

# 732A54: Lab 1

## Big Data Analytics

Carles Sans Fuentes

February 27, 2017

---

## Assignment 1

### Exercise 1

Part A code

```
1 ##1 maximum temperatures descending
2 #max
3 from pyspark import SparkContext
4 sc1 = SparkContext(appName = "1 max temperatures descending")
5 reading = sc1.textFile("/user/x_carsa/data/temperature-readings.csv")
6 separate = reading.map(lambda a: a.split(";"))
7 lines = separate.filter(lambda x: int(x[1][0:4]) >= 1950 and int(x[1][0:4]) <= 2014)
8 temperatures = lines.map(lambda x: (x[1][0:4], (x[0], float(x[3])))) ##added x[0] as it is the
    station column
9 maxTemperatures = temperatures.reduceByKey(lambda a,b: a if a[1]>b[1] else b)
10 maxTemperaturesTogether = maxTemperatures.repartition(1)
11 maxTemperaturesSorted = maxTemperaturesTogether.sortBy(ascending=False, keyfunc=lambda k: k
    [1][1]) ## the first [1] stands for k = length 1 and vector length 2, with the second [1,
    we acces to the second element of v]
12 maxTemperaturesSorted.saveAsTextFile("results/1resultsMax")
```

Output:

```
1
2 Row(_c0=u'1975', _c1=u'134520', value=36.1)
3 Row(_c0=u'1992', _c1=u'140360', value=35.4)
4 Row(_c0=u'1994', _c1=u'54300', value=34.7)
5 Row(_c0=u'2010', _c1=u'133250', value=34.4)
6 Row(_c0=u'2014', _c1=u'133250', value=34.4)
7 Row(_c0=u'1989', _c1=u'140200', value=33.9)
8 Row(_c0=u'1982', _c1=u'106100', value=33.8)
9 Row(_c0=u'1968', _c1=u'134520', value=33.7)
10 Row(_c0=u'1966', _c1=u'102190', value=33.5)
11 Row(_c0=u'1983', _c1=u'133260', value=33.3)
12 Row(_c0=u'2002', _c1=u'156770', value=33.3)
13 Row(_c0=u'1986', _c1=u'106100', value=33.2)
14 Row(_c0=u'1970', _c1=u'102190', value=33.2)
15 Row(_c0=u'2000', _c1=u'133470', value=33.0)
16 Row(_c0=u'1956', _c1=u'161670', value=33.0)
17 Row(_c0=u'1959', _c1=u'108640', value=32.8)
18 Row(_c0=u'1991', _c1=u'133260', value=32.7)
19 Row(_c0=u'2006', _c1=u'139120', value=32.7)
20 Row(_c0=u'1988', _c1=u'134460', value=32.6)
21 Row(_c0=u'2011', _c1=u'116430', value=32.5)
22 Row(_c0=u'1999', _c1=u'133470', value=32.4)
23 Row(_c0=u'2003', _c1=u'139120', value=32.2)
24 Row(_c0=u'2007', _c1=u'133250', value=32.2)
25 Row(_c0=u'2008', _c1=u'151280', value=32.2)
26 Row(_c0=u'1953', _c1=u'137030', value=32.2)
27 Row(_c0=u'1955', _c1=u'139340', value=32.2)
28 Row(_c0=u'1973', _c1=u'139200', value=32.2)
29 Row(_c0=u'2005', _c1=u'133470', value=32.1)
30 Row(_c0=u'1969', _c1=u'102190', value=32.0)
31 Row(_c0=u'1979', _c1=u'133470', value=32.0)
32 Row(_c0=u'2001', _c1=u'106160', value=31.9)
33 Row(_c0=u'1997', _c1=u'146070', value=31.8)
34 Row(_c0=u'1977', _c1=u'133470', value=31.8)
```

```

35 Row(_c0=u'2013', _c1=u'139120', value=31.6)
36 Row(_c0=u'2009', _c1=u'140460', value=31.5)
37 Row(_c0=u'2012', _c1=u'151240', value=31.3)
38 Row(_c0=u'1964', _c1=u'139200', value=31.2)
39 Row(_c0=u'1971', _c1=u'139200', value=31.2)
40 Row(_c0=u'1972', _c1=u'112080', value=31.2)
41 Row(_c0=u'1976', _c1=u'139200', value=31.1)
42 Row(_c0=u'1961', _c1=u'137030', value=31.0)
43 Row(_c0=u'1963', _c1=u'106270', value=31.0)
44 Row(_c0=u'1995', _c1=u'136420', value=30.8)
45 Row(_c0=u'1996', _c1=u'134410', value=30.8)
46 Row(_c0=u'1958', _c1=u'139340', value=30.8)
47 Row(_c0=u'1978', _c1=u'139200', value=30.8)
48 Row(_c0=u'1974', _c1=u'134520', value=30.6)
49 Row(_c0=u'1954', _c1=u'137030', value=30.5)
50 Row(_c0=u'1980', _c1=u'106100', value=30.4)
51 Row(_c0=u'1952', _c1=u'137030', value=30.4)
52 Row(_c0=u'1990', _c1=u'133260', value=30.2)
53 Row(_c0=u'2004', _c1=u'139120', value=30.2)
54 Row(_c0=u'1985', _c1=u'102200', value=29.8)
55 Row(_c0=u'1957', _c1=u'135440', value=29.8)
56 Row(_c0=u'1981', _c1=u'140480', value=29.7)
57 Row(_c0=u'1993', _c1=u'102210', value=29.7)
58 Row(_c0=u'1987', _c1=u'102200', value=29.6)
59 Row(_c0=u'1984', _c1=u'139200', value=29.5)
60 Row(_c0=u'1967', _c1=u'134520', value=29.5)
61 Row(_c0=u'1950', _c1=u'134110', value=29.4)
62 Row(_c0=u'1960', _c1=u'107400', value=29.4)
63 Row(_c0=u'1998', _c1=u'133260', value=29.2)
64 Row(_c0=u'1951', _c1=u'107440', value=28.5)
65 Row(_c0=u'1965', _c1=u'192710', value=28.5)
66 Row(_c0=u'1962', _c1=u'156730', value=27.4)

