# 211275026-陈畅-实验四

https://github.com/Cc17952/FBDP 环境安装 已有环境 Scala安装 Spark环境安装 使用pyspark遇到的问题 IP未设置 pyspark库无法导入 WARN: Setting default log level to "WARN". 实验内容 任务一 编写 Spark 程序,统计application\_data.csv中所有用户的贷款金额AMT\_CREDIT 的分布情况。 编写Spark程序,统计application\_data.csv中客户贷款金额AMT\_CREDIT 比客户收入AMT\_INCOME\_TOT... 任务二 统计所有男性客户(CODE\_GENDER=M)的小孩个数(CNT\_CHILDREN)类型占比情况。 统计每个客户出生以来每天的平均收入。 任务三 实验设计

选取变量

实验结果

实验存在的不足

## https://github.com/Cc17952/FBDP

## 环境安装

#### 已有环境

Scala版本	Spark 2.x版本	Spark 3.x版本	
2.11	✓		
2.12	✓	✓	

### Scala和Spark对应关系

Spark	版本	Hadoop版本
2.4.x		2.7.x
3.0.x		3.2.x

#### Spark和Hadoop对应关系

当前环境: Hadoop-2.9.2; Java-8

预下载版本: Scala-2.11.7; Spark-2.4.2

## Scala安装

下载链接: https://github.com/scala/scala/releases/tag/v2.11.7

参考教程: https://www.runoob.com/scala/scala-install.html

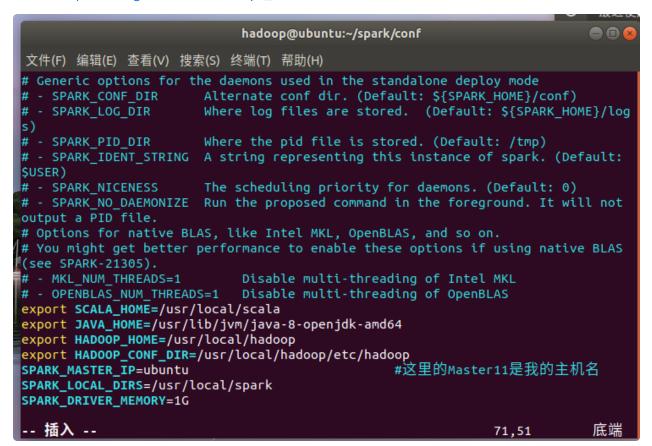
```
hadoop@ubuntu:/usr/local
 文件(F) 编辑(E) 查看(V) 搜索(S) 终端(T) 帮助(H)
export PATH=${JAVA_HOME}/bin:$PATH
export HADOOP_HOME=/usr/local/hadoop
export PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin
export HADOOP_COMMON_HOME=$HADOOP HOME
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
export HADOOP_HDFS_HOME=$HADOOP_HOME
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export HADOOP YARN HOME=$HADOOP HOME
export HADOOP_INSTALL=$HADOOP HOME
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native
export HADOOP_LIBEXEC_DIR=$HADOOP_HOME/libexec
export JAVA_LIBRARY_PATH=$HADOOP_HOME/lib/native:$JAVA_LIBRARY_PATH
export HDFS_DATANODE_USER=root
export HDFS_DATANODE_SECURE_USER=root
export HDFS_SECONDARYNAMENODE_USER=root
export HDFS_NAMENODE_USER=root
export YARN_RESOURCEMANAGER_USER=root
export YARN_NODEMANAGER_USER=root
export MAVEN_HOME=/usr/local/apache-maven-3.9.5
export PATH=${MAVEN_HOME}/bin:$PATH
export HBASE_HOME=/usr/local/hbase
export PATH=$PATH:$HBASE_HOME/bin:/$HBASE_HOME/sbin
export PATH=$PATH:/usr/local/scala/bin
                                                                   53,12
```

```
hadoop@ubuntu:~/Downloads$ sudo mv scala-2.11.7 /usr/local/scala [sudo] hadoop 的密码:
hadoop@ubuntu:~/Downloads$ cd ..
hadoop@ubuntu:~$ cd /usr/local/
hadoop@ubuntu:/usr/local$ sudo vim /etc/profile
hadoop@ubuntu:/usr/local$ source /etc/profile
hadoop@ubuntu:/usr/local$ source /etc/profile
hadoop@ubuntu:/usr/local$ scala
Welcome to Scala 2.11.12 (OpenJDK 64-Bit Server VM, Java 1.8.0_362).
Type in expressions for evaluation. Or try :help.
```

