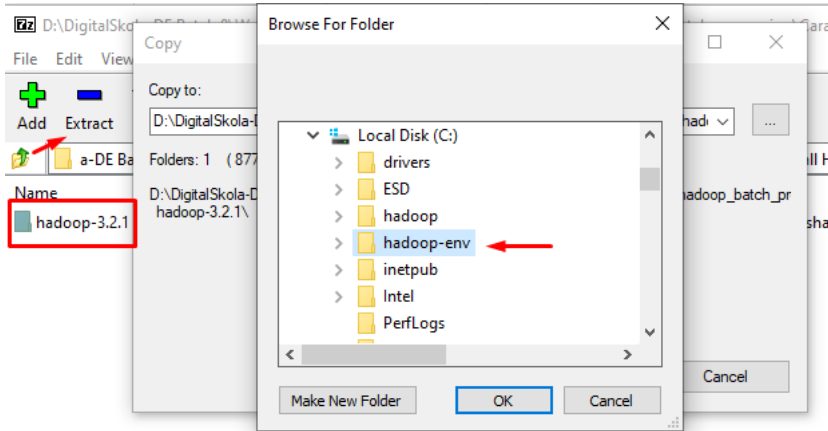
```
C:\> C:\WINDOWS\system32\cmd.exe
```

```
Microsoft Windows [Version 10.0.19043.2006]
(c) Microsoft Corporation. All rights reserved.

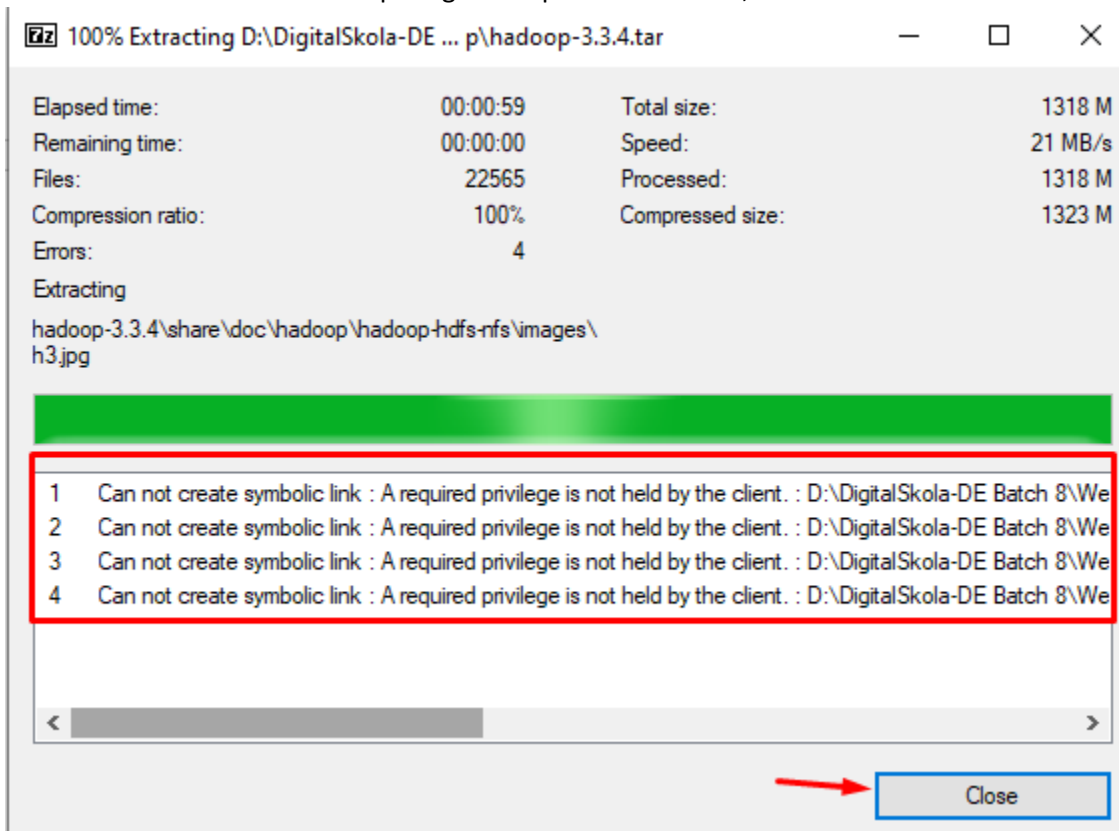
C:\Users\Galileo>java -version
java version "18.0.2.1" 2022-08-18
Java(TM) SE Runtime Environment (build 18.0.2.1+1-1)
Java HotSpot(TM) 64-Bit Server VM (build 18.0.2.1+1-1, mixed mode, sharing)
```

## Proses Instalasi Hadoop versi 3.2.1

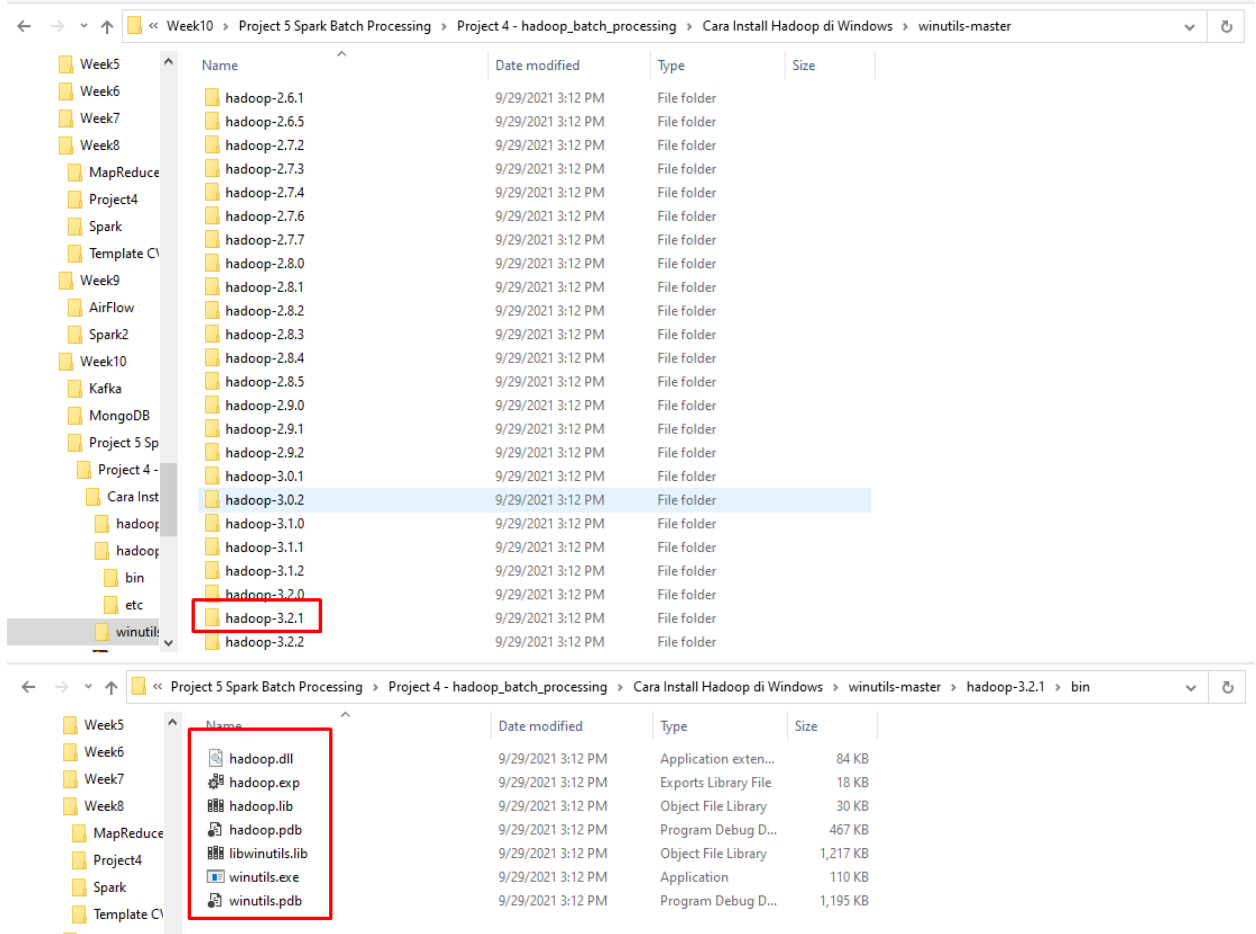
1. Download Hadoop versi 3.2.1 melalui link  
<https://archive.apache.org/dist/hadoop/common/hadoop-3.2.1/hadoop-3.2.1.tar.gz>
2. Copy file hadoop ke drive C, buat folder hadoop-env kemudian ekstrak file hadoop ke dalam folder hadoop-env tersebut. Sehingga file hasil ekstraksi ada di C:\hadoop-env



