

项目数据分析文档

● 实训背景

随着互联网技术的发展，网络信息呈暴涨式增长，丰富多彩的网络生活给人们带来更多体验的同时，也给高校图书馆带来更大的压力。如何让高校图书馆保持活跃，给广大师生提供优质的服务，是一个比较重要的研究课题，如何在数以十万甚至百万的书籍中给用户推荐最适合的图书是其中的关键。

HELLOWORLD2022接收到了大三的第一个订单——为一家书店制作一份推荐进货单。根据现在网络图书网站所拥有的图书数据进行图书推荐分析，了解大众喜好的同时根据不同读者的喜好，推荐给读者喜欢的图书类型、作者、出版社以及一个经济实惠的参考价格。最终交付给书店老板一张推荐进书单。

● 实训的目标

- 掌握数据仓库的概念并能设计数据仓库的库表结构。
- 熟练HDFS的shell命令的使用。
- 掌握Hive和Mysql的安装及配置。
- 掌握Hive数据仓库的内置数据类型。
- 掌握Hive的DDL语句和DML语句。
- 掌握如何建立事实表和维度表的方法。
- 掌握Sqoop的安装配置和导入导出命令。

● 实训的环境

- 项目开发环境:Pycharm2022，或者IDEA2022
- Linux操作系统: Centos7
- 分析框架: Hive3.1.3或Spark Sql
- 大数据平台: Hadoop3.3.2或2.9 Hive 3.1.3 Sqoop1.4.7开发语言:hql
- 分析框架: Hive3.1.3或Spark Sql、Spark Mlib

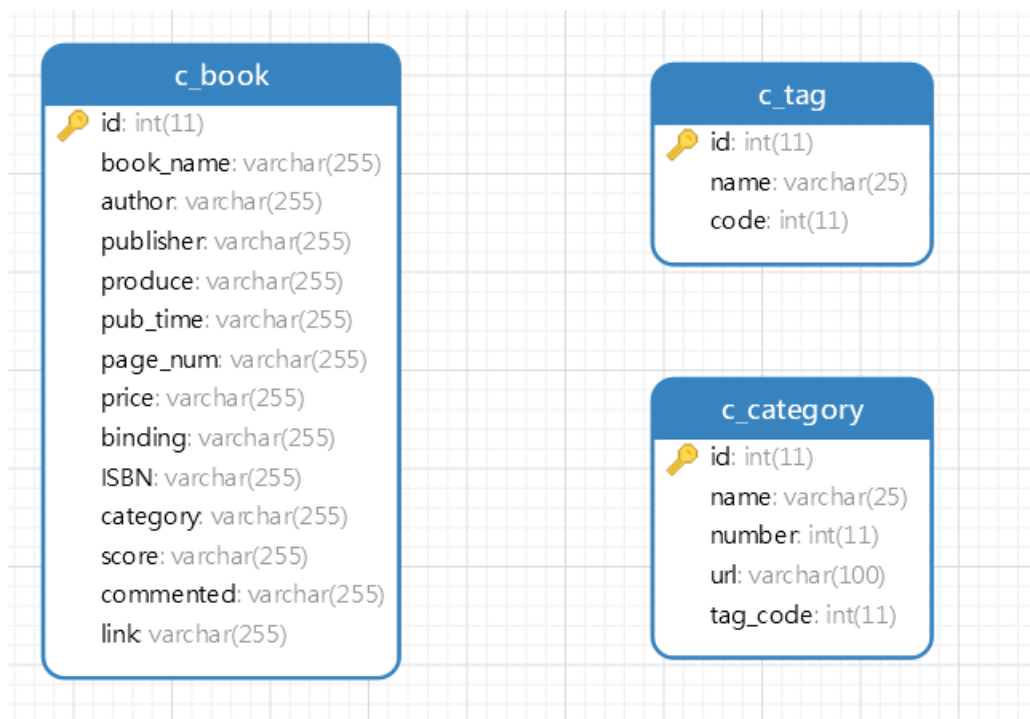
● 实训的过程

- 实训的过程包括:

