



南開大學
Nankai University

计算机学院

大数据计算与应用实验报告

推荐算法的实现及优化

姓名：齐明杰 李佳豪 纪潇洋

专业：信息安全

2023 年 6 月 8 日

1 实验目的

- 基于深度学习，借助 SVD 算法实现协同过滤的推荐算法
- 利用训练得到的权重预测 Test.txt 文件中 (u, i) 对的评分分数

2 数据集的处理

2.1 数据集介绍

- train.txt，包含不同用户对不同物品的打分数据，每个用户所评价的物品数各不相同，用于后续训练模型。
- test.txt，包含待评分的用户序号以及其需要评价的 6 个物品编号，用于测试所训练的模型。
- ItemAttribute.txt，包含了每个物品的两个属性值 (不是所有的物品都具有两个属性值，因此需要预处理)，可以通过这两个属性值计算物品之间的相似度，更加科学地进行基于物品的协同过滤，更加精准地训练模型。
- ResultForm.txt，保存预测结果的文档，将 Test.txt 文件中待测试物品评分输出到该文档中提交。

2.2 数据预处理

2.2.1 建立物品下标的有序映射

由于物品的序号不连续，需要建立有序的映射。

首先定义一个空集合 all_nodes，用于存储所有的节点；打开 train_path 文件，并循环读取每一行内容（使用 while ... readline() 实现）。每行内容包括两个部分：节点数量和节点关联信息；对于每一行，首先使用 strip().split('|') 方法切分出节点数量和节点关联信息，并将节点数量转化为整数类型。并提取其中的节点 ID，将其加入到 all_nodes 集合中；完成所有行的读取后，使用 sort 方法对 all_nodes 中的节点从小到大排序，为每个节点指定一个唯一的索引。

```
1 def get_idx(train_path):
2     all_nodes = set()
3     with open(train_path, 'r') as f:
4         while (line := f.readline()) != '':
5             _, num = map(int, line.strip().split('|'))
6             for _ in range(num):
7                 line = f.readline()
8                 item_id, _ = map(int, line.strip().split())
9                 all_nodes.add(item_id)
10    node_idx = {node: idx for idx, node in enumerate(sorted(all_nodes))}
11    return node_idx
```

2.2.2 处理训练集和测试集

具体分为以下几个步骤：

1. `get_train_data` 函数：读取训练集文件 `train.txt`，并将数据按照用户为键和物品为键分别存储到字典 `data_user` 和 `data_item` 中。在处理每条记录时，先读取用户 ID 和该用户评价的物品数量，再读取该用户评价的物品及其评分，对应地将物品 ID 和评分存储到 `data_user` 和 `data_item` 中。由于原始数据中物品 ID 不是连续的数字，因此需要一个 `node_idx` 字典来映射原始 ID 与连续 ID 之间的关系。
2. `get_attribute_data` 函数：读取物品属性文件 `itemAttribute.txt`，并将数据按照物品 ID 存储到字典 `attrs` 中。在处理每条记录时，先读取物品 ID 及其两个属性值 `attr1` 和 `attr2`，将它们转换成 0/1 格式后，将物品 ID 和属性列表存储到 `attrs` 中。只有在 `node_idx` 中出现的物品才会被存储，也就是说只有这些物品有相关的属性信息。
3. `get_data_by_attr` 函数：根据物品属性来划分数据。这个函数接受两个参数，第一个是之前得到的 `data_user` 字典，第二个是之前得到的 `attrs` 字典。返回两份数据，第一份是有属性 1 的，第二份是有属性 2 的。这里将属性值为 1 和 2 的物品对应的评价数据分别提取出来，存储到 `data1` 和 `data2` 字典中。
4. `get_test_data` 函数：读取测试集文件 `test.txt`，并按照用户 ID 存储到字典 `data` 中。在处理每条记录时，先读取用户 ID 及其评价的物品数量，再读取该用户评价的物品 ID，将这些物品 ID 存储到 `data` 中。
5. `split_data` 函数：将训练集数据划分为训练集和验证集。这个函数接受两个参数，第一个是之前得到的 `data_user` 字典，第二个是可选参数 `ratio`，表示训练集所占的比例，默认是 0.85。返回两份数据，第一份是训练集，第二份是验证集。

