BIG DATA LABORATORY

Course Code: ISL75 Credit:

0:0:1

Prerequisite: Nil Contact

Hours: 14L

Course Coordinator: Dr. Rajeshwari S B

PART-A

1. Compile all java files (driver.java mapper.java reducer.java)

javac -d . *.java

2. Set driver class in manifest

echo Main-Class: weather.driver > Manifest.txt

3. Create an executable jar file

jar cfm weather.jar Manifest.txt weather/*.class

- 4. input.txt is input file for Weather create Input File
- 5. Run the jar file

hadoop jar weather.jar input.txt output

6. To see the Output:

cat output/*

Write MapReduce programs for the following using Java:

1. To count the number of occurrences of each word in a given input text.

<Goodbye, 1> <Hadoop, 2> <Hello, 1>

2. To read N natural numbers and display the sum along with Odd and Even count.

```
ibmlab@ibmlab: ~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$ ls
driver.java input.txt Manifest.txt mapper.java oddeven oddeven.jar output reducer.java
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$ cat output/*
Sum of even Numbers
                       30
even Number count
Sum of odd Numbers
                       25
odd Number count
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/oddeven$
```

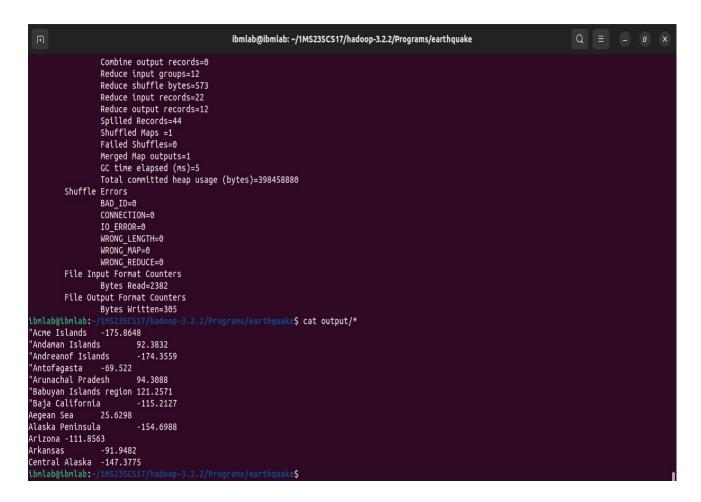
3. To analyse the given **Employee Data** and generate a statistics report with the total number of Female and Male Employees and their average Salary.

```
Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=4
                Total committed heap usage (bytes)=398458880
       Shuffle Errors
                BAD ID=0
                CONNECTION=0
               IO ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
       File Input Format Counters
                Bytes Read=8092
       File Output Format Counters
               Bytes Written=94
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/employee$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/employee$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/employee$
lbmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/employee$ cat output/*
 Average
               7117.073170731707
 Count 41.0
               6333.781194029851
 Average
M Count 67.0
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/employee$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/employee$
ibmlab@ibmlab:~/1MS23SCS17/hadoop-3.2.2/Programs/employee$
```

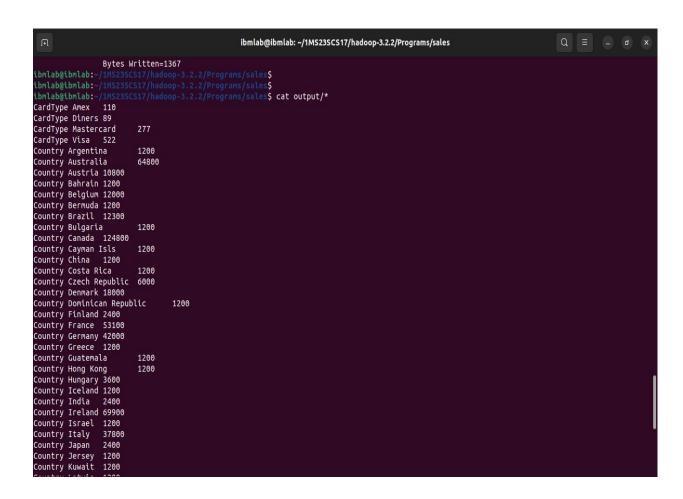
4. To analyse the **Titanic Ship Data** and find the average age of the people (both male and female) who died in the tragedy and also how many people are survived in each class.\

