



732A54 Big Data Analytics

Lab 1; pyspark

Group G6:

Hoda Fakharzadehjahromy(hodfa840)

Seyda Aqsa Iftekhar(syeif776)

1) What are the lowest and highest temperatures measured each year for the period 1950-2014. Provide the lists sorted in the descending order with respect to the maximum temperature. In this exercise you will use the *temperature-readings.csv* file.

The output should at least contain the following information (You can also include a Station column so that you may find multiple stations that record the highest (lowest) temperature.):

Year, temperature

Code:

```
from pyspark import SparkContext
import os
import sys

os.environ['PYSPARK_PYTHON'] = sys.executable
os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable
sc = SparkContext(appName="ex1" )
# MAP
temperature1 = sc.textFile("data/temperature-readings.csv")
lines = temperature1.map(lambda x: x.split(";"))
#print(lines.take(3))
# filter
year_temperature = lines.map(lambda x: (x[1][0:4],float(x[3])) )
year_temperature = year_temperature.filter(lambda x: (int(x[0]) >= 1950 and
int(x[0]) <= 2014))
#reduce
max_temperature = year_temperature.reduceByKey(lambda a,b: a if a >= b else b )
max_sorted_temperature = max_temperature.sortBy(ascending=False,keyfunc=lambda
k:k[1])

min_temperature = year_temperature.reduceByKey(lambda a,b: a if a < b else b )

min_sorted_temperature = min_temperature.sortBy(ascending=False,keyfunc=lambda
k:k[1])

#print(min_sorted_temperature.take( 10))
#print("-----")
#print(max_sorted_temperature.take( 10))
min_sorted_temperature.saveAsTextFile("res/ex1_min")
max_sorted_temperature.saveAsTextFile("res/ex1_max")
```

Output:

Max:

```
[('1992', 35.4), ('1994', 34.7), ('2014', 34.4), ('2010', 34.4), ('1989', 33.9),  
( '1982', 33.8), ('1968', 33.7), ('1966', 33.5), ('1983', 33.3), ('2002', 33.3),  
( '1986', 33.2), ('1970', 33.2), ('1956', 33.0), ('2000', 33.0), ('1959', 32.8),  
( '2006', 32.7), ('1991', 32.7), ('1988', 32.6), ('2011', 32.5)]
```

Min:

```
[('1999', -49.0), ('1978', -47.7), ('1987', -47.3), ('1967', -45.4), ('1956', -  
45.0), ('1980', -45.0), ('1971', -44.3), ('1986', -44.2), ('2001', -44.0),  
( '1965', -44.0), ('1981', -44.0), ('1979', -44.0), ('1959', -43.6), ('1985', -  
43.4), ('1958', -43.0), ('1998', -42.7), ('2012', -42.7), ('2014', -42.5),  
( '1977', -42.5)]
```

2) Count the number of readings for each month in the period of 1950-2014 which are higher than 10 degrees. Repeat the exercise, this time taking only distinct readings from each station. That is, if a station reported a reading above 10 degrees in some month, then it appears only once in the count for that month.

In this exercise you will use the *temperature-readings.csv* file.

The output should contain the following information:

Year, month, count

Code: 2a)

```
from pyspark import SparkContext  
from operator import add  
import os  
import sys  
  
os.environ['PYSPARK_PYTHON'] = sys.executable  
os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable  
  
sc = SparkContext(appName="ex2")  
  
temperature= sc.textFile("data/temperature-readings.csv")  
  
lines = temperature.map(lambda line: line.split(";"))
```

```

yearMontemp = lines.map(lambda x: ((x[1][0:4],x[1][5:7]), (x[0], float(x[3]))))
#print(year_month.take(5))
#print("-----")
#filtering
yearMontemp = yearMontemp.filter(lambda x: int(x[0][0])>=1950 and
int(x[0][0])<=2014 and x[1][1]>10)
#print(year_month.take(5))
#print("-----")
count_yearMontemp= yearMontemp.map(lambda x:(x[0],1)).reduceByKey(add).\
sortByKey(ascending=False)

count_yearMontemp.saveAsTextFile("output2/ex2a")