```

#### Part b code

```

1
2 #min
3 from pyspark import SparkContext
4 sc1 = SparkContext(appName = "1 min temperatures descending")
5 reading = sc1.textFile("/user/x_carsa/data/temperature-readings.csv")
6 separate = reading.map(lambda a: a.split(";"))
7 lines = separate.filter(lambda x: int(x[1][0:4]) >= 1950 and int(x[1][0:4]) <= 2014)
8 temperatures = lines.map(lambda x: (x[1][0:4], (x[0], float(x[3])))) ##added x[0] as it is the
    station column
9 maxTemperatures = temperatures.reduceByKey(lambda a,b: a if a[1]<b[1] else b)
10 maxTemperaturesTogether = maxTemperatures.repartition(1)
11 maxTemperaturesSorted = maxTemperaturesTogether.sortBy(ascending=True, keyfunc=lambda k: k
    [1][1]) ## the first [1] stands for k = length 1 and vector length 2, with the second [1,
    we acces to the second element of v]
12 maxTemperaturesSorted.saveAsTextFile("results/1resultsMin")

```

#### Output:

```

1
2 Row(_c0=u'1990', _c1=u'133260', value=-35.0)
3 Row(_c0=u'1952', _c1=u'144300', value=-35.5)
4 Row(_c0=u'1974', _c1=u'102190', value=-35.6)
5 Row(_c0=u'1954', _c1=u'107440', value=-36.0)
6 Row(_c0=u'1992', _c1=u'116430', value=-36.1)
7 Row(_c0=u'1975', _c1=u'136010', value=-37.0)
8 Row(_c0=u'1972', _c1=u'188800', value=-37.5)
9 Row(_c0=u'2000', _c1=u'146350', value=-37.6)
10 Row(_c0=u'1995', _c1=u'107400', value=-37.6)
11 Row(_c0=u'1957', _c1=u'108640', value=-37.8)
12 Row(_c0=u'1983', _c1=u'135520', value=-38.2)
13 Row(_c0=u'1989', _c1=u'116430', value=-38.2)
14 Row(_c0=u'1953', _c1=u'134110', value=-38.4)
15 Row(_c0=u'2009', _c1=u'116430', value=-38.5)
16 Row(_c0=u'1993', _c1=u'107400', value=-39.0)
17 Row(_c0=u'1984', _c1=u'140480', value=-39.2)
18 Row(_c0=u'1973', _c1=u'112080', value=-39.3)
19 Row(_c0=u'2008', _c1=u'116430', value=-39.3)
20 Row(_c0=u'1991', _c1=u'107400', value=-39.3)
21 Row(_c0=u'2005', _c1=u'116430', value=-39.4)
22 Row(_c0=u'1961', _c1=u'102190', value=-39.5)
23 Row(_c0=u'1964', _c1=u'188830', value=-39.5)
24 Row(_c0=u'1970', _c1=u'149160', value=-39.6)
25 Row(_c0=u'2004', _c1=u'135460', value=-39.7)
26 Row(_c0=u'1988', _c1=u'158750', value=-39.9)
27 Row(_c0=u'1960', _c1=u'151290', value=-40.0)
28 Row(_c0=u'1997', _c1=u'188790', value=-40.2)
29 Row(_c0=u'1994', _c1=u'116430', value=-40.5)

```

```

30 Row(_c0=u'2006', _c1=u'102190', value=-40.6)
31 Row(_c0=u'2007', _c1=u'158740', value=-40.7)
32 Row(_c0=u'2013', _c1=u'116430', value=-40.7)
33 Row(_c0=u'1963', _c1=u'133630', value=-41.0)
34 Row(_c0=u'1955', _c1=u'139340', value=-41.2)
35 Row(_c0=u'2003', _c1=u'158750', value=-41.5)
36 Row(_c0=u'1969', _c1=u'105230', value=-41.5)
37 Row(_c0=u'1996', _c1=u'146070', value=-41.7)
38 Row(_c0=u'2010', _c1=u'146050', value=-41.7)
39 Row(_c0=u'2011', _c1=u'146050', value=-42.0)
40 Row(_c0=u'1962', _c1=u'136010', value=-42.0)
41 Row(_c0=u'1950', _c1=u'180750', value=-42.0)
42 Row(_c0=u'1951', _c1=u'188830', value=-42.0)
43 Row(_c0=u'1968', _c1=u'133470', value=-42.0)
44 Row(_c0=u'1982', _c1=u'135520', value=-42.2)
45 Row(_c0=u'2002', _c1=u'102190', value=-42.2)
46 Row(_c0=u'1976', _c1=u'136010', value=-42.2)
47 Row(_c0=u'2014', _c1=u'158740', value=-42.5)
48 Row(_c0=u'1977', _c1=u'105230', value=-42.5)
49 Row(_c0=u'1998', _c1=u'133260', value=-42.7)
50 Row(_c0=u'2012', _c1=u'116490', value=-42.7)
51 Row(_c0=u'1958', _c1=u'102190', value=-43.0)
52 Row(_c0=u'1985', _c1=u'180940', value=-43.4)
53 Row(_c0=u'1959', _c1=u'134110', value=-43.6)
54 Row(_c0=u'1981', _c1=u'158750', value=-44.0)
55 Row(_c0=u'2001', _c1=u'116430', value=-44.0)
56 Row(_c0=u'1965', _c1=u'151290', value=-44.0)
57 Row(_c0=u'1979', _c1=u'173960', value=-44.0)
58 Row(_c0=u'1986', _c1=u'140480', value=-44.2)
59 Row(_c0=u'1971', _c1=u'112080', value=-44.3)
60 Row(_c0=u'1980', _c1=u'140480', value=-45.0)
61 Row(_c0=u'1956', _c1=u'144300', value=-45.0)
62 Row(_c0=u'1967', _c1=u'173960', value=-45.4)
63 Row(_c0=u'1987', _c1=u'158750', value=-47.3)
64 Row(_c0=u'1978', _c1=u'133470', value=-47.7)
65 Row(_c0=u'1999', _c1=u'158750', value=-49.0)
66 Row(_c0=u'1966', _c1=u'108640', value=-49.4)