由于2.11.7解压得到2.11.12、大版本没有出现问题所以采用2.11.12

## Spark环境安装

参考教程: https://blog.csdn.net/Nurbiya\_K/article/details/100982166



mv slaves.template slaves

vim slaves//修改节点名称(改为了h01、h02、h03)

```
LICENSE-antlr.txt 100% 1490 7.7MB/s 00:00
LICENSE-netlib.txt 100% 2300 10.8MB/s 00:00
LICENSE-paranamer.txt 100% 1661 8.4MB/s 00:00
spark-2.4.0-yarn-shuffle.jar 100% 9489KB 229.1MB/s 00:00
hadoop@ubuntu:/usr/local$
```

### scp -r spark/ h01:/usr/local/

```
hadoop@ubuntu:/usr/local$ scp -r scala/ h01:/usr/local/
Warning: Permanently added 'h01' (ECDSA) to the list of known hosts.
hadoop@h01's password:
```

#### scp -r scala/ h01:/usr/local/

```
9.2.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of
security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of furth
ive access operations
WARNING: All illegal access operations will be denied in a fut
localhost: Warning: Permanently added 'localhost' (ECDSA) to
osts.
hadoop@localhost's password:
localhost: starting nodemanager, logging to /usr/local/hadoop,
odemanager-ubuntu.out
hadoop@ubuntu:/usr/local/hadoop$ jps
8275 SecondaryNameNode
7732 NameNode
8872 NodeManager
9048 Jps
7980 DataNode
8476 ResourceManager
hadoop@ubuntu:/usr/local/hadoop$
```

#### 先启动hadoop

```
hadoop@ubuntu:/usr/local/spark$ sbin/start-all.sh
starting org.apache.spark.deploy.master.Master, logging to /usr/local/spark/logs
/spark-hadoop-org.apache.spark.deploy.master.Master-1-ubuntu.out
h01: Warning: Permanently added 'h01' (ECDSA) to the list of known hosts.
hadoop@h01's password: h03: ssh: Could not resolve hostname h03: Name or service
not known
h02: ssh: Could not resolve hostname h02: Name or service not known
```

再启动Spark

这里发现上述slaves文件配置错误,修改为h01。

```
hadoop@ubuntu:/usr/local/spark$ jps
9218 Jps
8275 SecondaryNameNode
7732 NameNode
8872 NodeManager
9098 Master
7980 DataNode
8476 ResourceManager
```

## 使用pyspark遇到的问题

#### IP未设置

vim spark-env.sh

来设置SPARK\_LOCAL\_IP=主机ip

### pyspark库无法导入

尝试了比较多的解决方法,例如pip install、修改环境变量

最终通过代码导入。

```
task1.py > ...
    import sys
    sys.path.append("/usr/local/spark/python/lib/py4j-0.10.7-src.zip")
    sys.path.append("/usr/local/spark/python/lib/pyspark.zip")
4
    from pyspark import SparkContext
```

WARN: Setting default log level to "WARN".

具体报错信息: Setting default log level to "WARN".To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).

```
sc = SparkContext( 'local', 'test')
# sc.setLogLevel("INFO")
```

解决方法: 注释行

```
2023-12-24 12:16:01 INFO SparkContext:54 - Invoking stop() from shutdown hook

| hadoop@ubuntu:-/lab4$ 2023-12-24 12:16:01 INFO AbstractConnector:318 - Stopped Spark@392ffb88{HTTP/1.1,[http/1.1]}{192.168.217.147:4040}
2023-12-24 12:16:01 INFO SparkUI:54 - Stopped Spark web UI at http://192.168.217.147:4040
2023-12-24 12:16:01 INFO MapOutputTrackerMasterEndpoint:54 - MapOutputTrackerMasterEndpoint stopped!
2023-12-24 12:16:01 INFO MemoryStore:54 - MemoryStore cleared
2023-12-24 12:16:01 INFO BlockManager:54 - BlockManager stopped
2023-12-24 12:16:01 INFO BlockManagerMaster:54 - BlockManagerMaster stopped
2023-12-24 12:16:01 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:54 - OutputCommitCoordinator stopped!
2023-12-24 12:16:01 INFO SparkContext:54 - Successfully stopped SparkContext
2023-12-24 12:16:01 INFO ShutdownHookManager:54 - Shutdown hook called
2023-12-24 12:16:01 INFO ShutdownHookManager:54 - Deleting directory /usr/local/spark/spark-9579c5e2-a646-48f6-ale8-db3ee7f5658a/pyspark-61abe4f9-6c7
2023-12-24 12:16:01 INFO ShutdownHookManager:54 - Deleting directory /usr/local/spark/spark-9579c5e2-a646-48f6-ale8-db3ee7f5658a
/usr/bin/python /home/hadoop/lab4/task1.py
```

