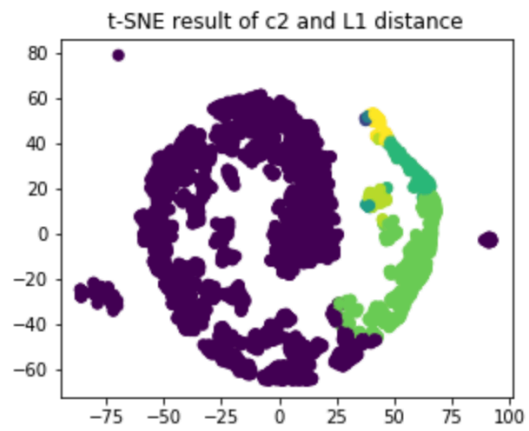
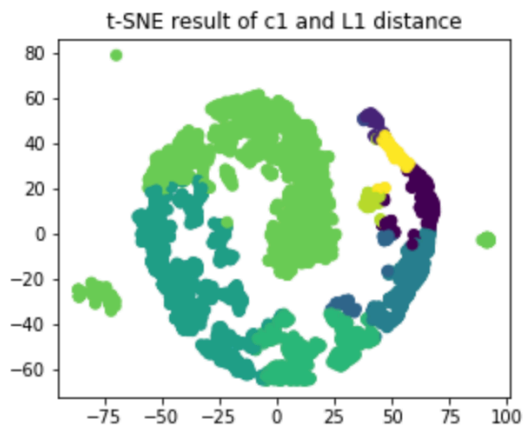
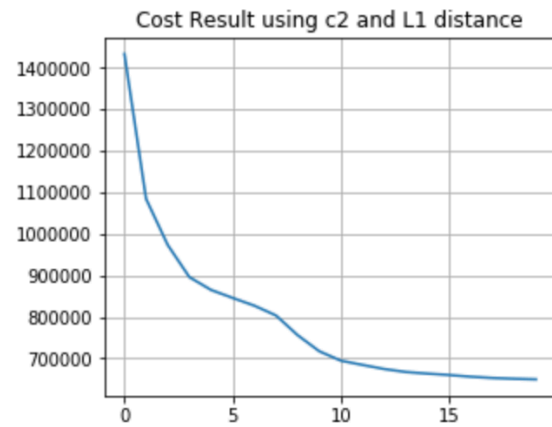