```

## Part c code

```

1 from functools import reduce
2 from itertools import groupby
3
4
5 file = open("/nfshome/x_carsa/Desktop/data/temperatures-big.csv", mode = "r")
6
7 working_data = []
8 for line in file:
9     temp = line.split(";")
10    if int(temp[1][0:4]) >= 1950 and int(temp[1][0:4]) <= 2014:
11        working_data.append((temp[1][0:4], (temp[0], float(temp[3]))))
12    else:
13        continue
14
15
16
17 def reduceByKey(func, iterable):
18     first = lambda p: p[0]
19     second = lambda p: p[1]
20     return map(lambda l: (l[0], reduce(func, map(second, l[1]))), groupby(sorted(iterable, key=
21         first), first))
22
23
24 maxData = reduceByKey(lambda x ,y: x if x[1] > y[1] else y, working_data)
25
26 for line in maxData:
27     print line

```

Output: It last for about 40 minutes whereas in the previous case it last much less (about a minute or two). Such differences are related to the non-usage of parallel execution of the algorithm rather than for a single computer with its nodes.

## Exercise 2

Count the number of readings for each month in the period of 1950-2014 which are higher than 10 degrees.

```

1 from pyspark import SparkContext
2 sc = SparkContext(appName = "2 number of readings each month")
3 reading = sc.textFile("/user/x_carsa/data/temperature-readings.csv")
4 separate = reading.map(lambda a: a.split(";"))
5 lines = separate.filter(lambda x: int(x[1][0:4]) >= 1950 and int(x[1][0:4]) <= 2014 and float(x
    [3]) > 10)
6 higher10temperatures = lines.map(lambda x: (x[1][0:7], 1)) ##added as it is the station column
7 counts = higher10temperatures.reduceByKey(lambda v1,v2: v1 + v2) ##summing counts of every one
8 maxTemperaturestogether = counts.repartition(1) ## putting everything in one partition
9
10 maxTemperaturestogether.sortBy(ascending = False, keyfunc= lambda x : x[1]) \
11     .saveAsTextFile("results/2resultshigher10")
12 print maxTemperaturestogether.take(20)

```

## Output:

```

1 (u'2014-07', 147681)
2 (u'2011-07', 146656)
3 (u'2010-07', 143419)
4 (u'2012-07', 137477)
5 (u'2013-07', 133657)
6 (u'2009-07', 133008)
7 (u'2011-08', 132734)
8 (u'2009-08', 128349)
9 (u'2013-08', 128235)
10 (u'2003-07', 128133)
11 (u'2002-07', 127956)
12 (u'2006-08', 127622)
13 (u'2008-07', 126973)
14 (u'2002-08', 126073)
15 (u'2005-07', 125294)
16 (u'2011-06', 125193)
17 (u'2012-08', 125037)
18 (u'2006-07', 124794)
19 (u'2010-08', 124417)
20 (u'2014-08', 124045)
21 (u'1997-07', 123496)
22 (u'2007-07', 123218)
23 (u'2013-06', 122181)
24 (u'1997-08', 121154)
25 (u'2001-07', 120529)
26 (u'1998-07', 120230)
27 (u'2000-07', 119769)
28 (u'2004-07', 119536)
29 (u'1999-07', 116385)
30 (u'2008-08', 114272)
31 (u'2004-08', 114168)
32 (u'2002-06', 114034)
33 (u'2005-08', 113950)
34 (u'2001-08', 113937)
35 (u'2007-08', 110428)
36 (u'2000-08', 109201)
37 (u'2003-08', 108501)
38 (u'1996-08', 107758)
39 (u'1997-06', 104696)
40 (u'1999-06', 103227)
41 (u'2007-06', 103046)
42 (u'2008-06', 102900)
43 (u'2010-06', 102716)
44 (u'2006-06', 102588)
45 (u'2014-06', 101711)
46 (u'1998-08', 101387)
47 (u'1996-07', 99916)
48 (u'2003-06', 99693)
49 (u'2011-09', 99335)
50 (u'1999-08', 97437)
51 (u'2006-09', 97181)
52 (u'2012-06', 94513)
53 (u'2001-06', 93375)
54 (u'2005-06', 90724)
55 (u'2004-06', 89628)
56 (u'1999-09', 89418)
57 (u'2009-09', 89106)
58 (u'2009-06', 87787)
59 (u'2000-06', 86592)
60 (u'2014-09', 86090)
61 (u'1998-06', 82608)
62 (u'2013-05', 81996)
63 (u'2013-09', 81960)
64 (u'1996-06', 80440)
65 (u'2001-09', 79657)
66 (u'1998-09', 76535)
67 (u'1988-07', 75521)
68 (u'2005-09', 75494)
69 (u'2010-09', 74816)

