

ASSIGNMENT 4.2

Dataset used:

television.txt

```
2017-12-02 01:32:03,624 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
(Onida,Lucid,18,Uttar Pradesh,232401,16200)
(Akai,Decent,16,Kerala,922401,12200)
(Lava,Attention,20,Assam,454601,24200)
(Zen,Super,14,Maharashtra,619082,9200)
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
(Onida,Lucid,18,Uttar Pradesh,232401,16200)
(Onida,Decent,14,Uttar Pradesh,232401,16200)
(Lava,Attention,20,Assam,454601,24200)
(Zen,Super,14,Maharashtra,619082,9200)
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
(Samsung,Decent,16,Kerala,922401,12200)
(Lava,Attention,20,Assam,454601,24200)
(Samsung,Super,14,Maharashtra,619082,9200)
(Samsung,Super,14,Maharashtra,619082,9200)
(Samsung,Super,14,Maharashtra,619082,9200)
grunt> DESCRIBE filtered_tv;
filtered_tv: {company_name: chararray,product_name: chararray,size_in_inches: int,state: chararray,pincode: int,price: int}
grunt>
```

Data in the text file had some invalid rows. So I have used filter operation to get rid of them. **filtered_tv** is the resulted relation/schema.

Apache Pig Commands:

1. **CONCAT()**: It is used to two or more expressions or columns of same type.

Example query: `concat_company_and_product = FOREACH filtered_tv GENERATE CONCAT (company_name, ' ', product_name) AS (company_with_product);`

Output: `DUMP concat_company_and_product;`

```
2017-12-02 02:05:03,561 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-02 02:05:03,562 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Samsung Optima)
(Onida Lucid)
(Akai Decent)
(Lava Attention)
(Zen Super)
(Samsung Optima)
(Onida Lucid)
(Onida Decent)
(Lava Attention)
(Zen Super)
(Samsung Optima)
(Samsung Decent)
(Lava Attention)
(Samsung Super)
(Samsung Super)
(Samsung Super)
grunt> DESCRIBE concat_company_and_product;
concat_company_and_product: {company_with_product: chararray}
grunt>
```

2. **TOKENIZE():** It is used to split a string having many words in a tuple. It returns a bag of token tuples of type same as input.

Example: tokenize_company_and_product = FOREACH concat_company_and_product
GENERATE TOKENIZE(company_with_product);

Output: DUMP tokenize_company_and_product;

```
2017-12-09 13:32:49,786 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 13:32:49,786 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
((Samsung),(Optima))
((Onida),(Lucid))
((Akai),(Decent))
((Lava),(Attention))
((Zen),(Super))
((Samsung),(Optima))
((Onida),(Lucid))
((Onida),(Decent))
((Lava),(Attention))
((Zen),(Super))
((Samsung),(Optima))
((Samsung),(Decent))
((Lava),(Attention))
((Samsung),(Super))
((Samsung),(Super))
((Samsung),(Super))
grunt> describe tokenize_company_and_product;
tokenize_company_and_product: {bag_of_tokenTuples_from_company_with_product: {tuple_of_tokens: (token: chararray)}}
grunt>
```

3. **SUM():** As the name implies, this function is used to get sum of all values of a specific column in a single column bag. We need to perform GROUP operation on a column or on all columns, after which we can apply SUM() to get global sum among the groups.

Example query to group data by company_name :

grouped_tv = GROUP filtered_tv BY company_name;

Output: DUMP grouped_tv;

```
2017-12-09 15:39:38,509 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 15:39:38,510 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen,{(Zen,Super,14,Maharashtra,619082,9200),(Zen,Super,14,Maharashtra,619082,9200)})
(Akai,{(Akai,Decent,16,Kerala,922401,12200)})
(Lava,{(Lava,Attention,20,Assam,454601,24200),(Lava,Attention,20,Assam,454601,24200),(Lava,Attention,20,Assam,454601,24200)})
(Onida,{(Onida,Lucid,18,Uttar Pradesh,232401,16200),(Onida,Lucid,18,Uttar Pradesh,232401,16200),(Onida,Decent,14,Uttar Pradesh,232401,16200)})
(Samsung,{(Samsung,Super,14,Maharashtra,619082,9200),(Samsung,Super,14,Maharashtra,619082,9200),(Samsung,Super,14,Maharashtra,619082,9200),(Samsung,Decent,16,Kerala,922401,12200),(Samsung,Optima,14,Madhy Pradesh,132401,14200),(Samsung,Optima,14,Madhy Pradesh,132401,14200)})
grunt>
```

