# Milestone 2 Management of Data

WQD7005 – DATA MINING

ZULKANAIN BIN HASAN

WQD180031
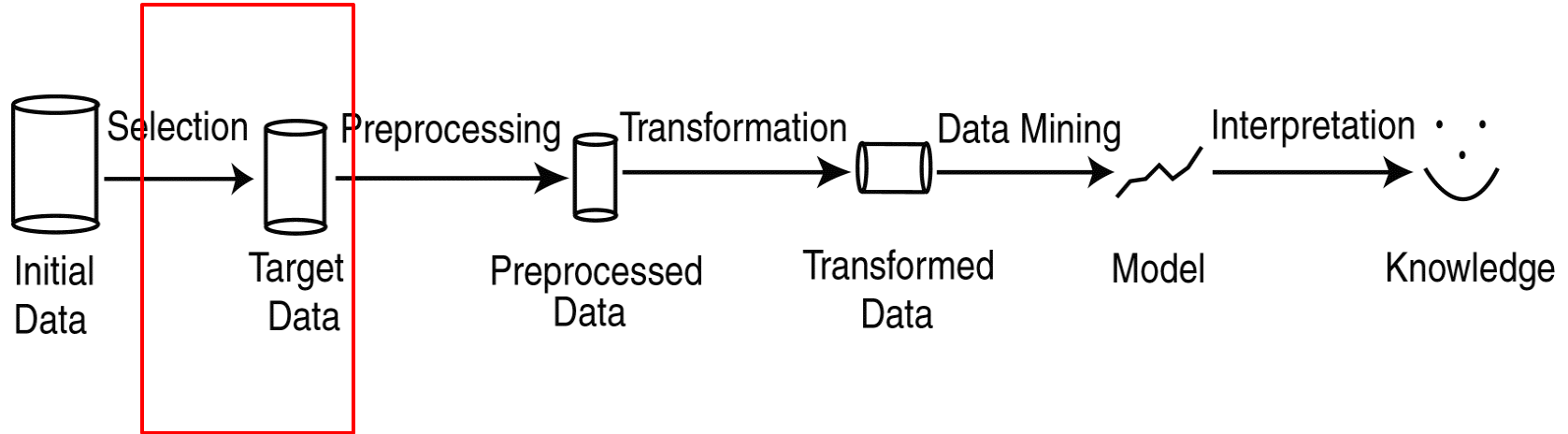
NUR ATHIRAH BT. RAZAK

WQD170072

**Milestone 2**



Initial Data → Selection → Target Data → Preprocessing → Preprocessed Data → Transformation → Transformed Data → Data Mining → Model → Interpretation → Knowledge

```
File  Edit  View  Search  Terminal  Help
(base) zulkanh@elitebook:~$ start-all.sh
This script is Deprecated. Instead use start-dfs.sh and start-yarn.sh
Starting namenodes on [localhost]
zulkanh@localhost's password:
localhost: starting namenode, logging to /home/zulkanh/bigdata/hadoop/logs/hadoo
p-zulkanh-namenode-elitebook.out
zulkanh@localhost's password:
localhost: starting datanode, logging to /home/zulkanh/bigdata/hadoop/logs/hadoo
p-zulkanh-datanode-elitebook.out
Starting secondary namenodes [0.0.0.0]
zulkanh@0.0.0.0's password:
0.0.0.0: starting secondarynamenode, logging to /home/zulkanh/bigdata/hadoop/log
s/hadoop-zulkanh-secondarynamenode-elitebook.out
starting yarn daemons
starting resourcemanager, logging to /home/zulkanh/bigdata/hadoop/logs/yarn-zulk
anh-resourcemanager-elitebook.out
zulkanh@localhost's password:
localhost: starting nodemanager, logging to /home/zulkanh/bigdata/hadoop/logs/ya
rn-zulkanh-nodemanager-elitebook.out
(base) zulkanh@elitebook:~$ jps
15889 DataNode
15729 NameNode
16083 SecondaryNameNode
16246 ResourceManager
```

```
File  Edit  View  Search  Terminal  Help
localhost: starting datanode, logging to /home/zulkanh/bigdata/hadoop/logs/hadoo
p-zulkanh-datanode-elitebook.out
Starting secondary namenodes [0.0.0.0]
zulkanh@0.0.0.0's password:
0.0.0.0: starting secondarynamenode, logging to /home/zulkanh/bigdata/hadoop/log
s/hadoop-zulkanh-secondarynamenode-elitebook.out
starting yarn daemons
starting resourcemanager, logging to /home/zulkanh/bigdata/hadoop/logs/yarn-zulk
anh-resourcemanager-elitebook.out
zulkanh@localhost's password:
localhost: starting nodemanager, logging to /home/zulkanh/bigdata/hadoop/logs/ya
rn-zulkanh-nodemanager-elitebook.out
(base) zulkanh@elitebook:~$ jps
15889 DataNode
15729 NameNode
16083 SecondaryNameNode
16246 ResourceManager
16872 Jps
16413 NodeManager
(base) zulkanh@elitebook:~$ hive

Logging initialized using configuration in jar:file:/home/zulkanh/bigdata/hive/l
ib/hive-common-1.2.2.jar!/hive-log4j.properties
hive>
```

