# Hadoop in Windows Server using PolyBase

## Requirements

Required Software: (Hadoop/ Hadoop bin files/Hadoop Configuration files/Java)

https://drive.google.com/open?id=1e6pYQxZrr1JQaAR5ywOWNxkEYjQKKpBL

SQL Server 2017 Developer Edition https://www.microsoft.com/en-us/sql-server/sql-server-downloads

SQL Server Edition **Standard, Web, Express with Advanced Services, Express** does not support head node they require head node

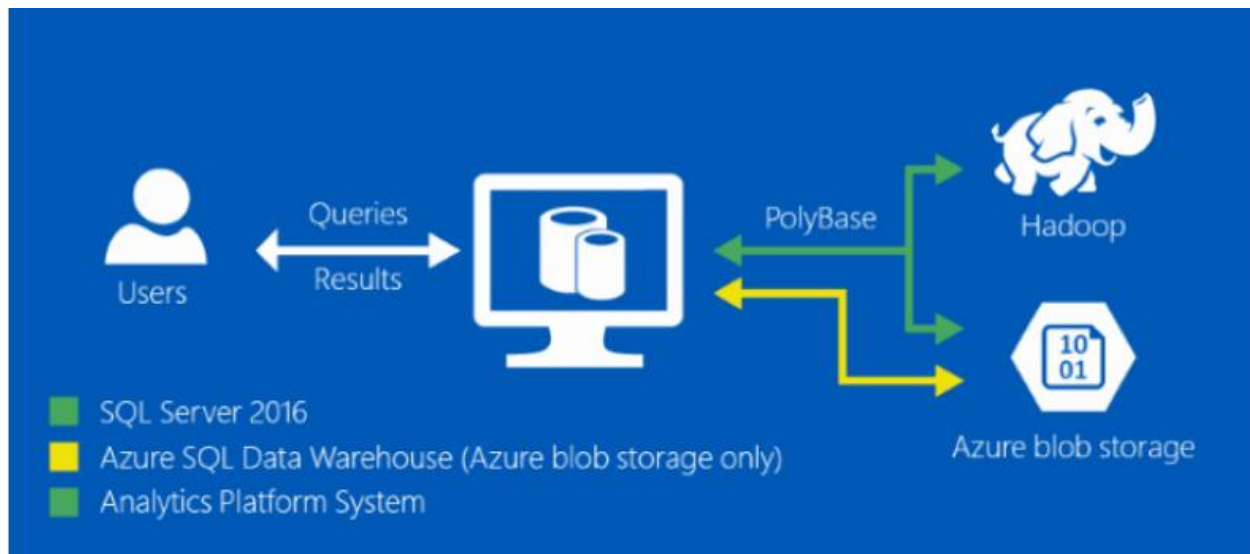**Enterprise** Edition only can be used for head node and other can be Scale out with multiple compute nodes.

Before we start Some of the difference between Hadoop and SQL

| HADOOP | SQL |
|---|---|
| **Schema on Read Approach** | **Schema on Write Approach** |
| When we write the data on Hadoop, has distributed file system so brings all the files without any rules. Then when we dictate the data we apply rules to the code that reads the data rather of preconfiguring the data ahead of time. | In Schema on Write while coping the data from database A to B the data is first checked before inserting to the database (checked the datatypes of the data and the data size) if not meet the requirement reject the data |
| **Data Storage Process Hadoop** | **Data Storage Process in SQL** |
| Data is stored in Compressed file of either text or others datatype. In the moment files or data stored are replicated in multiple nodes in Hadoop distributed filing system | Data is stored in a Logical form in Related Table and Columns |

## Installation of Java

Before installing Hadoop File Distributed System (HDFS) and **PolyBase** feature in SQL Server in 2016 or later requires Java. We need to first install Java in our system.

**PolyBase** enables your SQL Server 2016 and later instance to process Transact-SQL queries that read data from Hadoop. The same query can also access relational tables in your SQL Server. **PolyBase** enables the same query to also join the data from Hadoop and SQL Server. **PolyBase** is a mediator between SQL Server and Hadoop.