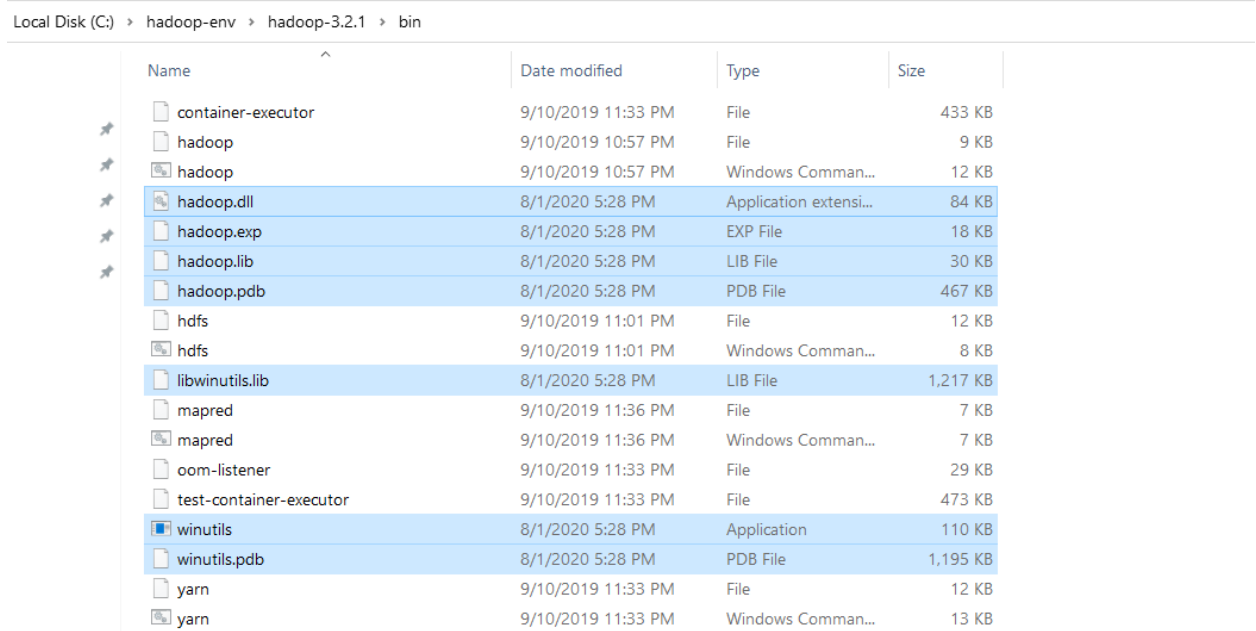
3. Ketika selesai dan menemukan peringatan seperti dibawah ini, maka klik close.



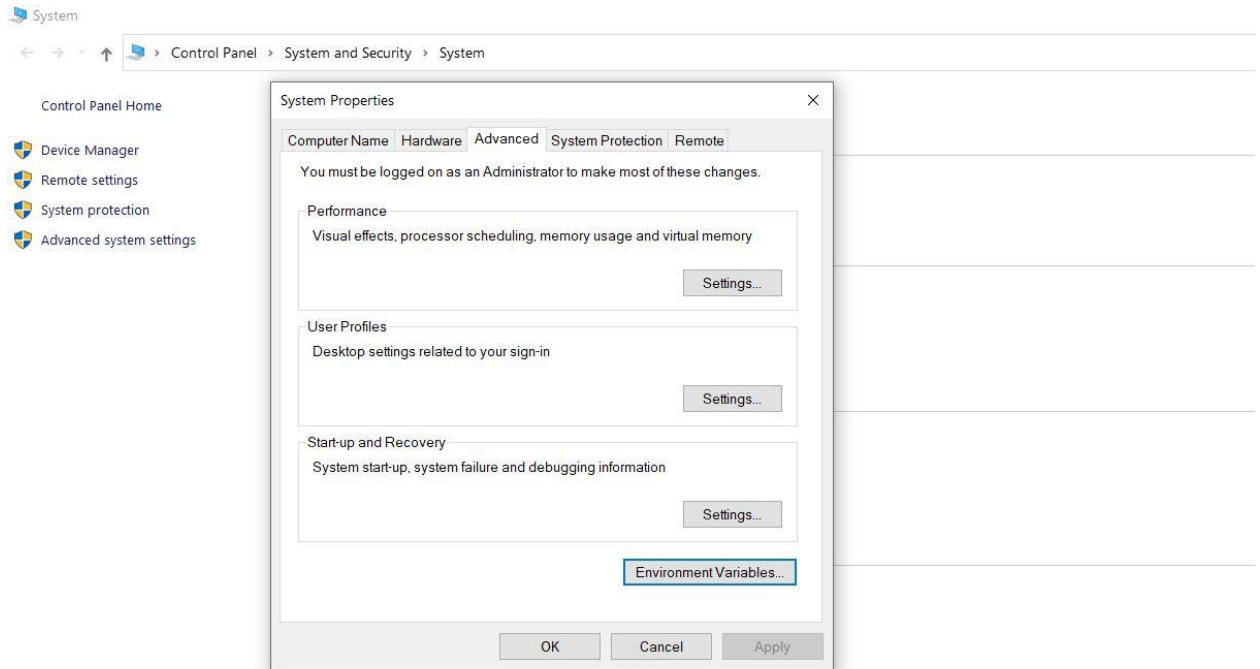
4. Agar Hadoop dapat dioperasikan pada OS Windows, diperlukan patch file Hadoop yang spesifik untuk Windows. Download file zip **bin** yang ada di link ini  
(<https://github.com/cdarlint/winutils/archive/refs/heads/master.zip>), kemudian extract, dan hasilnya seperti dalam gambar di bawah.



