

download

1. `wget https://archive.apache.org/dist/pig/pig-0.16.0/pig-0.16.0.tar.gz`
2. `tar xvf pig-0.16.0.tar.gz`

modify profile

1. `vim ~/.bashrc`

1. `export PIG_HOME=${HADOOP_BASE_PATH}/pig-0.16.0`
2. `export PATH=$PATH:${PIG_HOME}/bin`
- 3.
4. `# hadoop config setting`
5. `export PIG_CLASSPATH=${HADOOP_BASE_PATH}/hadoop-2.7.3/etc/hadoop`

start as localhost session

1. `$PIG_HOME/bin/pig -x local`
2. `# $HADOOP_HOME/bin/hdfs dfs -ls /input`
3. `grunt> ls /input`

start as service

1. `$PIG_HOME/bin/pig`
2. `# $HADOOP_HOME/bin/hdfs dfs -ls /input`
3. `grunt> ls /input`

WordCount

1. `// 加载输入文件，并按行分隔`
2. `grunt> a = LOAD '/input/immortals.txt' as (line:chararray);`
3. `// 将每行分割成单词`
4. `grunt> words = FOREACH a GENERATE flatten(TOKENIZE(line)) as w;`
5. `// 按单词分组`
6. `grunt> g = GROUP words by w;`
7. `// 单词记数`
8. `grunt> wordcount = FOREACH g GENERATE group,COUNT(words);`

WordCount2

带词频倒排序

```
1. a = LOAD '/input/immortals.txt' as (line:chararray);
2. words = FOREACH a GENERATE flatten(TOKENIZE(line)) as w;
3. g = GROUP words by w;
4. // 给单词数所在列加一个别名count
5. wordcount = FOREACH g GENERATE group,COUNT(words) as count;
6. // 将结果列交换，将变成{count, word}这种结构
7. r = foreach wordcount generate count,group;
8. // 按count分组
9. g2 = group r by count;
10. // 去掉无用的列
11. x = foreach g2 generate group,r.group;
12. // 按count倒排
13. y = order x by group desc;
```