**EXP NO: 3          MAP REDUCE PROGRAM TO PROCESS A WEATHER DATASET**
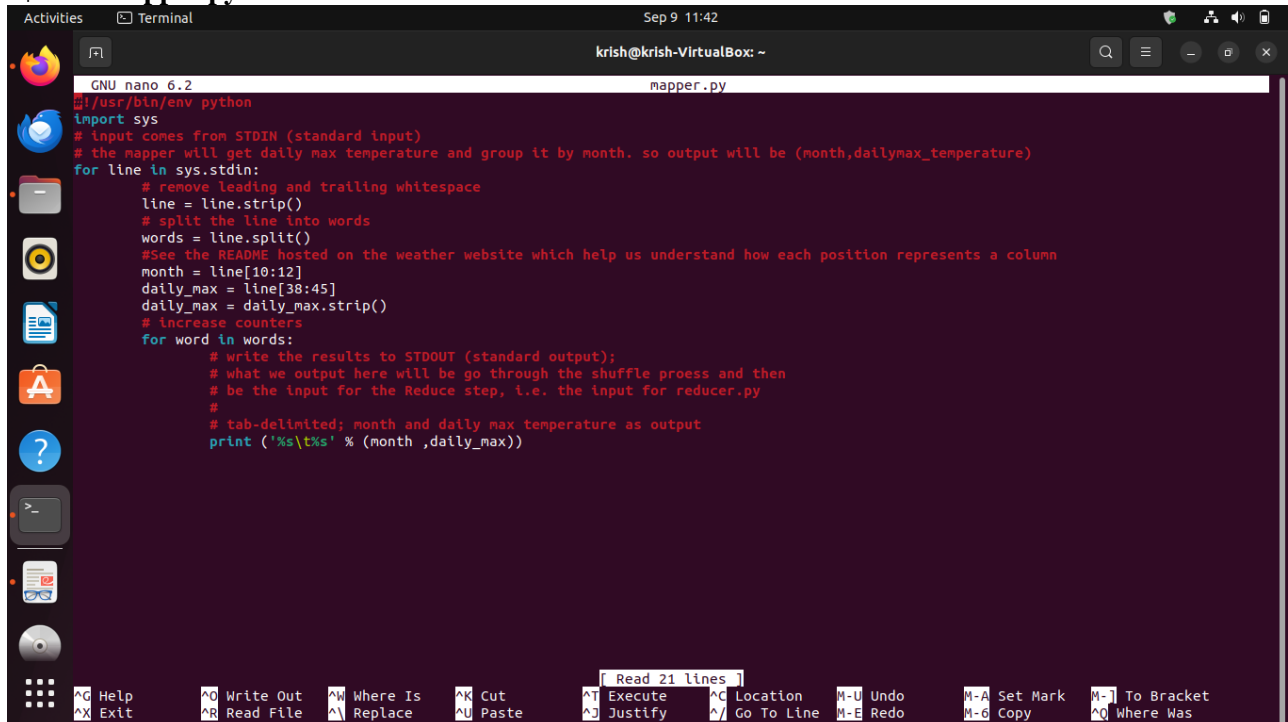
**$cd DA-Lab**
**$mkdir exp3**
**$cd exp3**
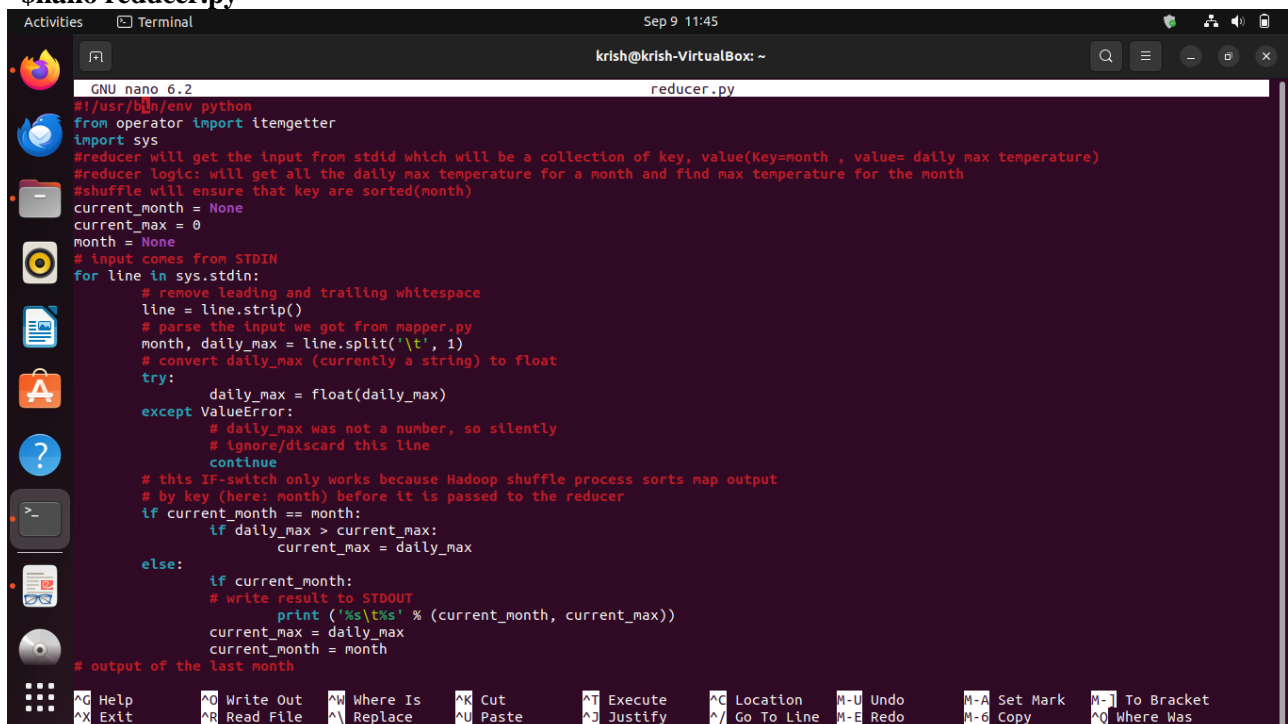
**$nano dataset.txt**

**$nano mapper.py**



**$nano reducer.py**

**$start-all.sh**

**$ jps**



**$hdfs dfs -mkdir /exp3**

**$hdfs dfs -copyFromLocal ~/DA-Lab/exp3/dataset.txt /exp2**



**$chmod 777 mapper.py reducer.py**

**$hadoop jar $HADOOP_STREAMING -input /exp3/dataset.txt -output /exp3/output -mapper ~/DA-Lab/exp3/mapper.py -reducer ~/DA-Lab/exp3/reducer.py**

**$hdfs dfs -cat /exp3/output/***