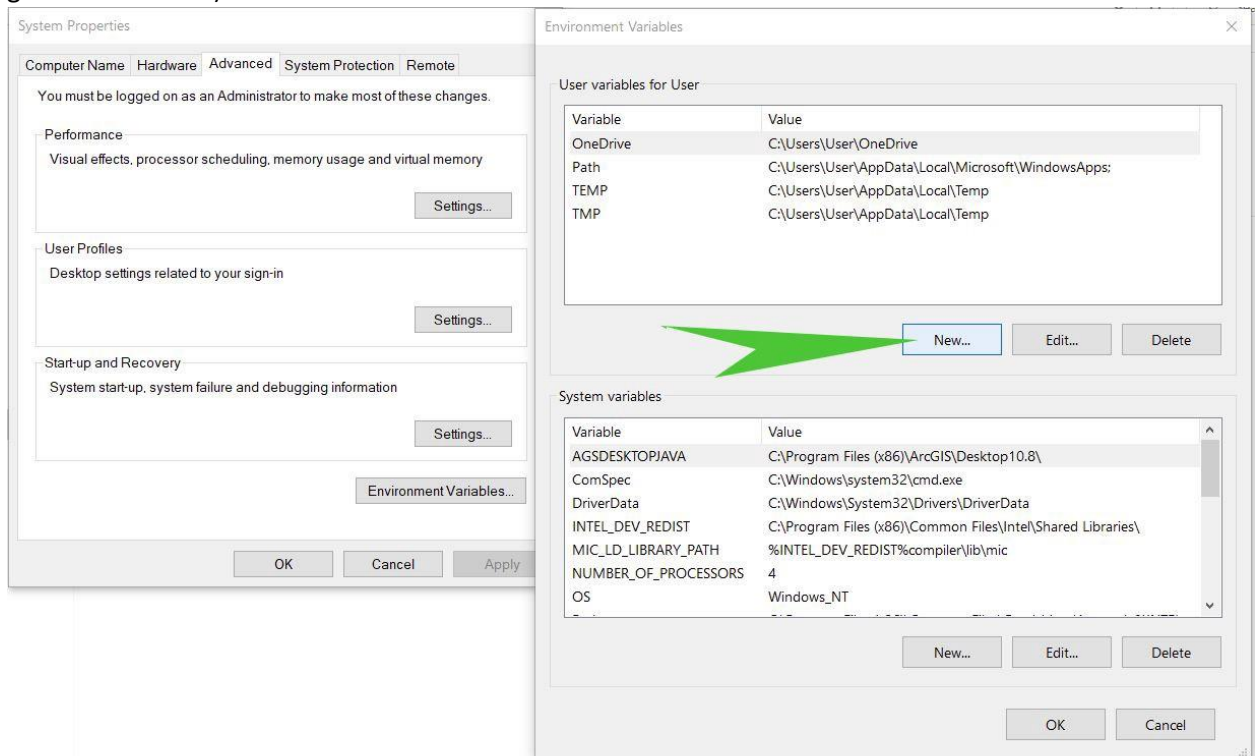
- Letakkan semua file-file tersebut di folder C:\hadoop-env\hadoop-3.2.1\bin, seperti pada file-file yang tersorot biru pada gambar di bawah.



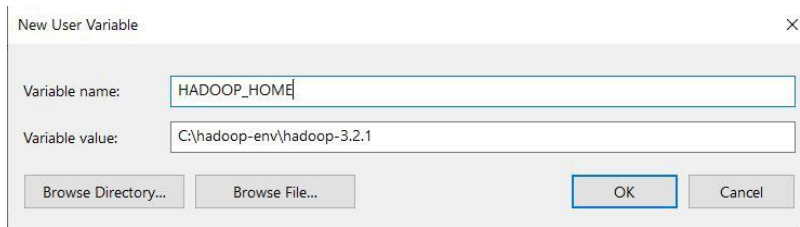
6. Setting pada environment variables Klik Control Panel > System and Security > System > Advanced system settings (lihat gambar di bawah)



7. Lalu, klik Environment Variables. Dan lanjutkan klik New pada user variables for User (lihat gambar di bawah)

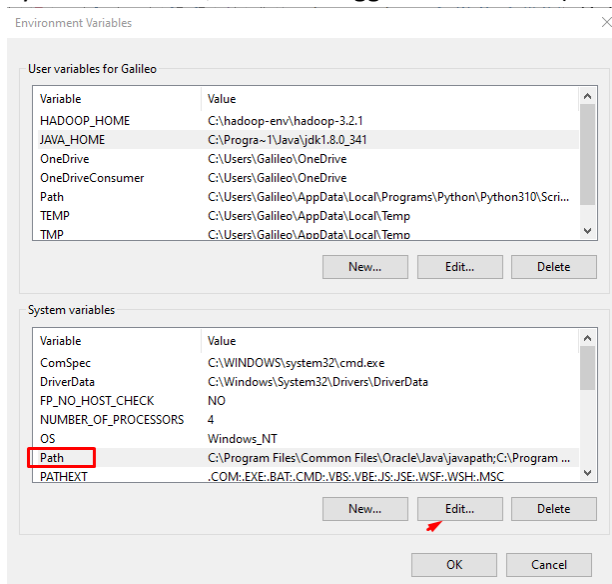


8. Setelah klik New, akan muncul popup Kemudian ketikkan HADOOP\_HOME pada form Variable name, dan C:\hadoop-env\hadoop-3.2.1 pada form Variable value, seperti pada gambar di bawah. Lalu klik OK.



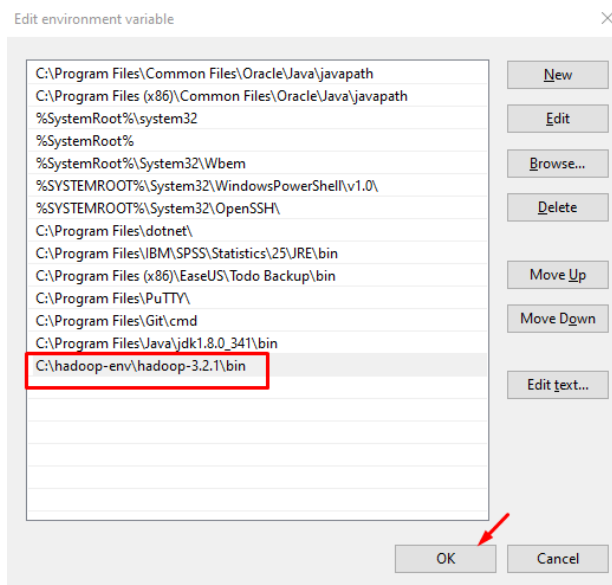
The 'New User Variable' dialog box is shown. It has two input fields: 'Variable name' with the text 'HADOOP\_HOME' and 'Variable value' with the text 'C:\hadoop-env\hadoop-3.2.1'. There are buttons for 'Browse Directory...', 'Browse File...', 'OK', and 'Cancel'.

9. Setelah klik OK di atas, akan muncul tampilan seperti pada gambar di bawah. Pada bagian System variables, scroll sehingga muncul Path (lihat gambar di bawah), select kemudian klik Edit.



The 'Environment Variables' dialog box is shown. It has two sections: 'User variables for Galileo' and 'System variables'. The 'System variables' section is expanded, and the 'Path' variable is selected and highlighted with a red box. The 'Edit...' button is also highlighted with a red arrow.

10. Klik New dan isi dengan direktori C:\hadoop-env\hadoop-3.2.1\bin. Lalu klik OK pada Edit environment variable, Environment Variables dan System Properties.