上述几个函数代码如下所示：

```
1 def get_train_data(train_path, node_idx):
2     data_user, data_item = defaultdict(list), defaultdict(list)
3     with open(train_path, 'r') as f:
4         while (line := f.readline()) != '':
5             user_id, num = map(int, line.strip().split('|'))
6             for _ in range(num):
7                 line = f.readline()
8                 item_id, score = line.strip().split()
9                 item_id, score = int(item_id), float(score)
10                # 把0-100的得分映射到0-10
11                score = score / 10
12                data_user[user_id].append([node_idx[item_id], score])
13                data_item[node_idx[item_id]].append([user_id, score])
14    return data_user, data_item
15
16
```

```
17 def get_attribute_data(attribute_path, node_idx):
18     attrs = defaultdict(list)
19     with open(attribute_path, 'r') as f:
20         while (line := f.readline()) != '':
21             item_id, attr1, attr2 = line.strip().split('|')
22             attr1 = 0 if attr1 == 'None' else 1
23             attr2 = 0 if attr2 == 'None' else 1
24             item_id = int(item_id)
25             if item_id in node_idx:
26                 attrs[node_idx[item_id]].extend([attr1, attr2])
27     return attrs
28 def get_data_by_attr(data_user, attrs):
29     # 返回两份数据，第一份是有属性1的，第二份是有属性2的
30     data1, data2 = defaultdict(list), defaultdict(list)
31     for user_id, items in data_user.items():
32         for item_id, score in items:
33             if attrs[item_id]:
34                 if attrs[item_id][0] == 1:
35                     data1[user_id].append([item_id, score])
36                 if attrs[item_id][1] == 1:
37                     data2[user_id].append([item_id, score])
38     return data1, data2
39 def get_test_data(test_path):
40     data = defaultdict(list)
41     with open(test_path, 'r') as f:
42         while (line := f.readline()) != '':
43             user_id, num = map(int, line.strip().split('|'))
44             for _ in range(num):
45                 line = f.readline()
46                 item_id = int(line.strip())
47                 data[user_id].append(item_id)
48     return data
49
50 def split_data(data_user, ratio=0.85, shuffle=True):
51     train_data, valid_data = defaultdict(list), defaultdict(list)
52     for user_id, items in data_user.items():
53         if shuffle:
54             np.random.shuffle(items)
55             train_data[user_id] = items[:int(len(items) * ratio)]
56             valid_data[user_id] = items[int(len(items) * ratio):]
57     return train_data, valid_data
```

2.2.3 统计数据特征

由于训练过程中需要用到全评分均值, 用户偏差, 物品偏差, 因此需要对处理好的数据进行数据统计, 方便后续进行偏差的计算。我们通过以下代码读取 train.txt 中的用户 ID, 物品 ID, 统计出 ID 的数量, 最大 ID 等信息。

```
1 with open(file_path, 'r') as f:
2     for line in f:
3         if '|' in line:
4             user_id, num_ratings = line.strip().split('|')
5             user_id = int(user_id)
6             user_set.add(user_id)
7             rating_count += int(num_ratings)
8
9             if user_id > max_user_id:
10                max_user_id = user_id
11        else:
12            item_id, _ = line.strip().split()
13            item_id = int(item_id)
14            item_set.add(item_id)
15
16            if item_id > max_item_id:
17                max_item_id = item_id
```

最终得到如下结果:

- Number of users: 19835
- Number of rated items: 455705
- Number of ratings: 5001507
- Max user ID: 19834
- Max item ID: 624960

2.2.4 获取均值与偏差

- 计算全局评分均值 μ : 遍历 train_data_user, 对每个用户的评分进行累加 sum, 将 sum 加到全局评分均值 μ 中, 将 μ 除以评分总数 ratings_num, 得到全局评分均值 μ
- 计算用户偏差 b_u : 遍历 train_data_user, 对每个用户的评分进行累加 sum, 将 sum 除以该用户的评分总数 len(train_data_user[user_id]), 得到该用户的平均评分, 将该用户的平均评分 $b_u[user_id]$ 减去全局评分均值 μ , 得到该用户的偏差 b_u
- 计算物品偏差 b_i : 遍历 train_data_item, 对每个物品的评分进行累加 sum, 将 sum 除以该物品的评分总数 len(train_data_item[item_id]), 得到该物品的平均评分, 将该物品的平均评分

