# EXP 3: Map Reduce program to process a weather dataset.

**AIM:** To implement MapReduce program to process a weather dataset.

**Procedure:**

Step 1: Create Data File
1. Log in with your Hadoop user.
2. Download the weather dataset and save it locally, for example, as `dataset.txt`.

Step 2: Mapper Logic
1. Create a file named `mapper.py`.
2. Implement the mapper logic:
   - The mapper processes each line of the dataset.
   - Extract the month and daily maximum temperature from each record and output them.

Step 3: Reducer Logic
1. Create a file named `reducer.py`.
2. Implement the reducer logic:
   - The reducer receives the output from the mapper, which contains the month and temperature data.
   - Aggregate the daily maximum temperatures by month and find the highest temperature for each month.

Step 4: Prepare Hadoop Environment
1. Start the necessary Hadoop services (daemons).
2. Create a directory in HDFS for storing the weather dataset.

Step 5: Upload Data to HDFS
1. Upload the dataset file to the HDFS directory created in the previous step.

Step 6: Make Python Files Executable
1. Provide executable permissions to the `mapper.py` and `reducer.py` files.

Step 7: Run the MapReduce Program Using Hadoop Streaming
1. Download the Hadoop Streaming JAR file if not already available.
2. Run the MapReduce job by specifying the input data (dataset), the output directory, and the mapper and reducer Python files using Hadoop Streaming.

Step 8: Check Output
1. View the results of the MapReduce job in the HDFS output directory.
2. If needed, you can copy the results to your local machine for further analysis.

**Commands:**

C:\hadoop\sbin>        start-all.cmd

C:\hadoop\sbin>        jps

C:\hadoop\sbin>        cd /
C:\>                   cd hadoop

C:\hadoop>                hadoop fs -mkdir /user/

C:\hadoop>                hadoop fs -put C:/DataAnalytics/sample_weather.csv /input

C:\hadoop>                hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.3.6.jar -input /user/sample_weather.csv -output /user/output-data -mapper "C:\Users\hp\Documents\weather\mapper.py"-reducer "C:\Users\hp\Documents\weather\reducer.py"

hadoop fs -cat /user/jayas/output/part-00000

**OUTPUT:**