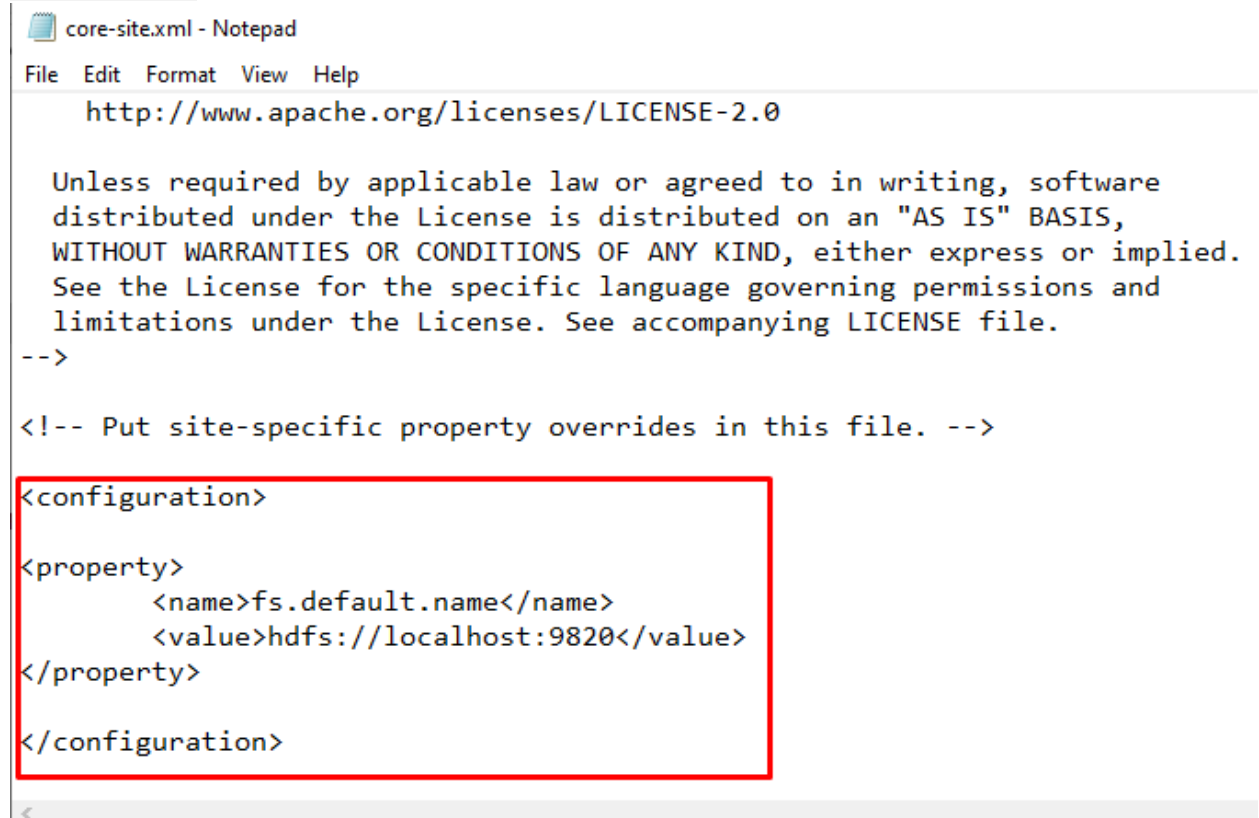
The 'Edit environment variable' dialog box is shown. It has a list of paths on the left, including 'C:\Program Files\Common Files\Oracle\Java\javapath', '%SystemRoot%\system32', and 'C:\hadoop-env\hadoop-3.2.1\bin'. The 'New' button is highlighted with a red arrow. The 'OK' button is also highlighted with a red arrow.

11. Sampai dengan tahap ini, semoga berjalan lancar. Lakukan testing untuk melihat apakah setting java dan hadoop di environment variables telah berhasil. Caranya: buka cmd, kemudian ketik `hadoop -version`. Hasilnya seperti pada gambar di bawah:

```
C:\Users\Galileo>hadoop -version
java version "1.8.0_341"
Java(TM) SE Runtime Environment (build 1.8.0_341-b10)
Java HotSpot(TM) 64-Bit Server VM (build 25.341-b10, mixed mode)
```

12. Kemudian Kita akan melakukan konfigurasi terhadap 4 file hadoop. Buka folder `C:\hadoop-env\hadoop-3.2.1\etc\hadoop`. Kemudian buka file `hdfs-site.xml`, `core-site.xml`, `mapred-site.xml`, `yarn-site.xml` di text editor.
13. Tambahkan code berikut pada file `core-site.xml`

```
<property>
<name>fs.default.name</name>
<value>hdfs://localhost:9820</value>
</property>
```



14. Tambahkan code berikut pada file `mapred-site.xml`

```
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
<description>MapReduce framework name</description>
</property>
```



mapred-site.xml - Notepad

File Edit Format View Help

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License. See accompanying LICENSE file.

-->

<!-- Put site-specific property overrides in this file. -->

```
<configuration>
  <property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
    <description>MapReduce framework name</description>
  </property>
</configuration>
```

15. Tambahkan code berikut pada file yarn-site.xml

```
<property>
<name>yarn.nodemanager.aux-
services</name><value>mapreduce_shuffle</value>
<description>Yarn Node Manager Aux Service</description>
</property>
```

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License. See accompanying LICENSE file.

-->

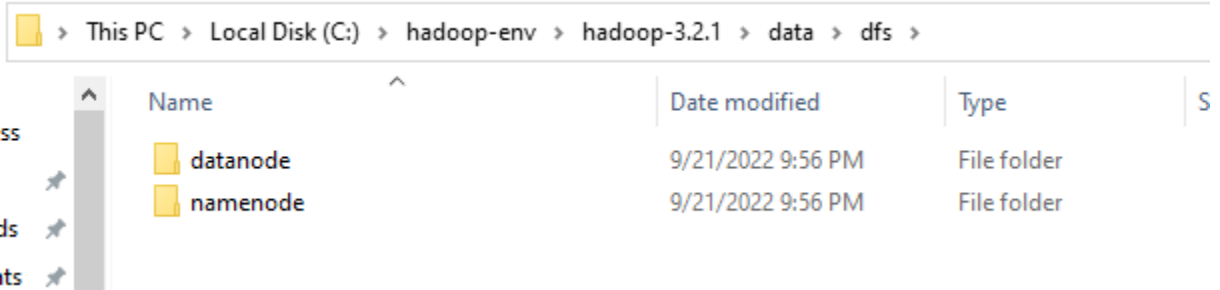
```
<configuration>

<!-- Site specific YARN configuration properties -->

<property>
    <name>yarn.nodemanager.aux-services</name><value>mapreduce_shuffle</value>
    <description>Yarn Node Manager Aux Service</description>
</property>

</configuration>
```

16. Pada direktori hadoop buatlah folder baru dengan nama **data**, didalam folder **data** buat folder baru dengan nama **dfs**. Pada folder **dfs** tersebut buat 2 folder baru dengan nama **datanode** dan **namenode**.



17. Tambahkan code berikut pada file hdfs-site.xml. Untuk tag value disesuaikan dengan direktori dimana folder namenode dan datanode dibuat.

```
<property>
<name>dfs.replication</name>
<value>1</value>
</property>
<property>
<name>dfs.namenode.name.dir</name>
<value>file:///C:/hadoop-env/hadoop-3.2.1/data/dfs/namenode</value>
</property>
<property>
<name>dfs.datanode.data.dir</name>
<value>file:///C:/hadoop-env/hadoop-3.2.1/data/dfs/datanode</value>
</property>
```

hdfs-site.xml - Notepad

File Edit Format View Help

```
<!-- Put site-specific property overrides in this file. -->

<configuration>

<property>
  <name>dfs.replication</name>
  <value>1</value>
</property>
<property>
  <name>dfs.namenode.name.dir</name>
  <value>file:///C:/hadoop-env/hadoop-3.2.1/data/dfs/namenode</value>
</property>
<property>
  <name>dfs.datanode.data.dir</name>
  <value>file:///C:/hadoop-env/hadoop-3.2.1/data/dfs/datanode</value>
</property>

</configuration>
```

Ln 1, Col 1 100% Unix (LF) UTF-8

18. [OPTIONAL] Pada file hadoop-env.cmd, sesuaikan direktori JAVA\_HOME dengan direktori java jdk.

