

```
In [ ]: import sys
from __future__ import division
from math import radians, cos, sin, asin, sqrt, exp
from datetime import datetime
from pyspark import SparkContext, StorageLevel

sc = SparkContext(appName = "Lab")
```

```
In [ ]: # sc.stop()
```

```
In [ ]: def haversine(lon1, lat1, lon2, lat2):
    """
    Calculate the great circle distance between two points
    on the earth (specified in decimal degrees)
    """
    # convert decimal degrees to radians
    lon1, lat1, lon2, lat2 = map(radians, [lon1, lat1, lon2, lat2])
    # haversine formula
    dlon = lon2 - lon1
    dlat = lat2 - lat1
    a = sin(dlat/2)**2 + cos(lat1) * cos(lat2) * sin(dlon/2)**2
    c = 2 * asin(sqrt(a))
    km = 6367 * c
    return km

def d_date(date1, date2):
    """
    Calculate the distance between two day
    """
    date_diff=abs((date2 - date1).days)
    return date_diff

def d_time(time1, time2):
    """
    Calculate the distance between two time
    """
    time_diff=abs((time2 - time1).total_seconds() / 3600)
    return time_diff

def sum_kernel(d_distance, d_date, d_time):
    """
    Calculate the sum of the kernel
    """
    distance_kernel=exp(-(d_distance**2)/(h_distance**2))
    time_kernel=exp(-(d_time**2)/(h_time**2))
    date_kernel=exp(-(d_date**2)/(h_date**2))

    return distance_kernel*time_kernel*date_kernel

def prod_kernel(d_distance, d_date, d_time):
    """
    Calculate the prod of the kernel
    """
    distance_kernel=exp(-(d_distance**2)/(h_distance**2))
    time_kernel=exp(-(d_time**2)/(h_time**2))
    date_kernel=exp(-(d_date**2)/(h_date**2))

    return distance_kernel*time_kernel*date_kernel
```

```

h_distance = 480
h_date = 6
h_time = 3
a = 58.4274 # Up to you
b = 14.826 # Up to you
date = "2014-11-04" # Up to you

```

```

In [ ]: DATE_TIMESTAMP=datetime.strptime(date, "%Y-%m-%d")

stations = sc.textFile("s3://bigdatalabs7051/data/stations/stations.csv")
temps = sc.textFile("s3://bigdatalabs7051/data/temperature-readings/temperature-read-
# lines = temps.map(lambda line: line.split(";")).map(lambda x:(str(x[0]), datetime.s

# Broadcast
# Station data: station_num, distance
stations_data=stations.map(lambda line: line.split(";")).map(lambda x:(str(x[0]), h
bc=sc.broadcast(stations_data.collectAsMap()) # Convert the rdd stations_data to

# Join
# intermediate: station_num, (date, time, temperature)
# joined_data: station_num, (geo distance, date distance, time, temperature)
joined_data = temps.sample(False, 0.1).map(lambda line: line.split(";"))\
    .map(lambda x:(str(x[0]), (datetime.strptime(x[1], "%Y-%m-%d"), dateti
    .filter(lambda x: (x[1][0]<=DATE_TIMESTAMP))\
    .map(lambda x: (x[0], (bc.value[x[0]], d_date(x[1][0], DATE_TIMESTAMP)
# date filtered here, time will be filtered when looping over times

# Persistence
joined_data.cache()

# print(joined_data.collect())

```

PythonRDD[5] at RDD at PythonRDD.scala:53

## Task 1: Use a kernel that is the sum of three Gaussian kernels

```

In [ ]: y_sum=[]

for time in ["0:00:00", "22:00:00", "20:00:00", "18:00:00", "16:00:00", "14:00:00",
"12:00:00", "10:00:00", "08:00:00", "06:00:00", "04:00:00"]:
    # Your code here

    # sum_kernel_rdd: 1, (sum_kernel, temperature)
    #sum_kernel_rdd=joined_data.map(lambda x:(1, (sum_kernel(x[1][0], x[1][1], d_time(t
    #sum_kernel_rdd.cache()
    #denominator=sum_kernel_rdd.groupByKey().mapValues(lambda values: sum([v[0] for
    #numerator=sum_kernel_rdd.groupByKey().mapValues(lambda values: sum([v[0]*v[1] f

    # Avoid using groupByKey and mapValues because it's much slower than reduceByKey

    # sum_kernel_rdd: 1, (sum_kernel, temperature)
    time=datetime.strptime(time, "%H:%M:%S")
    sum_kernel_rdd=joined_data.filter(lambda x:(x[1][1]>0 or (x[1][2] == 0 and x[
        .map(lambda x:(sum_kernel(x[1][0], x[1][1], d_time(time, x[1][2])), x[1][3]))
    sum_kernel_rdd=sum_kernel_rdd.map(lambda x:(x[0], x[0]*x[1])).reduce(lambda x, y
    #The result type returned by the reduce operation is consistent with the element
    y_sum.append(sum_kernel_rdd[1]/sum_kernel_rdd[0])