File   Edit   View   Search   Terminal   Help

```
starting yarn daemons
starting resourcemanager, logging to /home/zulkanh/bigdata/hadoop/logs/yarn-zulk
anh-resourcemanager-elitebook.out
zulkanh@localhost's password:
localhost: starting nodemanager, logging to /home/zulkanh/bigdata/hadoop/logs/ya
rn-zulkanh-nodemanager-elitebook.out
(base) zulkanh@elitebook:~$ jps
15889 DataNode
15729 NameNode
16083 SecondaryNameNode
16246 ResourceManager
16872 Jps
16413 NodeManager
(base) zulkanh@elitebook:~$ hive

Logging initialized using configuration in jar:file:/home/zulkanh/bigdata/hive/l
ib/hive-common-1.2.2.jar!/hive-log4j.properties
hive> show databases;
OK
default
wqd7005
wqd7007
Time taken: 0.836 seconds, Fetched: 3 row(s)
hive>
```

# Import dataset to HDFS



Hdfs dfs -put /home/students/Downloads/klse2.csv /user/hdfs

```
student@student-
File  Edit  View  Search  Terminal  Help
3144 DataNode
2747 NameNode
5134 HRegionServer
student@student-VirtualBox:~$ hdfs dfs -put /home/student/Downloads/klse-main-ma
rket-stock-rhbtradesmart-topvolume-data.csv /user/hdfs
student@student-VirtualBox:~$ hdfs dfs -ls /user/hdfs
Found 2 items
-rw-r--r--   1 student supergroup    6398990 2019-04-30 00:14 /user/hdfs/batting
.csv
-rw-r--r--   1 student supergroup    23304590 2019-10-13 19:10 /user/hdfs/klse-ma
in-market-stock-rhbtradesmart-topvolume-data.csv
student@student-VirtualBox:~$ hdfs dfs -rm /user/hdfs/batting.csv
19/10/13 19:15:44 INFO fs.TrashPolicyDefault: Namenode trash configuration: Dele
tion interval = 0 minutes, Emptier interval = 0 minutes.
Deleted /user/hdfs/batting.csv
student@student-VirtualBox:~$ hdfs dfs -ls /user