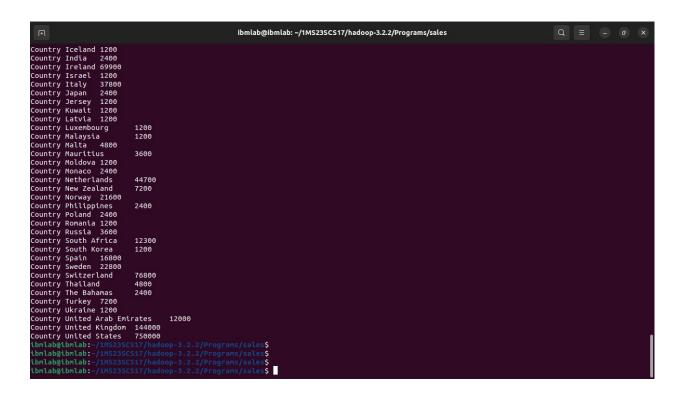
```
6b@iselab1-ThinkCentre-M720t:~/dd/hadoop-3.2.2/titanic$ cat output/*
SurvivorClass 1 1.0
SurvivorClass 2 1.0
SurvivorClass 3 2.0
Average Age of Died female 34.0
Average Age of Died male 37.75
```

5. To analyse the **Earthquake Data** and generate statistics with region and magnitude/region and depth/region and latitude/region and longitude.



6. To analyse the given **Sales Records** over a period of time and generate data about the country's total sales and the total number of products. Country's total sales and the frequency of the payment mode.





PART-B

7. Write a **Spark program using Python**, to analyse the given **Weather Report Data** and to generate a report with cities having maximum and minimum temperature for a particular year.

```
ibmlab@ibmlab: ~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather
        at java.base/jdk.internal.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
        at java.base/java.lang.reflect.Method.invoke(Method.java:568)
        at py4j.reflection.MethodInvoker.invoke(MethodInvoker.java:244)
        at py4j.reflection.ReflectionEngine.invoke(ReflectionEngine.java:374)
        at py4j.Gateway.invoke(Gateway.java:282)
        at py4j.commands.AbstractCommand.invokeMethod(AbstractCommand.java:132)
        at py4j.commands.CallCommand.execute(CallCommand.java:79)
        at py4j.ClientServerConnection.waitForCommands(ClientServerConnection.java:182)
        at py4j.ClientServerConnection.run(ClientServerConnection.java:106)
        at java.base/java.lang.Thread.run(Thread.java:840)
24/07/23 15:19:53 INFO SparkContext: Invoking stop() from shutdown hook
24/07/23 15:19:53 INFO SparkContext: SparkContext is stopping with exitCode 0.
24/07/23 15:19:53 INFO SparkUI: Stopped Spark web UI at http://172-1-31-216.lightspeed.toldoh.sbcglobal.net:4040
24/07/23 15:19:53 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
24/07/23 15:19:53 INFO MemoryStore: MemoryStore cleared
24/07/23 15:19:53 INFO BlockManager: BlockManager stopped
24/07/23 15:19:53 INFO BlockManagerMaster: BlockManagerMaster stopped
24/07/23 15:19:53 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
24/07/23 15:19:53 INFO SparkContext: Successfully stopped SparkContext
24/07/23 15:19:53 INFO ShutdownHookManager: Shutdown hook called
24/07/23 15:19:53 INFO ShutdownHookManager: Deleting directory /tmp/spark-b8de80c5-cf68-4841-bae9-4592d992cd8b
24/07/23 15:19:53 INFO ShutdownHookManager: Deleting directory /tmp/spark-ff950943-8a48-4008-8ae9-82a8236cfab6
24/07/23 15:19:53 INFO ShutdownHookManager: Deleting directory /tmp/spark-ff950943-8a48-4008-8ae9-82a8236cfab6/pyspark-2cc48d7b-88cd-43ac-8fac
-060d13265bf6
lbmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$
l<mark>bmlab@ibmlab:~/1</mark>MS23SCS17/spark-3.5.1-bin-hadoop3/weather$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$ cat minimum/*
(1950, -11)
(1949, 78)
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$
lbmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$ cat maximum/*
(1949, 111)
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$
 .bmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/weather$
```