```

```

70 (u'1997-09', 74472)
71 (u'1991-07', 73385)
72 (u'2004-09', 73334)
73 (u'1973-07', 71522)
74 (u'1991-08', 71185)
75 (u'2003-09', 70459)
76 (u'2012-09', 70427)
77 (u'1990-07', 70031)
78 (u'1988-08', 69913)
79 (u'1987-07', 68135)
80 (u'1989-07', 67880)
81 (u'1989-08', 67793)
82 (u'1990-08', 67604)
83 (u'1995-08', 66920)
84 (u'1974-07', 66277)
85 (u'2002-05', 66116)
86 (u'2002-09', 65928)
87 (u'1974-08', 64470)
88 (u'1975-07', 64408)
89 (u'1976-07', 64109)
90 (u'2000-09', 63837)
91 (u'1988-06', 63572)
92 (u'1992-07', 62911)
93 (u'1975-08', 62565)
94 (u'2007-09', 61346)
95 (u'1978-07', 60998)
96 (u'2008-09', 60989)
97 (u'1976-08', 60898)
98 (u'2009-05', 60867)
99 (u'1989-06', 60822)
100 (u'1979-07', 60719)
101 (u'1994-07', 60691)

```

Part b: 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.

```

1 from pyspark import SparkContext
2 sc = SparkContext(appName = "2 number of readings each month one per station max")
3 reading = sc.textFile("/user/x_carsa/data/temperature-readings.csv")
4 separate = reading.map(lambda a: a.split(";"))
5 lines = separate.filter(lambda x: int(x[1][0:4]) >= 1950 and int(x[1][0:4]) <= 2014 and float(x
    [3]) > 10)
6 higher10temperatures = lines.map(lambda x: (x[1][0:7], (x[0], 1))).distinct() ##DISTINCT = one
    count for each station
7
8 counts = higher10temperatures.reduceByKey(lambda v1, v2: (v1[0], v1[1] + v2[1])) ##summing counts
    of every one
9 counts = counts.map(lambda x: (x[0], x[1][1]))
10 maxTemperaturestogether = counts.repartition(1) ## putting everything in one partition
11 ## putting everything in one partition
12
13 maxTemperaturestogether.sortBy(ascending= False, keyfunc= lambda x: x[1])\
14     .saveAsTextFile("results/2resultshigher10one_eachmonth")
15 print maxTemperaturestogether.take(20)

```

Output:

```

1 (u'1972-10', 378)
2 (u'1973-06', 377)
3 (u'1973-05', 377)
4 (u'1973-09', 376)
5 (u'1972-08', 376)
6 (u'1972-05', 375)
7 (u'1971-08', 375)
8 (u'1972-06', 375)
9 (u'1972-09', 375)
10 (u'1971-09', 374)
11 (u'1972-07', 374)
12 (u'1971-06', 374)
13 (u'1973-08', 373)
14 (u'1971-05', 373)
15 (u'1974-06', 372)
16 (u'1974-08', 372)
17 (u'1974-05', 370)
18 (u'1970-08', 370)
19 (u'1971-07', 370)
20 (u'1973-07', 370)
21 (u'1974-09', 370)
22 (u'1975-09', 369)
23 (u'1970-09', 369)