File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?

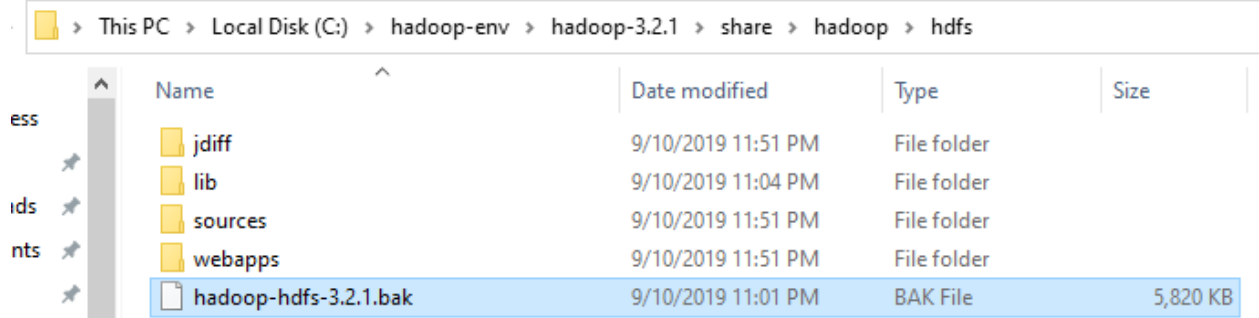
core-site.xml hdfs-site.xml hadoop-env.cmd yam-site.xml mapred-site.xml template

```
1 @echo off
2 @rem Licensed to the Apache Software Foundation (ASF) under one or more
3 @rem contributor license agreements. See the NOTICE file distributed with
4 @rem this work for additional information regarding copyright ownership.
5 @rem The ASF licenses this file to You under the Apache License, Version 2.0
6 @rem (the "License"); you may not use this file except in compliance with
7 @rem the License. You may obtain a copy of the License at
8 @rem
9 @rem http://www.apache.org/licenses/LICENSE-2.0
10 @rem
11 @rem Unless required by applicable law or agreed to in writing, software
12 @rem distributed under the License is distributed on an "AS IS" BASIS,
13 @rem WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
14 @rem See the License for the specific language governing permissions and
15 @rem limitations under the License.
16
17 @rem Set Hadoop-specific environment variables here.
18
19 @rem The only required environment variable is JAVA_HOME. All others are
20 @rem optional. When running a distributed configuration it is best to
21 @rem set JAVA_HOME in this file, so that it is correctly defined on
22 @rem remote nodes.
23
24 @rem The java implementation to use. Required.
25 set JAVA_HOME=C:\Program Files\Java\jdk1.8.0_181
26
27 @rem The jsvc implementation to use. Jsvc is required to run secure datanodes.
28 @rem set JSVC_HOME=%JSVC_HOME%
29
30 @rem set HADOOP_CONF_DIR=
31
32 @rem Extra Java CLASSPATH elements. Automatically insert capacity-scheduler.
33 if exist %HADOOP_HOME%\contrib\capacity-scheduler (
34   if not defined HADOOP_CLASSPATH (
35     set HADOOP_CLASSPATH=%HADOOP_HOME%\contrib\capacity-scheduler\*.jar
36   ) else (
37     set HADOOP_CLASSPATH=%HADOOP_CLASSPATH%;%HADOOP_HOME%\contrib\capacity-sch
```

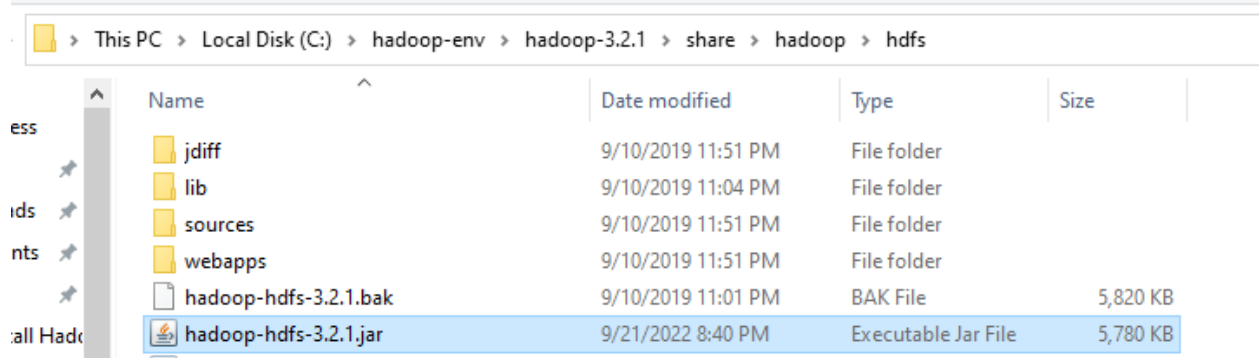
Penulisan direktori JAVA\_HOME sebenarnya tidak boleh terdapat folder yang mengandung spasi, namun apabila ingin tetap terdapat folder dengan spasi, maka ditambahkan "...." agar tidak terjadi error.

19. Melakukan format pada Name node. Di tahap ini cukup **krusial** karena jika tidak dilakukan akan ada error saat menjalankan formatting. Sebelum melakukan format, lakukan 3 tahapan ini:

- Download file `hadoop-hdfs-3.2.1.jar` pada link <https://github.com/FahaoTang/big-data/raw/master/hadoop-hdfs-3.2.1.jar>
- Buka folder `C:\hadoop-env\hadoop-3.2.1\share\hadoop\hdfs` kemudian Rename nama file `hadoop-hdfs-3.2.1.jar` menjadi `hadoop-hdfs-3.2.1.bak`



- Copy file `hadoop-hdfs-3.2.1.jar` yang sebelumnya telah di download pada langkah 9.a. ke dalam folder `C:\hadoop-env\hadoop-3.2.1\share\hadoop\hdfs`



