

# lab1

Zhixuan\_Duan(zhidu838)

2020\_07\_01

## Assignment 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.

### spark code

```
##### lab1-1 #####
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
result1 = "BDA/output/max_temperature"
result2 = "BDA/output/min_temperature"

sc = SparkContext(appName="lab1-1")
lines = sc.textFile(data1)
lines01 = lines.map(lambda x: x.split(";"))
lines02 = lines01.filter(lambda x: int(x[1][0:4]) >= 1950 and
                           int(x[1][0:4]) <= 2014)
dateTemp = lines02.map(lambda x: (x[1][0:4], (x[0], float(x[3]))))

min = (lambda x, y: x if x[1] < y[1] else y)
max = (lambda x, y: x if x[1] > y[1] else y)

minTemp = dateTemp.reduceByKey(min)
maxTemp = dateTemp.reduceByKey(max)

minMaxTemp = minTemp.union(maxTemp).reduceByKey(lambda x,y: (x[0],x[1],y[0],y[1]))

# sort by highest temp on descending order
sortedMinMaxTemp1 = minMaxTemp.sortBy(keyfunc=lambda x: x[1][3], ascending=False)
# sort by lowest temp on descending order
sortedMinMaxTemp2 = minMaxTemp.sortBy(keyfunc=lambda x: x[1][1], ascending=False)

# list_year-station-maxTemp
output1 = sortedMinMaxTemp1.map(lambda x: '%s,%s,%s' % (x[0], x[1][2], x[1][3]))
# list_year-station-minTemp
output2 = sortedMinMaxTemp2.map(lambda x: '%s,%s,%s' % (x[0], x[1][0], x[1][1]))

output1.saveAsTextFile(result1)
output2.saveAsTextFile(result2)
```

## Max\_Temperature

```
year temp station
1975 36.1 86200
1992 35.4 63600
1994 34.7 117160
2014 34.4 96560
2010 34.4 75250
1989 33.9 63050
1982 33.8 94050
1968 33.7 137100
1966 33.5 151640
1983 33.3 98210
...
1987 29.6 96310
1984 29.5 105370
1967 29.5 75040
1960 29.4 173810
1950 29.4 75040
1998 29.2 63600
1965 28.5 116500
1951 28.5 75040
1962 27.4 76380
1962 27.4 86200
```

## Min\_Temperature

```
year temp station
1990 -35.0 147270
1990 -35.0 166870
1952 -35.5 192830
1974 -35.6 166870
1974 -35.6 179950
1954 -36.0 113410
1992 -36.1 179960
1975 -37.0 157860
1972 -37.5 167860
1995 -37.6 182910
...
1971 -44.3 166870
1956 -45.0 160790
1980 -45.0 191900
1980 -45.0 191900
1967 -45.4 166870
1987 -47.3 123480
1978 -47.7 155940
1999 -49.0 192830
1999 -49.0 192830
1966 -49.4 179950
```

## Assignment 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

### spark code

```
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
result2_1 = 'BDA/output/over_ten_mth_temp_counts'
result2_2 = 'BDA/output/over_ten_temp_distinct_counts'
startYear = 1950
endYear = 2014
standardTemp = 10

sc = SparkContext(appName="lab1-2")

lines = sc.textFile(data1)
lines01 = lines.map(lambda x: x.split(";"))

lines02 = lines01.filter(lambda x: (int(x[1][0:4]) >= startYear and
                                     int(x[1][0:4]) <= endYear))

# 1. year, month, counts
temperatures = lines02.map(lambda x:
                           (x[1][0:7], (float(x[3]), 1))) \
    .filter(lambda (x, (y, z)): y > standardTemp)

reading_counts = temperatures.reduceByKey(lambda (temp1, count1), (temp2, count2):
                                           (temp1, count1 + count2)) \
    .map(lambda (x, (y, z)): (x, z))

reading_counts.saveAsTextFile(result2_1)

# 2. year, month, distinct counts
station_temperatures = lines02.map(lambda x:
                                   (x[1][0:7], (x[0], float(x[3])))) \
    .filter(lambda (x, (y, z)): z > standardTemp)

year_station = station_temperatures.map(lambda (x, (y, z)):
                                         (x, (y, 1))).distinct()

reading_counts = year_station.reduceByKey(lambda (station1, count1), (station2, count2):
                                           (station1, count1 + count2)) \
    .map(lambda (x, (y, z)): (x, y))

reading_counts.saveAsTextFile(result2_2)
```