设计数据仓库的库表结构、加载数据仓库需要的数据、首页界面的图书类别排名的数据分析、图书的各个类别里分数排行的类别得分排行、不同作者的图书得分排行分析、同一本书不同出版社的得分和价格区别分析、图书类别的词云分析、导出分析结果到Mysql数据库。

• 实训任务一:设计数据仓库的库表结构

举例:原型设计中招聘大数据首页界面功能点所需要的库表结构设计如下, 学生可参考这个功能的库表设计出其他功能点所需要的库表结构。



• 实训任务二:加载数据仓库的事实表数据

在slave1启动hive metastore服务:

```
hive --service metastore
```

• 实训任务三: 首页图书可视化大屏展示

1、展示图书总本数的SQL:

创建一个采集图书总数的统计表book_num:

```
bookNum Scala 复制代码
1 val BookNum = spark.sql("select count(*) as `图书总数` from bigdata.book")
```

图书总数
14683

2、展示图书大类统计的SQL:

创建一个图书大类统计的统计表book_category:

```
1 val BookCategory = spark.sql(  
2   ""  
3   |SELECT  
4   |t.`name` as `类别名称`,  
5   |count(*) as `图书数量`  
6   |FROM  
7   |bigdata.tag t  
8   |JOIN bigdata.category c ON t.`code`=c.tag_code  
9   |JOIN bigdata.book b ON c.`name`=b.category  
10  |GROUP BY t.`name`  
11  |"".stripMargin  
12  )
```

```
=====2.图书大类统计=====  
+-----+-----+  
|category| num|  
+-----+-----+  
|    文学|8653|  
|    科技| 343|  
|    生活| 597|  
|    经管| 432|  
|    流行|3238|  
|    文化|1420|  
+-----+-----+
```

3、展示得分等宽分布的SQL（运用了SparkMlib中的等宽分布函数Bucketizer）：

创建一个得分分布的统计表book_score：

```

1  val df1 = spark.sql("select score from bigdata.book order by score")
2  val splits = Array(Double.NegativeInfinity, 2.0, 4.0, 6.0, 8.0, 10.0, Double.PositiveInfinity)
3  new Bucketizer()
4  .setInputCol("score")
5  .setOutputCol("scorepartithon")
6  .setSplits(splits)
7  .transform(df1)
8  .createTempView("scorebuck")
9  val bookScore = spark.sql(
10  """
11  |SELECT
12  |scorepartithon ,
13  |concat("[", min(score), "-", max(score), "]") as `partitionsplit`,
14  |COUNT(*) as `num`
15  |FROM
16  |scorebuck
17  |GROUP BY scorepartithon
18  |ORDER BY scorepartithon
19  |""".stripMargin
20  )

```

```

=====3.得分分布统计=====
+-----+-----+-----+
|scorepartithon|partitionsplit| num|
+-----+-----+-----+
|          0.0|    [0.0-0.0]|1294|
|          1.0|    [2.6-3.5]|   2|
|          2.0|    [4.0-5.9]| 109|
|          3.0|    [6.0-7.9]|4029|
|          4.0|    [8.0-9.9]|9239|
|          5.0|   [10.0-10.0]|  10|
+-----+-----+-----+

```

4、展示图书价格等频分布的统计表book_price的SQL（运用了SparkMlib中的等频分布函数QuantileDiscretizer）：

```

1  val df = spark.sql("select price from bigdata.book where price > 1")
2  val discretizer = new QuantileDiscretizer()
3  .setInputCol("price")
4  .setOutputCol("partition")
5  .setNumBuckets(6).fit(df).transform(df)
6  discretizer.createTempView("qd")
7  val bookPrice = spark.sql(
8  """
9  |select
10 |partition,
11 |concat("[", min(price), "-", max(price), "]") as `partitionsplit`,
12 |count(1) as `num`
13 |from
14 |qd
15 |group by partition
16 |order by partition
17 |""".stripMargin)