$bi[item_id]$ 减去全局评分均值 miu ，得到该物品的偏差 b_i 返回三个变量 miu 、 b_x 和 b_i ，用于后续的评分预测

```

1 def get_bias(train_data_user, train_data_item):
2     """
3     :param train_data_user: 用户-[物品, 评分]字典
4     :param train_data_item: 物品-[用户, 评分]字典
5     :return: 全评分均值, 用户偏差, 物品偏差
6     """
7     miu = 0.0
8     bx = np.zeros(user_num, dtype=np.float64)
9     bi = np.zeros(item_num, dtype=np.float64)
10    for user_id in train_data_user:
11        sum = 0.0
12        for item_id, score in train_data_user[user_id]:
13            miu += score
14            sum += score
15        bx[user_id] = sum / len(train_data_user[user_id])
16    miu /= ratings_num
17    for item_id in train_data_item:
18        sum = 0.0
19        for user_id, score in train_data_item[item_id]:
20            sum += score
21        bi[item_id] = sum / len(train_data_item[item_id])
22    bx -= miu
23    bi -= miu
24    return miu, bx, bi

```

3 实验原理

3.1 基于物品 (item) 的协同过滤算法

该方法的基本思想是使用物品之间的相似度来进行预测，即对于物品 i 找到其他相似的对象，然后基于对相似对象们的评分 (Ratings)，估计用户对 i 的评分。这里的相似度计算和预测函数可以沿用 User-user 中的方法。

$$r_{xi} = \frac{\sum_{j \in N(i;x)} s_{ij} \cdot r_{xj}}{\sum_{j \in N(i;x)} s_{ij}}$$

这里的 s_{ij} 表示对象 i 和 j 的相似度， r_{xj} 表示用户 x 对对象 j 的评分， $N(i; x)$ 一系列与 i 相似且被用户 x 评过分的对象。

具体到应用场景上，User-user 因为是基于相似用户，所以对社交性有更高的要求，适用于新闻推荐场景。因为新闻本身的兴趣点往往是分散的，相比用户对不同新闻的兴趣偏好，新闻的及时性、热点性

往往是其更重要的属性，而 User-user 正适用于发现热点，以及跟踪热点的趋势。另一方面，Item-item 更适用于兴趣变化较为稳定的应用。比如在 Amazon 的电商场景中，用户在一个时间段内更倾向于寻找一类商品，这时利用物品相似度为其推荐相关物品是契合用户动机的。

在实际情况中，Item-item 协同推荐比 User-user 协同推荐效果更好，因为用户的口味可能是多样的，对象相对简单。

3.2 局部和全局效应建模

在之前的模型中，存在一些问题

- 相似度衡量标准是“任意的”
- 成对相似性忽略用户之间的相互依赖关系
- 采用加权平均可能会有限制

因此我们提出解决方案：可以直接用机器学习的方法从数据中估计出 w_{ij} ，而不是使用 s_{ij} 来衡量相似度。

$$r_{xi} = b_{xi} + \frac{\sum_{j \in \mathcal{N}(i;x)} s_{ij} \cdot (x_{xj} - b_{xj})}{\sum_{j \in \mathcal{N}(i;x)} s_{ij}}$$

其中：

$$b_{xi} = \mu + b_x + b_i$$

整体平均评分为 μ ，用户 x 的评分偏差为 b_x = 用户 x 的平均评分 μ ，电影 i 的评分偏差为 b_i = 电影 i 的平均评分 μ 。

所以这里我们用 w_{ij} 计算加权总和，从数据中进行评估，而不直接使用相似度计算，公式如下：

$$\hat{r}_{xi} = b_{xi} + \sum_{j \in \mathcal{N}(i;x)} w_{ij} (r_{xj} - b_{xj})$$

3.3 矩阵分解

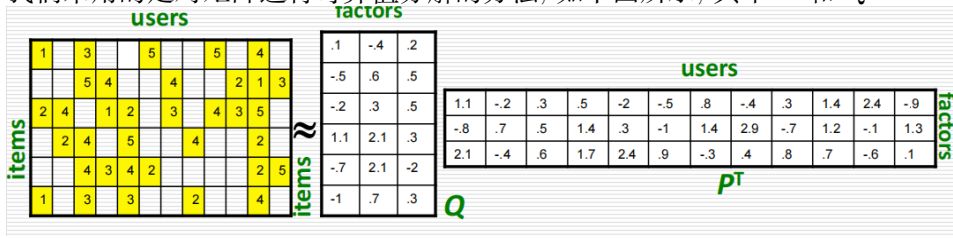
我们之前使用权值 w_{ij} 得出了一个预测公式，权重 w_{ij} 根据其作用得出，明确说明相邻电影之间的相互关系。接下来，要关注潜在因子模型 (Latent factor model)，以提取出“区域”相关性 (“regional” correlations)。