## Set up

1. Check either Java 1.8.0 is already installed on your system or not, use **"Java -version"** to check.
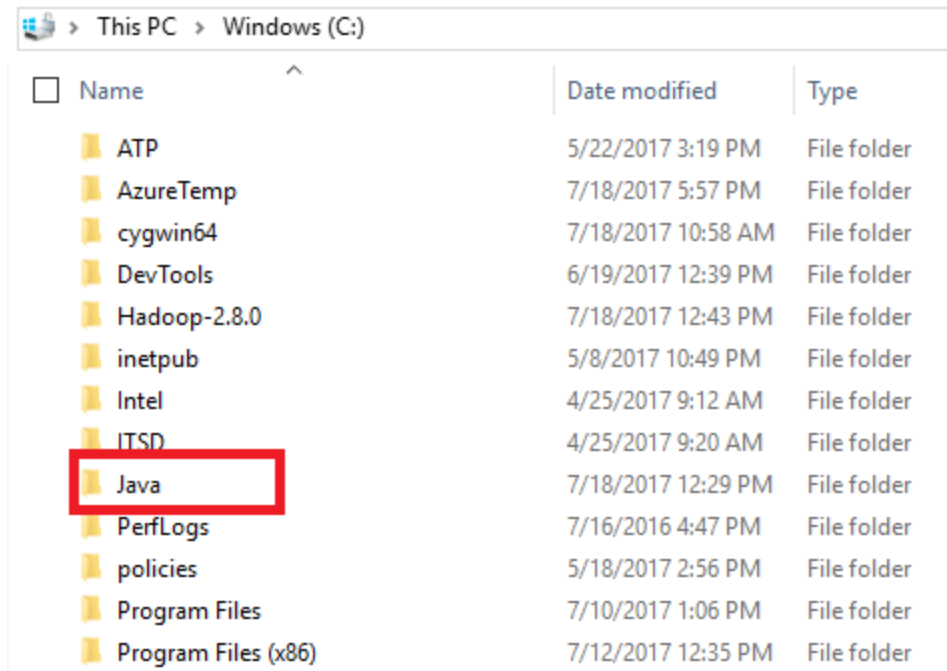
2. If Java is not installed on your system then first install java under **"C:\JAVA"**

| Name | Date modified | Type |
|------|---------------|------|
| ATP | 5/22/2017 3:19 PM | File folder |
| AzureTemp | 7/18/2017 5:57 PM | File folder |
| cygwin64 | 7/18/2017 10:58 AM | File folder |
| DevTools | 6/19/2017 12:39 PM | File folder |
| Hadoop-2.8.0 | 7/18/2017 12:43 PM | File folder |
| inetpub | 5/8/2017 10:49 PM | File folder |
| Intel | 4/25/2017 9:12 AM | File folder |
| ITSD | 4/25/2017 9:20 AM | File folder |
| Java | 7/18/2017 12:29 PM | File folder |
| PerfLogs | 7/16/2016 4:47 PM | File folder |
| policies | 5/18/2017 2:56 PM | File folder |
| Program Files | 7/10/2017 1:06 PM | File folder |
| Program Files (x86) | 7/12/2017 12:35 PM | File folder |

3. Extract file Hadoop 2.8.0.tar.gz or Hadoop-2.8.0.zip and place under **"C:\Hadoop-2.8.0"**.

| Name | Date modified | Type |
|------|---------------|------|
| ATP | 5/22/2017 3:19 PM | File folder |
| AzureTemp | 7/18/2017 5:57 PM | File folder |
| cygwin64 | 7/18/2017 10:58 AM | File folder |
| DevTools | 6/19/2017 12:39 PM | File folder |
| Hadoop-2.8.0 | 7/18/2017 12:43 PM | File folder |
| inetpub | 5/8/2017 10:49 PM | File folder |
| Intel | 4/25/2017 9:12 AM | File folder |
| ITSD | 4/25/2017 9:20 AM | File folder |
| Java | 7/18/2017 12:29 PM | File folder |
| PerfLogs | 7/16/2016 4:47 PM | File folder |
| policies | 5/18/2017 2:56 PM | File folder |
| Program Files | 7/10/2017 1:06 PM | File folder |
| Program Files (x86) | 7/12/2017 12:35 PM | File folder |