```

24 (u'1976-05', 369)  
25 (u'1970-06', 369)  
26 (u'1976-06', 368)  
27 (u'1975-06', 368)  
28 (u'1975-08', 367)  
29 (u'1975-05', 367)  
30 (u'1970-05', 366)  
31 (u'1976-09', 365)  
32 (u'1977-06', 364)  
33 (u'1967-05', 363)  
34 (u'1976-08', 363)  
35 (u'1974-07', 362)  
36 (u'1970-07', 362)  
37 (u'1967-09', 361)  
38 (u'1966-09', 360)  
39 (u'1966-06', 360)  
40 (u'1966-08', 359)  
41 (u'1969-09', 359)  
42 (u'1967-06', 359)  
43 (u'1965-09', 358)  
44 (u'1978-09', 358)  
45 (u'1967-08', 358)  
46 (u'1975-07', 358)  
47 (u'1969-08', 357)  
48 (u'1968-06', 357)  
49 (u'1968-08', 357)  
50 (u'1976-07', 356)  
51 (u'1968-09', 356)  
52 (u'1968-05', 355)  
53 (u'1965-06', 355)  
54 (u'1979-05', 354)  
55 (u'1978-06', 354)  
56 (u'1965-08', 354)  
57 (u'1966-05', 354)  
58 (u'1977-08', 354)  
59 (u'1968-07', 353)  
60 (u'1977-09', 353)  
61 (u'1978-05', 352)  
62 (u'1969-06', 352)  
63 (u'1966-07', 352)  
64 (u'1967-07', 351)  
65 (u'1979-06', 351)  
66 (u'1977-05', 351)  
67 (u'1979-09', 351)  
68 (u'1977-07', 350)  
69 (u'1978-08', 350)  
70 (u'1965-07', 349)  
71 (u'1973-10', 349)  
72 (u'1969-07', 349)  
73 (u'1971-10', 347)  
74 (u'1969-10', 346)  
75 (u'1979-07', 345)  
76 (u'1996-06', 345)  
77 (u'1970-10', 345)  
78 (u'1974-04', 344)  
79 (u'1965-05', 344)  
80 (u'1978-07', 343)  
81 (u'1996-07', 342)  
82 (u'1996-05', 342)  
83 (u'1996-08', 341)  
84 (u'1978-10', 340)  
85 (u'1996-09', 340)  
86 (u'1975-10', 340)  
87 (u'1979-08', 340)  
88 (u'1997-09', 340)  
89 (u'1982-06', 339)  
90 (u'1997-06', 338)  
91 (u'1980-09', 338)  
92 (u'1980-05', 337)  
93 (u'1981-05', 337)  
94 (u'1997-08', 337)  
95 (u'1983-06', 337)  
96 (u'1983-05', 336)  
97 (u'1965-10', 335)  
98 (u'1981-09', 335)  
99 (u'1969-05', 335)  
100 (u'1981-08', 334)  
101 (u'1982-09', 334)  
102 (u'1997-07', 333)  
103 (u'1984-05', 333)  
104 (u'1983-09', 332)  
105 (u'1980-06', 332)  
106 (u'1981-06', 331)  
107 (u'1999-06', 330)  
108 (u'1983-08', 330)  
109 (u'1982-05', 330)  
110 (u'1980-08', 330)

```

111 (u'1999-07', 329)
112 (u'1981-07', 329)
113 (u'1999-09', 328)
114 (u'1985-09', 327)
115 (u'1984-09', 327)
116 (u'1999-08', 327)
117 (u'1998-09', 326)
118 (u'1998-08', 326)
119 (u'2002-06', 326)
120 (u'1998-07', 326)
121 (u'1982-08', 326)
122 (u'1998-06', 326)
123 (u'1981-10', 325)
124 (u'1999-05', 325)
125 (u'2000-08', 325)
126 (u'1985-05', 325)
127 (u'1980-07', 324)
128 (u'1967-10', 324)
129 (u'1984-06', 324)
130 (u'2001-07', 324)
131 (u'2002-07', 324)
132 (u'2001-06', 324)
133 (u'1985-06', 324)
134 (u'2002-05', 324)
135 (u'1987-06', 323)
136 (u'2003-06', 323)
137 (u'2000-05', 323)
138 (u'2002-09', 323)
139 (u'2001-08', 323)
140 (u'1986-09', 323)
141 (u'1987-09', 323)
142 (u'2002-08', 322)
143 (u'2001-09', 322)
144 (u'1968-04', 322)
145 (u'1998-05', 322)
146 (u'2000-09', 322)
147 (u'1988-06', 322)
148 (u'2003-05', 321)
149 (u'2004-05', 321)
150 (u'2003-07', 321)
151 (u'1984-10', 321)
152 (u'1982-07', 321)
153 (u'2000-06', 321)
154 (u'1991-06', 321)
155 (u'2004-09', 321)
156 (u'1987-05', 320)
157 (u'2010-06', 320)
158 (u'2000-07', 320)
159 (u'1988-05', 320)
160 (u'2003-09', 320)
161 (u'2004-08', 320)
162 (u'1987-08', 320)
163 (u'2003-08', 320)
164 (u'1997-05', 319)
165 (u'1987-07', 319)
166 (u'2004-06', 319)
167 (u'2004-07', 319)
168 (u'2010-05', 319)
169 (u'2011-07', 319)
170 (u'1983-07', 319)
171 (u'2010-07', 318)

```

## Exercise 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.

```

1 from pyspark import SparkContext
2
3 sc = SparkContext(appName = "3 average temperature for each month")
4 reading = sc.textFile("/user/x_carsa/data/temperature-readings.csv")
5 separate = reading.map(lambda a: a.split(";"))
6 lines = separate.filter(lambda x: int(x[1][0:4]) >= 1960 and int(x[1][0:4]) <= 2014)
7 temperatures = lines.map(lambda x: ((x[1], x[0]), (float(x[3]), float(x[3])))) ##added x[0] as
   it is the station column, and one to divide it at then end by the averages
8 minmax = temperatures.reduceByKey(lambda v1,v2: (v1[0] if v1[0]>v2[0] else v2[0], v1[1] if v1
   [1]<v2[1] else v2[1])) # first is max, second is min, done just to get the max and min per
   day in case there is more than one

```

```

9 calculations= minmax.map(lambda x: ((x[0][0][0:7], x[0][1]), (x[1][0], x[1][1], 1))) # getting
    together by index and adding 1 to be able to divide by counts
10 averageperday = calculations.reduceByKey(lambda v1,v2: (v1[0] + v2[0],v1[1] + v2[1], v1[2]+ v2
    [2])) #
11 averageperday = averageperday.map(lambda x: ((x[0]), ((x[1][0]+x[1][1])/(2*x[1][2]))))
12 temperaturestogether = averageperday.repartition(1)
13
14 print temperaturestogether.take(20)
15
16 temperaturestogether.sortBy(ascending = False, keyfunc= lambda x : x[1]) \
17     .saveAsTextFile("results/3average_temperature_eachmonth")