```

```
In [ ]: time_list=[“0:00:00”, “22:00:00”, “20:00:00”, “18:00:00”, “16:00:00”, “14:00:00”,  
“12:00:00”, “10:00:00”, “08:00:00”, “06:00:00”, “04:00:00”]  
dictionary = {k: v for k, v in zip(time_list, y_sum)}  
print(“latitude,longitude:”, a, ”, b)  
print(“date:”, date)  
print(dictionary)  
print(“HyperParams:”)  
print(“h_distance:”, h_distance)  
print(“h_date:”, h_date)  
print(“h_time:”, h_time)
```

output:

```
latitude,longitude: 58.4274 , 14.826  
date: 2014-11-04  
{'0:00:00': 5.570851580782074, '22:00:00': 5.741075353218582, '20:00:00':  
5.729431354715042, '18:00:00': 5.8618844900509455, '16:00:00': 6.076963659810907,  
'14:00:00': 6.247533933860337, '12:00:00': 6.249859464055124, '10:00:00':  
6.048901637752291, '08:00:00': 5.679243570182346, '06:00:00': 5.359907996425515,  
'04:00:00': 5.2696688195185315}  
HyperParams:  
h_distance: 480  
h_date: 6  
h_time: 3
```

**Q1:** 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.

**Answer:** The parameter of the Gaussian kernel function is the bandwidth, which determines the decay speed and smoothness of the kernel function. The smaller the bandwidth, the narrower the scope of the kernel function, the greater the impact of the current parameter change (in other words: more consideration of the current parameter), and vice versa.

According to common sense, time has the greatest impact on temperature, and date also has a certain impact on temperature. For example, winter is colder and summer is hotter, and the impact of distance is the least obvious (the distance from Oslo to Copenhagen is 563km, but the temperature difference is sometimes less than 1 degree Celsius). . Therefore the choice for time bandwidth should be the smallest and for the distance bandwidth the choice should be the largest.

We choose h\_distance=483 because the distance from Uppsala to Tampere is 483km, and the average temperature difference between the two places on May 18 is 5 degrees Celsius.

Choose h\_date=6, because Linkoping's average daily temperature increased by 5 degrees Celsius from May 18 to May 24.

Choose h\_time=3, because Linkoping May 18 from 7:00 to 10:00, the temperature increased by 5 degrees Celsius

## Task 2: Use a kernel that is the prod of three Gaussian kernels

```
In [ ]: y_prod=[]
```

```
for time in ["0:00:00", "22:00:00", "20:00:00", "18:00:00", "16:00:00", "14:00:00",
"12:00:00", "10:00:00", "08:00:00", "06:00:00", "04:00:00"]:
    # Your code here
    # prod_kernel_rdd: 1, (prod_kernel, temperature)
    time=datetime.strptime(time,"%H:%M:%S")
    prod_kernel_rdd=joined_data.filter(lambda x:(x[1][1]>0 or (x[1][2] == 0 and x
        .map(lambda x:(prod_kernel(x[1][0],x[1][1],d_time(time,x[1][2])),x[1][3]))
    prod_kernel_rdd=prod_kernel_rdd.map(lambda x:(x[0],x[0]*x[1])).reduce(lambda x
    #The result type returned by the reduce operation is consistent with the element
    y_prod.append(prod_kernel_rdd[1]/prod_kernel_rdd[0])
```

```
In [ ]: time_list=["0:00:00", "22:00:00", "20:00:00", "18:00:00", "16:00:00", "14:00:00",
"12:00:00", "10:00:00", "08:00:00", "06:00:00", "04:00:00"]
dictionary = {k: v for k, v in zip(time_list, y_prod)}
print("latitude,longitude:",a,",",b)
print("date:",date)
print(dictionary)
print("HyperParams:")
print("h_distance:",h_distance)
print("h_date:",h_date)
print("h_time:",h_time)
```