```

Output 2a):

```

(('2003', '09'), 70459)
(('2003', '08'), 108501)
(('2003', '07'), 128133)
(('2003', '06'), 99693)
(('2003', '05'), 48264)
(('2003', '04'), 11786)
(('2003', '03'), 3581)
(('2003', '02'), 3)
(('2002', '10'), 4939)
(('2002', '09'), 65928)

```

Code 2b) distinct reading for each station:

```

from pyspark import SparkContext
from operator import add
import os
import sys

os.environ['PYSPARK_PYTHON'] = sys.executable
os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable

```

```

sc = SparkContext(appName="ex2")

temperature= sc.textFile("./data/temperature-readings.csv")

lines = temperature.map(lambda line: line.split(";"))

year_month = lines.map(lambda x: ((x[1][0:4],x[1][5:7]), (x[0], float(x[3]))))
print(year_month.take(20))
print('-----')
-----')
#filtering
year_month = year_month.filter(lambda x: int(x[0][0])>=1950 and
int(x[0][0])<=2014 and x[1][1]>10)
print(year_month.take(20))

print('-----')
-----')

#selecting distinct key
year_month_unique = year_month. map(lambda x: (x[0], (x[1][0], 1)))
print(year_month_unique.take(20))

year_month_unique = year_month. map(lambda x: (x[0], (x[1][0], 1))).distinct()
print(year_month_unique.take(20))
station_month_counts = year_month_unique. map(lambda x: (x[0],
1)).reduceByKey(add).sortByKey(ascending=False)
station_month_counts.saveAsTextFile("./res/ex2b_updated")
print(station_month_counts.collect())

```

Output:

```

(('1951', '02'), 1)
(('1950', '12'), 1)
(('1950', '11'), 2)
(('1950', '10'), 46)
(('1950', '09'), 50)
(('1950', '08'), 49)
(('1950', '07'), 49)
(('1950', '06'), 47)
(('1950', '05'), 46)
(('1950', '04'), 36)
(('1950', '03'), 26)

```

3) Find the average monthly temperature for each available station in Sweden. Your result should include average temperature for each station for each month in the period of 1960-2014. Bear in mind that not every station has the readings for each month in this timeframe.

In this exercise you will use the *temperature-readings.csv* file.

The output should contain the following information:

Year, month, station number, average monthly temperature

```
from pyspark import SparkContext
import os
import sys

os.environ['PYSPARK_PYTHON'] = sys.executable
os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable

sc = SparkContext(appName='ex')

temperature = sc.textFile("data/temperature-readings.csv").cache()

lines = temperature.map(lambda x: x.split(";"))

#temperature = temperature.sample(False,0.01)

temperature = lines.map(lambda x: ((int(x[1][0:4]),int(x[1][5:7]),int(x[1][8:10])), (x[0], float(x[3]))))

# filter
select_temperature = temperature.filter(lambda x: x[0][0] >= 1960 and x[0][0] <= 2014)

select_temperature = select_temperature.map(lambda x: ((x[0][0],x[0][1],x[0][2],x[1][0]), (x[1][1],x[1][1])))
# #writing a function to calculate max and min together
#
def min_and_max(a,b):

    minimum = a[0] if a[0]<b[0] else b[0]
    maximum = b[1] if a[1]<b[1] else b[1]
    return (minimum,maximum)

select_temperature = select_temperature.reduceByKey(min_and_max)
```

```
##sum the daily values for each month
```

```
select_temperature1 = select_temperature.map(lambda x:  
((x[0][0],x[0][1],x[0][3]),(sum(x[1]),2)))
```

```
ave = select_temperature1.reduceByKey(lambda a, b: (a[0]+b[0],  
a[1]+b[1])).map(lambda x: (x[0], round(x[1][0]/x[1][1],  
ndigits=3))).sortByKey(ascending=False)
```

```
ave.saveAsTextFile("assignment3_updated")  
print(ave.collect())
```

