

SHUANGLN

工作汇报

屠晓涵

2016.10.12

本周工作

1, 明确选题背景

2, 构建基于矩阵分解的推荐引擎

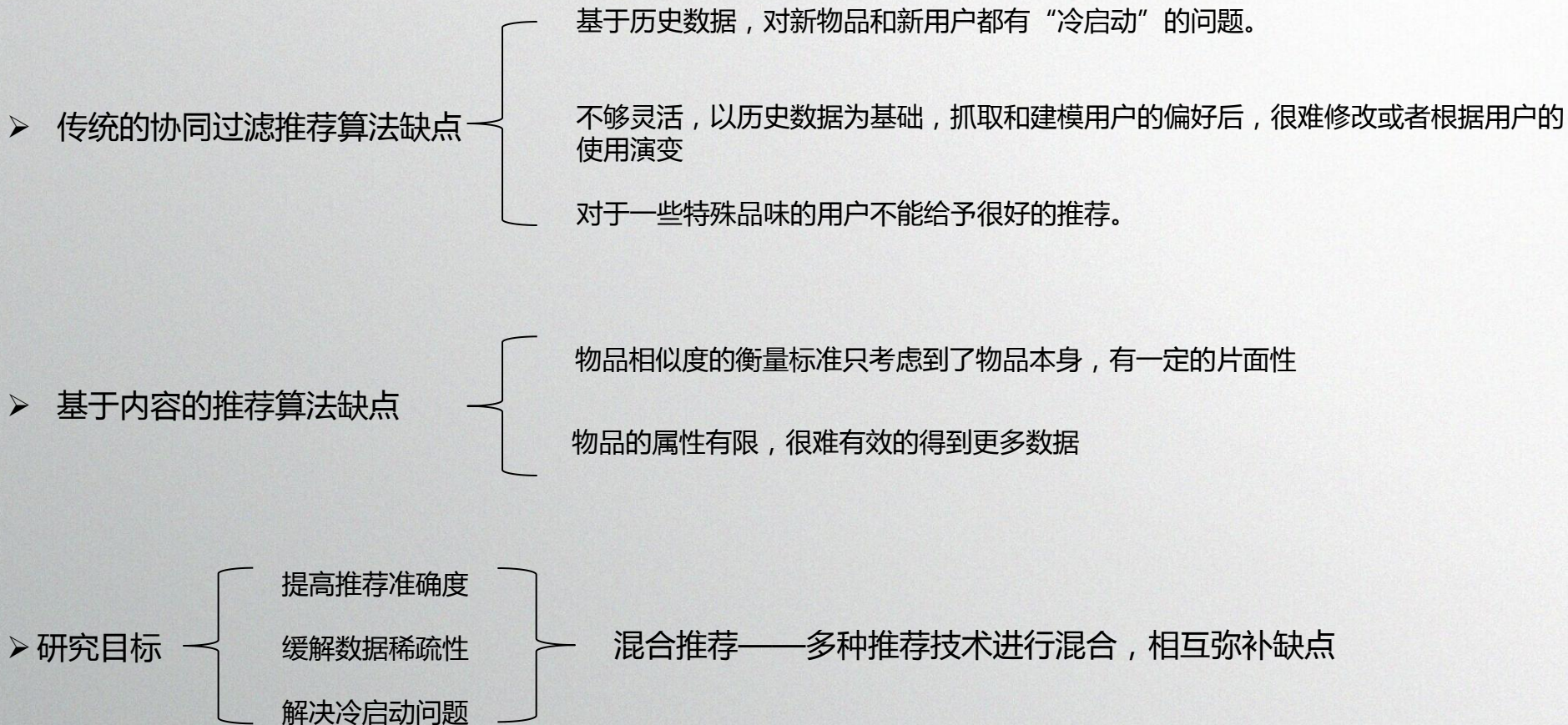
——提取有效特征

——训练推荐模型

——使用推荐模型

——对用户推荐

选题背景



● 构建基于矩阵分解推荐引擎—SVD算法原理

基于SVD的推荐：根据已有的评分情况，分析电影包含各个因子的程度和用户对于各个因子的喜爱程度，再反过来分析数据得到预测结果。

SVD原理：

$$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T p_u$$

为了防止过拟合，加入参数，得到优化函数：

$$\min \sum (r_{ui} - \mu + b_i + b_u + q_i^T p_u)^2 + \lambda(b_i^2 + b_u^2 + \|q_i\|^2 + \|p_u\|^2)$$

效果评测：

$$\text{RMSE}(P, A) = \sqrt{\frac{\sum_{i=1}^n \sum_{j=1}^m J_{ij} (A_{ij} - P_{ij})^2}{\sum_{i=1}^n \sum_{j=1}^m J_{ij}}}$$