```

## Output:

```

1 ((u'2014-07', u'96000'), 26.3)
2 ((u'1994-07', u'96550'), 23.071052631578944)
3 ((u'1983-08', u'54550'), 23.0)
4 ((u'1994-07', u'78140'), 22.970967741935482)
5 ((u'1994-07', u'85280'), 22.872580645161293)
6 ((u'1994-07', u'75120'), 22.85806451612903)
7 ((u'1994-07', u'65450'), 22.856451612903225)
8 ((u'1994-07', u'96000'), 22.80806451612903)
9 ((u'1994-07', u'95160'), 22.76451612903226)
10 ((u'1994-07', u'86200'), 22.711290322580645)
11 ((u'2002-08', u'78140'), 22.700000000000003)
12 ((u'1994-07', u'76000'), 22.698387096774198)
13 ((u'1997-08', u'78140'), 22.666129032258066)
14 ((u'1994-07', u'105260'), 22.65967741935484)
15 ((u'1975-08', u'54550'), 22.642857142857142)
16 ((u'2006-07', u'76530'), 22.598387096774193)
17 ((u'1994-07', u'86330'), 22.54838709677419)
18 ((u'2006-07', u'75120'), 22.52741935483871)
19 ((u'1994-07', u'54300'), 22.469354838709677)
20 ((u'2006-07', u'78140'), 22.45806451612903)
21 ((u'2001-07', u'96550'), 22.408333333333335)
22 ((u'2010-07', u'98180'), 22.37903225806452)
23 ((u'2006-07', u'65450'), 22.37741935483871)
24 ((u'1994-07', u'85210'), 22.375806451612902)
25 ((u'1994-07', u'98180'), 22.367741935483874)
26 ((u'2014-07', u'98180'), 22.367741935483874)
27 ((u'2002-08', u'98180'), 22.366129032258062)
28 ((u'1994-07', u'92100'), 22.31774193548387)
29 ((u'1994-07', u'86470'), 22.30806451612903)
30 ((u'1994-07', u'83230'), 22.272580645161288)
31 ((u'1994-07', u'64290'), 22.259677419354837)
32 ((u'1994-07', u'97490'), 22.258064516129032)
33 ((u'1994-07', u'94180'), 22.25322580645161)
34 ((u'1972-07', u'173960'), 22.244999999999997)
35 ((u'1994-07', u'74080'), 22.241935483870968)
36 ((u'2006-07', u'54300'), 22.23709677419355)
37 ((u'2002-08', u'98210'), 22.23548387096774)
38 ((u'1994-07', u'106070'), 22.232258064516135)
39 ((u'1994-07', u'75100'), 22.229032258064517)
40 ((u'1994-07', u'53440'), 22.197499999999998)
41 ((u'1994-07', u'83270'), 22.177419354838705)
42 ((u'1994-07', u'103080'), 22.16451612903226)
43 ((u'1994-07', u'82110'), 22.161290322580644)
44 ((u'1994-07', u'97120'), 22.135483870967743)
45 ((u'2010-07', u'98210'), 22.111290322580647)
46 ((u'1994-07', u'53430'), 22.096774193548388)
47 ((u'1997-08', u'86330'), 22.07903225806452)
48 ((u'2006-07', u'66500'), 22.05483870967742)
49 ((u'1994-07', u'76530'), 22.033870967741933)
50 ((u'1997-08', u'98210'), 21.983870967741936)
51 ((u'2014-07', u'98210'), 21.962903225806453)
52 ((u'1994-07', u'62400'), 21.951612903225808)
53 ((u'1997-08', u'62400'), 21.938709677419357)
54 ((u'1994-07', u'108110'), 21.90806451612903)
55 ((u'1994-07', u'83130'), 21.900000000000002)
56 ((u'1997-08', u'98180'), 21.887096774193544)
57 ((u'2006-07', u'98210'), 21.872580645161293)
58 ((u'1997-08', u'98290'), 21.864516129032257)
59 ((u'1991-08', u'78040'), 21.85)
60 ((u'2010-07', u'78140'), 21.830645161290324)
61 ((u'1994-07', u'63340'), 21.7758064516129)
62 ((u'1994-07', u'91130'), 21.76290322580645)
63 ((u'1994-07', u'105370'), 21.761290322580642)
64 ((u'2008-07', u'83420'), 21.75)
65 ((u'1994-07', u'64130'), 21.72741935483871)
66 ((u'1997-08', u'96550'), 21.725806451612904)
67 ((u'1994-07', u'83440'), 21.716129032258067)
68 ((u'2006-07', u'85210'), 21.706451612903226)
69 ((u'1994-07', u'74420'), 21.690322580645162)
70 ((u'2003-07', u'98180'), 21.68548387096774)

```



```
71 ((u'1997-08', u'52230'), 21.68548387096774)
```

## Exercise 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 200 mm. In this exercise you will use the temperature-readings.csv and precipitation-readings.csv file.

```
1 from pyspark import SparkContext
2
3 sc = SparkContext(appName = "4 average temperature for each month")
4 filetemp = sc.textFile("/user/x_carsa/data/temperature-readings.csv")
5
6 temp = filetemp.map(lambda a: a.split(";"))
7
8
9 calculatedTemp = temp.map(lambda x: ((x[0]), (x[3])))\
10     .filter(lambda (station_number, temp): (temp >= 25 and temp <=30))\
11     .reduceByKey(max)\
12     .repartition(1)\
13     .sortByKey()
14
15
16 fileprec = sc.textFile("/user/x_carsa/data/precipitation-readings.csv")
17 prec = fileprec.map(lambda a: a.split(";"))
18
19 calculatedPrec = prec.map(lambda x: ((x[0],x[1]), float(x[3])))\
20     .reduceByKey(lambda x1, x2: (x1+x2))\
21     .filter(lambda x: (x >= 100 and x <=200))\
22     .reduceByKey(max)\
23     .map(lambda ((station_number, date), prec): (station_number, prec))\
24     .repartition(1)\
25     .sortByKey()
26 print calculatedPrec.take(20)
27
28 together = calculatedPrec.join(calculatedTemp)\
29     .repartition(1)\
30     .sortBy(ascending= False, keyfunc= lambda x: (x[0]))
31
32 print together.take(20)
33 together.saveAsTextFile("results/4stations_with_max_temperatures_precipitations")
```