4. Set the path HADOOP_HOME Environment variable on windows 10(see Step 1,2,3 and 4 below).
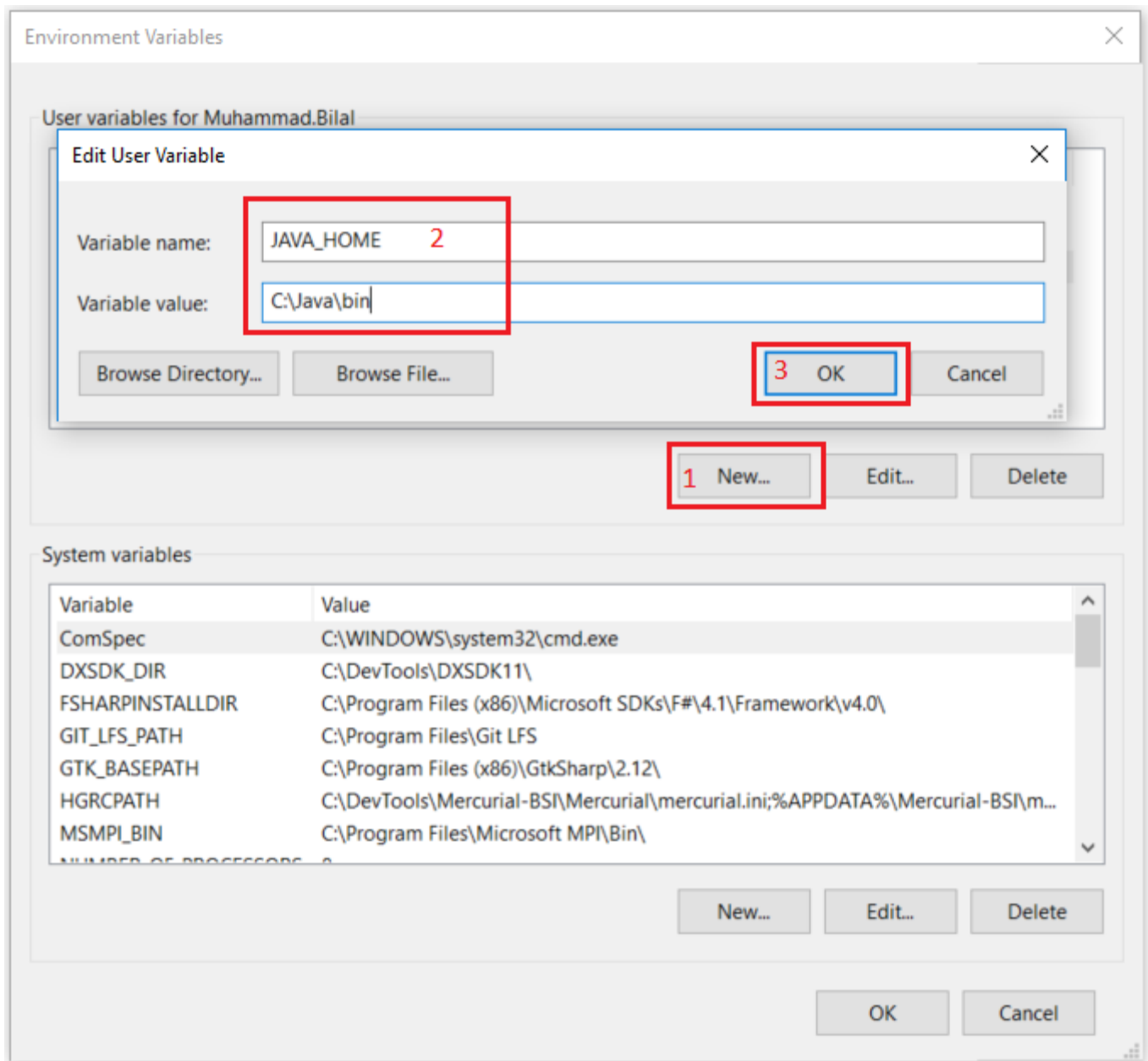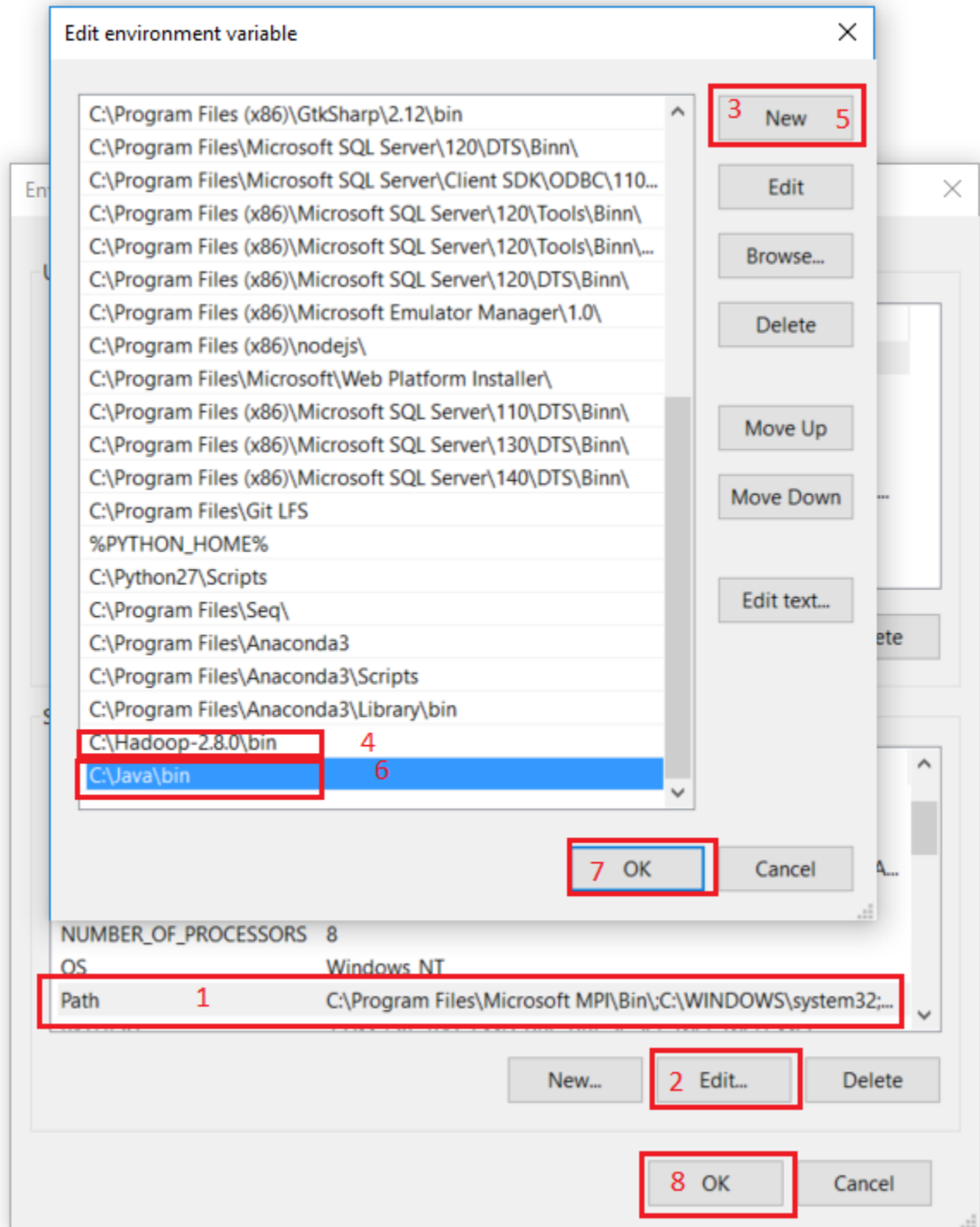
5. Set the path JAVA_HOME Environment