8. Write a **Spark program using Python**, to analyse the given **Insurance Data** and generate a statistics report with the construction building name and the count of building/ county name and its frequency.

```
tbmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
cat construction/*
('Wood', 17)
('Reinforced Masonry', 2)
('Reinforced Concrete', 3)
('Masonry', 2)
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
cat county/*
('ALACHUA COUNTY', 24)
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-hadoop3/insurance$
ibmlab@ibmlab:~/1MS23SCS17/spark-3.5.1-bin-ha
```

Pig

1. Write Pig Latin scripts for Crop Production Dataset.

crop_prod = LOAD 'crop_production.csv' USING PigStorage(',') AS (State_Name:chararray,
District_Name:chararray, Crop_Year:int, Season:chararray, Crop:chararray, Area:float,
Production:float);

DESCRIBE crop_prod;

```
grunt> crop_prod = LOAD 'crop_production.csv' USING PigStorage(',') AS (State_Na
me:chararray, District_Name:chararray, Crop_Year:int, Season:chararray, Crop:cha
rarray, Area:float, Production:float);
grunt> DESCRIBE crop_prod;
crop_prod: {State_Name: chararray,District_Name: chararray,Crop_Year: int,Season_
: chararray,Crop: chararray,Area: float,Production: float}
```

a. Calculate total production of each crop.

```
total_production = GROUP crop_prod BY Crop;
sum_production = FOREACH total_production GENERATE group AS Crop,
SUM(crop_prod.Production) AS Total_Production;
DUMP sum_production;
```

```
(Cowpea(Lobia),240638.0)
(Small millets,5630375.600000024)
(Oilseeds total, 4.386756E7)
(Rajmash Kholar,18590.0)
(other oilseeds, 4769908.810001001)
(Cond-spcs other, 2260.4000107347965)
(Other Dry Fruit, 0.0)
(Total foodgrain, 4.3270757E7)
(Other Vegetables,963212.0)
(Moong(Green Gram), 1.830318779909374E7)
(Peas (vegetable),0.0)
(Rapeseed &Mustard,9.086926582185636E7)
(Other Rabi pulses, 4805261.400009513)
(Other Citrus Fruit, 0.0)
(Other Fresh Fruits, 399443.0)
(Ricebean (nagadal),5230.0)
(other misc. pulses, 9704.219999995083)
(Arcanut (Processed),192831.0)
(Cashewnut Processed, 8121.0)
(Other Kharif pulses, 4349860.380000353)
(Peas & beans (Pulses),8752955.239967406)
(Other Cereals & Millets,1186045.6999955177)
(Beans & Mutter(Vegetable),211359.0)
```

b. Find the average production per year for each crop.

grouped_by_crop_year = GROUP crop_prod BY (Crop, Crop_Year);

average_production = FOREACH grouped_by_crop_year GENERATE group.Crop AS Crop, group.Crop_Year AS Crop_Year, AVG(crop_prod.Production) AS Avg_Production;

DUMP average production;

```
(Peas & beans (Pulses),2010,1770.018237082067)
(Peas & beans (Pulses),2011,2126.6575757575756)
(Peas & beans (Pulses),2012,2333.253472222222)
(Peas & beans (Pulses),2013,2209.551211006501)
(Peas & beans (Pulses),2014,1828.0802575073528)
(Other Cereals & Millets, 1998, 2142.553191489362)
(Other Cereals & Millets, 1999, 471. 42857142857144)
(Other Cereals & Millets, 2000, 1797, 4130434782608)
(Other Cereals & Millets, 2001, 199. 23529411764707)
(Other Cereals & Millets, 2002, 164.88888888888889)
(Other Cereals & Millets, 2003, 4401.6)
(Other Cereals & Millets, 2004, 1745. 9638554216867)
(Other Cereals & Millets, 2005, 2311.50000001395)
(Other Cereals & Millets, 2006, 1694.3488372093022)
(Other Cereals & Millets, 2007, 2920.691379307159)
(Other Cereals & Millets, 2008, 2936.808510638298)
(Other Cereals & Millets, 2009, 1051. 128205128205)
(Other Cereals & Millets, 2010, 2934.9285714285716)
(Other Cereals & Millets, 2011, 1716.5405405405406)
(Other Cereals & Millets, 2012, 1106.1052631578948)
(Other Cereals & Millets, 2014, 70.00526289720284)
(Beans & Mutter(Vegetable),2002,1020.9277108433735)
(Beans & Mutter(Vegetable),2003,1507.404761904762)
```

c. Filter all crops grown in 'Karnataka'
specific_state = FILTER crop_prod BY State_Name == 'Karnataka';
unique_crops = GROUP specific_state BY Crop;
DUMP unique_crops;