实际上这个报错信息可以不用理会

## 实验内容

## 任务一

编写 Spark 程序,统计application\_data.csv中所有用户的贷款金额AMT\_CREDIT 的分布情况。

```
task1-1.py
                                                                       Python
    \# -*- coding: utf-8 -*-
 1
2
    import sys
    import math
3
     from pyspark import SparkContext
4
5
     from pyspark.sql import SparkSession
6
     from pyspark.sql.functions import col
7
8
     sys.path.append("/usr/local/spark/python/lib/py4j-0.10.7-src.zip")
9
     sys.path.append("/usr/local/spark/python/lib/pyspark.zip")
     spark = SparkSession.builder.appName("LoanAmountDistribution").getOrCreate
10
     ()
11
     df = spark.read.csv("file:///home/hadoop/lab4/data/application data.csv",
     header=True, inferSchema=True)
12
    df = df.withColumn("AMT_CREDIT_BIN", ((col("AMT_CREDIT") / 10000).cast("in
13
     t") * 10000).alias("AMT CREDIT BIN"))
     credit_distribution = df.groupBy("AMT_CREDIT_BIN").count().sort("AMT_CREDI
14
    T_BIN")
15
16
     credit distribution.show(truncate=False)
     result = credit_distribution.rdd.map(lambda x: "(({},{}) , {})".format(x
17
     ['AMT_CREDIT_BIN'], x['AMT_CREDIT_BIN'] + 10000, x['count'])).collect()
18
19 with open('./result/task1-1.txt', 'w') as file:
         for line in result:
20 -
21
             file.write(line + '\n')
```

```
资源管理器
                       task1.py 3
                                        ≣ task1-1.txt ×
        中に甘り自
                       result > ≣ task1-1.txt
                              ((40000,50000) , 561)
                               ((50000,60000),
{} launch.json
                               ((60000,70000),
∨ data
                                                , 1226)
                               ((70000,80000)
application_data.csv
                              ((80000,90000),
                                                   668)
                              ((90000, 100000) , 1939)
≣ task1-1.txt
                               ((100000,110000)
                               ((110000,120000)
task1.py
                               ((120000,130000)
                               ((130000,140000)
                                                    , 4792)
                               ((140000,150000)
                               ((150000,160000)
                               ((160000,170000)
                               ((170000,180000)
                               ((180000,190000)
                                                     8745)
                               ((190000,200000)
                               ((200000,210000)
                                                     4017)
                               ((210000,220000)
                                                   , 10013)
                               ((220000,230000)
                               ((230000,240000)
                                                     3343)
                               ((240000,250000)
                               ((250000,260000)
                                                   , 5186)
                               ((260000,270000)
                        问题 ③ 输出 调试控制台
                                         |1930
|1323
|4792
|2239
|3653
|1919
                        |110000
|120000
                         130000
                         |140000
|150000
|160000
                         |170000
|170000
|180000
|190000
                                          |2131
|8745
|1537
                         |200000
|210000
                                          |4017
|1475
                         220000
                                          10013
                        only showing top 20 rows
```