● 构建基于矩阵分解推荐引擎

实验数据集：

电影评分数据集：movielens，有1k,1m,10m,20m数据量

<http://www.grouplens.org/datasets/movielens>

1) 用户数据文件 (6040个用户) 每个用户至少有20个评分

用户ID::性别::年龄::职业编号::邮编

2) 电影数据文件 (3883行记录) 3706个电影

电影ID::电影名称::电影种类

3) 评
用户

```
[root@centos ml-1m]# tail movies.dat
3943::Bamboozled (2000)::Comedy
3944::Bootmen (2000)::Comedy|Drama
3945::Digimon: The Movie (2000)::Adventure|Animation|Children's
3946::Get Carter (2000)::Action|Drama|Thriller
3947::Get Carter (1971)::Thriller
3948::Meet the Parents (2000)::Comedy
3949::Requiem for a Dream (2000)::Drama
3950::Tigerland (2000)::Drama
3951::Two Family House (2000)::Drama
3952::Contender, The (2000)::Drama|Thriller
```

```
[root@centos ml-1m]# tail users.dat
6031::F::18::0::45123
6032::M::45::7::55108
6033::M::50::13::78232
6034::M::25::14::94117
6035::F::25::1::78734
6036::F::25::15::32603
6037::F::45::1::76006
6038::F::56::1::14706
6039::F::45::0::01060
6040::M::25::6::11106
```

```
[root@centos ml-1m]# tail ratings.dat
6040::2022::5::956716207
6040::2028::5::956704519
6040::1080::4::957717322
6040::1089::4::956704996
6040::1090::3::956715518
6040::1091::1::956716541
6040::1094::5::956704887
6040::562::5::956704746
6040::1096::4::956715648
6040::1097::4::956715569
```

● 构建基于矩阵分解推荐引擎

提取有效特征：

- 加载数据

```
val rawData = sc.textFile("/movielens/ml-100k/u.data")
```

- 抽取 用户id,电影id,评分 3种数据：

```
val rawRatings = rawData.map(_.split("\t").take(3))
```

- 转换为 Rating对象

```
val ratings = rawRatings.map { case Array(user, movie, rating) =>  
    Rating(user.toInt, movie.toInt, rating.toDouble) }
```

```
res5: org.apache.spark.mllib.recommendation.Rating = Rating(196,242,3.0)
```

● 构建基于矩阵分解推荐引擎

训练推荐模型：

- `val model = SVD.train(ratings, 50, 10, 0.01)`
- `model.userFeatures`
- `model.userFeatures.count`

由于spark操作数延迟性的转化操作。只有当物品因子或用户因子RDD调用了执行操作时，计算才发生。

使用推荐模型：

预测：通常有2种，为某用户推荐物品，或找出与某物品相关或相似的其他物品

预测单个用户对单个电影的评分

```
val predictedRating = model.predict(789, 123)
```

● 构建基于矩阵分解推荐引擎

为推荐推荐电影：

```
scala> val userId = 789
userId: Int = 789

scala> val K = 10
K: Int = 10

scala> val topKRecs = r
```

```
Rating(789,56,5.793640
Rating(789,1098,5.3348
Rating(789,789,5.30942
Rating(789,614,5.21985
Rating(789,12,5.196173
Rating(789,42,5.194872
Rating(789,708,5.09558
Rating(789,187,5.07636
Rating(789,272,5.07128
Rating(789,129,5.02280
```



主页

信息

用户信息

展示交互多的用户top-n

图表展示

电影信息

搜索服务

作业管理

作业执行列表

提交作业

程序管理

主页 > 分析 > 用户详细信息

用户详细信息：

用户ID：1， 年龄：24， 性别：M， 职业：technician， 邮编：85711

总共评分的电影数量：272

推荐电影10部：

337, House of Yes, The (1997), 5.6866927440369

362, Blues Brothers 2000 (1998), 5.5925973691194

390, Fear of a Black Hat (1993), 6.1546333725552

516, Local Hero (1983), 6.1976937181776

618, Picnic (1955), 5.6303375527675

811, Thirty-Two Short Films About Glenn Gould (1993), 5.8691777171141

980, Mother Night (1996), 5.6797919485961

1126, Old Man and the Sea, The (1958), 6.2573181664471

1129, Chungking Express (1994), 6.0838361147073

1192, Boys of St. Vincent, The (1993), 5.7909167182614

相似用户：

评分电影记录：

电影ID

电影名称

● 构建基于矩阵分解推荐引擎

推荐效果评价指标：

(1) 用户满意度

指从用户角度来看，对结果的认可程度。测评有两种方法，一是进行用户调查，二是
在线评测

(2) 预测准确度

- ◆ 评分预测：均方根误差RMSE、平均绝对误差MAE
- ◆ TopN推荐：准确率/召回率

```
import org.apache.spark.mllib.evaluation.RegressionMetrics
val predictedAndTrue = ratingsAndPredictions.map { case ((user, product), (actual, predicted)) => (actual, predicted)
val regressionMetrics = new RegressionMetrics(predictedAndTrue)
println("Mean Squared Error = " + regressionMetrics.meanSquaredError)
println("Root Mean Squared Error = " + regressionMetrics.rootMeanSquaredError)
// Mean Squared Error = 0.08231947642632852
// Root Mean Squared Error = 0.2869137090247319
```

● 下周任务

下周任务： 算法改进

相关功能实现

SHUANGLN

2016

感谢您的观看