Karnataka,TUMKUR,2011,Kharif ,Peas & beans (Pulses),85.0,358.0),(Karnataka,HASSAN,2010 Peas & beans (Pulses),167.0,1471.0),(Karnataka, MANDYA,2005,Rabi ,Peas & beans (Pulses),32.0,225.0),(Karnataka,HASSAN,2010,Rabi Peas & beans (Pulses),320.0,2421.0),(Karnataka,MANDYA,20, beans (Pulses),36.0,182.0),(Karnataka,MANDYA,2010,Kharif ,Peas & beans (Pulses),318.0,2128.0),(Karnataka,MANDYA,2010,Rabi ,Peas & Peas & beans (Pulses),176.0,596.0),(Karnataka,HASSAN,2010,Kharif ,Peas & beans (Pulses),515.0,4586.0),(Karnataka,BANGALORE RURAL,2001,Whole Year ,Peas & beans (Pulses),411.0,2526.0),(Karnataka,CHIKBALLAPUR,2010,Kharif ,Peas & beans (Pulses),488.0,6545.0),(Karnataka,BENGALURU URBAN,1998,Kharif ,Peas & beans (Pulses),131.0,119.0),(Karnataka,BENGALURU URBAN,1998,Rharif ,Peas & beans (Pulses),131. Peas & beans (Pulses),488.0,6545.0),(Karnataka,BÉNGALURÚ URBÁN,1998,Kharif Peas & beans (Pulses),19.0,168.0),(Karnataka,RAMANAGARA,2010,Kharif,Peas & beans (Pulses),24.0,213.0),(Karnataka,CHIKBALLAPUR,2010,Rabi ,Peas & beans (Pulses),456.0,2887.0),(Karnataka,BANGALORE RURAL ,Rabi ,Peas & beans (Pulses),27.0,272.0),(Karnataka,MANDYA,1998,Summer ,Peas & beans (Pulses),66.0,652.0),(Karnataka,CHIKBALLAPUR,2010,Summer , r ,Peas & beans (Pulses),49.0,182.0),(Karnataka,MANDYA,1998,Rabi ,Peas & beans (Pulses),88.0,664.0),(Karnataka,GADAG,1998,Summer Peas & beans (Pulses),11.0,95.0),(Karnataka,MANDYA,1998,Kharif ,Peas & beans (Pulses),88.0,664.0),(Karnataka,GADAG,1998,Summer ,P. (Pulses),3.0,27.0),(Karnataka,BIJAPUR,2010,Kharif ,Peas & beans (Pulses),85.0,780.0),(Karnataka,BIJAPUR,2011,Summer ,Peas & beans (Pulses),34.0,303.0),(Karnataka,BENGALURU URBAN,2001,Whole Year ,Peas & beans (Pulses),11.0,942.0),(Karnataka,CHIKBALLAPUR,2011,Kharif ,Peas & beans (Pulses),508.0,5295.0),(Karnataka,RAMANAGARA,2010 Rabi ,Peas & beans (Pulses),508.0,5295.0),(Karnataka,RAMANAGARA,2010 Rabi ,Peas & beans (Pulses),508.0,5295.0),(Karnataka,RAMANAGARA,2010,Summer ,Peas & beans (Pulses),508.0,5295.0),(Karnataka,RAMANAGARA,2010 Rabi ,Peas & beans (Pulses),508.0,5295.0),(Karnataka,RAMANAGARA,2010,Summer ,Peas & beans (Pulses), peas & beans (Pulses),23.0,116.0),(Karnataka,HASSAN,2005,Rabi,Peas & beans (Pulses),178.0,1657.0),(Karnataka,MANDYA,2011,Rabi ,2878.0),(Karnataka,HASSAN,2005,Summer ,Peas & beans (Pulses),12.0,84.0),(Ka nataka,CHIKBALLAPUR,2011,Rabi ,Peas & beans (Pulses),151.0,482.0)}) Beans & Mutter(Vegetable),{(Karnataka,CHAMARAJANAGAR,2002,Whole Year ,Beans & Mutter(Vegetable),328.0,686.0),(Karnataka,BANGALORE RURAL,2003,Whole r "Beans & Mutter(Vegetable),462.0,2708.0),(Karnataka,KOLAR,2003,Whole Year "Beans & Mutter(Vegetable),1335.0,13505.0),(Karnataka,GULBARGA,2003,Whole Year "Beans & Mutter(Vegetable),525.0,3761.0),(Karnataka,CHITRADURGA,2002,Whol , Beans & Mutter(Vegetable),21.0,157.0),(Karnataka,MANDYA,2003,Whole Year ,Beans & Mutter(Vegetable),523.0,1054.0),(Karnataka,KOLAR,2002,Whole , Beans & Mutter(Vegetable),2575.0,23457.0),(Karnataka,MYSORÉ,2002,Whole Year ,Beans & Mutter(Vegetable),104.0,950.0),(Karnataka,CHIKMÁGALUR,2002, ole Year ,Beans & Mutter(Vegetable),187.0,1006.0),(Karnataka,HASSAN,2003,Whole Year ,Beans & Mutter(Vegetable),193.0,1321.0),(Karnataka,CHITRADURGA,2 03,Whole Year ,Beans & Mutter(Vegetable),26.0,156.0),(Karnataka,BELGAUM,2003,Whole Year ,Beans & Mutter(Vegetable),800.0,5124.0),(Karnataka,KODAGU,20 2,Whole Year ,Beans & Mutter(Vegetable),25.0,198.0),(Karnataka,CHIKMAGALUR,2003,Whole Year ,Beans & Mutter(Vegetable),2282.0,15472.0),(Karnataka,SHIM PAGA, 2002, Whole Year , Beans & Mutter(Vegetable), 25.0, 198.0), (Karnataka, CHIKMAUALUK, 2003, Whole Year , Beans & Mutter(Vegetable), 545.0, 2086.0), (Karnataka, SHIM)