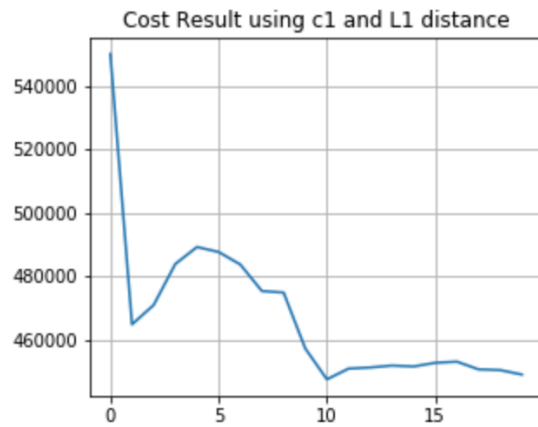
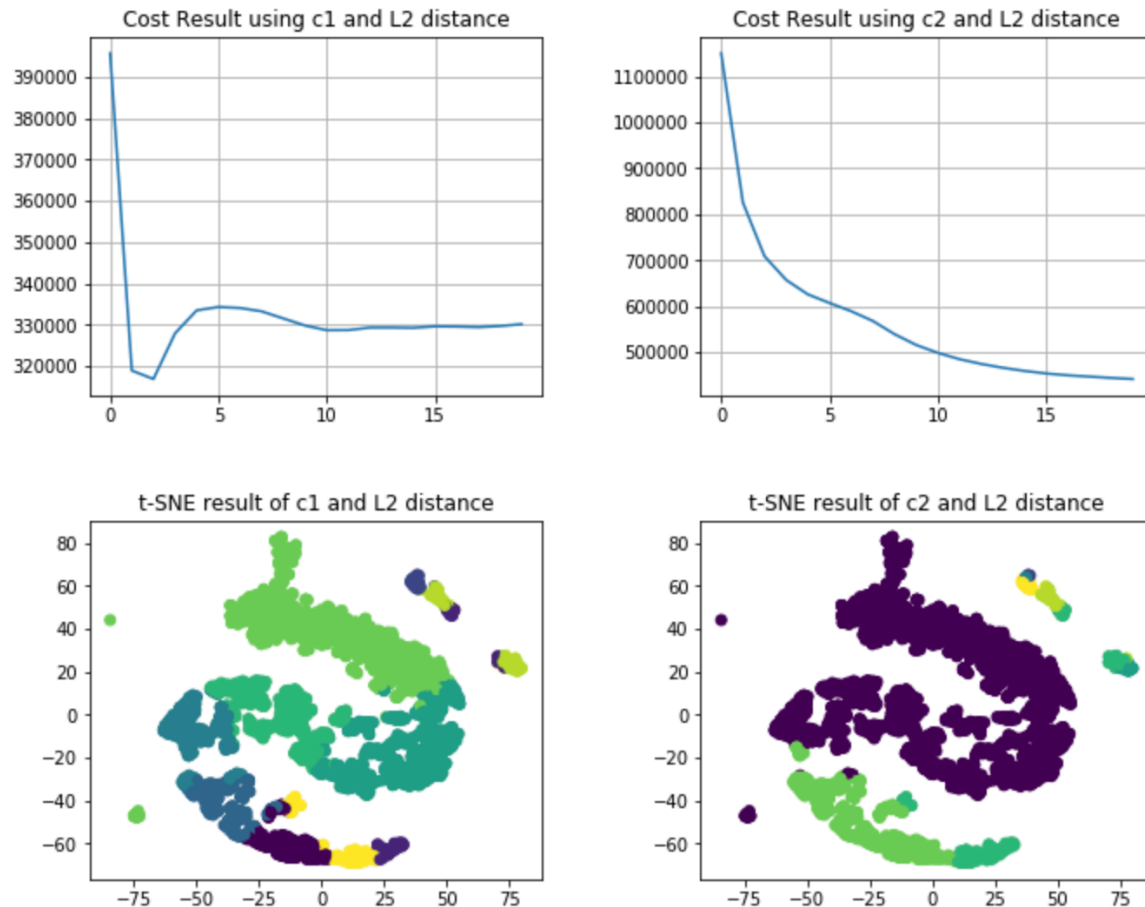