运行结果

编写Spark程序,统计application\_data.csv中客户贷款金额AMT\_CREDIT 比客户收入AMT\_INCOME\_TOTAL差值最高和最低的各十条记录。

```
task1-2.py
                                                                       Python
    \# -*- coding: utf-8 -*-
 1
2
    import sys
    import math
3
     from pyspark import SparkContext
4
     from pyspark.sql import SparkSession
5
6
     from pyspark.sql.functions import col
7
     sys.path.append("/usr/local/spark/python/lib/py4j-0.10.7-src.zip")
     sys.path.append("/usr/local/spark/python/lib/pyspark.zip")
8
9
     spark = SparkSession.builder.appName("LoanAmountDistribution").getOrCreate
     ()
10
     df = spark.read.csv("file:///home/hadoop/lab4/data/application data.csv",
     header=True, inferSchema=True)
11
     df = df.withColumn("DIFF", col("AMT CREDIT") - col("AMT INCOME TOTAL"))
12
     top records = df.orderBy(col("DIFF").desc()).limit(10)
13
     bottom records = df.orderBy(col("DIFF").asc()).limit(10)
14
15
16
     output_path = "./result/task1-2.txt"
17 with open(output_path, 'w') as file:
         file.write("Top 10 records with highest difference:\n")
18
         file.write(top_records.select("SK_ID_CURR", "NAME_CONTRACT_TYPE", "AMT
19
    _CREDIT", "AMT_INCOME_TOTAL", "DIFF").toPandas().to_string(index=False))
         file.write("\n\nBottom 10 records with lowest difference:\n")
20
         file.write(bottom_records.select("SK_ID_CURR", "NAME_CONTRACT_TYPE",
21
     "AMT_CREDIT", "AMT_INCOME_TOTAL", "DIFF").toPandas().to_string(index=Fals
     e))
```

	🕏 task1	-2.py 3	≣ task1-2.txt ×							
	result >	≣ task1-2.tx	t-							
	1 Top 10 records with highest difference:									
						AMT INCOME TOTAL	DIFF			
		43329		loans			3645000.0			
	4	21095		loans			3600382.5			
data.csv	5	43417		loans			3600000.0			
		31589		loans			3569130.0			
		23843		loans			3567969.0			
		24000	7 Cash	loans		587250.0				
		11733	7 Cash	loans	4050000.0	760846.5				
3	10	12092	6 Cash	loans	4050000.0	783000.0	3267000.0			
3	11	11708	5 Cash	loans	3956274.0	749331.0	3206943.0			
	12	22813	5 Cash	loans	4050000.0	864900.0	3185100.0			
	13	Bottom 10	records with	highest	difference:					
	14	SK_ID_CUR	R NAME_CONTRAC	T_TYPE	AMT_CREDIT	AMT_INCOME_TOTAL	DIFF			
	15	11496	7 Cash	loans	562491.0	117000000.0	-116437509.0			
	16	33614	7 Cash	loans	675000.0	18000090.0	-17325090.0			
	17	38567	4 Cash	loans	1400503.5	13500000.0	-12099496.5			
	18	19016	0 Cash	loans	1431531.0	9000000.0	-7568469.0			
	19	25208	4 Cash	loans	790830.0	6750000.0	-5959170.0			
	20	33715		loans		4500000.0	-4050000.0			
	21	31774		loans		4500000.0	-3664620.0			
	22	31060		loans		3950059.5				
	23	43298		loans	1755000.0	4500000.0				
		15747	1 Cash	loans	953460.0	3600000.0	-2646540.0			

运行结果

## 任务二

基于Hive或者Spark SQL对application\_data.csv进行如下统计:

```
task2.py
                                                                        Python
    \# -*- coding: utf-8 -*-
 1
 2
 3
    import sys
4
     sys.path.append("/usr/local/spark/python/lib/py4j-0.10.7-src.zip")
     sys.path.append("/usr/local/spark/python/lib/pyspark.zip")
5
     import math
6
7
     from pyspark import SparkContext
8
     from pyspark.sql import SparkSession
9
     from pyspark.sql.functions import col
     from pyspark.sql import functions as F
10
11
12
     spark = SparkSession.builder.appName("GenderChildIncomeAnalysis").getOrCre
     ate()
     df = spark.read.csv("file:///home/hadoop/lab4/data/application_data.csv",
13
     header=True, inferSchema=True)
14
    # task2-1
15
     male_customers = df.filter(col("CODE_GENDER") == "M")
16
17
     child_count_stats = male_customers.groupBy("CNT_CHILDREN").count()
18
19
     total male customers = male customers.count()
     child_count_stats = child_count_stats.withColumn("TYPE_RATIO", col("coun")
20
     t") / total male customers)
21
     # child_count_stats.select("CNT_CHILDREN", "TYPE_RATIO").show(truncate=Fal
22
     se)
     outputpath task1 = "./result/task2-1.txt"
23
24 • with open(outputpath task1, 'w') as file:
25
         file.write("CNT CHILDREN, TYPE RATIO\n")
         for row in child count stats.collect():
26 =
             file.write("{},{}\n".format(row['CNT_CHILDREN'], row['TYPE_RATI
27
     0'1))
28
29
    # task2-2
     income_stats = df.withColumn("avg_income", F.abs(col("AMT_INCOME_TOTAL")
30
     / col("DAYS BIRTH")))
31
     filtered_income_stats = income_stats.filter(col("avg_income") > 1).orderBy
     (col("avg_income").desc())
32
33
     outputpath task2 = "./result/task2-2.csv"
34
     pandas_filtered_income_stats = filtered_income_stats.select("SK_ID_CURR",
     "avg income").toPandas()
35
     pandas_filtered_income_stats.to_csv(outputpath_task2, header=True, index=F
     alse)
```

```
36
37 spark.ston()
```