## over\_ten\_mth\_temp\_counts

```
year-month counts
1997-04 3756
1974-07 66277
2003-05 48264
1981-10 9882
1981-03 395
1987-05 17191
2009-07 133008
1986-11 1198
1983-09 38692
1952-03 1
...
2013-10 37839
1973-05 27886
1960-12 2
1972-11 993
1959-05 12445
1971-10 13326
1960-05 14333
1966-07 55544
1976-02 19
1983-06 49265
```

## over\_ten\_temp\_distinct\_counts

```
year-month counts
1997-04 91130
1974-07 77220
2003-05 182930
1981-10 91130
1983-09 157860
1987-05 122330
1979-04 74180
2009-07 106040
1986-11 81540
1966-08 72120
...
2013-10 136360
1973-05 105450
1960-12 134110
1972-11 65360
1959-05 135440
1971-10 65360
1960-05 97620
1966-07 85160
1976-02 75040
1955-01 84310
```

### Assignment 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

#### spark code

```
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
result3 = "BDA/output/avg_station_month_temp"
startYear = 1960
endYear = 2014

sc = SparkContext(appName="lab1-3")

lines = sc.textFile(data1)
lines01 = lines.map(lambda a: a.split(";"))
lines02 = lines01.filter(lambda x: (int(x[1][0:4]) >= startYear and
                                     int(x[1][0:4]) <= endYear))

stationTemp = lines02.map(lambda x: ((x[1], x[0]), (float(x[3]), float(x[3]))))

# get the daily min and max Temps
stationMinMaxTemp = stationTemp.reduceByKey(lambda (mintemp1, maxtemp1), (mintemp2, maxtemp2):
                                              (min(mintemp1, mintemp2), max(maxtemp1, maxtemp2)))

# get the monthly avg_Temp
stationMonthlyAvgTemps = stationMinMaxTemp.map(lambda ((day, station), (mintemp, maxtemp)):
                                              ((day[0:7], station), (sum((mintemp, maxtemp)), 2))) \
    .reduceByKey(lambda (temp1, count1), (temp2, count2):
                  (temp1 + temp2, count1 + count2)) \
    .map(lambda ((month, station), (temp, count)):
          ((month, station), temp / float(count)))

stationMonthlyAvgTemps.saveAsTextFile(result3)
```

#### avg\_station\_month\_temp

```
year-month station temp
2001-01 63280 0.0032258064516128343
1969-03 83620 -4.780645161290323
1978-09 156730 5.5566666666666675
1991-11 106500 -0.05166666666666667
1995-06 112080 12.001666666666667
1998-05 172770 5.282258064516129
2003-06 177930 8.711666666666668
1978-09 95160 10.040000000000001
```

```
1976-10 89240 6.138709677419354
2001-02 159770 -14.35
...
1976-08 54400 16.81129032258065
1990-05 181900 6.77258064516129
1964-12 87150 0.04516129032258059
2003-08 134590 12.309677419354841
1965-04 158750 1.6566666666666665
1972-07 94140 19.853225806451615
2002-10 83440 4.762903225806452
1970-12 163920 -5.527419354838709
1965-01 74470 -0.7516129032258065
1983-12 161780 -8.535483870967743
```

## Assignment 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 100mm 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