```

```

=====4.图书价格分布统计=====
+-----+-----+-----+
|partition|partitionsplit| num|
+-----+-----+-----+
|      0.0|  [1.04-20.99]|2290|
|      1.0|  [21.0-29.99]|2280|
|      2.0|  [30.0-41.7]|2251|
|      3.0|  [42.0-57.3]|2344|
|      4.0|  [58.0-88.0]|2422|
|      5.0|[89.0-13670.0]|2345|
+-----+-----+-----+

```

5、展示图书出版社图书数量top10的统计表book_publisher的SQL:

bookPublisher

Scala

复制代码

```
1 val bookPublisher = spark.sql(  
2   ""  
3   |select  
4   |publisher,  
5   |ROUND(AVG(score),1) as `avgscore`,  
6   |count(1) as `num`  
7   |from  
8   |bigdata.book  
9   |group by publisher  
10  |order by count(1) desc limit 10  
11  |"".stripMargin)
```

publisher	avgscore	num
人民文学出版社	8.2	854
上海译文出版社	8.2	689
北京联合出版公司	7.7	453
新星出版社	7.8	367
译林出版社	8.4	351
南海出版公司	8.1	313
广西师范大学出版社	7.7	310
中信出版社	8.1	309
江苏凤凰文艺出版社	7.7	282
生活·读书·新知三联书店	8.3	279

6、展示图书发售时间顺序的统计表book_time的SQL:

bookTime

Scala

复制代码

```
1 val bookTime = spark.sql(  
2   ""  
3   |select  
4   |pub_time,  
5   |count(1) as `num`  
6   |from  
7   |bigdata.book  
8   |where pub_time != '其他'  
9   |group by pub_time  
10  |order by pub_time desc  
11  |"".stripMargin)
```

```

=====6.图书出版时间统计=====
+-----+----+
|pub_time|num|
+-----+----+
| 2022-9|107|
| 2022-8|239|
| 2022-7|300|
| 2022-6|265|
| 2022-5|206|
| 2022-4|162|
| 2022-3|148|
| 2022-2| 82|
| 2022-12| 1|
| 2022-10| 29|
| 2022-1|129|
| 2021-9| 92|
| 2021-8|107|
| 2021-7| 87|
| 2021-6| 71|

```

7、展示作者出版的图书数量的统计表book_author的SQL:

▼ bookAuthor

Scala

📄 复制代码

```

1  val bookAuthor = spark.sql(
2  """
3  |select
4  |author,
5  |count(1) as `num`
6  |from
7  |bigdata.book
8  |where author != '其他'
9  |group by author
10 |order by num desc limit 10
11 |""".stripMargin)

```

```

=====7.统计作者出版的图书数量=====
+-----+-----+
|          author|num|
+-----+-----+
|          鲁迅|161|
|          余华| 95|
|        王小波| 87|
|        余秋雨| 79|
|[英] 阿加莎·克里斯蒂| 77|
|          金庸| 77|
|          韩寒| 74|
|[日] 村上春树| 67|
|        村上春樹| 58|
|          古龙| 56|
|        荒木飛呂彦| 54|
|          王朔| 53|

```

8、展示评论数量多的top15的统计表book_comment的SQL:

bookComment

Scala | 复制代码

```

1  val bookComment = spark.sql(
2  """
3  |select
4  |book_name,
5  |commented
6  |from
7  |bigdata.book
8  |ORDER BY commented desc limit 15
9  |""".stripMargin)

```



```

=====8.最热图书的top15=====
+-----+-----+
|          book_name | commented |
+-----+-----+
|      追风筝的人 |    752525 |
|      解忧杂货店 |    726255 |
|          活着 |    716460 |
|        小王子 |    701191 |
|        白夜行 |    497365 |
|  嫌疑人X的献身 |    489689 |
|          三体 |    446737 |
|          围城 |    426057 |
|        白夜行 |    407002 |
|      百年孤独 |    383237 |
|      红楼梦 |    380856 |
|      挪威的森林 |    340927 |
|  房思琪的初恋乐园 |    324633 |
|          看见 |    320965 |
|  平凡的世界（全三部） |    301872 |
+-----+-----+