由于 CF 本身存在泛化能力较弱的缺点，无法将两个物品相似这一信息推广到其他物品的相似性计算之上。这会造成一个比较严重的后果，那就是很强的“头部效应”，容易与大量对象有高相似性；而尾部对象（冷门对象）因为特征向量比较稀疏，与其他对象的相似性就会比较低，因而很少被推荐。如下所示由 4 个对象组成的共现矩阵。

$$\begin{matrix} A[0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 & 1] \\ B[0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0] \\ C[0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0] \\ D[1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1] \end{matrix} \Rightarrow \begin{matrix} A & B & C & D \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \begin{bmatrix} - & 0.00 & 0.00 & 0.71 \\ 0.00 & - & 0.00 & 0.18 \\ 0.00 & 0.00 & - & 0.18 \\ 0.71 & 0.18 & 0.18 & - \end{bmatrix} \end{matrix}$$

可以看到 A, B, C 之间毫无相似性，但 D 却与 A, B, C 都有相似性。然而 D 的这些相似性仅仅因为它比较热门，这就是 CF 最大的缺陷：头部效应明显，无法处理稀疏矩阵。为解决上述问题，同时增加模型的泛化能力，矩阵分解技术被提出。该方法在协同过滤共现矩阵的基础上，使用更稠密的隐向量表示用户和物品，挖掘用户和物品的隐含兴趣和隐含特征，在一定程度上弥补了协同过滤模型处理稀疏矩阵能力不足的问题。

我们采用的是对矩阵进行奇异值分解的方法, 如下图所示, 其中 P 和 Q



这里的 factors 是一个超参数, 即隐向量的维度, 由用户决定, factor 数量越大, 隐向量表达能力更强, 效果会越好, 但也会增大计算量, 降低泛化性, 因此这需要权衡。在该方法中, 我们希望重构出的矩阵中“已知”评分的误差尽可能小, 而未知评分的部分我们暂时忽略。我们想知道用户 i 对对象 j 的评分, 只需要用对应的行列相乘即可得到结果。

我们采用最小化误差的方法来求 P 和 Q:

$$\min_{P,Q} \sum_{(i,x) \in R} (r_{xi} - q_i \cdot p_x)^2$$

该目标函数的目的是让原始评分 r_{ui} 与用户向量和物品向量之积 $q_i^T p_u$ 的差尽量小, 这样才能最大限度地保存共现矩阵的原始信息。这里的 K 是所有用户评分样本的集合。为了减少过拟合现象, 加入正则化项后的目标函数如下所示:

$$\min_{q^*, p^*} \sum_{(u,i) \in K} (r_{u,i} - q_i^T p_u)^2 + \lambda(\|q_i\| + \|p_u\|)^2$$

4 核心算法解析

4.1 基于物品属性的协同过滤

4.1.1 余弦相似度

d1, d2 是两个由用户评分组成的两个物品属性的向量, self.matrix_item 存储着 item_id:[(user_)], 的键值对。

```

1 def CosineCorrelation(self, id1, id2):
2     dic1 = {id1:self.matrix_item.get(id1)}
3     dic2 = {id2:self.matrix_item.get(id2)}
4     n1 = np.zeros(self.NumOfItem)
5     n2 = np.zeros(self.NumOfItem)
6     //对numpy进行初始化, 因为使用numpy进行数据计算比较快
7     for id, v1 in dic1.items():
8         for it in v1:
9             n1[it[0]] = it[1] - self.itemBia[id]
10    for id, v2 in dic2.items():
11        for it in v2:
12            n2[it[0]] = it[1] - self.itemBia[id]
13    return np.multiply(n1,n2).sum() / (np.linalg.norm(n1)* np.linalg.norm(n2))