20. Buka Command prompt dan ketikkan perintah **`hdfs namenode -format`**. maka lihatlah hasilnya, pada `Startup_msg` akan menampilkan Starting NameNode.

```
Administrator: Command Prompt - hdfs namenode -format
Microsoft Windows [Version 10.0.18363.1440]
(c) 2019 Microsoft Corporation. All rights reserved.

C:\WINDOWS\system32>hdfs namenode -format
21/05/10 16:57:49 INFO namenode.NameNode: STARTUP_MSG:
/*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = Aziz/192.168.56.1
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 2.9.2
STARTUP_MSG: classpath = C:\hadoop\etc\hadoop;c:\hadoop\share\hadoop\common\lib\activation-1.1.jar;c:\hadoop\share\hadoop\common\lib\apacheds-118n-2.0.0-M15.jar;c:\hadoop\share\hadoop\common\lib\apacheds-kerberos-codec-2.0.0-M15.jar;c:\hadoop\share\hadoop\common\lib\api-asn1-api-1.0.0-M20.jar;c:\hadoop\share\hadoop\common\lib\api-util-1.0.0-M20.jar;c:\hadoop\share\hadoop\common\lib\asm-3.2.jar;c:\hadoop\share\hadoop\common\lib\avro-1.7.7.jar;c:\hadoop\share\hadoop\common\lib\commons-beanutils-1.7.0.jar;c:\hadoop\share\hadoop\common\lib\commons-beanutils-core-1.8.0.jar;c:\hadoop\share\hadoop\common\lib\commons-cli-1.2.jar;c:\hadoop\share\hadoop\common\lib\commons-codec-1.4.jar;c:\hadoop\share\hadoop\common\lib\commons-collections-3.2.2.jar;c:\hadoop\share\hadoop\common\lib\commons-compress-1.4.1.jar;c:\hadoop\share\hadoop\common\lib\commons-configuration-1.6.jar;c:\hadoop\share\hadoop\common\lib\commons-digester-1.8.jar;c:\hadoop\share\hadoop\common\lib\commons-io-2.4.jar;c:\hadoop\share\hadoop\common\lib\commons-lang-2.6.jar;c:\hadoop\share\hadoop\common\lib\commons-lang3-3.4.jar;c:\hadoop\share\hadoop\common\lib\commons-logging-1.1.3.jar;c:\hadoop\share\hadoop\common\lib\commons-math3-3.1.1.jar;c:\hadoop\share\hadoop\common\lib\commons-net-3.1.jar;c:\hadoop\share\hadoop\common\lib\curator-client-2.7.1.jar;c:\hadoop\share\hadoop\common\lib\curator-framework-2.7.1.jar;c:\hadoop\share\hadoop\common\lib\curator-recipes-2.7.1.jar;c:\hadoop\share\hadoop\common\lib\gson-2.2.4.jar;c:\hadoop\share\hadoop\common\lib\guava-11.0.2.jar;c:\hadoop\share\hadoop\common\lib\hadoop-annotations-2.9.2.jar;c:\hadoop\share\hadoop\common\lib\hadoop-auth-2.9.2.jar;c:\hadoop\share\hadoop\common\lib\hadoop-htrace-core-4.1.0-incubating.jar;c:\hadoop\share\hadoop\common\lib\httpclient-4.5.2.jar;c:\hadoop\share\hadoop\common\lib\httpcore-4.4.4.jar;c:\hadoop\share\hadoop\common\lib\jackson-core-asl-1.9.13.jar;c:\hadoop\share\hadoop\common\lib\jackson-jaxrs-1.9.13.jar;c:\hadoop\share\hadoop\common\lib\jackson-mapper-asl-1.9.13.jar;c:\hadoop\share\hadoop\common\lib\jackson-xc-1.9.13.jar;c:\hadoop\share\hadoop\common\lib\java-xmlbuilder-0.4.jar;c:\hadoop\share\hadoop\common\lib\jaxb-api-2.2.2.jar;c:\hadoop\share\hadoop\common\lib\jaxb-impl-2.2.3-1.jar;c:\hadoop\share\hadoop\common\lib\jcip-annotations-1.0-1.jar;c:\hadoop\share\hadoop\commo
```

```
Select Command Prompt
2022-09-21 20:21:40,254 INFO util.GSet: Computing capacity for map cachedBlocks
2022-09-21 20:21:40,254 INFO util.GSet: VM type = 64-bit
2022-09-21 20:21:40,255 INFO util.GSet: 0.25% max memory 889 MB = 2.2 MB
2022-09-21 20:21:40,256 INFO util.GSet: capacity = 2^18 = 262144 entries
2022-09-21 20:21:40,280 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.window.num.buckets = 10
2022-09-21 20:21:40,280 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.num.users = 10
2022-09-21 20:21:40,282 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.windows.minutes = 1,5,25
2022-09-21 20:21:40,292 INFO namenode.FSNamesystem: Retry cache on namenode is enabled
2022-09-21 20:21:40,293 INFO namenode.FSNamesystem: Retry cache will use 0.03 of total heap and retry cache entry expiry
time is 600000 millis
2022-09-21 20:21:40,300 INFO util.GSet: Computing capacity for map NameNodeRetryCache
2022-09-21 20:21:40,300 INFO util.GSet: VM type = 64-bit
2022-09-21 20:21:40,302 INFO util.GSet: 0.029999999329447746% max memory 889 MB = 273.1 KB
2022-09-21 20:21:40,302 INFO util.GSet: capacity = 2^15 = 32768 entries
2022-09-21 20:21:40,380 INFO namenode.FSImage: Allocated new BlockPoolId: BP-212676949-192.168.88.102-1663766580352
2022-09-21 20:21:40,444 INFO common.Storage: Storage directory C:\hadoop\data\namenode has been successfully formatted.
2022-09-21 20:21:40,504 INFO namenode.FSImageFormatProtobuf: Saving image file C:\hadoop\data\namenode\current\fsimage.ckpt_000000000000000000 using no compression
2022-09-21 20:21:40,779 INFO namenode.FSImageFormatProtobuf: Image file C:\hadoop\data\namenode\current\fsimage.ckpt_000000000000000000 of size 402 bytes saved in 0 seconds .
2022-09-21 20:21:40,808 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
2022-09-21 20:21:40,852 INFO namenode.FSNamesystem: Stopping services started for active state
2022-09-21 20:21:40,853 INFO namenode.FSNamesystem: Stopping services started for standby state
2022-09-21 20:21:40,863 INFO namenode.FSImage: FSImageSaver clean checkpoint: txid=0 when meet shutdown.
2022-09-21 20:21:40,864 INFO namenode.NameNode: SHUTDOWN_MSG:
/*****
SHUTDOWN_MSG: Shutting down NameNode at User/192.168.88.102
*****/
C:\hadoop\sbin>
```

21. Lalu aktifkan Hadoop Services dengan buka kembali command prompt, dan arahkan ke C:\hadoop-env\hadoop-3.2.1\sbin.

```
Command Prompt
2022-09-21 21:55:42,711 INFO namenode.FSNamesystem: Retry cache will use 0.03 of total heap and retry cache entry expiry
time is 600000 millis
2022-09-21 21:55:42,719 INFO util.GSet: Computing capacity for map NameNodeRetryCache
2022-09-21 21:55:42,719 INFO util.GSet: VM type = 64-bit
2022-09-21 21:55:42,723 INFO util.GSet: 0.029999999329447746% max memory 889 MB = 273.1 KB
2022-09-21 21:55:42,723 INFO util.GSet: capacity = 2^15 = 32768 entries
2022-09-21 21:55:42,827 INFO namenode.FSImage: Allocated new BlockPoolId: BP-1217908123-192.168.88.102-1663772142804
2022-09-21 21:55:42,898 INFO common.Storage: Storage directory C:\hadoop-env\hadoop-3.2.1\data\dfs\namenode has been suc
cessfully formatted.
2022-09-21 21:55:42,984 INFO namenode.FSImageFormatProtobuf: Saving image file C:\hadoop-env\hadoop-3.2.1\data\dfs\name
ode\current\fsimage.ckpt_000000000000000000 using no compression
2022-09-21 21:55:43,251 INFO namenode.FSImageFormatProtobuf: Image file C:\hadoop-env\hadoop-3.2.1\data\dfs\name
ode\current\fsimage.ckpt_000000000000000000 of size 402 bytes saved in 0 seconds .
2022-09-21 21:55:43,280 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
2022-09-21 21:55:43,296 INFO namenode.FSImage: FSImageSaver clean checkpoint: txid=0 when meet shutdown.
2022-09-21 21:55:43,297 INFO namenode.NameNode: SHUTDOWN_MSG:
/*****
SHUTDOWN_MSG: Shutting down NameNode at User/192.168.88.102
*****/
C:\Users\Galileo>cd C:\hadoop-env\hadoop-3.2.1\sbin
```

22. Jalankan command .\start-dfs.cmd.

```
C:\hadoop-env\hadoop-3.2.1\sbin>.\start-dfs.cmd
```

Dan muncullah seperti pada gambar di bawah, yang menandakan service Hadoop dfs telah aktif.