```

9、展示图书大类在2000年以来数量变化的统计表category_time的SQL:

▼ categoryTime

Scala

📄 复制代码

```

1  val categoryTime = spark.sql(
2  ""
3  |select
4  |t.`name`,
5  |`year`,
6  |COUNT(1) as `num`
7  |from
8  |bigdata.book b
9  |JOIN bigdata.category c ON b.category=c.`name`
10 |JOIN bigdata.tag t on c.tag_code = t.`code`
11 |WHERE pub_time != '其他' and year BETWEEN 2000 AND 2022
12 |GROUP BY t.`name`,`year`
13 |ORDER BY t.`name` DESC , `year` desc
14 |"".stripMargin)

```

```

22/10/07 14:34:38 INFO BlockManager: Initialized BlockManag
=====9.图书大类在2000年以来数量变化=====
+----+----+----+
|name|year|num|
+----+----+----+
|经管|2022| 67|
|经管|2021| 30|
|经管|2020| 39|
|经管|2019| 32|
|经管|2018| 39|
|经管|2017| 34|
|经管|2016| 18|
|经管|2015| 18|
|经管|2014| 19|
|经管|2013| 12|
|经管|2012| 16|
|经管|2011| 14|
|经管|2010| 17|
|经管|2009| 6|
|经管|2008| 7|
|经管|2007| 7|
|经管|2006| 3|
|经管|2005| 6|
|经管|2004| 4|
|经管|2003| 4|
+----+----+----+
only showing top 20 rows

```

10、展示每个图书大类最热图书top3的统计表category_top3的SQL:

```
1  val tempTable = spark.sql(  
2  ""  
3  |SELECT  
4  |t.`name` as `categoty_name`,  
5  |b.book_name,  
6  |b.commented  
7  |FROM  
8  |bigdata.book b  
9  |JOIN bigdata.category c ON b.category=c.`name`  
10 |JOIN bigdata.tag t on c.tag_code = t.`code`  
11 |GROUP BY book_name,t.`name`,commented  
12 |HAVING COUNT(t.`name`) < 3  
13 |ORDER BY t.`name`,commented DESC  
14 |"".stripMargin)  
15 tempTable.createTempView("book_temp")  
16  
17 val categoryTop3 = spark.sql(  
18 ""  
19 |select  
20 |a.categoty_name,  
21 |a.book_name,  
22 |a.commented  
23 |from book_temp as a  
24 |left join book_temp b on a.categoty_name = b.categoty_name and a.commente  
25 |group by a.categoty_name,a.book_name,a.commented  
26 |having count(b.categoty_name) < 3  
27 |order by a.categoty_name asc,a.commented desc  
28 |"".stripMargin)
```

```

=====10.每个图书大类最热图书top3=====
+-----+-----+-----+
|category_name|book_name|commented|
+-----+-----+-----+
|文化|万历十五年|182612|
|文化|穆斯林的葬礼|156079|
|文化|明朝那些事儿(1-9)|144189|
|文学|追风筝的人|752525|
|文学|活着|716460|
|文学|小王子|701191|
|流行|解忧杂货店|726255|
|流行|白夜行|497365|
|流行|嫌疑人X的献身|489689|
|生活|傲慢与偏见|219734|
|生活|窗边的小豆豆|194547|
|生活|杀死一只知更鸟|124960|
|科技|浪潮之巅|31707|
|科技|思考,快与慢|19515|
|科技|时间的秩序|18559|
|经管|富爸爸,穷爸爸|72664|
|经管|小狗钱钱|50117|
|经管|影响力|48225|

```

11、展示最热门的作者排行top10的统计表author_top10的SQL:

▼ authorTop10

Scala

复制代码

```

1  val authorTop10 = spark.sql(
2  ""
3  |SELECT
4  |author,
5  |sum(commented) as `commenteds`
6  |FROM
7  |bigdata.book
8  |WHERE author != '其他'
9  |GROUP BY author
10 |ORDER BY sum(commented) DESC limit 10
11 |"".stripMargin)