### spark code

```
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
data2 = "BDA/input/precipitation-readings.csv"
result4 = "BDA/output/max_temp_prec"

sc = SparkContext(appName="lab1-4")

# max_temp
temp = sc.textFile(data1)
temp = temp.map(lambda a: a.split(";"))
temp = temp.map(lambda x: (x[0], float(x[3])))

maxTemp = temp.reduceByKey(max)
maxTemp = maxTemp.filter(lambda a: a[1] > 25 and a[1] < 30)

# max_prec
prec = sc.textFile(data2)
prec = prec.map(lambda a: a.split(";"))
prec = prec.map(lambda x: (x[0]+'+', '+x[1]', float(x[3])))

maxPrec = prec.reduceByKey(lambda v1, v2: v1+v2)
maxPrec = maxPrec.filter(lambda a: a[1] >= 100 and a[1] <= 200)

# merge them together
maxTempPrec = maxTemp.leftOuterJoin(maxPrec)
maxTempPrec = maxTempPrec.map(lambda a: '%s,%s,%s' % (a[0], a[1][0], a[1][1]))

maxTempPrec.saveAsTextFile(result4)
```

### max\_temp\_prec

```
station temp prec
99280 25.5 34.00000000000001
87440 26.2 45.1
82360 29.9 84.30000000000001
78550 29.9 49.199999999999996
71190 28.3 80.39999999999999
68560 28.5 52.99999999999999
66110 26.3 72.60000000000002
```

```
52240 29.0 79.6
188790 26.6 34.199999999999996
183750 29.4 38.1
...
158740 29.2 61.1
157870 29.9 61.299999999999999
144310 28.4 53.5
139260 29.4 86.300000000000001
139120 27.9 95.4
135460 29.2 62.1
122260 28.5 47.6
178860 27.4 44.900000000000006
133500 27.2 50.100000000000001
132170 28.0 39.6
```



## Assignment 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. HINT (not for the SparkSQL lab): Avoid using joins here! stations-Ostergotland.csv is small and if distributed will cause a number of unnecessary shuffles when joined with precipitationRDD. If you distribute precipitation-readings.csv then either repartition your stations RDD to 1 partition or make use of the collect to acquire a python list and broadcast function to broadcast the list to all nodes. The output should contain the following information: Year, month, average monthly precipitation

### spark code

```
from pyspark import SparkContext

data1 = "BDA/input/precipitation-readings.csv"
data2 = "BDA/input/stations-Ostergotland.csv"
result5 = "BDA/output/avg_prec_ost"

sc = SparkContext(appName="lab1-5")

prec = sc.textFile(data1)
stations = sc.textFile(data2)

# stations data
stations = stations.map(lambda x: x.split(";"))
stations = stations.map(lambda x: x[0])
stations = stations.coalesce(1)
stations = stations.collect()
stations = sc.broadcast(stations)

# prec data
prec = prec.map(lambda x: x.split(";"))
prec = prec.map(lambda x: (x[0], x[1][0:4], x[1][5:7], float(x[3])))
prec = prec.filter(lambda x: x[0][0] in stations.value)

# avg_prec
prec_2016 = prec.filter(lambda x: int(x[0][1]) >= 1993 and
                        int(x[0][1]) <= 2016)
prec_2016 = prec_2016.reduceByKey(lambda x, y: x + y)
prec_2016 = prec_2016.map(lambda line: ((line[0][1], line[0][2]), (line[1], 1)))
prec_2016 = prec_2016.reduceByKey(lambda x, y: (x[0] + y[0], x[1] + y[1]))
prec_2016 = prec_2016.map(lambda line: (line[0][0], line[0][1], (line[1][0] / line[1][1])))
prec_2016 = prec_2016.coalesce(1)

prec_2016.saveAsTextFile(result5)
```

avg\_prec\_ost

```
year month prec
2016 05 29.250000000000007
2016 02 21.5625
2015 12 28.925
2015 11 63.887500000000002
2015 09 101.29999999999998
2015 08 26.9875
2015 07 119.09999999999995
2015 06 78.662500000000001
2015 05 93.225
2015 04 15.3375
...
2013 07 54.562499999999999
2013 05 47.925000000000001
2013 04 38.2875000000000016
2013 02 25.525000000000002
2012 11 68.65
2012 10 65.583333333333334
2012 09 72.75
2012 07 59.066666666666667
2012 04 62.783333333333335
2012 01 43.550000000000003
```