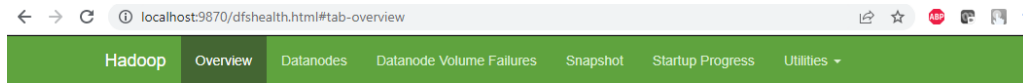




24. Memastikan semua services telah aktif, ketikkan perintah **jps**

```
C:\hadoop-env\hadoop-3.2.1\sbin>jps
15616 NameNode
4992 Jps
14500 ResourceManager
10040 DataNode
5404 NodeManager
```

25. Membuka Hadoop Web UI dengan cara jalankan browser kemudian ketikan alamat URL Name node: <http://localhost:9870/dfshealth.html>



## Overview 'localhost:9820' (active)

Started:	Wed Sep 21 21:56:41 +0700 2022
Version:	3.2.1, rb3cbbb467e22ea829b3808f4b7b01d07e0bf3842
Compiled:	Tue Sep 10 22:56:00 +0700 2019 by rohithsharmaks from branch-3.2.1
Cluster ID:	CID-3c1422bc-cc68-43cf-bffb-de0d6c08120a
Block Pool ID:	BP-1217908123-192.168.88.102-1663772142804

## Summary

Security is off.

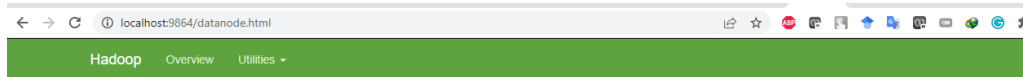
Safemode is off.

1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).

Heap Memory used 78.61 MB of 155 MB Heap Memory. Max Heap Memory is 889 MB.

Non Heap Memory used 51.09 MB of 52.15 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

Data node: <http://localhost:9864/datanode.html>



## DataNode on User:9866

Cluster ID:	CID-3c1422bc-cc68-43cf-bffb-de0d6c08120a
Version:	3.2.1, rb3cbbb467e22ea829b3808f4b7b01d07e0bf3842


## Block Pools

Namenode Address	Block Pool ID	Actor State	Last Heartbeat	Last Block Report	Last Block Report Size (Max Size)
localhost:9820	BP-1217908123-192.168.88.102-1663772142804	RUNNING	1s	an hour	0 B (64 MB)

## Volume Information

Directory	Storage Type	Capacity Used	Capacity Left	Capacity Reserved	Reserved Space for Replicas	Blocks
C:\hadoop-env\hadoop-3.2.1\data\dfs\datanode	DISK	323 B	20.7 GB	0 B	0 B	0

Yarn: <http://localhost:8088/cluster>



Cluster

About

Nodes

Node Labels

Applications

NEW

NEW SAVING

SUBMITTED

ACCEPTED

RUNNING

FINISHED

FAILED

KILLED

Scheduler

Tools

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores
0	0	0	0	0	0 B	8 GB	0 B	0

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy Nodes	Reb
1	0	0	0	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation
Capacity Scheduler	[memory-mb (unit=Mi), vcores]	<memory:1024, vCores:1>	<memory:8192, vCores:4>

Show 20 entries

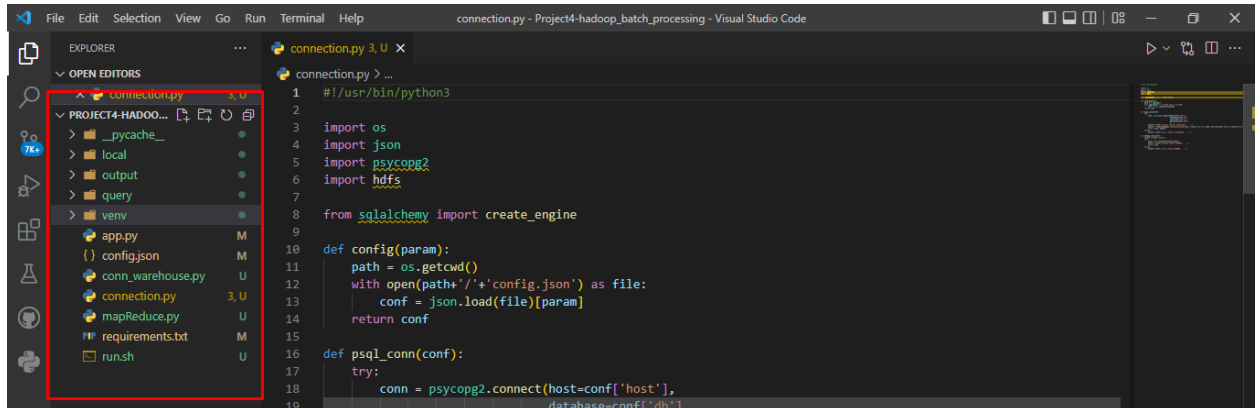
ID	User	Name	Application Type	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU VCo	Allocated Memory MB	Reserved CPU VCo	Reserved Memory MB
No data available in table															

Showing 0 to 0 of 0 entries

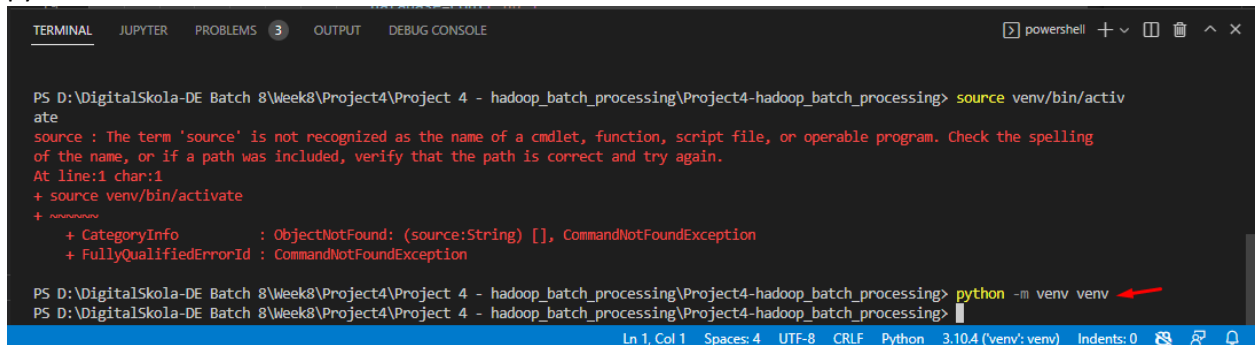


## Data Analytic dengan Hadoop

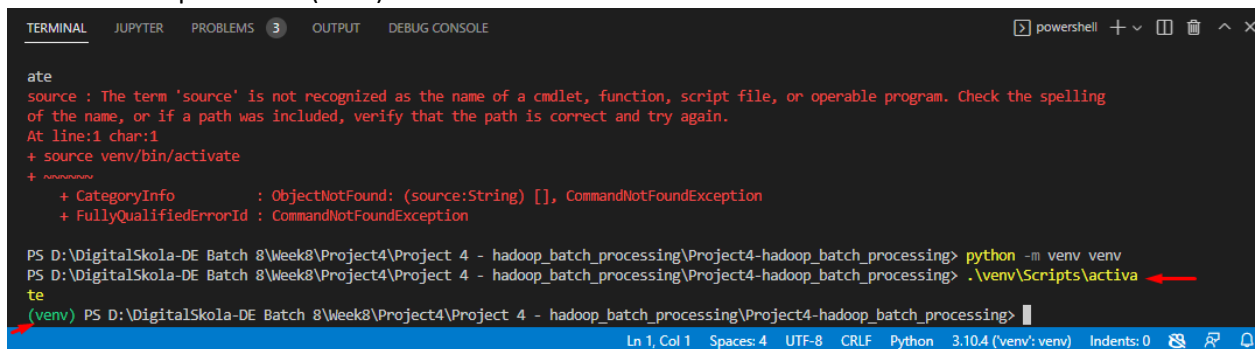
1. Buka visual studio code. Open folder dan arah ke folder Project4-hadoop\_batch\_processing