SUM query on each company group based on price:

```
sum_of_prices = FOREACH grouped_tv GENERATE group, SUM(filtered_tv.price);
```

Output: DUMP sum_of_prices;

```
2017-12-09 15:41:17,915 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2017-12-09 15:41:17,923 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2017-12-09 15:41:17,926 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2017-12-09 15:41:17,933 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2017-12-09 15:41:17,934 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-09 15:41:17,934 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2017-12-09 15:41:17,934 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-09 15:41:17,934 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2017-12-09 15:41:17,967 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 15:41:17,967 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen,18400)
(Akai,12200)
(Lava,72600)
(Onida,48600)
(Samsung,82400)
grunt>
```

4. MIN(): This function is used to get minimum value for a particular column in a single column bag. Like SUM(), this function is also applied on grouped data to find lowest value of a column from each group.

Example query:

```
min_of_prices = FOREACH grouped_tv GENERATE group, MIN(filtered_tv.price);
```

Output: DUMP min_of_prices;

```
2017-12-09 16:16:17,007 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-09 16:16:17,012 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2017-12-09 16:16:17,012 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-09 16:16:17,012 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2017-12-09 16:16:17,053 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 16:16:17,053 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen,9200)
(Akai,12200)
(Lava,24200)
(Onida,16200)
(Samsung,9200)
grunt> DESCRIBE min_of_prices;
min_of_prices: {group: chararray,int}
grunt>
```

5. **MAX():** This function is used to get maximum value for a particular column in a single column bag. This function is applied on grouped data to find highest value of a column from each group.

Example query:

```
max_of_prices = FOREACH grouped_tv GENERATE group, MAX(filtered_tv.price);
```

Output: DUMP max_of_prices;

```
2017-12-09 16:20:04,566 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated.
Instead, use dfs.bytes-per-checksum
2017-12-09 16:20:04,570 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2017-12-09 16:20:04,571 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-09 16:20:04,571 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2017-12-09 16:20:04,581 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 16:20:04,581 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen,9200)
(Akai,12200)
(Lava,24200)
(Onida,16200)
(Samsung,14200)
grunt> DESCRIBE max_of_prices;
max_of_prices: {group: chararray,int}
grunt>
```

6. **LIMIT:** This operator is used to get limited number of tuples from a relation. The resulting relation will have same schema as that of source relation.

Example query:

```
limited_tuples_from_tv = LIMIT tv 10;
```

Output: DUMP limited_tuples_from_tv;

```
2017-12-09 16:38:34,817 [main] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_0001_m_000001_1' to file:/tmp/temp1952696468/tmp-1940080738/_temporary/0/task_0001_m_000001
2017-12-09 16:38:34,831 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2017-12-09 16:38:34,855 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 16:38:34,855 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
(Onida,Lucid,18,Uttar Pradesh,232401,16200)
(Akai,Decent,16,Kerala,922401,12200)
(Lava,Attention,20,Assam,454601,24200)
(Zen,Super,14,Maharashtra,619082,9200)
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
(Onida,Lucid,18,Uttar Pradesh,232401,16200)
(Onida,Decent,14,Uttar Pradesh,232401,16200)
(Onida,NA,16,Kerala,922401,12200)
(Lava,Attention,20,Assam,454601,24200)
grunt> DESCRIBE limited_tuples_from_tv;
limited_tuples_from_tv: {company_name: chararray,product_name: chararray,size_in_inches: int,state: chararray,pincode: int,price: int}
grunt>
```

7. STORE: This operator is used to store data of a relation into local file system or to HDFS. Output directory should not exist as STORE command creates one and places the output file into it.

Example query to store data into HDFS:

```
STORE limited_tuples_from_tv INTO 'hdfs://localhost:9000/pig_data/' USING PigStorage(',');
```

