

TDDE31 - Lab 3 - Improved

Sindre Jonsson Wold (sinwo612),
Jonathan Vikbladh (jonvi252)

VT 2021

Code

```
1 from __future__ import division
2 from math import radians, cos, sin, asin, sqrt, exp
3 from datetime import datetime
4 from pyspark import SparkContext
5
6 sc = SparkContext(appName="lab_kernel")
7
8 def haversine(lon1, lat1, lon2, lat2):
9     """
10     Calculate the great circle distance between two points
11     on the earth (specified in decimal degrees)
12     """
13
14     # convert decimal degrees to radians
15
16     lon1, lat1, lon2, lat2 = map(radians, [lon1, lat1, lon2, lat2])
17
18     # haversine formula
19
20     dlon = lon2 - lon1
21     dlat = lat2 - lat1
22     a = sin(dlat/2)**2 + cos(lat1)*cos(lat2)*sin(dlon/2)**2
23     c = 2*asin(sqrt(a))
24     km = 6367*c
25     return km
26
27 def date_diff(date1, date2):
28     date1 = datetime.strptime(date1, '%Y-%m-%d')
29     date2 = datetime.strptime(date2, '%Y-%m-%d')
30
31     diff = (date1 - date2).days % 365
32
33     if diff > 182:
34         diff = (date2 - date1).days % 365
35
36     return diff
37
38 def time_diff(time1, time2):
39     time1 = datetime.strptime(time1, '%H:%M:%S')
40     time2 = datetime.strptime(time2, '%H:%M:%S')
41
42     diff_day = (time1-time2).days
43
44     mod_number = 60*60*24/2
45
46     if diff_day < 0:
47         diff=(time2 - time1).seconds % mod_number
48     else:
49         diff=(time1 - time2).seconds % mod_number
50
```

```

51     return diff
52
53 def kernel(pred, sample):
54
55     kernel_dist = exp(-(haversine(pred[0][0], pred[0][1], sample[0][0], sample[0][1])**2)/
56         h_distance)
57     kernel_date = exp(-(date_diff(pred[1], sample[1])**2)/h_date)
58     kernel_time = exp(-(time_diff(pred[2], sample[2])**2)/h_time)
59
60     #kernel_sum = kernel_dist + kernel_date + kernel_time
61     kernel_sum = kernel_dist * kernel_date * kernel_time
62
63     return kernel_sum
64
65 h_distance = 50# Up to you
66 h_date = 100# Up to you
67 h_time = 700000# Up to you
68 a = 58.4274 # Up to you
69 b = 14.826 # Up to you
70 date = "1998-09-21" # Up to you
71
72 stations = sc.textFile("BDA/input/stations.csv")
73 temps = sc.textFile("BDA/input/temperature-readings.csv")
74
75 station_lines = stations.map(lambda line: line.split(";")).cache()
76 temp_lines = temps.map(lambda line: line.split(";")).cache()
77
78 station_lat_long = station_lines.map(lambda p: (p[0], (float(p[3]), float(p[4]))))
79 station_lat_long.cache()
80 station_date_time_temp = temp_lines.map(lambda p: (p[0], (p[1], p[2], float(p[3]))))
81 station_date_time_temp.cache()
82
83 station_date_time_temp = station_date_time_temp.filter(lambda x: datetime.strptime(x[1][0],
84     '%Y-%m-%d') < datetime.strptime(date, '%Y-%m-%d'))
85 station_date_time_temp.cache()
86
87 station_lat_long = station_lat_long.collectAsMap()
88 bc = sc.broadcast(station_lat_long)
89
90 station_lat_long_date_time_temp = station_date_time_temp.map(lambda x: (x[0], (bc.value[x
91     [0]], x[1][0], x[1][1], x[1][2])))
92 station_lat_long_date_time_temp.cache()
93
94 for i, time in enumerate(["00:00:00", "22:00:00", "20:00:00", "18:00:00", "16:00:00", "
95     14:00:00", "12:00:00", "10:00:00", "08:00:00", "06:00:00", "04:00:00"]):
96
97     kernels = station_lat_long_date_time_temp.map(lambda x: (kernel(((a, b), date, time), (x
98         [1][0], x[1][1], x[1][2])), x[1][3]))
99     kernels.cache()
100     kernels = kernels.map(lambda x: (x[0]*x[1], x[0]))
101     kernels.cache()
102     kernels = kernels.reduce(lambda a, b: (a[0] + b[0], a[1] + b[1]))
103
104     temp = kernels[0]/kernels[1]
105     print(time, temp)