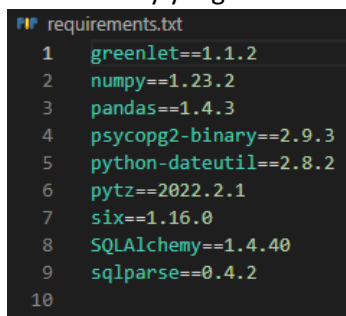
2. Pada Visual Studio Code aktifkan terminal kemudian buat virtual environment dengan perintah `python -m venv venv`



3. Aktifkan virtual environment tersebut dengan perintah `.\venv\Scripts\activate`  
Pastikan terdapat tulisan (venv) di sebelah kiri



4. Install library yang dibutuhkan sesuai list pada file requirements.txt



dengan perintah pip install -r .\requirements.txt

```
TERMINAL JUPYTER PROBLEMS 3 OUTPUT DEBUG CONSOLE
(venv) PS D:\DigitalSkola-DE Batch 8\Week8\Project4\Project 4 - hadoop_batch_processing\Project4-hadoop_batch_processing> pip install -r .\requirements.txt
Collecting greenlet==1.1.2
  Using cached greenlet-1.1.2-cp310-cp310-win_amd64.whl (101 kB)
Collecting numpy==1.23.2
  Using cached numpy-1.23.2-cp310-cp310-win_amd64.whl (14.6 MB)
Collecting pandas==1.4.3
  Using cached pandas-1.4.3-cp310-cp310-win_amd64.whl (10.5 MB)
Collecting psycpg2-binary==2.9.3
  Using cached psycpg2-binary-2.9.3-cp310-cp310-win_amd64.whl (1.2 MB)
Collecting python-dateutil==2.8.2
  Using cached python-dateutil-2.8.2-py2.py3-none-any.whl (247 kB)
```

Tunggu sampai proses instalasi library selesai

```
TERMINAL JUPYTER PROBLEMS 1 OUTPUT DEBUG CONSOLE
Collecting six==1.16.0
  Using cached six-1.16.0-py2.py3-none-any.whl (11 kB)
Collecting SQLAlchemy==1.4.40
  Using cached SQLAlchemy-1.4.40-cp310-cp310-win_amd64.whl (1.6 MB)
Collecting sqlparse==0.4.2
  Using cached sqlparse-0.4.2-py3-none-any.whl (42 kB)
Installing collected packages: pytz, sqlparse, six, psycpg2-binary, numpy, greenlet, SQLAlchemy, python-dateutil, pandas
Successfully installed SQLAlchemy-1.4.40 greenlet-1.1.2 numpy-1.23.2 pandas-1.4.3 psycpg2-binary-2.9.3 python-dateutil-2.8.2 pytz-2022.2.1 six-1.16.0 sqlparse-0.4.2
WARNING: You are using pip version 22.0.4; however, version 22.2.2 is available.
You should consider upgrading via the 'D:\DigitalSkola-DE Batch 8\Week8\Project4\Project 4 - hadoop_batch_processing\Project4-hadoop_batch_processing\venv\Scripts\python.exe -m pip install --upgrade pip' command.
(venv) PS D:\DigitalSkola-DE Batch 8\Week8\Project4\Project 4 - hadoop_batch_processing\Project4-hadoop_batch_processing>
```

5. Jalankan program proses python hadoop batch processing dengan perintah **python app.py**




```
(venv) PS D:\DigitalSkola-DE Batch 8\Week8\Project4\Project 4 - hadoop_batch_processing\Project4-hadoop_batch_processing> python .\app.py
[INFO] Service ETL is Starting .....
[INFO] Success connect Warehouse .....
[INFO] Success connect PostgreSQL .....
[INFO] Success connect HADOOP .....
[INFO] Service ETL is Running .....
[INFO] Upload Data in HADOOP Success .....
[INFO] Upload Data in LOCAL Success .....
[INFO] Update WDH Success .....
[INFO] Service ETL is Success .....
(venv) PS D:\DigitalSkola-DE Batch 8\Week8\Project4\Project 4 - hadoop_batch_processing\Project4-hadoop_batch_processing>
```

6. Bila program berjalan dengan baik tanpa error, maka di hadoop akan muncul folder bernama digitalskola

localhost:9870/explorer.html#/

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

## Browse Directory

/ Go!   

Show 25 entries Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<input type="checkbox"/>	drwxr-xr-x	Galileo	supergroup	0 B	Sep 22 15:57	0	0 B	digitalskola

Showing 1 to 1 of 1 entries

Previous 1 Next




Hadoop, 2019.

7. Bila di-explore isi dari folder digitalskola, maka akan ada folder bernama project dan di dalam folder project akan ada file dim\_order\_20220922.csv

localhost:9870/explorer.html#/digitalskola/project

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

## Browse Directory

/digitalskola/project Go!   

Show 25 entries Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<input type="checkbox"/>	-rw-r--r--	Galileo	supergroup	1.22 KB	Sep 22 15:57	1	128 MB	dim_order_20220922.csv

Showing 1 to 1 of 1 entries

Previous 1 Next

Hadoop, 2019.

8. Setelahnya jalankan program python mapreduce untuk mengolah data dim\_order\_20220922 yang telah ada di hadoop dengan perintah **python .\mapReduce.py -r hadoop hdfs:///digitalskola/project/dim\_order\_20220922.csv**  
Dan tunggu sampai proses mapreduce selesai running

```
(venv) PS D:\DigitalSkola-DE Batch 8\Week8\Project4\Project 4 - hadoop_batch_processing\Project4-hadoop_batch_processing> python .\mapReduce.py -r hadoop hdfs:///digitalskola/project/dim_order_20220922.csv
No configs found; falling back on auto-configuration
No configs specified for hadoop runner
Looking for hadoop binary in C:\hadoop-env\hadoop-3.2.1\bin...
Found hadoop binary: C:\hadoop-env\hadoop-3.2.1\bin\hadoop.cmd
Using Hadoop version 3.2.1
Looking for Hadoop streaming jar in C:\hadoop-env\hadoop-3.2.1...
Looking for Hadoop streaming jar in D:\DigitalSkola-DE Batch 8\Week8\Project4\Project 4 - hadoop_batch_processing\Project4-hadoop_batch_processing...
Looking for Hadoop streaming jar in /home/hadoop/contrib...
Looking for Hadoop streaming jar in /usr/lib/hadoop-mapreduce...
Hadoop streaming jar not found. Use --hadoop-streaming-jar
Creating temp directory C:\Users\Galileo\AppData\Local\Temp\mapReduce.Galileo.20220926.142237.165822
uploading working dir files to hdfs:///user/Galileo/tmp/mrjob/mapReduce.Galileo.20220926.142237.165822/files/wd...
Copying other local files to hdfs:///user/Galileo/tmp/mrjob/mapReduce.Galileo.20220926.142237.165822/files/
```

9. Hasil output mapreduce yang dihasilkan adalah untuk total pendapatan (counting) transaksi setiap bulannya di tahun 2022.

[illegible]

 ordercount\_output\_local\_map.txt - Notepad

File Edit Format View Help

"2022-01"	2
"2022-02"	1
"2022-03"	1
"2022-04"	1
"2022-05"	2
"2022-06"	2
"2022-07"	3

## Hasil