```


注：一开始考虑到计算量，先使用计算量较小的余弦相似度而不是皮尔逊卡方进行计算相似度。

4.1.2 协同过滤核心代码

- self.matrix_user 存储了 user: [(item_id, score),] 的键值对，key 为用户 id，value 是包含了改用户打分的列表。
- 遍历用户，_score 为将来该用户对所有电影评分的数组，要是用户已经打过的电影，则无需使用协同过滤的方式进行预测；若是用户未评分的电影，对所有电影两两使用余弦相似度函数求相似度 Sij，然后取其前 neighbor 个值进行进一步计算
- self.itemBia 和 self.userBia 存储了 item 和 user 的偏置，self.GlobalMean 则是分数全局平均，根据前 neighbor 个相似物品的相似度求加权和，以及加上物品偏置和用户偏置、全局平均。

```
1 def fulfill_table(self, neighbor=2):
2     # 协同过滤
3     for user, score_lis in self.matrix_user.items():
4         _score = np.zeros(self.NumOfItem) # 分数
5         _temp = [] # 已经打分的movie_id
6         for it in score_lis:
7             _score[it[0]] = it[1]
8             _temp.append(it[0])
9         for movie in trange(self.NumOfItem):
10            # 遍历所有电影计算相似度
11            if _score[movie] == 0:
12                # 取两个计算
13                _neighbor = [(0,-1) for _ in range(neighbor)]
14                for v in _temp:
15                    Corr = self.CosineCorrelation(movie, v)
16                    if v != movie and Corr>_neighbor[0][1]:
17                        _neighbor.pop(0)
18                        _neighbor.append((v,Corr))
19                        _neighbor.sort(key=lambda x:x[1])
20                _tsum=0.0
21                for ne in _neighbor:
22                    _tsum += ne[1]
23                for ne in _neighbor:
24                    _score[movie] += _score[movie]*ne[1]/_tsum
25                _score[movie] +=
                    self.GlobalMean+self.itemBia[movie]+self.userBia[user]
```

注：由于物品数量 50w 个，使用上述方法直接计算物品两两之间相似度实际上会花费大量时间，就算将原本物品之间的相似度提前计算存入本地，共 $\binom{500000}{2}$ 条目，若使用 numpy 进行存储，将会产生近

千 GB 的空间。因此直接使用上算法进行 $O(n^3)$ 的计算, 大约要耗费 7 天的时间, 而解决该问题则希望采取并行的方式加速, 而且计算两两之间相似度时可以考虑使用随机算法寻找 neighbor。

不过, SVD 实际解决了协同过滤空间、时间的复杂度的难题, 于是下面直接讲解我们小组的主要实验。

4.2 基于全局与局部效应和 SVD 的协同过滤

下面介绍进行分数预测的三个模型;

- Base 模型: 求完整训练集的 GlobalMean, 用户 Bia 和物品偏置
- SVD 模型: 利用矩阵分解求出 P 和 Q 矩阵, 结合 base 模型进行评分
- attr_SVD 模型: 只使用对应属性的值做 svd, 再结合 base、svd 共三个模型进行评分

4.2.1 Base 模型

模型参数如下:

表 1: 模型输入输出

模型输入	模型输出
完整数据集	GlobalMean; user_bia: vector with shape (number_of_users,) item_bia: vector with shape (number_of_items)

- 初始化函数: 定义了一些必要的变量和数据结构, 包括用户偏置、物品偏置、训练集数据、测试集数据和全局平均分等。
- 计算全局平均分函数: 遍历训练集数据, 计算所有评分的平均值, 作为全局平均分。
- 预测评分函数: 根据用户偏置、物品偏置和全局平均分, 计算出给定用户对给定物品的预测评分。
- 测试函数: 遍历测试集数据, 对每个用户和物品组合进行预测评分, 并将结果保存在一个字典中。如果物品不在训练集数据中, 则将预测评分设置为全局平均分的 10 倍, 这是由于, 在读取训练分数时已经将分数除 10 以避免训练时出现数字过大溢出的情况。最后将预测结果写入到指定的文件中, 并返回预测结果字典。

Baseline 代码如下:

```
1 class Baseline:
2     def __init__(self, data_path='./data/train_user.pkl'):
3         self.bx = load_pkl(bx_pkl) # 用户偏置
4         self.bi = load_pkl(bi_pkl) # 物品偏置
5         self.idx = load_pkl(idx_pkl)
6         self.train_user_data = load_pkl(data_path)
7         self.test_data = load_pkl(test_pkl)
```

```
8         self.globalmean = self.get_globalmean() # 全局平均分
9     def get_globalmean(self):
10         score_sum, count = 0.0, 0
11         for user_id, items in self.train_user_data.items():
12             for item_id, score in items:
13                 score_sum += score
14                 count += 1
15         return score_sum / count
16     def predict(self, user_id, item_id):
17         pre_score = self.globalmean + \
18             self.bx[user_id] + \
19             self.bi[item_id]
20         return pre_score
21     def test(self, write_path='./result/result.txt'):
22         print('Start testing...')
23         predict_score = defaultdict(list)
24         for user_id, item_list in self.test_data.items():
25             for item_id in item_list:
26                 if item_id not in self.idx:
27                     pre_score = self.globalmean * 10
28                 else:
29                     new_id = self.idx[item_id]
30                     pre_score = self.predict(user_id, new_id) * 10
31                     if pre_score > 100.0:
32                         pre_score = 100.0
33                     elif pre_score < 0.0:
34                         pre_score = 0.0
35                 predict_score[user_id].append((item_id, pre_score))
36         print('Testing finished.')
37     def write_result(predict_score, write_path):
38         print('Start writing...')
39         with open(write_path, 'w') as f:
40             for user_id, items in predict_score.items():
41                 f.write(f'{user_id}|6\n')
42                 for item_id, score in items:
43                     f.write(f'{item_id} {score}\n')
44         print('Writing finished.')
45     if write_path:
46         write_result(predict_score, write_path)
47     return predict_score
```