output:

```
latitude,longitude: 58.4274 , 14.826
date: 2014-11-04
{'0:00:00': 7.303006795891909, '22:00:00': 7.982068570228097, '20:00:00':
8.12157139414748, '18:00:00': 8.22400511629453, '16:00:00': 8.40477703134592, '14:00:00':
8.868385461565387, '12:00:00': 9.113797702887448, '10:00:00': 8.690957561793367,
'08:00:00': 7.987008400145532, '06:00:00': 7.51869787901216, '04:00:00':
7.335325908719835}
HyperParams:
h_distance: 480
h_date: 6
h_time: 3
```

**Q2:** 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.

**Answer:** Using the product of the three Gaussian kernels, overall, the predicted air temperature is higher than using the kernel that is the sum of three Gaussian kernels.

The product of Gaussian kernel functions can enhance the similarity between samples, that is, when the values of the three Gaussian kernel functions are larger at the same time, the similarity is higher. Conversely, if one of the Gaussian kernels has a lower value, then the kernel value will also be lower. Therefore, it pays more attention to the common characteristics of the three samples in the feature space. The sum of Gaussian kernel functions can comprehensively consider the similarity between samples. It pays more attention to the respective characteristics of the three samples in the feature space, and considers them comprehensively by adding.

# Task 3: Using at least two MLlib library models to predict the hourly temperatures for a date and place in Sweden.

## Linear Regression

```
In [ ]: from datetime import date

# Transform time to float, range from 0-100
def time_to_float(time):
    hour, minute, sec = map(int, time.split(":"))
    total_minutes = hour * 60 + minute
    relative_minutes = total_minutes / (24 * 60)
    return relative_minutes * 100

# Transform date to float, range from 0-100
def date_to_float(date_str):
    year, month, day = map(int, date_str.split("-"))
    date_obj = date(year, month, day)
    year_start = date(year, 1, 1)
    year_end = date(year, 12, 31)
    total_days = (date_obj - year_start).days
    relative_days = total_days / (year_end - year_start).days
    return relative_days * 100
```

```
In [ ]: from pyspark.mllib.regression import LabeledPoint
from pyspark.mllib.regression import LinearRegressionWithSGD

# Broadcast
# station_num, longitude, latitude
stations_data=stations.map(lambda line: line.split(";")).map(lambda x:(str(x[0]),(x[1],x[2])))
bc=sc.broadcast(stations_data.collectAsMap()) # Convert the rdd stations_data to broadcast variable

# Predict!
a = 58.4274 # latitude
b = 14.826 # longitude
date1 = "2014-11-04" # Up to you

y_linreg=[]
predict_data=[]
for time1 in ["0:00:00", "22:00:00", "20:00:00", "18:00:00", "16:00:00", "14:00:00", "12:00:00", "10:00:00", "08:00:00", "06:00:00", "04:00:00"]:
    # intermediate: (longitude, latitude), day, time, temp
    # LabeledPoint(labeled_data) : temperature, longitude, latitude, date, time
    labeled_data = temps.sample(False, 0.1).map(lambda line: line.split(";"))\
        .map(lambda x:(bc.value[str(x[0])],date_to_float(x[1]),time_to_float(x[2]),float(x[3])))\
        .filter(lambda x: x[1]<date_to_float(date1) or (x[1]==date_to_float(date1) and x[2]>time1))\
        .map(lambda x: LabeledPoint(x[3], [x[0][0], x[0][1], x[1], x[2]]))

    train_data, test_data = labeled_data.randomSplit([0.7, 0.3])

    # Train the linear regression model
    model = LinearRegressionWithSGD.train(train_data, iterations=100, step=0.0005)
```

```
# Your code here
predict_data.append([b, a, date_to_float(date1), time_to_float(time1)])

rdd1 = sc.parallelize(predict_data)
print("Training data:")
print(rdd1.collect())
predictions = model.predict(rdd1)
y_linreg=predictions.collect()

print("Predict result:")
print(y_linreg)
```