PAGA, 2002, Whole Year , Beans & Mutter(Vegetable), 28.0, 221.0), (Karnataka, BANGALORE RURAL, 2002, Whole Year , Beans & Mutter(Vegetable), 545.0, 2086.0), (Karnataka, CADAG, 2002, Whole Year , Beans & Mutter(Vegetable), 855.0, 5834.0), (Karnataka, CADAG, 2003, Whole Year , Beans & Mutter(Vegetable), 202.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegetable), 229.0, 1597.0), (Karnataka, SHIMOGA, 2003, Whole Year , Beans & Mutter(Vegeta Mataka,BIDAPUR,2002,Whole Year ,Beans & Mutter(Vegetable),28.0,215.0),(Karnataka,MYSORE,2003,Whole Year ,Beans & Mutter(Vegetable),2644.0),(Karnataka,MYSORE,2003,Whole Year ,Beans & Mutter(Vegetable),2644.0),(Karnataka,MYSORE,2003,Whole Year ,Beans & Mutter(Vegetable),130.0,533.0),(Karnataka,DAVANGERE,2002,Whole Year ,Beans & Mutter(Vegetable),130.0,533.0),(Karnataka,DAVANGERE,2003,Whole Year ,Beans & Mutter(Vegetable),431.0,3005.0),(Karnataka,HAVERI,2002,Whole Year ,Beans & Mutter(Vegetable),3.0,24.0),(Karnataka,BIJAPUR,2003,Whole Year ,Beans & Mutter(Vegetable),4.0,28.0),(Karnataka,BIJAPUR,2003,Whole Year ,Beans & Mutter(Vegetable),524.0,562.0),(Karnataka,BIJAPUR,2003,Whole Year ,BIJAPUR,2003,Whole Year ,BIJAP 5.0),(Karnataka,BELGAUM,2002,Whole Year ,Beans & Mutter(Vegetable),99.0,626.0),(Karnataka,BIDAR,2003,Whole Year ,Beans & Mutter(Vegetable),185.0,1200 0),(Karnataka,HAVERI,2003,Whole Year ,Beans & Mutter(Vegetable),237.0,1653.0),(Karnataka,KODAGU,2003,Whole Year

d. Calculate the total area used for each crop in the year 2010.