## Appendix

```
knitr::opts_chunk$set(echo = TRUE,
                      eval = FALSE,
                      warning = FALSE,
                      message = FALSE,
                      fig.width = 7,
                      fig.height = 4,
                      fig.align = 'center')

##### lab1-1 #####
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
result1 = "BDA/output/max_temperature"
result2 = "BDA/output/min_temperature"

sc = SparkContext(appName="lab1-1")
lines = sc.textFile(data1)
lines01 = lines.map(lambda x: x.split(";"))
lines02 = lines01.filter(lambda x: int(x[1][0:4]) >= 1950 and
                          int(x[1][0:4]) <= 2014)
dateTemp = lines02.map(lambda x: (x[1][0:4], (x[0], float(x[3]))))

min = (lambda x, y: x if x[1] < y[1] else y)
max = (lambda x, y: x if x[1] > y[1] else y)

minTemp = dateTemp.reduceByKey(min)
maxTemp = dateTemp.reduceByKey(max)

minMaxTemp = minTemp.union(maxTemp).reduceByKey(lambda x,y: (x[0],x[1],y[0],y[1]))

# sort by highest temp on descending order
sortedMinMaxTemp1 = minMaxTemp.sortBy(keyfunc=lambda x: x[1][3], ascending=False)
# sort by lowest temp on descending order
sortedMinMaxTemp2 = minMaxTemp.sortBy(keyfunc=lambda x: x[1][1], ascending=False)

# list_year-station-maxTemp
output1 = sortedMinMaxTemp1.map(lambda x: '%s,%s,%s' % (x[0], x[1][2], x[1][3]))
# list_year-station-minTemp
output2 = sortedMinMaxTemp2.map(lambda x: '%s,%s,%s' % (x[0], x[1][0], x[1][1]))

output1.saveAsTextFile(result1)
output2.saveAsTextFile(result2)
year temp station
1975 36.1 86200
1992 35.4 63600
1994 34.7 117160
2014 34.4 96560
2010 34.4 75250
1989 33.9 63050
1982 33.8 94050
1968 33.7 137100
1966 33.5 151640
```

```

1983 33.3 98210
...
1987 29.6 96310
1984 29.5 105370
1967 29.5 75040
1960 29.4 173810
1950 29.4 75040
1998 29.2 63600
1965 28.5 116500
1951 28.5 75040
1962 27.4 76380
1962 27.4 86200
year  temp station
1990 -35.0 147270
1990 -35.0 166870
1952 -35.5 192830
1974 -35.6 166870
1974 -35.6 179950
1954 -36.0 113410
1992 -36.1 179960
1975 -37.0 157860
1972 -37.5 167860
1995 -37.6 182910
...
1971 -44.3 166870
1956 -45.0 160790
1980 -45.0 191900
1980 -45.0 191900
1967 -45.4 166870
1987 -47.3 123480
1978 -47.7 155940
1999 -49.0 192830
1999 -49.0 192830
1966 -49.4 179950
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
result2_1 = 'BDA/output/over_ten_mth_temp_counts'
result2_2 = 'BDA/output/over_ten_temp_distinct_counts'
startYear = 1950
endYear = 2014
standardTemp = 10

sc = SparkContext(appName="lab1-2")

lines = sc.textFile(data1)
lines01 = lines.map(lambda x: x.split(";"))

lines02 = lines01.filter(lambda x: (int(x[1][0:4]) >= startYear and
                                   int(x[1][0:4]) <= endYear))

# 1. year, month, counts
temperatures = lines02.map(lambda x:

```

```

        (x[1][0:7], (float(x[3]), 1))) \
    .filter(lambda (x, (y, z)): y > standardTemp)

reading_counts = temperatures.reduceByKey(lambda (temp1, count1), (temp2, count2):
                                           (temp1, count1 + count2)) \
    .map(lambda (x, (y, z)): (x, z))

reading_counts.saveAsTextFile(result2_1)

# 2. year, month, distinct counts
station_temperatures = lines02.map(lambda x:
                                   (x[1][0:7], (x[0], float(x[3])))) \
    .filter(lambda (x, (y, z)): z > standardTemp)

year_station = station_temperatures.map(lambda (x, (y, z)):
                                         (x, (y, 1))).distinct()

reading_counts = year_station.reduceByKey(lambda (station1, count1), (station2, count2):
                                           (station1, count1 + count2)) \
    .map(lambda (x, (y, z)): (x, y))

reading_counts.saveAsTextFile(result2_2)
year-month counts
1997-04 3756
1974-07 66277
2003-05 48264
1981-10 9882
1981-03 395
1987-05 17191
2009-07 133008
1986-11 1198
1983-09 38692
1952-03 1
...

2013-10 37839
1973-05 27886
1960-12 2
1972-11 993
1959-05 12445
1971-10 13326
1960-05 14333
1966-07 55544
1976-02 19
1983-06 49265
year-month counts
1997-04 91130
1974-07 77220
2003-05 182930
1981-10 91130
1983-09 157860
1987-05 122330
1979-04 74180

```

```

2009-07 106040
1986-11 81540
1966-08 72120
...

2013-10 136360
1973-05 105450
1960-12 134110
1972-11 65360
1959-05 135440
1971-10 65360
1960-05 97620
1966-07 85160
1976-02 75040
1955-01 84310
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
result3 = "BDA/output/avg_station_month_temp"
startYear = 1960
endYear = 2014

sc = SparkContext(appName="lab1-3")

lines = sc.textFile(data1)
lines01 = lines.map(lambda a: a.split(";"))
lines02 = lines01.filter(lambda x: (int(x[1][0:4]) >= startYear and
                                     int(x[1][0:4]) <= endYear))

stationTemp = lines02.map(lambda x:((x[1], x[0]), (float(x[3]), float(x[3]))))

# get the daily min and max Temps
stationMinMaxTemp = stationTemp.reduceByKey(lambda (mintemp1, maxtemp1), (mintemp2, maxtemp2):
                                              (min(mintemp1, mintemp2), max(maxtemp1, maxtemp2)))

# get the monthly avg_Temp
stationMonthlyAvgTemps = stationMinMaxTemp.map(lambda ((day, station), (mintemp, maxtemp)):
                                              ((day[0:7], station), (sum((mintemp, maxtemp)), 2))) \
    .reduceByKey(lambda (temp1, count1), (temp2, count2):
                  (temp1 + temp2, count1 + count2)) \
    .map(lambda ((month, station), (temp, count)):
          ((month, station), temp / float(count)))

stationMonthlyAvgTemps.saveAsTextFile(result3)
year-month station temp
2001-01 63280 0.0032258064516128343
1969-03 83620 -4.780645161290323
1978-09 156730 5.5566666666666675
1991-11 106500 -0.05166666666666667
1995-06 112080 12.001666666666667
1998-05 172770 5.282258064516129
2003-06 177930 8.711666666666668
1978-09 95160 10.040000000000001
1976-10 89240 6.138709677419354