variable on windows 10(see Step 1,2,3 and 4 below).
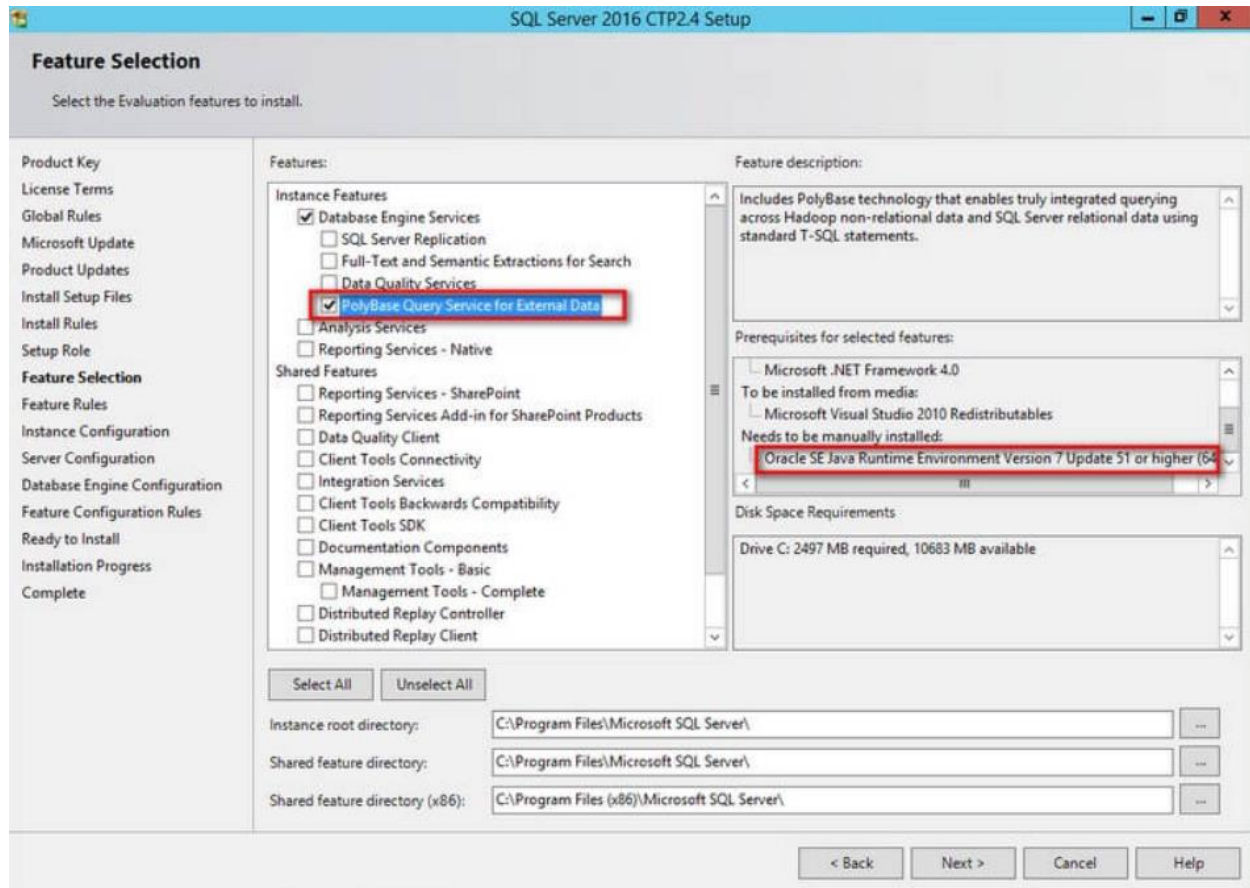