Homework1

1. Iterative K-means clustering on Spark (1)



(2) and (3)



(4)

No, although from the above plots we could see that using c1 will get a better cost than c2 in the end, but we only got 20 iterations. If we increase the 'MAX_ITER', the cost by using c2 will eventually be lower than by using c1.

Because random initialization sometimes can result in creating centroids in such a way that they are clumped together in space. Then in the end we'll have several clusters cramped tightly, and make other clusters more sparse, which will result in a higher cost.

(5)

Assume we have:

`k` clusters

`p` data points, with `d` dimensions

`n` maximum iterations

Then:

for calculating the distance between each points and centroids: $O(n) = pdk$

for choosing the closest centroids by using quick sort: $O(n) = p * k \log(k)$

for recomputing the centroids: $O(n) = pd$

So altogether, $O(n) = n(pdk + pk\log(k) + pd) = np(dk + k\log(k) + d)$

Code:

```

1  import operator
2  import sys
3  from pyspark import SparkConf, SparkContext
4  import numpy as np
5  import matplotlib.pyplot as plt
6  from scipy import linalg # calculating L1 and L2 distance
7  import time # calculate computing time
8  import pandas as pd
9  from sklearn.manifold import TSNE # for dimensionality reduction
10 import random
11
12 # Macros.
13 MAX_ITER = 20
14 DATA_PATH = "gs://big_data_hw/hw1/data.txt"
15 C1_PATH = "gs://big_data_hw/hw1/c1.txt"
16 C2_PATH = "gs://big_data_hw/hw1/c2.txt"
17
18 # Helper functions.
19 def closest(p, centroids, norm):
20     closest_c = min([(i, linalg.norm(p - c, norm))
21                     for i, c in enumerate(centroids)],
22                     key=operator.itemgetter(1))[0]
23     return closest_c
24
25 # K-means clustering
26 def kmeans(data, centroids, norm=2):
27     cost = []
28
29     for t in range(MAX_ITER):
30         combo = data.map(lambda point: (point, 1)).cache()
31
32         #within cluster cost and running time
33         costById = combo.map(lambda x: (x[0], linalg.norm(x[1][0] - centroids[x[0]], norm))) \
34             .reduceByKey(lambda x, y: x + y).collect()
35         cost.append(sum([costById[i][1] for i in range(len(costById))]))
36
37         reduce1 = combo.reduceByKey(lambda a, b: ([a[0][i]+b[0][i] for i in range(len(a[0]))], a[1]+b[1]))
38
39         # Average the points for each centroid: divide sum of points by count
40         map1 = reduce1.sortByKey().map(lambda x: [x[1][0][i] / x[1][1] for i in range(len(x[1][0]))])
41
42         # Use collect() to turn RDD into list
43         centroids = map1.collect()
44
45         centroIndex = data.map(lambda point: closest(point, centroids, norm)).collect()
46
47     return cost, centroIndex
48
49 def main(norm):
50     random.seed(0)
51     # Spark settings
52     conf = SparkConf()
53     sc = SparkContext.getOrCreate(conf=conf)

```

```

55 # Load the data, cache this since we're accessing this each iteration
56 data = sc.textFile(DATA_PATH).map(
57     lambda line: np.array([float(x) for x in line.split(' ')])
58 ).cache()
59 # Load the initial centroids c1, split into a list of np arrays
60 centroids1 = sc.textFile(C1_PATH).map(
61     lambda line: np.array([float(x) for x in line.split(' ')])
62 ).collect()
63 # Load the initial centroids c2, split into a list of np arrays
64 centroids2 = sc.textFile(C2_PATH).map(
65     lambda line: np.array([float(x) for x in line.split(' ')])
66 ).collect()
67
68 # calculate the cost and the final central point index
69 cost1, centroIndex1 = kmeans(data, centroids1, norm)
70 cost2, centroIndex2 = kmeans(data, centroids2, norm)
71
72 # dimensionality reduction
73 data_np = np.array(data.collect())
74 data_embedded = TSNE(n_c, init_transform(data_np)
75 vis_x = data_embedded[:, References:
76 vis_y = data_embedded[:, temp.py:34
77 temp.py:43
78 plt.figure(figsize=(10,8 temp.py:45
79 # plot for cost temp.py:62
80 plt.subplot(2,2,1) temp.py:66
81 plt.plot(range(20), cost_)
82 plt.title('Cost Result using c1 and L%d distance' % (norm))
83 plt.grid(True)
84 plt.subplot(2,2,2)
85 plt.plot(range(20), cost2)
86 plt.title('Cost Result using c2 and L%d distance' % (norm))
87 plt.grid(True)
88
89 # plot for clustering result
90 plt.subplot(2,2,3)
91 plt.scatter(vis_x, vis_y, c = centroIndex1)
92 plt.title('t-SNE result of c1 and L%d distance' % (norm))
93 plt.subplot(2,2,4)
94 plt.scatter(vis_x, vis_y, c = centroIndex2)
95 plt.title('t-SNE result of c2 and L%d distance' % (norm))
96
97 plt.subplots_adjust(top=0.92, bottom=0.08, left=0.10,
98                     right=0.95, hspace=0.35, wspace=0.35)
99
100 plt.show()
101
102 sc.stop()

```

2. Binary classification with Spark MLlib

(1) data loading

```

In [1]: # create SparkSession and load data
from pyspark.sql import SparkSession
ADULT_PATH = "gs://big_data_hw/hw1/adult.data.csv"

spark = SparkSession \
    .builder \
    .appName("Python Spark SQL basic example") \
    .config("spark.some.config.option", "some-value") \
    .getOrCreate()

df = spark.read.csv(ADULT_PATH, inferSchema = True)
# rename each column
oldColumns = df.columns
newColumns = ("age", "workclass", "fnlwgt", "education",
              "education_num", "marital_status", "occupation",
              "relationship", "race", "sex", "capital_gain",
              "capital_loss", "hours_per_week", "native_country", "income")

df = reduce(lambda data, idx: data.withColumnRenamed(oldColumns[idx], newColumns[idx]), xrange(len(oldColumns)), df)