Output:

```
1 No output
```

## Exercise 5

Calculate the average monthly precipitation for the Östergötland region (list of stations is provided in the separate file). In order to do this, you will first need to calculate the total daily precipitation before calculating the monthly average. In this exercise you will use the precipitation-readings.csv and stations-Ostergotland.csv files.

```
1 from pyspark import SparkContext
2 sc = SparkContext(appName = "5 average monthly precipitation for the Ostergotland")
3
4 reading2 = sc.textFile("/user/x_carsa/data/stations-Ostergotland.csv")
5 separate2 = reading2.map(lambda a: a.split(";"))
6 stations = separate2.map(lambda observation: int(observation[0]))
7 stations = stations.distinct().collect() #collect transforms the rdd to a python list object
8 stations = {station: True for station in stations}
9
10
11 reading1 = sc.textFile("/user/x_carsa/data/precipitation-readings.csv")#x[0] = station number
12
13
14 percip_daily = reading1.map(lambda line: line.split(";")) \
```

```

15         .filter(lambda obs: stations.get(int(obs[0]), False)) \
16         .map(lambda obs: (obs[1], float(obs[3]))) \
17         .reduceByKey(lambda value1, value2: (value1 + value2)) #total per day
18
19 counting = percip_daily.map(lambda x: (x[0][0:7], (float(x[1]), 1)))
20 average = counting.reduceByKey(lambda v1, v2: (v1[0]+v2[0], v1[1] + v2[1])) #total
21     precipitation and number of days
22 average = average.map(lambda x: ((x[0]), (x[1][0]/x[1][1])))
23 print average.take(20)
24 results = average.repartition(1) \
25     .sortBy(ascending= False, keyfunc= lambda x: x[0]) \
26     .saveAsTextFile("results/5_average_monthly_precipitation")

```

## Output:

```

1 (u'2016-07', 0.0)
2 (u'2016-06', 12.710000000000004)
3 (u'2016-05', 7.548387096774194)
4 (u'2016-04', 7.173333333333335)
5 (u'2016-03', 5.151612903225806)
6 (u'2016-02', 5.948275862068965)
7 (u'2016-01', 5.761290322580644)
8 (u'2015-12', 7.4645161290322575)
9 (u'2015-11', 17.036666666666665)
10 (u'2015-10', 0.5838709677419355)
11 (u'2015-09', 27.013333333333335)
12 (u'2015-08', 6.9645161290322575)
13 (u'2015-07', 30.73548387096774)
14 (u'2015-06', 20.976666666666667)
15 (u'2015-05', 24.058064516129033)
16 (u'2015-04', 4.09)
17 (u'2015-03', 10.99677419354839)
18 (u'2015-02', 7.092857142857143)
19 (u'2015-01', 15.254838709677419)
20 (u'2014-12', 9.151612903225804)
21 (u'2014-11', 13.979999999999999)
22 (u'2014-10', 18.616129032258065)
23 (u'2014-09', 12.920000000000003)
24 (u'2014-08', 23.43548387096774)
25 (u'2014-07', 5.932258064516129)
26 (u'2014-06', 20.036666666666667)
27 (u'2014-05', 14.96774193548387)
28 (u'2014-04', 8.469999999999999)
29 (u'2014-03', 9.435483870967742)
30 (u'2014-02', 12.489285714285716)
31 (u'2014-01', 16.148387096774194)
32 (u'2013-12', 10.90645161290323)
33 (u'2013-11', 12.366666666666667)
34 (u'2013-10', 13.903225806451612)
35 (u'2013-09', 6.983333333333336)
36 (u'2013-08', 13.954838709677418)
37 (u'2013-07', 14.080645161290322)
38 (u'2013-06', 16.353333333333335)
39 (u'2013-05', 12.36774193548387)
40 (u'2013-04', 10.209999999999999)
41 (u'2013-03', 1.9064516129032258)
42 (u'2013-02', 7.292857142857143)
43 (u'2013-01', 5.554838709677419)
44 (u'2012-12', 12.954838709677418)
45 (u'2012-11', 13.73)
46 (u'2012-10', 12.693548387096774)
47 (u'2012-09', 14.55)
48 (u'2012-08', 13.319354838709677)
49 (u'2012-07', 11.432258064516128)
50 (u'2012-06', 26.44)
51 (u'2012-05', 4.445161290322581)
52 (u'2012-04', 12.556666666666668)
53 (u'2012-03', 1.6548387096774195)
54 (u'2012-02', 5.931034482758621)
55 (u'2012-01', 8.429032258064517)
56 (u'2011-12', 8.154838709677417)
57 (u'2011-11', 2.6933333333333334)
58 (u'2011-10', 8.46774193548387)
59 (u'2011-09', 10.513333333333334)
60 (u'2011-08', 16.69677419354839)
61 (u'2011-07', 18.370967741935484)
62 (u'2011-06', 17.669999999999998)
63 (u'2011-05', 7.325806451612904)
64 (u'2011-04', 2.983333333333333)
65 (u'2011-03', 3.838709677419355)
66 (u'2011-02', 5.253571428571428)
67 (u'2011-01', 6.800000000000001)
68 (u'2010-12', 7.196774193548386)
69 (u'2010-11', 18.71)
70 (u'2010-10', 10.167741935483871)