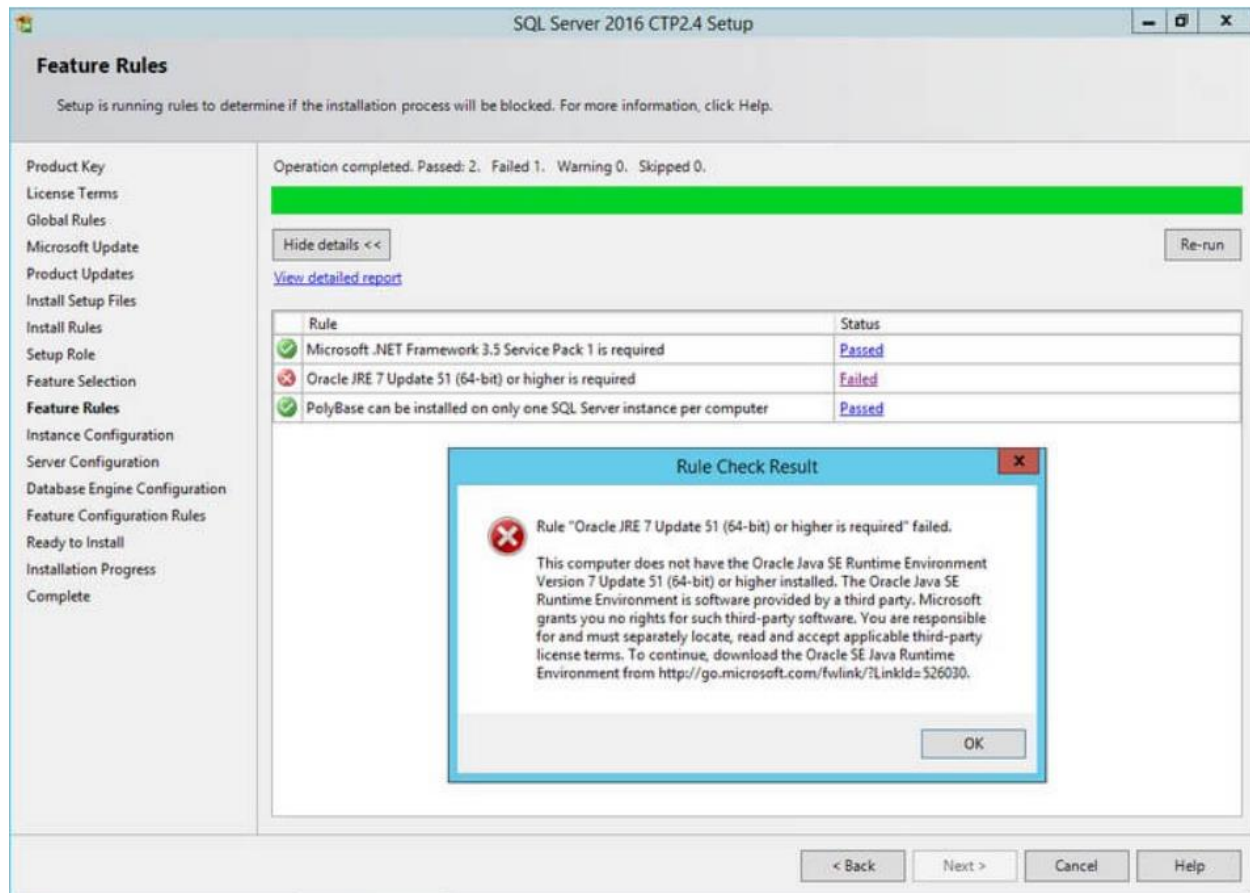6. Next, we set the Hadoop bin directory path and JAVA bin directory path.