```

(2) preprocessing

```

In [35]: # preprocessing and model building
from pyspark.ml import Pipeline
from pyspark.ml.feature import *
from pyspark.ml.classification import LogisticRegression
from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
from pyspark.ml.evaluation import BinaryClassificationEvaluator, MulticlassClassificationEvaluator
from pyspark.mllib.evaluation import MulticlassMetrics
import matplotlib.pyplot as plt

In [5]: # piping the preprocess and regression
categorical_variables = ['workclass', 'education', 'marital_status',
                        'occupation', 'relationship', 'race', 'sex', 'native_country']
continuous_variables = ['age', 'fnlwgt', 'education_num', 'capital_gain', 'capital_loss', 'hours_per_week']

indexers = [StringIndexer(inputCol=column,
                          outputCol=column+"_index") for column in categorical_variables]
encoder = OneHotEncoderEstimator(
    inputCols=[indexer.getOutputCol() for indexer in indexers],
    outputCols=["{0}_encoded".format(indexer.getOutputCol()) for indexer in indexers]
)
assembler = VectorAssembler(
    inputCols=encoder.getOutputCols() + continuous_variables,
    outputCol="features"
)
response = StringIndexer(inputCol='income', outputCol='label')
preprocess = Pipeline(stages=indexers + [encoder, assembler, response])

# split data into 70% training and 30% testing and set the seed to 100
train, test = df.randomSplit([0.7, 0.3], seed = 100)
pre_model = preprocess.fit(train)
train = pre_model.transform(train)
test = pre_model.transform(test)

```

(3) modeling

```

lr = LogisticRegression(maxIter=10, featuresCol='features', labelCol='label')
model = lr.fit(train)
prediction = model.transform(test)

```

```

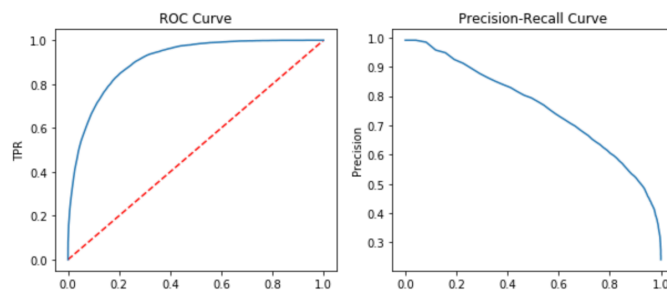
In [14]: # ROC for training data
plt.figure(figsize=(10,4))

plt.subplot(1,2,1)
plt.plot([0, 1], [0, 1], 'r--')
plt.plot(model.summary.roc.select('FPR').collect(),
        model.summary.roc.select('TPR').collect())
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC Curve')

plt.subplot(1,2,2)
pr = model.summary.pr.toPandas()
plt.plot(pr['recall'], pr['precision'])
plt.ylabel('Precision')
plt.xlabel('Recall')
plt.title('Precision-Recall Curve')

plt.show()

```



(4) evaluation

```
In [49]: # testing AUC and accuracy and confusion matrix
evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")
evaluator.evaluate(prediction)
print("%s is: %.6f \n" % (evaluator.getMetricName(), evaluator.evaluate(prediction)))

acc = MulticlassClassificationEvaluator(labelCol='label', predictionCol='prediction', metricName='accuracy')
print("%s is: %.6f \n" % (acc.getMetricName(), acc.evaluate(prediction)))

predictionAndLabels = prediction.select('prediction', 'label').rdd
metrics = MulticlassMetrics(predictionAndLabels)
print("confusion matrix is: ")
print(metrics.confusionMatrix())

areaUnderROC is: 0.902738

accuracy is: 0.848401

confusion matrix is:
DenseMatrix([[6860.,  530.],
              [ 944., 1389.]])
```