Output:

```
[acadgild@localhost ~]$ hadoop fs -ls /
17/12/09 17:03:26 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
Found 5 items
drwxr-xr-x - acadgild supergroup          0 2015-11-09 19:21 /hbasestorage
drwxr-xr-x - acadgild supergroup          0 2017-12-09 17:01 /pig_data
drwxrwxr-x - acadgild supergroup          0 2015-11-05 13:46 /tmp
drwxr-xr-x - acadgild supergroup          0 2015-11-17 01:56 /user
drwxr-xr-x - acadgild supergroup          0 2015-11-05 12:56 /zookeeper
[acadgild@localhost ~]$ hadoop fs -ls /pig_data
17/12/09 17:04:12 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
Found 2 items
-rw-r--r--  3 acadgild supergroup          0 2017-12-09 17:01 /pig_data/_SUCCESS
-rw-r--r--  3 acadgild supergroup        397 2017-12-09 17:01 /pig_data/part-r-00000
[acadgild@localhost ~]$ hadoop fs -cat /pig_data/part-r-00000
17/12/09 17:04:35 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
Zen,Super,14,Maharashtra,619082,9200
Akai,Decent,16,Kerala,922401,12200
Lava,Attention,20,Assam,454601,24200
Lava,Attention,20,Assam,454601,24200
Onida,NA,16,Kerala,922401,12200
Onida,Lucid,18,Uttar Pradesh,232401,16200
Onida,Lucid,18,Uttar Pradesh,232401,16200
Onida,Decent,14,Uttar Pradesh,232401,16200
Samsung,Optima,14,Madhya Pradesh,132401,14200
Samsung,Optima,14,Madhya Pradesh,132401,14200
[acadgild@localhost ~]$ █
```

Example query to store data into local file system:

```
STORE limited_tuples_from_tv INTO 'file:///home/acadgild/pig_data' USING PigStorage(',');
```

Output:

```
[acadgild@localhost ~]$ pwd
/home/acadgild
[acadgild@localhost ~]$ ls
assignment_stuff  Downloads                hive-site.xml  pig_1512803760957.log  television.txt
derby.log         eclipse                  metastore_db   pig_1512805790632.log  Templates
Desktop           eclipse-jee-neon-M3-linux-gtk-x86_64.tar.gz  Music          pig_data               Videos
Documents         hadoop                  Pictures        Public                 workspace
[acadgild@localhost ~]$ cd pig_data/
[acadgild@localhost pig_data]$ ls
part-r-00000 _SUCCESS
[acadgild@localhost pig_data]$ cat part-r-00000
Zen,Super,14,Maharashtra,619082,9200
Akai,Decent,16,Kerala,922401,12200
Lava,Attention,20,Assam,454601,24200
Lava,Attention,20,Assam,454601,24200
Onida,NA,16,Kerala,922401,12200
Onida,Lucid,18,Uttar Pradesh,232401,16200
Onida,Lucid,18,Uttar Pradesh,232401,16200
Onida,Decent,14,Uttar Pradesh,232401,16200
Samsung,Optima,14,Madhya Pradesh,132401,14200
Samsung,Optima,14,Madhya Pradesh,132401,14200
[acadgild@localhost pig_data]$ █
```

8. **DISTINCT:** This operator returns unique tuples from a relation removing the redundant ones.

We can see from STORE query output that it has 3 redundant tuples.

Example query:

```
distinct_tuples_from_tv = DISTINCT limited_tuples_from_tv;
```

Output:

```
2017-12-09 17:37:28,337 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 17:37:28,337 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen,Super,14,Maharashtra,619082,9200)
(Akai,Decent,16,Kerala,922401,12200)
(Lava,Attention,20,Assam,454601,24200)
(Onida,NA,16,Kerala,922401,12200)
(Onida,Lucid,18,Uttar Pradesh,232401,16200)
(Onida,Decent,14,Uttar Pradesh,232401,16200)
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
grunt> DESCRIBE distinct_tuples_from_tv;
distinct_tuples_from_tv: {company_name: chararray,product_name: chararray,size_in_inches: int,state: chararray,pincode: int,price: int}
grunt>
```

9. **FLATTEN:** This operator changes the structure of tuples and bags in a relation by un-nesting them. FLATTEN replaces a bag by tuples inside that bag and a tuple by fields of a tuple.

Example query on a bag:

```
flattened_bag = FOREACH grouped_tv GENERATE FLATTEN(filtered_tv);
```

Output: DUMP flattened_bag;

```
2017-12-09 19:33:31,049 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 19:33:31,049 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen,Super,14,Maharashtra,619082,9200)
(Zen,Super,14,Maharashtra,619082,9200)
(Akai,Decent,16,Kerala,922401,12200)
(Lava,Attention,20,Assam,454601,24200)
(Lava,Attention,20,Assam,454601,24200)
(Lava,Attention,20,Assam,454601,24200)
(Onida,Lucid,18,Uttar Pradesh,232401,16200)
(Onida,Lucid,18,Uttar Pradesh,232401,16200)
(Onida,Decent,14,Uttar Pradesh,232401,16200)
(Samsung,Super,14,Maharashtra,619082,9200)
(Samsung,Super,14,Maharashtra,619082,9200)
(Samsung,Super,14,Maharashtra,619082,9200)
(Samsung,Decent,16,Kerala,922401,12200)
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
(Samsung,Optima,14,Madhya Pradesh,132401,14200)
grunt> DESCRIBE flattened_bag;
flattened_bag: {filtered_tv::company_name: chararray,filtered_tv::product_name: chararray,filtered_tv::size_in_inches: int,filtered_tv::state: chararray,filtered_tv::pincode: int,filtered_tv::price: int}
grunt>
```