Install **Polybase** SQL Server 2016 or Later

Since **Polybase** is now part of SQL Server, we can use the SQL Server 2016 installation media to do the installation. And because it was designed to interact with Hadoop, we will need to install the Oracle Java SE Runtime Environment (JRE) 7.51 (x64) or higher prior to running the SQL Server 2016 installation media.



If You don't have JRE installed then the installation will fail and screen as below appears.

Install JRE and set the **Environment Variable** after installation which is mention in **<u>Set Up</u>** steps above.

And Continue the installation of Polybase and completed successfully.

## Configuration

1. Edit file **C:/Hadoop-2.8.0/etc/hadoop/core-site.xml**, paste below xml paragraph and save this file.

```
<configuration>
   <property>
       <name>fs.defaultFS</name>
       <value>hdfs://localhost:9000</value>
   </property>
</configuration>
```

2. Rename "mapred-site.xml.template" to "mapred-site.xml" and edit this file **C:/Hadoop-2.8.0/etc/hadoop/mapred-site.xml**, paste below xml paragraph and save this file.

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

3. Create folder **"data"** under **"C:\Hadoop-2.8.0"**

- Create folder **"datanode"** under **"C:\Hadoop-2.8.0\data"**
- Create folder **"namenode"** under **"C:\Hadoop-2.8.0\data"**

4. Edit file **C:\Hadoop-2.8.0/etc/hadoop/hdfs-site.xml**, paste below xml paragraph and save this file.

```
<configuration>
   <property>
       <name>dfs.replication</name>
       <value>1</value>
   </property>
   <property>
       <name>dfs.namenode.name.dir</name>
       <value>C:\hadoop-2.8.0\data\namenode</value>
   </property>
   <property>
       <name>dfs.datanode.data.dir</name>
       <value>C:\hadoop-2.8.0\data\datanode</value>
   </property>
</configuration>
```

5. Edit file **C:/Hadoop-2.8.0/etc/hadoop/yarn-site.xml**, paste below xml paragraph and save this file.

```
<configuration>
   <property>
       <name>yarn.nodemanager.aux-services</name>
       <value>mapreduce_shuffle</value>
   </property>
   <property>
       <name>yarn.nodemanager.auxservices.mapreduce.shuffle.class</name>
       <value>org.apache.hadoop.mapred.ShuffleHandler</value>
   </property>
</configuration>
```

6. Edit file **C:/Hadoop-2.8.0/etc/hadoop/hadoop-env.cmd** by closing the command line**"JAVA_HOME=%JAVA_HOME%"** instead of set **"JAVA_HOME=C:\Java"** (On C:\java this is path to file jdk.18.0)

```
@rem The java implementation to use.  Required.
@rem set JAVA_HOME=%JAVA_HOME%
set JAVA_HOME=C:\Java\jdk1.8.0_181
@rem The jsvc implementation to use. Jsvc is required to run secure datanodes.
@rem set JSVC_HOME=%JSVC_HOME%
```

**Hadoop Configuration**

1.  Download                file                Hadoop                Configuration.zip
    (Link: https://drive.google.com/open?id=1e6pYQxZrr1JQaAR5ywOWNxkEYjQKKp
    BL)
2.  Delete file bin on C:\Hadoop-2.8.0\bin, replaced by file bin on file just download
    (from Hadoop Configuration.zip).
3.  Open cmd and typing command **"hdfs namenode –format"**. You will see