```

lab3.py

Results for date 1998-09-21

Sum kernel

```

('00:00:00', 5.996665147375697)
('22:00:00', 7.0872167104208055)
('20:00:00', 6.854560383883817)

```

('18:00:00', 4.510200963718676)
('16:00:00', 6.962538658555636)
('14:00:00', 7.137336495761705)
('12:00:00', 5.99699280172776)
('10:00:00', 7.086651582283834)
('08:00:00', 6.8541295124527)
('06:00:00', 4.50999808600833)
('04:00:00', 6.961718165730446)

Product kernel

('00:00:00', 11.579258846340272)
('22:00:00', 10.97037686654864)
('20:00:00', 9.804128770217604)
('18:00:00', 10.091994279981927)
('16:00:00', 11.321951202300768)
('14:00:00', 13.83763881521831)
('12:00:00', 11.579258846340272)
('10:00:00', 12.112140424475651)
('08:00:00', 10.964781068592096)
('06:00:00', 10.091994279981927)
('04:00:00', 8.791166547029102)

Answers to questions

Show that your choice for the kernels' width is sensible, i.e. it gives more weight to closer points. Discuss why your definition of closeness is reasonable.

We choose the 50 for the distance kernel, as a reasonable radius for relevant predictions is about 10 km, which this width achieves (see figure 1). Stations further away than about 10 km are assigned almost zero weight, and do not contribute to the prediction.

Similarly we choose the weight 100 for the date measure, as this assigns almost zero weight for measurements more than 15 days from the day of interest (see figure 2), which we think is reasonable considering the rate of which the temperature changes during the year.

Lastly, we assign the weight 700000 for the time measure, as this includes measurements inside a 1800 second window (see figure 3), corresponding to 30 minutes, which we think is reasonable considering the rate of temperature change during the day.

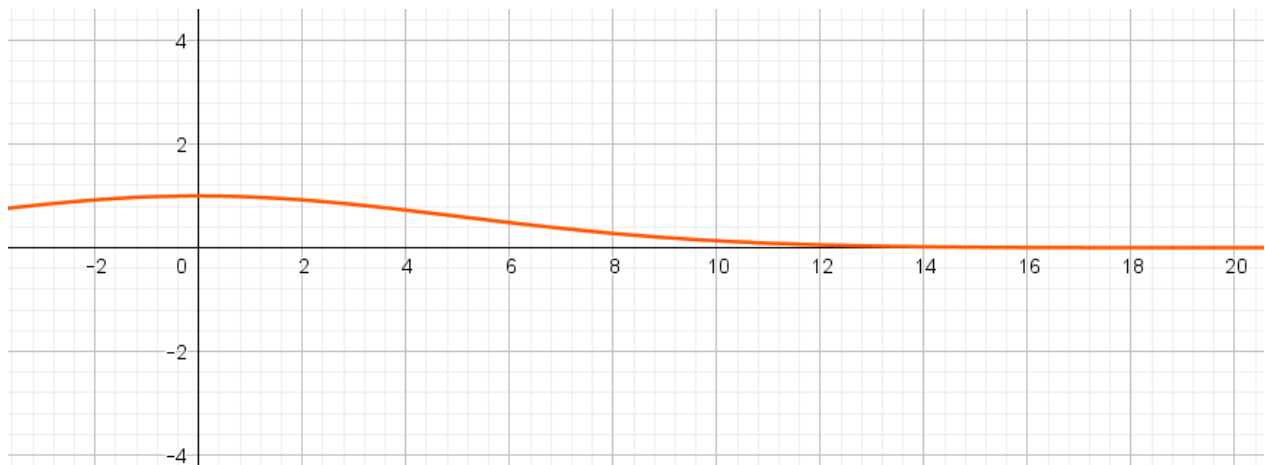


Figure 1: Kernel function for the distance measure.