```

```

2001-02 159770 -14.35
...
1976-08 54400 16.81129032258065
1990-05 181900 6.77258064516129
1964-12 87150 0.04516129032258059
2003-08 134590 12.309677419354841
1965-04 158750 1.6566666666666665
1972-07 94140 19.853225806451615
2002-10 83440 4.762903225806452
1970-12 163920 -5.527419354838709
1965-01 74470 -0.7516129032258065
1983-12 161780 -8.535483870967743
from pyspark import SparkContext

data1 = "BDA/input/temperature-readings.csv"
data2 = "BDA/input/precipitation-readings.csv"
result4 = "BDA/output/max_temp_prec"

sc = SparkContext(appName="lab1-4")

# max_temp
temp = sc.textFile(data1)
temp = temp.map(lambda a: a.split(";"))
temp = temp.map(lambda x: (x[0], float(x[3])))

maxTemp = temp.reduceByKey(max)
maxTemp = maxTemp.filter(lambda a: a[1] > 25 and a[1] < 30)

# max_prec
prec = sc.textFile(data2)
prec = prec.map(lambda a: a.split(";"))
prec = prec.map(lambda x: (x[0]+'+', '+x[1]', float(x[3])))

maxPrec = prec.reduceByKey(lambda v1, v2: v1+v2)
maxPrec = maxPrec.filter(lambda a: a[1] >= 100 and a[1] <= 200)

# merge them together
maxTempPrec = maxTemp.leftOuterJoin(maxPrec)
maxTempPrec = maxTempPrec.map(lambda a: '%s,%s,%s' % (a[0], a[1][0], a[1][1]))

maxTempPrec.saveAsTextFile(result4)
station temp prec
99280 25.5 34.000000000000001
87440 26.2 45.1
82360 29.9 84.300000000000001
78550 29.9 49.199999999999996
71190 28.3 80.39999999999999
68560 28.5 52.99999999999999
66110 26.3 72.600000000000002
52240 29.0 79.6
188790 26.6 34.199999999999996
183750 29.4 38.1
...

```

```

158740 29.2 61.1
157870 29.9 61.299999999999999
144310 28.4 53.5
139260 29.4 86.300000000000001
139120 27.9 95.4
135460 29.2 62.1
122260 28.5 47.6
178860 27.4 44.900000000000006
133500 27.2 50.100000000000001
132170 28.0 39.6
from pyspark import SparkContext

data1 = "BDA/input/precipitation-readings.csv"
data2 = "BDA/input/stations-Ostergotland.csv"
result5 = "BDA/output/avg_prec_ost"

sc = SparkContext(appName="lab1-5")

prec = sc.textFile(data1)
stations = sc.textFile(data2)

# stations data
stations = stations.map(lambda x: x.split(";"))
stations = stations.map(lambda x: x[0])
stations = stations.coalesce(1)
stations = stations.collect()
stations = sc.broadcast(stations)

# prec data
prec = prec.map(lambda x: x.split(";"))
prec = prec.map(lambda x: (x[0],x[1][0:4],x[1][5:7], float(x[3])))
prec = prec.filter(lambda x: x[0][0] in stations.value)

# avg_prec
prec_2016 = prec.filter(lambda x: int(x[0][1]) >= 1993 and
                        int(x[0][1]) <= 2016)
prec_2016 = prec_2016.reduceByKey(lambda x, y: x + y)
prec_2016 = prec_2016.map(lambda line: ((line[0][1], line[0][2]), (line[1], 1)))
prec_2016 = prec_2016.reduceByKey(lambda x, y: (x[0] + y[0], x[1] + y[1]))
prec_2016 = prec_2016.map(lambda line: (line[0][0], line[0][1], (line[1][0] / line[1][1])))
prec_2016 = prec_2016.coalesce(1)

prec_2016.saveAsTextFile(result5)
year month prec
2016 05 29.250000000000007
2016 02 21.5625
2015 12 28.925
2015 11 63.887500000000002
2015 09 101.299999999999998
2015 08 26.9875
2015 07 119.099999999999995
2015 06 78.662500000000001
2015 05 93.225

```



```
2015 04 15.3375
...
2013 07 54.562499999999999
2013 05 47.925000000000001
2013 04 38.287500000000016
2013 02 25.525000000000002
2012 11 68.65
2012 10 65.583333333333334
2012 09 72.75
2012 07 59.066666666666667
2012 04 62.783333333333335
2012 01 43.550000000000003
```