```

```

=====11.最热门的作者排行top10=====
+-----+-----+
|                author|commenteds|
+-----+-----+
|          [日] 东野圭吾|    1945203|
|              余华|    1447532|
|          [日] 村上春树|    1370741|
|              刘慈欣|    1335178|
|[哥伦比亚] 加西亚·马尔克斯|    860820|
|              三毛|    858186|
|              金庸|    766822|
|            王小波|    724032|
|            张爱玲|    579536|
|            钱锺书|    570090|
+-----+-----+

```

12、展示作者书的包装统计表author_blininding的SQL:

authorTop10 Scala [复制代码](#)

```

1  val bookauthorTemp = readMysql("book_author",spark).createTempView("book_a
2  val authorBlinding = spark.sql(
3  ""
4  |SELECT
5  | author,
6  | binding,
7  | COUNT( 1 ) as `num`
8  |FROM
9  | bigdata.book
10 |WHERE
11 | author IN (SELECT author FROM book_author) and binding in ('平装','精
12 |GROUP BY
13 | author,
14 | binding
15 |ORDER BY
16 | author DESC
17 |"".stripMargin)

```

author	binding	num
古龙	平装	52
鲁迅	平装	121
[英] 阿加莎	平装	65
王小波	精装	40
鲁迅	精装	18
古龙	其他	3
[英] 阿加莎	精装	9
余秋雨	精装	24
[英] 阿加莎	其他	3
[日] 村上春	精装	23
余秋雨	平装	38
[日] 村上春	平装	43
余秋雨	其他	10
[日] 村上春	其他	1
鲁迅	其他	12
韩寒	平装	56
余华	精装	32
王小波	平装	41
韩寒	其他	1
韩寒	精装	17
余华	平装	50
王小波	其他	6
余华	平装	66

13、展示图书评论分布统计表book_commented的SQL:

▼ bookCommented

Scala

📄 复制代码

```

1  val bookCommented = spark.sql(
2  ""
3  |select
4  |partition,
5  |concat("[", min(commented), "-", max(commented), "]") as `partitionsplit`
6  |count(1) as `num`
7  |from
8  |pricedis
9  |group by partition
10 |order by partition
11 |"".stripMargin)

```

	partition	partitionsplit	num
>	3	[639-2414]	2787
	4	[2418-752525]	2795
	0	[0-53]	2753
	2	[213-638]	2793
	1	[54-212]	2804

14、展示最终推荐书单统计表book_list的SQL:

bookList

Scala

复制代码

```

1  val finaltemp = spark.sql(
2  ""
3  |select distinct
4  |book_name,
5  |author,
6  |publisher,
7  |isbn,
8  |CASE
9  | WHEN score>=2.6 AND score<=3.5 THEN 1
10 | WHEN score>=4.0 AND score<=5.9 THEN 2
11 | WHEN score>=6.0 AND score<=7.9 THEN 3
12 | WHEN score>=8.0 AND score<=9.9 THEN 4
13 | ELSE 5
14 |END AS `score_weight`,
15 | CASE
16 | WHEN price>=1.04 AND price<=23.99 THEN 1
17 | WHEN price>=24 AND price<35 THEN 2
18 | WHEN price>=35 AND price<49 THEN 3
19 | WHEN price>=49 AND price<79 THEN 4
20 | ELSE 5
21 |END AS `price_weight`,
22 | CASE
23 |   WHEN commented>=0 AND commented<=53 THEN 1
24 |   WHEN commented>=54 AND commented<=212 THEN 2
25 |   WHEN commented>=213 AND commented<=638 THEN 3
26 |   WHEN commented>=639 AND commented<=2414 THEN 4
27 |   ELSE 5
28 | END AS `commented_weight`
29 |from
30 |bigdata.book
31 |where score != 0 and price != 0 and commented != 0
32 |"".stripMargin)

```

partition	partitionsplit	num
3	[639-2414]	2787
4	[2418-752525]	2795
0	[0-53]	2753
2	[213-638]	2793
1	[54-212]	2804

实训项目分析步骤配置