```

```

71 (u'2010-09', 8.616666666666667)
72 (u'2010-08', 20.91290322580645)
73 (u'2010-07', 17.883870967741938)
74 (u'2010-06', 9.729999999999999)
75 (u'2010-05', 12.999999999999998)
76 (u'2010-04', 4.756666666666666)
77 (u'2010-03', 4.622580645161291)
78 (u'2010-02', 11.30357142857143)
79 (u'2010-01', 6.964516129032259)
80 (u'2009-12', 10.345161290322583)
81 (u'2009-11', 12.843333333333334)
82 (u'2009-10', 11.0)
83 (u'2009-09', 5.989999999999999)
84 (u'2009-08', 11.916129032258064)
85 (u'2009-07', 21.903225806451616)
86 (u'2009-06', 9.953333333333335)
87 (u'2009-05', 10.483870967741936)
88 (u'2009-04', 0.56)
89 (u'2009-03', 6.674193548387096)
90 (u'2009-02', 5.310714285714288)
91 (u'2009-01', 3.074193548387097)
92 (u'2008-12', 8.416129032258064)
93 (u'2008-11', 9.350000000000001)
94 (u'2008-10', 11.52903225806452)
95 (u'2008-09', 9.473333333333333)
96 (u'2008-08', 26.809677419354838)
97 (u'2008-07', 16.490322580645163)
98 (u'2008-06', 8.586666666666668)
99 (u'2008-05', 4.47741935483871)
100 (u'2008-04', 4.05)
101 (u'2008-03', 8.167741935483871)
102 (u'2008-02', 5.844827586206897)
103 (u'2008-01', 8.703225806451613)
104 (u'2007-12', 10.587096774193547)
105 (u'2007-11', 10.136666666666667)
106 (u'2007-10', 5.441935483870968)
107 (u'2007-09', 12.373333333333335)
108 (u'2007-08', 10.483870967741936)
109 (u'2007-07', 18.5741935483871)
110 (u'2007-06', 21.790000000000003)
111 (u'2007-05', 7.8419354838709685)
112 (u'2007-04', 4.250000000000001)
113 (u'2007-03', 7.8419354838709685)
114 (u'2007-02', 7.085714285714286)
115 (u'2007-01', 13.283870967741937)
116 (u'2006-12', 5.75483870967742)
117 (u'2006-11', 14.343333333333334)
118 (u'2006-10', 22.870967741935484)
119 (u'2006-09', 3.8533333333333335)
120 (u'2006-08', 28.661290322580644)
121 (u'2006-07', 5.6096774193548375)
122 (u'2006-06', 6.2266666666666675)
123 (u'2006-05', 10.129032258064516)
124 (u'2006-04', 8.873333333333333)
125 (u'2006-03', 5.393548387096776)
126 (u'2006-02', 7.44642857142857)
127 (u'2006-01', 4.08076923076923)
128 (u'2005-12', 10.961290322580645)
129 (u'2005-11', 6.52)
130 (u'2005-10', 7.3645161290322605)
131 (u'2005-09', 2.79)

```

## Exercise 6

Compare the average monthly temperature (find the difference) in the period 1950-2014 for each station in Östergötland with long-term monthly averages in the period of 1950-1980. Make a plot of your results.

```

1 from pyspark import SparkContext
2 sc = SparkContext(appName = "6 Comparison Östergötland average monthly temperature 1950-2014
   with long-term monthly averages in the period of 1950-1980 ")
3
4 reading2 = sc.textFile("/user/x_carsa/data/stations-Östergötland.csv")
5
6 separate2 = reading2.map(lambda a: a.split(";"))
7 stations = separate2.map(lambda x: int(x[0]))
8 stations = stations.distinct().collect() # Taken different stations as in previous step
9 stations = {station: True for station in stations}
10
11 #Extracting temperatures

```

```

12
13
14 reading2 = sc.textFile("/user/x_carsa/data/temperature-readings.csv")
15 myfilter = reading2.map(lambda a: a.split(";")) \
16     .filter(lambda x: (int(x[1][0:4]) >= 1950 and int(x[1][0:4]) <= 2014 and stations
17         .get(int(x[0]), False)))
18 average_short_by_station = myfilter.map(lambda x: ((x[1], int(x[0])), (float(x[3]), float(x[3]))
19     )) \
20     .reduceByKey(lambda (min1, max1), (min2, max2): (min(min1, min2), max(
21         max1, max2))) \
22     .map(lambda ((date, station), (tmin, tmax)): ((date[0:7], station), (
23         tmin+tmax, 2))) \
24     .reduceByKey(lambda (temp1, count1), (temp2, count2): (temp1 + temp2,
25         count1 + count2)) \
26     .map(lambda ((date, station), (ttemp, tcount)): (date, station, ttemp/
27         float(tcount)))
28 average_short = average_short_by_station.map(lambda (date, station, avgtemp): (date, (avgtemp,
29     1))) \
30     .reduceByKey(lambda (temp1, count1), (temp2, count2): (temp1 + temp2,
31         count1 + count2)) \
32     .map(lambda (date, (ttemp, tcount)): (date, ttemp/float(tcount)))
33
34 mylong = avglong.collect()
35 avgMonth = {month: temp for (month, temp) in mylong} # first part (month: temp, = structure type
36     , for = old structure type)
37 together = average_short.map(lambda x: (x[0], abs(x[1]) - abs(avgMonth.get(x[0][5:7], 0))))
38
39 together.sortBy(ascending= False, keyfunc= lambda x: x[0], numPartitions = 1)\
40     .saveAsTextFile("results/6_average_diff")

```

Output:

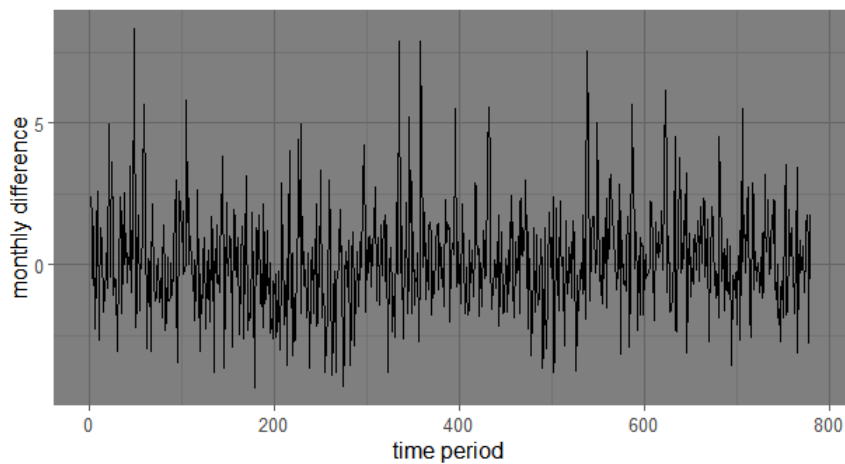


Figure 1: Average monthly difference on temperature