4.2.2 SVD 模型

Input:

UserBiasMatrix: User bias matrix; ItemBiasMatrix: Item bias matrix ; P: Matrix P with shape (item, factor); Q: Matrix Q with shape (user, factor)

11.2 表 2 模型输入输出

模型输入	模型输出
UserBiasMatrix: User bias matrix	scores
ItemBiasMatrix: Item bias matrix	
P: Matrix P with shape (item, factor)	
Q: Matrix Q with shape (user, factor)	

我们将模型四个矩阵分别使用 `np.random.normal(0, 0.1, (matrix))` 进行初始化, 模型参数如下:

SVD 模型参数	
参数	值
学习率	5e-3
P 正则化系数	0.05
Q 正则化系数	0.05
b_x 正则化系数	0.05
b_i 正则化系数	0.05
隐向量维度 $factor$	50

4.2.2.1 模型训练

SVD 模型代码如下:

```
1 # 有偏置, 有正则化(四个正则化项)的矩阵分解
2 class SVD:
3     def __init__(self, model_path='./model',
4                   data_path='./data/train_user.pkl', lr=5e-3,
5                   lamda1=1e-2, lamda2=1e-2, lamda3=1e-2, lamda4=1e-2,
6                   factor=50):
7         self.bx = load_pkl(bx_pkl) # 用户偏置
8         self.bi = load_pkl(bi_pkl) # 物品偏置
9         self.lr = lr # 学习率
10        self.lamda1 = lamda1 # 正则化系数1, 乘以P
```

```
11     self.lamda2 = lamda2 # 正则化系数2, 乘以Q
12     self.lamda3 = lamda3 # 正则化系数3, 乘以bx
13     self.lamda4 = lamda4 # 正则化系数4, 乘以bi
14     self.factor = factor # 隐向量维度
15     self.idx = load_pkl(idx_pkl)
16     self.train_user_data = load_pkl(data_path)
17     self.train_data, self.valid_data = split_data(self.train_user_data)
18     self.test_data = load_pkl(test_pkl)
19     self.globalmean = self.get_globalmean() # 全局平均分
20     # 随机初始化矩阵Q(物品)和P(用户)
21     self.Q = np.random.normal(0, 0.1, (self.factor, len(self.bi)))
22     self.P = np.random.normal(0, 0.1, (self.factor, len(self.bx)))
23     self.model_path = model_path
24
25
26 def get_globalmean(self):
27     score_sum, count = 0.0, 0
28     for user_id, items in self.train_user_data.items():
29         for item_id, score in items:
30             score_sum += score
31             count += 1
32     return score_sum / count
33
34
35 def predict(self, user_id, item_id):
36     pre_score = self.globalmean + \
37         self.bx[user_id] + \
38         self.bi[item_id] + \
39         np.dot(self.P[:, user_id], self.Q[:, item_id])
40     return pre_score
41
42
43 def loss(self, is_valid=False):
44     loss, count = 0.0, 0
45     data = self.valid_data if is_valid else self.train_data
46     for user_id, items in data.items():
47         for item_id, score in items:
48             loss += (score - self.predict(user_id, item_id)) ** 2
49             count += 1
50     # 如果是训练集, 计算正则化项
51     if not is_valid:
52         loss += self.lamda1 * np.sum(self.P ** 2)
```

```
53         loss += self.lamda2 * np.sum(self.Q ** 2)
54         loss += self.lamda3 * np.sum(self.bx ** 2)
55         loss += self.lamda4 * np.sum(self.bi ** 2)
56     loss /= count
57     return loss
58
59
60     def rmse(self):
61         rmse, count = 0.0, 0
62         for user_id, items in self.train_user_data.items():
63             for item_id, score in items:
64                 rmse += (score - self.predict(user_id, item_id)) ** 2
65                 count += 1
66         rmse /= count
67         rmse = np.sqrt(rmse)
68         return rmse
69
70     def train(self, epochs=10, save=False, load=False):
71         if load:
72             self.load_weight()
73         print('Start training...')
74         for epoch in range(epochs):
75             for user_id, items in tqdm(self.train_data.items(), desc=f'Epoch {epoch + 1}'):
76                 for item_id, score in items:
77                     error = score - self.predict(user_id, item_id)
78                     self.bx[user_id] += self.lr * (error - self.lamda3 *
79                                                         self.bx[user_id])
79                     self.bi[item_id] += self.lr * (error - self.lamda4 *
80                                                         self.bi[item_id])
80                     self.P[:, user_id] += self.lr * (error * self.Q[:, item_id] -
81                                                         self.lamda1 * self.P[:, user_id])
81                     self.Q[:, item_id] += self.lr * (error * self.P[:, user_id] -
82                                                         self.lamda2 * self.Q[:, item_id])
82             print(f'Epoch {epoch + 1} train loss: {self.loss():.6f} valid loss:
83                   {self.loss(is_valid=True):.6f}')
83         print('Training finished.')
84
85         if save:
86             self.save_weight()
87
88
```

```
89 def test(self, write_path='./result/result.txt', load=True):
90     if load:
91         self.load_weight()
92         print('Start testing...')
93         predict_score = defaultdict(list)
94         for user_id, item_list in self.test_data.items():
95             for item_id in item_list:
96                 if item_id not in self.idx:
97                     pre_score = self.globalmean * 10
98                 else:
99                     new_id = self.idx[item_id]
100                     pre_score = self.predict(user_id, new_id) * 10
101                     if pre_score > 100.0:
102                         pre_score = 100.0
103                     elif pre_score < 0.0:
104                         pre_score = 0.0
105
106                 predict_score[user_id].append((item_id, pre_score))
107         print('Testing finished.')
```