Part of the Output:

```
((1969, 4, '74420'), 2.533), ((1969, 4, '74390'), 2.465), ((1969, 4, '74350'), 2.552),  
((1969, 4, '74240'), 2.548), ((1969, 4, '74180'), 0.41), ((1969, 4, '73660'), 2.7),  
((1969, 4, '73470'), 2.677), ((1969, 4, '73430'), 2.617), ((1969, 4, '73250'), 2.81),  
((1969, 4, '73090'), 3.578), ((1969, 4, '72450'), 3.288), ((1969, 4, '72400'), 4.093),  
((1969, 4, '72300'), 4.008), ((1969, 4, '72290'), 4.142), ((1969, 4, '72150'), 3.942),  
((1969, 4, '72120'), 3.098), ((1969, 4, '72080'), 4.23), ((1969, 4, '72070'), 4.332),  
((1969, 4, '71470'), 2.56), ((1969, 4, '71430'), 3.117), ((1969, 4, '71420'), 5.315),  
((1969, 4, '71380'), 4.067), ((1969, 4, '68550'), 2.07), ((1969, 4, '66600'), 2.763),  
((1969, 4, '66470'), 3.743), ((1969, 4, '66410'), 1.327), ((1969, 4, '66120'), 2.11),  
((1969, 4, '65450'), 3.65), ((1969, 4, '65160'), 0.92), ((1969, 4, '65130'), 4.088),  
((1969, 4, '64520'), 3.605), ((1969, 4, '64400'), 2.14), ((1969, 4, '64130'), 4.428),  
((1969, 4, '64030'), 4.488), ((1969, 4, '64020'), 2.858), ((1969, 4, '63630'), 3.952),  
((1969, 4, '63530'), 3.537), ((1969, 4, '63500'), 4.1), ((1969, 4, '63450'), 2.995),  
((1969, 4, '63440'), 3.535), ((1969, 4, '63340'), 4.358), ((1969, 4, '63230'), 3.717),  
((1969, 4, '63220'), 1.022), ((1969, 4, '63050'), 1.148), ((1969, 4, '62560'), 3.953),  
((1969, 4, '62410'), 2.1), ((1969, 4, '62190'), 3.535), ((1969, 4, '62180'), 1.728),  
((1969, 4, '62080'), 5.23), ((1969, 4, '62030'), 5.44), ((1969, 4, '55570'), 2.645),  
((1969, 4, '54550'), 2.82), ((1969, 4, '54390'), 4.545), ((1969, 4, '54300'), 3.648),  
((1969, 4, '54230'), 3.217), ((1969, 4, '53650'), 4.213), ((1969, 4, '53560'), 4.525),  
((1969, 4, '53540'), 4.793), ((1969, 4, '53470'), 3.95), ((1969, 4, '53430'), 4.78),  
((1969, 4, '53360'), 5.087), ((1969, 4, '53260'), 2.482), ((1969, 4, '53250'), 4.177),  
((1969, 4, '53200'), 3.27), ((1969, 4, '52550'), 4.562), ((1969, 4, '52230'), 3.853),  
((1969, 4, '192830'), -6.093), ((1969, 4, '192710'), -4.953), ((1969, 4, '191900'), -6.907),  
((1969, 4, '189780'), -5.92), ((1969, 4, '188830'), -5.313), ((1969, 4, '188800'), -3.82),  
((1969, 4, '183980'), -3.862), ((1969, 4, '183760'), -6.555), ((1969, 4, '182930'), -4.732),  
((1969, 4, '182720'), -2.742), ((1969, 4, '181900'), -3.628), ((1969, 4, '180940'), -8.917),  
((1969, 4, '180750'), -3.297), ((1969, 4, '179950'), -6.647), ((1969, 4, '178740'), -3.49),  
((1969, 4, '173810'), -1.665), ((1969, 4, '172790'), -3.84), ((1969, 4, '171920'), -3.685),  
((1969, 4, '171810'), -1.733), ((1969, 4, '170670'), -0.687), ((1969, 4, '169880'), -2.0),  
((1969, 4, '167980'), -3.038), ((1969, 4, '167860'), -4.317), ((1969, 4, '167710'), -2.798),  
((1969, 4, '166870'), -4.36), ((1969, 4, '166810'), -3.66), ((1969, 4, '163950'), -3.885),  
((1969, 4, '163690'), -0.882), ((1969, 4, '162980'), -0.988), ((1969, 4, '162970'), -0.817),  
((1969, 4, '162880'), -0.503), ((1969, 4, '162860'), -4.043), ((1969, 4, '161940'), -2.418),  
((1969, 4, '161790'), -0.825), ((1969, 4, '161770'), -1.077)
```

4) Provide a list of stations with their associated maximum measured temperatures and maximum measured daily precipitation. Show only those stations where the maximum temperature is between 25 and 30 degrees and maximum daily precipitation is between 100 mm and 200mm.