specific_year = FILTER crop_prod BY Crop_Year == 2010;

total_area = GROUP specific_year BY Crop;

sum_area = FOREACH total_area GENERATE group AS Crop, SUM(specific_year.Area) AS Total_Area;

DUMP sum area;

```
Wheat, 2.9837101E7)
(Banana,319391.1396789551)
Barley,677806.0)
(Garlic,124884.0)
(Ginger,4554.0)
Masoor,1618613.0)
(Potato,1339266.0)
Khesari,512130.0)
Linseed,257083.0)
Sesamum, 2091602.8100008965)
(Tapioca,96568.09997558594)
(Tobacco,482107.0)
Arecanut, 399306.49981689453)
(Cardamom,19081.0)
Coconut ,1255481.125)
Sannhamp,9221.0)
(Soyabean,9558319.0)
Turmeric, 106831.7799949646)
(Arhar/Tur,4289565.0)
Cashewnut, 255319.7510986328)
(Coriander,401461.0)
(Groundnut,5861460.25)
Guar seed,3126187.0)
(Safflower,242802.0)
Sugarcane,4795159.140007984)
(Sunflower,889231.9000000954)
(Dry ginger,101069.92992973328)
Horse-gram,470891.0)
(Niger seed,240962.0)
Castor seed,889033.0)
(Black pepper,197783.55114746094)
(Cotton(lint),1.0929723E7)
(Dry chillies,430337.0)
(Sweet potato,38264.0)
Cowpea(Lobia),100769.0)
(Small millets,554661.0)
(Oilseeds total,31854.0)
(other oilseeds,1300390.0)
(Moong(Green Gram),3506872.9999752045)
Rapeseed &Mustard,5482169.0)
(Other Rabi pulses,69217.0)
(Other Kharif pulses,278997.0)
Peas & beans (Pulses),656383.0)
(Other Cereals & Millets,108567.0)
```

2. Write Pig Latin scripts for **Olympic Athletes and Hosts Datasets**.

```
athletes = LOAD 'olympic_athletes.csv' USING PigStorage(',') AS (athlete_url: chararray, athlete_full_name: chararray, games_participations: int, first_game: chararray, athlete_year_birth: float, athlete_medals: chararray, bio: chararray);

hosts = LOAD 'olympic_hosts.csv' USING PigStorage(',') AS (game_slug: chararray, game_end_date: chararray, game_start_date: chararray, game_location: chararray, game_name: chararray, game_season: chararray, game_year: int);

DESCRIBE athletes;
```

DESCRIBE hosts;

```
grunt> athletes = LOAD 'olympic_athletes.csv' USING PigStorage(',') AS (athlete_url: chararray, athlete_full_name: chararray, games_participations: in t, first_game: chararray, athlete_year_birth: float, athlete_medals: chararray, bio: chararray);
grunt> hosts = LOAD 'olympic_hosts.csv' USING PigStorage(',') AS (game_slug: chararray, game_end_date: chararray, game_start_date: chararray, game_loc ation: chararray, game_name: chararray, game_season: chararray, game_year: int);
grunt> DESCRIBE athletes;
athletes: {athlete_url: chararray,athlete_full_name: chararray,games_participations: int,first_game: chararray,athlete_year_birth: float,athlete_medal s: chararray,bio: chararray,game_start_date: chararray,game_location: chararray,game_name: chararray,game_season: chararray,game_year: int}
grunt> DESCRIBE hosts;
hosts: {game_slug: chararray,game_end_date: chararray,game_start_date: chararray,game_location: chararray,game_name: chararray,game_season: chararray,game_year: int}
```