统计所有男性客户(CODE\_GENDER=M)的小孩个数(CNT\_CHILDREN)类型占比情况。

```
🕏 task2.py 3
              ≡ task2-1.txt ×
result > ≡ task2-1.txt
     CNT CHILDREN, TYPE RATIO
  2 1,0.215688327511
  3 6,0.000104703071607
     3,0.0137636946858
  5 5,0.000314109214822
  6 9,9.51846105522e-06
     4,0.00161813837939
  8 8,9.51846105522e-06
  9 7,3.80738442209e-05
 10 11,9.51846105522e-06
     14,9.51846105522e-06
 11
     2,0.099115734968
 12
 13 0,0.669319144481
 14
```

统计每个客户出生以来每天的平均收入。

```
🕏 task2.py 4
              ■ task2-2.csv ×
result > III task2-2.csv
         SK ID CURR, avg income
         114967,9274.673008323425
         336147,1146.2105196128375
         385674,996.2364401151207
         190160,547.945205479452
         219563,417.51716459454445
         310601, 373.63408059023834
         157471,360.4325190228274
         252084,348.9995346672871
         199821,269.6548418024928
         337151,243.75710958236283
    11
         141198,243.62367661212704
    12
         429258,241.65939450896153
    13
         196091, 240.76187758596092
         317748,240.44883783061715
         432980,239.56558773424192
         217276,235.32048408785298
         445335,234.30843510366373
         387126,230.46532045654084
         304300,223.5839682013912
         123587,207.24967490247073
    21
         399467,195.92192148610405
         441639,192.4557351809084
    23
    24
         440768,192.04096873999788
         225210,188.75388133737036
         206341,186.00165334802975
         134526,183.68846436443792
         214063,180.08271655144367
         336135, 177. 46478873239437
         111903,174.13512885999535
         269498,172.6027397260274
         194130,172.0051983793288
         431111, 166.75931072818233
          251262 150 32023366067604
```

## 任务三

基于Spark MLlib 或者Spark ML编写程序对贷款是否违约进行分类,并评估实验结果的准确率。

## 实验设计

选用朴素贝叶斯分类器和决策树算法分别对贷款违约进行预测,比较两个模型的优劣。

#### 选取变量

#### 从5C角度考虑:

- 1. 个人品质: "NAME\_EDUCATION\_TYPE", "FLAG\_CONT\_MOBILE"
- 2. 还款能力: "AMT\_INCOME\_TOTAL","AMT\_CREDIT",
  "NAME\_INCOME\_TYPE", "CNT\_CHILDREN"
- 3. 资本实力: "FLAG\_OWN\_CAR","FLAG\_OWN\_REALTY","NAME\_HOUSING\_TYPE"
- 4. 担保:未找到明确的质押信息
- 5. 外部经营(信誉)环境: "OBS\_30\_CNT\_SOCIAL\_CIRCLE"

#### 实验结果

```
dt = DecisionTreeClassifier(featuresCol="features", labelCol="TARGET")
dt_model = dt.fit(trainingData)
dt_predictions = dt_model.transform(testData)
dt_evaluator = MulticlassClassificationEvaluator(labelCol="TARGET", pred
dt_accuracy = dt_evaluator.evaluate(dt_predictions)
print("DecisionTree Accuracy:", dt_accuracy)

hadoop@ubuntu:~/lab4$ /usr/bin/python /home/hadoop/lab4/task3.py
DecisionTree Accuracy:0.917754
NaiveBayes Accuracy:0.863485
```

从准确率的角度来说,该实验中决策树的精度高于朴素贝叶斯。

#### 实验存在的不足

由于时间关系,没有对变量做筛选再选择,并比较不同变量的效果等,而是直接采取了部分特征变量。在变量特征提取阶段,对于"CNT\_CHILDREN"这样类型变量也没有做归一化处理,可能导致模型训练精度有误。最后的比较只参考了accuracy一个指标。