## Testing

1.  Open cmd and change directory to "C:\Hadoop-2.8.0\sbin" and type **"start-all.cmd"** to
    start
    apache.



2.  Make sure these apps are running

- Hadoop Namenode
- Hadoop datanode
- YARN Resource Manager
- YARN Node Manager

Open: http://localhost:8088

Hadoop In Windows

Open: http://localhost:50070



## Create Directory in Hadoop



## Import file to MMESDataFiles directory



Prepared By: Suman Pantha

Simple procedure to configure database for **PolyBase**, Hadoop Connectivity,
And Create External Data Source, Create External File Format and Create External Table
also create Statistics on external table for query optimization.

```sql
CREATE DATABASE PolybaseDB
use PolybaseDB

SELECT SERVERPROPERTY ('IsPolybaseInstalled') AS IsPolybaseInstalled;

--Prestep: Configuring Hadoop flavor
exec sp_configure 'hadoop connectivity',7
Reconfigure

exec sp_configure 'allow polybase export',1
Reconfigure

--(a) Creating external data source --Hadoop  HDP Cluster

--DROP EXTERNAL DATA SOURCE [HadoopCluster]
CREATE EXTERNAL DATA SOURCE HadoopCluster
WITH (TYPE = Hadoop,
     LOCATION = N'hdfs://localhost:9000')

select * from sys.external_data_sources;

--(b) Creating external file formates -delimited text

--DROP EXTERNAL FILE FORMAT [TextFile]
CREATE EXTERNAL FILE FORMAT TextFile
WITH (FORMAT_TYPE = DelimitedText,
     FORMAT_OPTIONS (FIELD_TERMINATOR = N',',
     USE_TYPE_DEFAULT = True));


select * from sys.external_file_formats;
```

```sql
-- (c) Creating external tables refering to data in external Hadoop Cluster

--DROP EXTERNAL TABLE [dbo].[SensorDataHDP]
CREATE EXTERNAL TABLE [dbo].[SensorDataHDP]
(
    Id int NOT NULL,
    BeneAccountCreditId int NOT NULL,
    BeneficiaryId int NULL,
    SendAmount varchar(50) NOT NULL,
    TypeofTraxId int NOT NULL
)
WITH (LOCATION = '/MMESDataFiles/',
        DATA_SOURCE = HadoopCluster,
        FILE_FORMAT = TextFile,
        REJECT_TYPE = Value,
        REJECT_VALUE = 0
);

SELECT * FROM SensorDataHDP
INSERT INTO SensorDataHDP
SELECT
        [Id]
        ,[BeneAccountCreditId]
        ,[BeneficiaryId]
        ,[SendAmount]
        ,[TypeofTraxId]
FROM [PolybaseDB].[dbo].[EremitData]
where Id < 100

SELECT * FROM SensorDataHDP

CREATE STATISTICS STAT_SensorDataHDPSendAmount on SensorDataHDP(SendAmount)

SELECT
        hdp.*
FROM SensorDataHDP hdp
INNER JOIN [EremitData] ed on hdp.Id = ed.Id
WHERE hdp.SendAmount between 100 and 200
```

```sql
SELECT
        hdp.*
FROM SensorDataHDP hdp
INNER JOIN [EremitData] ed on hdp.Id = ed.Id
WHERE hdp.SendAmount between 100 and 200
```
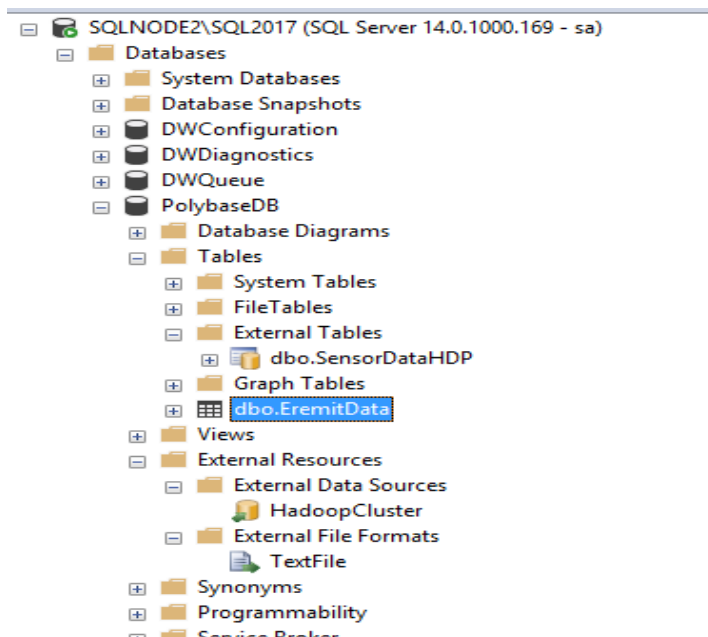
150 %

Results    Messages

|   | Id | BeneAccountCreditId | BeneficiaryId | SendAmount | TypeofTraxId |
|---|----|---------------------|---------------|------------|--------------|
| 1 | 3  | 4035                | 4035          | 125        | 1            |
| 2 | 44 | 827                 | 827           | 100        | 1            |
| 3 | 50 | 827                 | 827           | 100        | 1            |
| 4 | 82 | 536                 | 536           | 190        | 1            |

Structure of External Data Source, External File Format and External Table in SSMS.



File Resides in HDFS