I. Filter athletes participated in the "Tokyo 2020" games.

tokyo_2020_athletes = FILTER athletes BY first_game == 'Tokyo 2020';
DUMP tokyo_2020_athletes;

```
https://olympics.com/en/athletes/xinxin-wang,Xinxin WANG,1,Tokyo 2020,1998.0,,)
https://olympics.com/en/athletes/lidianny-echevarria-benitez,Lidianny ECHEVARRIA BENITEZ,1,Tokyo 2020,1996.0,,)
https://olympics.com/en/athletes/leila-consuelo-martinez-ortega,Leila Consuelo MARTINEZ ORTEGA,1,Tokyo 2020,1994.0,,)
 https://olympics.com/en/athletes/nadezda-makroguzova,Nadezda MAKROGUZOVA,1,Tokyo 2020,1997.0,,)
https://olympics.com/en/athletes/svetlana-kholomina,śvetlana KHOLOMINA,1,Tókyo 2020,1997.0,,)
https://olympics.com/en/athletes/tanja-huberli,Tanja HUBERLI,1,Tokyo 2020,1992.0,,)
https://olympics.com/en/athletes/nina-betschart,Nina BETSCHART,1,Tokyo 2020,1995.0,,)
https://olympics.com/en/athletes/kelly-claes,Kelly CLAES,1,Tokyo 2020,1995.0,,)
https://olympics.com/en/athletes/sarah-sponcil,Sarah SPONCIL,1,Tokyo 2020,1996.0,,)
 https://olympics.com/en/athletes/miki-ishii,Miki ISHII,1,Tokyo 2020,1989.0,,)
 https://olympics.com/en/athletes/megumi-murakami,Megumi MURAKAMI,1,Tokyo 2020,1985.0,,)
 https://olympics.com/en/athletes/katja-stam,Katja STAM,1,Tokyo 2020,1998.0,,
https://olympics.com/en/athletes/raisa-schoon,Raisa SCHOON,1,Tokyo 2020,,,)´
https://olympics.com/en/athletes/fernanda-pereyra,Fernanda PEREYRA,1,Tokyo 2020,1991.0,,)
https://olympics.com/en/athletes/julia-sude,Julia SUDE,1,Tokyo 2020,1987.0,,)
https://olympics.com/en/athletes/viktoria-orsi-toth,Viktoria ORSI TOTH,1,Tokyo 2020,1990.0,,)
 https://olympics.com/en/athletes/gaudencia-makokha,Gaudencia MAKOKHA,1,Tokyo 2020,1992.0,,
 https://olympics.com/en/athletes/brackcides-khadambi,Brackcides KHADAMBI,1,Tokyo 2020,1984.0,,)
https://olympics.com/en/athletes/anders-berntsen-mol,Anders Berntsen MOL,1,Tokyo 2020,1997.0,",)
https://olympics.com/en/athletes/christian-sandlie-sorum,Christian Sandlie SORUM,1,Tokyo 2020,1995.0,",)
https://olympics.com/en/athletes/oleg-stoyanovskiy,oleg STOYANOVSKIY,1,Tokyo 2020,1996.0,",)
https://olympics.com/en/athletes/cherif-younousse,Cherif YOUNOUSSE,1,Tokyo 2020,1995.0,",)
https://olympics.com/en/athletes/ahmed-tijan,Ahmed TIJAN,1,Tokyo 2020,1995.0,",)
https://olympics.com/en/athletes/edgars-tocs,Edgars TOCS,1,Tokyo 2020,1988.0,,)
https://olympics.com/en/athletes/alvaro-morais-filho,Alvaro MORAIS FILHO,1,Tokyo 2020,1990.0,,)
 https://olympics.com/en/athletes/julius-thole,Julius THOLE,1,Tokyo 2020,1997.0,,)
https://olympics.com/en/athletes/clemens-wickler,Clemens WICKLER,1,Tokyo 2020,1995.0,,)
https://olympics.com/en/athletes/ilya-leshukov,Ilya LESHUKOV,1,Tokyo 2020,1995.0,,)
https://olympics.com/en/athletes/pablo-herrera-allepuz,Pablo HERRERA ALLEPUZ,1,Tokyo 2020,,,)
 https://olympics.com/en/athletes/josue-gaston-gaxiola-leyva,Josue Gaston GAXIOLA LEYVA,1,Tokyo 2020,1997.0,,)
https://olympics.com/en/athletes/jose-luis-rubio-camargo,Jose Luis RUBIO CAMARGO,1,Tokyo 2020,1996.0,,)
 https://olympics.com/en/athletes/michal-bryl,Michal BRYL,1,Tokyo 2020,1994.0,,)
 https://olympics.com/en/athletes/tri-bourne,Tri BOURNE,1,Tokyo 2020,1989.0,,
 https://olympics.com/en/athletes/adrian-heidrich,Adrian HEIDRICH,1,Tokyo 2020,1994.0,,)
 https://olympics.com/en/athletes/mirco-gerson,Mirco GERSON,1,Tokyo 2020,1992.0,,)
https://olympics.com/en/athletes/nicolas-capogrosso,Nicolas CAPOGROSSO,1,Tokyo 2020,1995.0,,)
https://olympics.com/en/athletes/julian-amado-azaad, Julian Amado AZAAD,1,Tokyo 2020,1990.0,,)
https://olympics.com/en/athletes/christopher-mchugh,Christopher MCHUGH,1,Tokyo 2020,1989.0,,)
https://olympics.com/en/athletes/damien-schumann,Damien SCHUMANN,1,Tokyo 2020,1987.0,,)
https://olympics.com/en/athletes/ondrej-perusic,Ondrej PERUSIC,1,Tokyo 2020,1994.0,,)
https://olympics.com/en/athletes/david-schweiner,David SCHUKEINER,1,Tokyo 2020,1994.0,,)
 https://olympics.com/en/athletes/enrico-rossi,Enrico ROSSI,1,Tokyo 2020,1993.0,,)
https://olympics.com/en/athletes/mohammed-abicha,Mohammed ABICHA,1,Tokyo 2020,1980.0,,)
(https://olympics.com/en/athletes/zouheir-elgraoui,Zouheir ELGRAOUI,1,Tokyo 2020,1994.0,,)
```