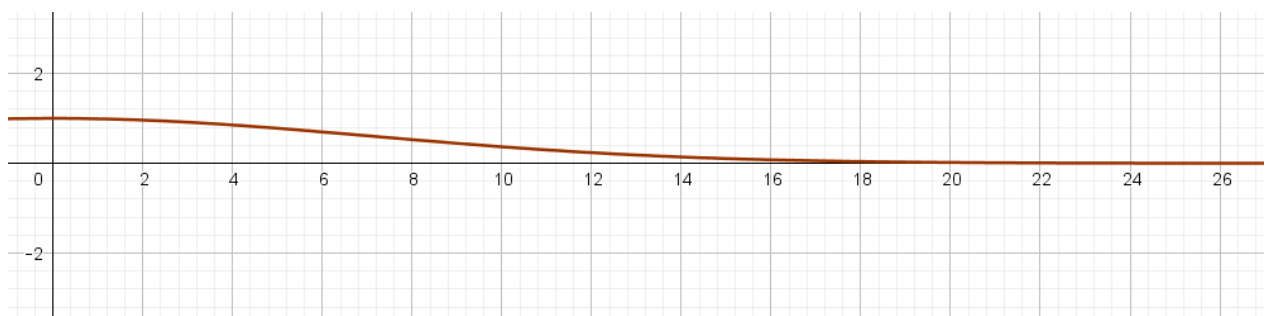


Figure 2: Kernel function for the date measure.

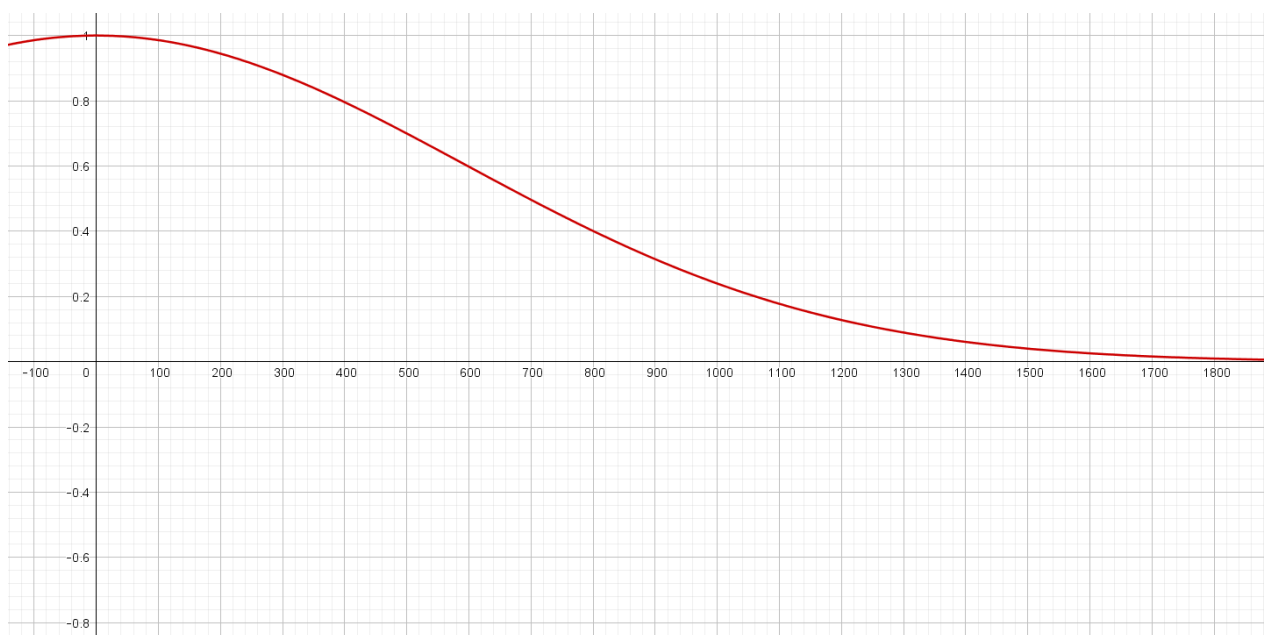


Figure 3: Kernel function for the time measure.

Repeat the exercise using a kernel that is the product of the three Gaussian kernels above. Compare the results with those obtained for the additive kernel. If they differ, explain why.

We saw that when multiplying the kernels the weights were now dependent on each of the three explanatory variables being significant, rather than when we added the kernels in which case the weights were determined individually. Therefore the multiplicative kernel is to prefer since intuitively one would require all of the explanatory variables to be significant to make a good prediction.