训练思路

- 首先将数据集打乱后划分为训练集和测试集，比例设置为 0.85
- loss 计算：对于每个打分元组 (user_id, item_id, score)，通过 predict 函数计算预测值 \hat{score} ，score 减去预测值，计算它们的平方和。
- 梯度传导，根据求导公式，对四个输入进行更新。
- 输出每一个 epoch 在训练集的 loss
- 值得注意的是，模型两个偏置量 b_u 和 b_i 是需要在训练过程中优化的，而其初值我们将其设置为 Baseline 模型的偏置值，以达到“起跑线领先”的效果。

4.2.3 SVD_attr 模型

这个模型在上述 SVD 模型的基础上，结合了 ItemAttribute.txt 文件的使用，力求充分利用已知数据来做出更精准的预测。

- 首先，通过继承 Solution 类创建 Solution_Attr 类，需要包括模型路径、数据路径、学习率、正则化权重等参数。
- loss 函数的实现：如果 is_valid 参数为 False，表示当前处理的是训练集数据，那么就对每个用户和物品评分计算预测结果与真实结果的平方误差，并对所有误差求和并除以总数得到平均误差；如果 is_valid 参数为 True，则表示当前处理的是验证集数据。在这种情况下，不需要进行正则化，而是通过结合当前模型的预测结果和另外两个属性模型的预测结果来计算损失。这里采用了

简单的加权平均方法，其中权重分别是 0.8、0.1 和 0.1，分别表示当前模型、属性 1 模型和属性 2 模型的贡献比例。

- 将属性值运用到推荐系统 SVD 中：首先需要根据属性值构建相应的模型，接下来就是如何将不同模型的预测结果结合起来。常用的方法有加权平均、矩阵分解等。在这段代码中，我们使用了加权平均的方法，将当前模型、属性 1 模型和属性 2 模型的预测结果进行加权平均得到最终的预测结果，从而提高推荐系统的精度和覆盖率。

