**Experiment 10: Week 11:**

**10. Install and Run Pig then write Pig Latin scripts to sort, group, join, project, and filter your data.**

**ALGORITHM**

**STEPS FOR INSTALLING APACHE PIG**

**1) Extract the pig-0.15.0.tar.gz and move to home directory**

**2) Set the environment of PIG in bashrc file.**

**3) Pig can run in two modes**

   **Local Mode and Hadoop Mode**
   **Pig –x local and pig**

**4) Grunt Shell**

   **Grunt >**

**5) LOADING Data into Grunt Shell**

   **DATA = LOAD <CLASSPATH> USING PigStorage (DELIMITER) as (ATTRIBUTE : DataType1, ATTRIBUTE : DataType2…..)**

**6) Describe Data**

   **Describe DATA;**

**7) DUMP Data**

   **Dump DATA;**

**8) FILTER Data**

   **FDATA = FILTER DATA by ATTRIBUTE = VALUE;**

**9) GROUP Data**

   **GDATA = GROUP DATA by ATTRIBUTE;**

**10) Iterating Data**

   **FOR_DATA = FOREACH DATA GENERATE GROUP AS GROUP_FUN, ATTRIBUTE = <VALUE>**

**11) Sorting Data**

**SORT_DATA = ORDER DATA BY ATTRIBUTE WITH CONDITION;**

**12) LIMIT Data**

**LIMIT_DATA = LIMIT DATA COUNT;**

**13) JOIN Data**

**JOIN DATA1 BY (ATTRIBUTE1,ATTRIBUTE2….) , DATA2 BY (ATTRIBUTE3,ATTRIBUTE….N)**

**INPUT:**
    **Input as Website Click Count Data**

**sOUTPUT:**