Found 2 items
drwxr-xr-x   - student supergroup    0 2019-10-13 19:15 /user/hdfs
drwxr-xr-x   - student supergroup    0 2019-04-30 00:05 /user/hive
student@student-VirtualBox:~$ hdfs dfs -ls /user/hdfs
Found 1 items
-rw-r--r--   1 student supergroup    23304590 2019-10-13 19:10 /user/hdfs/klse-ma
in-market-stock-rhbtradesmart-topvolume-data.csv
student@student-VirtualBox:~$
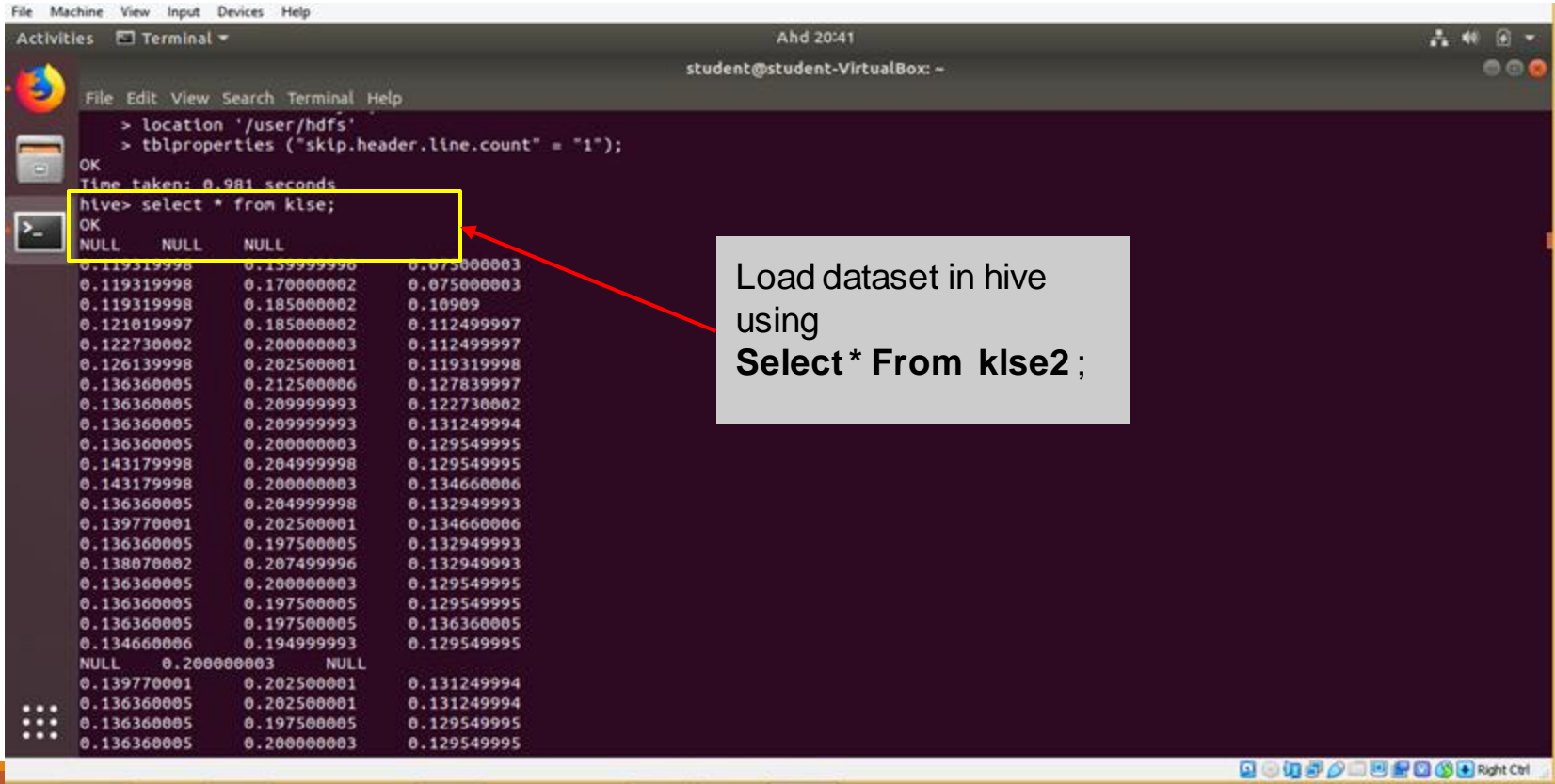```

Manage to import dataset to HDFS

Remove all irrelevent document

# Create **EXTERNAL TABLE** to Hive

Create Table is a statement used to create a table in hive

```
student@student-VirtualBox: ~
File  Edit  View  Search  Terminal  Help
java:62)
        at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAcces
sorImpl.java:43)
        at java.lang.reflect.Method.invoke(Method.java:498)
        at org.apache.hadoop.util.RunJar.run(RunJar.java:226)
        at org.apache.hadoop.util.RunJar.main(RunJar.java:141)
FAILED: ParseException line 1:46 cannot recognize input near ',' 'High' ',' in c
olumn type
hive> CREATE EXTERNAL TABLE IF NOT EXISTS klse (High double, High double, Low do
uble, Adj double, Adj double )
    > row format delimited
    > fields terminated by ","
    > location '/user/hdfs'
    > tblproperties ("skip.header.line.count" = "1");
FAILED: SemanticException [Error 10036]: Duplicate column name: high
hive> CREATE EXTERNAL TABLE IF NOT EXISTS klse (High double, Low double, Adj dou
ble)
    > row format delimited
    > fields terminated by ","
    > location '/user/hdfs'
    > tblproperties ("skip.header.line.count" = "1");
OK
Time taken: 0.981 seconds
hive>
```

# LOAD DATASET