II. Filter the games held in "China".

```
games_in_china = FILTER hosts BY game_location == 'China';
DUMP games_in_china;
```

```
(beijing-2022,2022-02-20T12:00:00Z,2022-02-04T15:00:00Z,China,Beijing 2022,Winter,2022)
(beijing-2008,2008-08-24T12:00:00Z,2008-08-08T00:00:00Z,China,Beijing 2008,Summer,2008)
```

III. Group games by season and count the number of games in each session. grouped_by_season = GROUP hosts BY game_season;

counted_by_season = FOREACH grouped_by_season GENERATE group AS game_season, COUNT(hosts) AS num_games;

DUMP counted_by_season;

```
(Summer, 28)
(Winter, 24)
(game_season, 1)
(Melbourne 1956, 1)
```

iv. Filter games that occurred after the year 2000.

games_after_2000 = FILTER hosts BY game_year > 2000; DUMP games_after_2000;

```
(beijing-2022,2022-02-20T12:00:00Z,2022-02-04T15:00:00Z,China,Beijing 2022,Winter,2022)
(tokyo-2020,2021-08-08T14:00:00Z,2021-07-23T11:00:00Z,Japan,Tokyo 2020,Summer,2020)
(pyeongchang-2018,2018-02-25T08:00:00Z,2018-02-08T23:00:00Z,Republic of Korea,PyeongChang 2018,Winter,2018)
(rio-2016,2016-08-21T21:00:00Z,2016-08-05T12:00:00Z,Brazil,Rio 2016,Summer,2016)
(sochi-2014,2014-02-23T16:00:00Z,2014-02-07T04:00:00Z,Russian Federation,Sochi 2014,Winter,2014)
(london-2012,2012-08-12T19:00:00Z,2012-07-27T07:00:00Z,Great Britain,London 2012,Summer,2012)
(vancouver-2010,2010-02-28T04:00:00Z,2010-02-12T16:00:00Z,Canada,Vancouver 2010,Winter,2010)
(beijing-2008,2008-08-24T12:00:00Z,2008-08-08T00:00:00Z,China,Beijing 2008,Summer,2008)
(turin-2006,2006-02-26T19:00:00Z,2006-02-10T07:00:00Z,Italy,Turin 2006,Winter,2006)
(athens-2004,2004-08-29T18:00:00Z,2004-08-13T06:00:00Z,Greece,Athens 2004,Summer,2004)
(salt-lake-city-2002,2002-02-24T08:00:00Z,2002-02-08T15:00:00Z,United States,Salt Lake City 2002,Winter,2002)
orunt>
```

9. Write **Pig Latin scripts** for the following queries on **Olympic Athletes** and **Hosts Dataset**:

- a. Filter athletes participated in the "Tokyo 2020" games.
- b. Filter the games held in "China".
- **C.** Group games by season and count the number of games in each session.
- d. Filter games that occurred after the year 2000.