In this exercise you will use the *temperature-readings.csv* and *precipitation-readings.csv* files. The output should contain the following information:

Station number, maximum measured temperature, maximum daily precipitation

Code :

```
from pyspark import SparkContext
import os
import sys
from operator import add

os.environ['PYSPARK_PYTHON'] = sys.executable
os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable
sc = SparkContext(appName='ex')

temperature = sc.textFile("data/temperature-readings.csv").cache()
temperature = temperature.map(lambda a: a.split(';'))
temperature = temperature.map(lambda x: (x[0],float(x[3])))
temperature = temperature.reduceByKey(max)
temperature = temperature.filter(lambda x: x[1] >= 25 and x[1] <= 30 )

#print(temperature.collect())

precipitation = sc.textFile("data/precipitation-readings.csv").cache()
precipitation = precipitation.map(lambda x: x.split(';') )

precipitation = precipitation.map(lambda x: ((x[0],x[1]),float(x[3])))
#calculating daily precipitation
precipitation1 = precipitation.reduceByKey(add)
precipitation1=precipitation1.map(lambda x: (x[0][0],x[1]))
#finding daily precipitation
precipitation1 = precipitation1.reduceByKey(max)
precipitation1 = precipitation1.filter(lambda x: x[1] >= 100 and x[1] <= 200)

#print(precipitation1.collect())
```



```
# join the RDDs
stations = temperature.join(precipitation1)
# #the results
stations.saveAsTextFile("assignment4_updated")
print(stations.collect())
```

Output:

```
[ ]
```

5) Calculate the average monthly precipitation for the Östergötland region (list of stations is provided in the separate file) for the period 1993-2016. In order to do this, you will first need to calculate the total monthly precipitation for each station before calculating the monthly average (by averaging over stations).

In this exercise you will use the *precipitation-readings.csv* and *stations-Ostergotland.csv* files

The output should contain the following information:

Year, month, average monthly precipitation

Code:

```
from pyspark import SparkContext
import os
import sys
from operator import add

os.environ['PYSPARK_PYTHON'] = sys.executable
os.environ['PYSPARK_DRIVER_PYTHON'] = sys.executable

# create the spark application
sc = SparkContext(appName="exe")

ostergotland = sc.textFile('data/stations-Ostergotland.csv')
stations = ostergotland.map(lambda line: line.split(";")[0]).collect()
stations = sc.broadcast(stations)
#print(stations.value)
precipitation = sc.textFile("data/precipitation-readings.csv").cache()
precipitation = precipitation.map(lambda x: x.split(';') )
precipitation1 = precipitation.filter(lambda x: x[0] in stations.value)
```

```

precipitation1 = precipitation1.filter(lambda x: int(x[1][0:4]) >= 1993 and int(x[1][0:4])
<= 2016)
precipitation1 = precipitation1.map(lambda x: ((x[0],x[1][0:7]),float(x[3])))

#calculate the monthly precipitation for each station
precipitation1 = precipitation1.reduceByKey(add)

tmp = precipitation1.map(lambda x: (x[0][1], (x[1], 1)))
tmp = tmp.reduceByKey(lambda x,y: (x[0]+y[0],x[1]+y[1]))
#print(tmp.collect())
precip_ave = tmp.map(lambda x: (x[0],x[1][0]/x[1][1]))
# #precip_ave.saveAsTextFile("assignment5_updated")
print(precip_ave.collect())