Training data:

```
[[14.826, 58.4274, 84.34065934065934, 0.0], [14.826, 58.4274, 84.34065934065934, 91.66666666666666], [14.826, 58.4274, 84.34065934065934, 83.33333333333334], [14.826, 58.4274, 84.34065934065934, 75.0], [14.826, 58.4274, 84.34065934065934, 66.66666666666666], [14.826, 58.4274, 84.34065934065934, 58.33333333333336], [14.826, 58.4274, 84.34065934065934, 50.0], [14.826, 58.4274, 84.34065934065934, 41.66666666666667], [14.826, 58.4274, 84.34065934065934, 33.33333333333333], [14.826, 58.4274, 84.34065934065934, 25.0], [14.826, 58.4274, 84.34065934065934, 16.66666666666664]]
```

Predict result:

```
[13.604115554952378, 15.996997551067711, 15.779462824148135, 15.56192809722856, 15.344393370308984, 15.126858643389408, 14.909323916469832, 14.691789189550256, 14.47425446263068, 14.256719735711105, 14.039185008791529]
```

## ridge regression

```
In [ ]: from pyspark.mllib.regression import LabeledPoint
from pyspark.mllib.regression import RidgeRegressionWithSGD

# Broadcast
# station_num, longitude, latitude
stations_data=stations.map(lambda line: line.split(",")).map(lambda x:(str(x[0]),(bc.broadcast(stations_data.collectAsMap()) # Convert the rdd stations_data to

# Predict!
a = 58.4274 # latitude
b = 14.826 # longitude
date1 = "2014-11-04" # Up to you

y_linreg=[]
predict_data=[]
for time1 in ["0:00:00", "22:00:00", "20:00:00", "18:00:00", "16:00:00", "14:00:00", "12:00:00", "10:00:00", "08:00:00", "06:00:00", "04:00:00"]:
    # intermediate: (longitude, latitude), day, time, temp
    # LabeledPoint(labeled_data) : temperature, longitude, latitude, date, time
    labeled_data = temps.sample(False, 0.1).map(lambda line: line.split(","))
        .map(lambda x:(bc.value[str(x[0])],date_to_float(x[1]),time_to_
        .filter(lambda x: x[1]<date_to_float(date1) or (x[1]==date_to_.
        .map(lambda x: LabeledPoint(x[3], [x[0][0], x[0][1], x[1], x[2]]
```

```

train_data, test_data = labeled_data.randomSplit([0.7, 0.3])

# Train the linear regression model
model = RidgeRegressionWithSGD.train(train_data, iterations=100, step=0.0005)

# Your code here
predict_data.append([b, a, date_to_float(date1), time_to_float(time1)])

rdd1 = sc.parallelize(predict_data)
print("Training data:")
print(rdd1.collect())
predictions = model.predict(rdd1)
y_linreg=predictions.collect()

print("Predict result:")
print(y_linreg)

```

Training data:

```

[[14.826, 58.4274, 84.34065934065934, 0.0], [14.826, 58.4274, 84.34065934065934,
91.66666666666666], [14.826, 58.4274, 84.34065934065934, 83.33333333333334], [14.826,
58.4274, 84.34065934065934, 75.0], [14.826, 58.4274, 84.34065934065934,
66.66666666666666], [14.826, 58.4274, 84.34065934065934, 58.33333333333336], [14.826,
58.4274, 84.34065934065934, 50.0], [14.826, 58.4274, 84.34065934065934,
41.66666666666667], [14.826, 58.4274, 84.34065934065934, 33.33333333333333], [14.826,
58.4274, 84.34065934065934, 25.0], [14.826, 58.4274, 84.34065934065934,
16.66666666666664]]

```

Predict result:

```

[13.57574968051529, 15.985517198023985, 15.766447423705014, 15.547377649386041,
15.328307875067068, 15.109238100748097, 14.890168326429125, 14.671098552110152,
14.45202877779118, 14.232959003472207, 14.013889229153236]

```

**Q3:** Repeat the exercise using at least two MLlib library models to predict the hourly temperatures for a date and place in Sweden. Compare the results with two Gaussian kernels. If they differ, explain why

**Answer:** we can see that the predicted result from LinearRegressionWithSGD and RidgeRegressionWithSGD are all similar to each other(most critical hyperparameters are set to the same, there isn't too much features and the regularizer parameter is set as 0.01 by default which didn't make too significant impact). And the regression results are not to kernel results. It might be caused by the hyperparameters(such as learning rate) are still needs to be optimized to get a better prediction.