Example query on a tuple:

```
flattened_tuple = FOREACH distinct_tuples_from_tv GENERATE FLATTEN($0);
```

Output: DUMP flattened_tuple;

```
2017-12-09 19:31:47,956 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 19:31:47,956 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen)
(Akai)
(Lava)
(Onida)
(Onida)
(Onida)
(Samsung)
grunt> DESCRIBE flattened_tuple;
flattened_tuple: {company_name: chararray}
grunt>
```

10. isEmpty(): This function is used to check if a bag or map is empty. If a relation has at least one empty bag or map, it will be stored to an assigned relation.

Using COGROUP operator to group two relations on different columns we can obtain a relation that contains some empty bags.

Example query for COGROUP:

```
cogroup_data = COGROUP limited_tuples_from_tv BY company_name,
distinct_tuples_from_tv by state;
```

Output: DUMP cogroup_data;

```
2017-12-09 19:51:56,513 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 19:51:56,513 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Zen,{(Zen,Super,14,Maharashtra,619082,9200)},{})
(Akai,{(Akai,Decent,16,Kerala,922401,12200)},{})
(Lava,{(Lava,Attention,20,Assam,454601,24200),(Lava,Attention,20,Assam,454601,24200)},{})
(Assam,{},{(Lava,Attention,20,Assam,454601,24200)})
(Onida,{(Onida,Decent,14,Uttar Pradesh,232401,16200),(Onida,Lucid,18,Uttar Pradesh,232401,16200),(Onida,Lucid,18,Uttar Pradesh,232401,16200),(Onida,NA,16,Kerala,922401,12200)},{})
(Kerala,{},{(Onida,NA,16,Kerala,922401,12200),(Akai,Decent,16,Kerala,922401,12200)})
(Samsung,{(Samsung,Optima,14,Madhya Pradesh,132401,14200),(Samsung,Optima,14,Madhya Pradesh,132401,14200)},{})
(Maharashtra,{},{(Zen,Super,14,Maharashtra,619082,9200)})
(Uttar Pradesh,{},{(Onida,Decent,14,Uttar Pradesh,232401,16200),(Onida,Lucid,18,Uttar Pradesh,232401,16200)})
(Madhya Pradesh,{},{(Samsung,Optima,14,Madhya Pradesh,132401,14200)})
grunt> DESCRIBE cogroup_data;
cogroup_data: {group: chararray,limited tuples_from_tv: {(company_name: chararray,product_name: chararray,size_in_inches: int,state: chararray,pincode: int,price: int)},distinct tuples_from_tv: {(company_name: chararray,product_name: chararray,size_in_inches: int,state: chararray,pincode: int,price: int)}}
grunt>
```


Example query for IsEmpty() taking cogroup_data as input relation:

isempty_result = FILTER cogroup_data BY IsEmpty(limited_tuples_from_tv);

Output: DUMP isempty_result;

```
2017-12-09 19:57:49,357 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-09 19:57:49,357 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Assam,{},{(Lava,Attention,20,Assam,454601,24200)})
(Kerala,{},{(Onida,NA,16,Kerala,922401,12200),(Akai,Decent,16,Kerala,922401,12200)})
(Maharashtra,{},{(Zen,Super,14,Maharashtra,619082,9200)})
(Uttar Pradesh,{},{(Onida,Decent,14,Uttar Pradesh,232401,16200),(Onida,Lucid,18,Uttar Pradesh,232401,16200)})
(Madhya Pradesh,{},{(Samsung,Optima,14,Madhya Pradesh,132401,14200)})
grunt> DESCRIBE isempty_result;
isempty_result: {group: chararray,limited_tuples_from_tv: {(company_name: chararray,product_name: chararray,size_in_inches: int,state: chararray,pincode: int,price: int)},distinct_tuples_from_tv: {(company_name: chararray,product_name: chararray,size_in_inches: int,state: chararray,pincode: int,price: int)}}
grunt> █
```