```

Output:

```

[('2014-09', 48.45), ('2009-05', 54.167), ('2009-08', 61.567), ('2016-04', 26.9), ('1998-
05', 38.367), ('2002-02', 47.583), ('2016-02', 21.562), ('1997-03', 9.55), ('1999-01',
61.933), ('2009-03', 34.483), ('2011-12', 42.133), ('2015-09', 101.3), ('2015-10', 2.263),
('2006-12', 29.733), ('2008-11', 46.75), ('1994-06', 45.1), ('1999-05', 27.383), ('2004-01',
26.4), ('2004-09', 37.2), ('2006-08', 148.083), ('1995-02', 31.8), ('1996-02', 15.733),
('1998-10', 53.417), ('2007-06', 108.95), ('2013-01', 21.525), ('2014-08', 90.812), ('2008-
02', 28.25), ('2004-11', 54.133), ('2012-02', 28.667), ('2000-02', 24.05), ('2008-12',
43.483), ('2003-07', 113.467), ('1997-01', 5.783), ('2010-06', 48.65), ('2003-08', 55.283),
('1997-09', 43.567), ('2001-03', 30.267), ('2004-02', 27.35), ('2007-08', 54.167), ('2006-
06', 31.133), ('2006-07', 28.983), ('1994-09', 94.6), ('2008-10', 59.567), ('1996-06',
51.667), ('2010-03', 23.883), ('2016-03', 19.963), ('2010-09', 43.083), ('1995-03', 34.4),
('2010-08', 108.05), ('1997-11', 64.45), ('2007-04', 21.25), ('2015-11', 63.888), ('2000-
05', 25.317), ('2007-10', 28.117), ('2000-09', 27.517), ('2004-10', 78.183), ('2005-07',
104.35), ('2012-12', 66.933), ('2005-10', 38.05), ('2005-08', 76.967), ('2016-05', 29.25),
('2010-05', 67.167), ('2012-09', 72.75), ('1998-03', 33.9), ('1994-12', 50.9), ('2015-12',
28.925), ('2002-05', 72.133), ('2008-08', 138.517), ('2010-01', 35.983), ('2002-11',
53.683), ('2015-02', 24.825), ('1995-01', 26.0), ('2005-03', 23.467), ('2013-02', 25.525),
('2011-02', 24.517), ('1993-10', 43.2), ('1993-09', 40.6), ('2001-04', 46.067), ('2002-06',
98.783), ('2013-04', 38.288), ('2015-01', 59.113), ('2013-10', 53.875), ('2014-07',
22.988), ('2000-06', 62.017), ('1999-02', 25.15), ('2007-12', 54.7), ('2008-03', 42.2),
('2000-11', 108.117), ('2007-03', 40.517), ('1997-02', 35.433), ('2011-01', 35.133),
('1995-05', 26.0), ('2007-01', 68.633), ('1998-11', 28.967), ('2006-01', 17.683), ('2008-
06', 42.933), ('2009-10', 56.833), ('1995-11', 31.6), ('2001-01', 36.417), ('2001-08',
69.967), ('2003-09', 8.883), ('2007-11', 50.683), ('2005-02', 33.5), ('2012-01', 43.55),
('2000-08', 54.85), ('1999-04', 54.55), ('2009-06', 49.767), ('2008-01', 44.967), ('2008-
07', 85.2), ('2003-11', 54.45), ('1994-02', 22.5), ('1994-05', 25.1), ('1993-07', 95.4),
('1995-07', 43.6), ('1999-12', 66.0), ('2002-08', 8.25), ('2007-05', 40.517), ('2000-12',
63.267), ('2001-02', 36.767), ('2016-06', 47.663), ('2011-09', 52.567), ('2015-08',
26.987), ('1998-06', 88.883), ('2007-09', 61.867), ('2015-05', 93.225), ('2012-04',
62.783), ('1996-04', 8.1), ('2000-04', 36.167), ('2006-11', 71.717), ('2009-02', 24.783),

```

('1994-08', 58.8), ('1999-11', 18.45), ('2008-04', 20.25), ('2011-05', 37.85), ('2012-06', 132.2), ('2012-10', 65.583), ('2013-08', 54.075), ('2015-04', 15.337), ('1997-04', 25.95), ('2014-12', 35.463), ('2012-05', 22.967), ('1997-06', 86.983), ('2006-04', 44.367), ('2007-02', 33.067), ('1999-08', 54.8), ('1995-06', 97.2), ('2012-08', 68.817), ('1998-07', 85.167), ('2010-07', 92.4), ('2013-11', 46.375), ('2003-12', 52.117), ('2001-11', 26.383), ('2002-07', 80.517), ('1995-10', 20.867), ('2010-04', 23.783), ('2005-11', 32.6), ('1994-04', 23.1), ('2014-06', 75.138), ('2001-10', 60.483), ('2009-01', 15.883), ('2014-05', 58.0), ('2016-07', 0.0), ('2005-01', 18.05), ('2013-03', 7.387), ('2001-07', 40.283), ('2002-01', 55.0), ('2003-06', 66.667), ('2004-06', 56.85), ('2006-09', 19.267), ('2011-06', 88.35), ('2014-03', 36.563), ('2013-09', 26.188), ('1993-06', 56.5), ('1996-01', 10.417), ('1999-09', 55.517), ('1997-12', 69.6), ('2010-12', 37.183), ('2007-07', 95.967), ('1996-11', 67.117), ('2010-10', 52.533), ('2011-08', 86.267), ('1994-01', 22.1), ('2010-02', 52.75), ('1998-04', 44.417), ('2015-07', 119.1), ('2002-09', 16.067), ('2012-03', 8.55), ('2005-12', 56.633), ('1997-10', 58.15), ('1993-05', 21.1), ('2016-01', 22.325), ('1993-08', 80.7), ('1995-08', 16.05), ('1997-08', 24.617), ('1998-12', 59.383), ('2004-05', 39.7), ('2009-11', 64.217), ('2001-06', 26.767), ('2009-04', 2.8), ('2009-09', 29.95), ('2011-03', 19.833), ('2011-10', 43.75), ('1995-12', 5.117), ('2004-08', 75.483), ('1997-05', 60.8), ('2004-12', 24.25), ('2006-02', 34.75), ('2003-03', 6.9), ('1995-04', 61.7), ('2011-07', 94.917), ('1994-03', 37.6), ('2002-10', 60.5), ('2013-12', 42.263), ('1996-05', 63.233), ('2003-02', 9.117), ('2008-05', 23.133), ('2009-12', 53.45), ('2013-07', 54.562), ('2008-09', 47.367), ('1999-10', 18.55), ('1993-04', 0.0), ('2014-10', 72.137), ('1996-10', 22.45), ('2014-01', 62.575), ('2003-01', 17.717), ('1996-08', 37.717), ('2000-01', 18.617), ('2000-07', 135.867), ('2002-12', 20.917), ('2003-04', 51.417), ('2003-05', 68.45), ('2005-05', 55.383), ('2014-02', 43.713), ('2001-12', 35.183), ('1995-09', 134.55), ('1994-11', 15.7), ('1996-07', 84.033), ('1998-09', 56.833), ('2006-03', 27.867), ('1994-10', 33.2), ('1998-01', 44.317), ('1998-08', 86.467), ('2004-03', 28.483), ('2006-05', 52.333), ('2012-07', 59.067), ('2013-06', 61.325), ('2014-04', 31.763), ('2005-04', 11.65), ('2005-09', 13.95), ('1993-12', 37.1), ('1999-06', 50.25), ('1996-12', 39.55), ('2002-04', 29.917), ('1996-03', 10.033), ('2011-04', 14.917), ('2000-10', 110.3), ('1998-02', 50.033), ('2006-10', 118.167), ('2015-06', 78.663), ('2001-09', 110.633), ('1993-11', 42.8), ('1999-03', 42.233), ('2004-07', 96.0), ('2012-11', 68.65), ('2015-03', 42.613), ('2000-03', 21.683), ('2009-07', 113.167), ('2013-05', 47.925), ('2014-11', 52.425), ('2002-03', 26.933), ('2010-11', 93.55), ('1996-09', 57.467), ('1997-07', 41.967), ('2001-05', 33.983), ('2011-11', 13.467), ('1999-07', 29.083), ('2003-10', 45.583), ('2004-04', 20.617), ('1994-07', 0.0), ('2005-06', 67.967)]