SVD_attr 模型代码如下：

```
1 class SVD_Attr(SVD):
2     def __init__(self, model_path='./model',
3                 data_path='./data/train_user.pkl',
4                 lr=5e-3,
5                 lamda1=1e-2, lamda2=1e-2, lamda3=1e-2, lamda4=1e-3,
6                 factor=50):
7         super().__init__(model_path, data_path, lr, lamda1, lamda2, lamda3, lamda4,
8                          factor)
9         self.attr1 = SVD(model_path='./model_attr1')
10        self.attr2 = SVD(model_path='./model_attr2')
11        self.attr1.load_weight()
12        self.attr2.load_weight()
13
14    def loss(self, is_valid=False):
15        loss, count = 0.0, 0
16        data = self.valid_data if is_valid else self.train_data
17        # 如果是训练集
18        if not is_valid:
19            for user_id, items in data.items():
20                for item_id, score in items:
21                    loss += (score - self.predict(user_id, item_id)) ** 2
22                    count += 1
23                loss += self.lamda1 * np.sum(self.P ** 2)
24                loss += self.lamda2 * np.sum(self.Q ** 2)
25                loss += self.lamda3 * np.sum(self.bx ** 2)
26                loss += self.lamda4 * np.sum(self.bi ** 2)
27            loss /= count
28            # 如果是验证集，不需要正则化，同时结合另外两个属性模型的预测结果
29            # 按0.8*本模型+0.1*属性1模型+0.1*属性2模型的比例计算loss
30        else:
31            for user_id, items in data.items():
32                for item_id, score in items:
33                    loss += (score - self.predict(user_id, item_id)) ** 2
34                    count += 1
```

```
34         loss *= 0.8
35         loss += 0.1 * self.attr1.loss(is_valid=True)
36         loss += 0.1 * self.attr2.loss(is_valid=True)
37         loss /= count
38     return loss
39
40     def predict(self, user_id, item_id):
41         return 0.8 * super().predict(user_id, item_id) + \
42             0.1 * self.attr1.predict(user_id, item_id) + 0.1 *
                self.attr2.predict(user_id, item_id)
```

5 推荐算法实验结果

5.1 Baseline 计算结果

表 4: Baseline 计算结果

用户 ID	物品 ID	评分值
0	208031	89.0074548913919
0	193714	99.3189492922157
0	393064	93.41849680352792
0	207030	100.0
0	112040	73.495419880451
0	464229	85.39440977944089
1	55971	100.0
1	583090	92.74407389873193
1	180171	100.0
1	617646	100.0
1	175835	100.0
1	553890	100.0

5.2 SVD 计算结果

表 5: SVD 计算结果

用户 ID	物品 ID	评分值
0	208031	91.56452132477179
0	193714	65.18324435919928
0	393064	37.28245349410778
0	207030	79.5223011677439
0	112040	51.451792136036715
0	464229	77.17008874354663
1	55971	89.14414602588442
1	583090	87.30347518892893
1	180171	87.8170572411186
1	617646	91.58699804481857
1	175835	89.11096494711528
1	553890	89.7784090003257

5.3 带属性的 SVD 计算结果

表 6: 带属性的 SVD 计算结果

用户 ID	物品 ID	评分值
0	208031	98.86593463013403
0	193714	68.80643768416341
0	393064	84.07892229673476
0	207030	91.38121457103834
0	112040	59.83807203259182
0	464229	78.60939763535667
1	55971	89.88532431945927
1	583090	83.57154890059475
1	180171	90.84206666471499
1	617646	91.65473111878568
1	175835	89.76072049399745
1	553890	90.0442265908409

6 属性的利用

题目的 itemAttribute.txt 文件中给出物品的两个特征, 直观上看, 这两个属性一定程度决定了物品之间的相似度, 即将两个属性组成的元组 (Attr1, Attr2), 看作物品投影到二维空间的向量, 那么使用余弦相似度 (皮尔逊卡方相似度) 可以得到二者之间的相似度。而我们在使用协同过滤计算物品之间相似度时, 每一个物品的维数都是 user 的数量 (本实验中为 19w), 维数过高导致计算量巨大。那么使用这个二维向量计算其相似度可以减少计算量。

同时第二个问题是物品数量规模在 50w, 计算两两之间相似度依旧是耗费时间; 我们小组的解决方案是利用 SVD, 对**单独含有属性 1 或者属性 2 的三元组** (user, item, score) 进行训练, 最终将导致同时含有属性 1 或者含有属性 2 的物品有了更高的相似度, 为其增加了权重。

但缺点便是我们没有利用属性的具体大小, 使得物品属性的条件减弱。

7 SVD 过程中的问题和思考

7.1 时间与空间

- **时间:** svd 模型每一个 epoch 大概需要 50s 的训练时间。属性提取后的两个 svd 属性模型训练每个 epoch 大概 25s 的时间, 带两个 svd 属性模型的主模型训练仍需 50s/epoch。
- **磁盘:** 共 500 万条数据, 我们以稀疏的形式存储两次 (item_matrix 和 user_matrix), 每个 37mb, PQ 矩阵大小分别为 [item_number:factor],[user_number:factor], P 矩阵在磁盘存储占 7mb, Q 占 170mb
- **内存:** 训练可选择读取权重继续训练 or 不读取权重开始训练, 我们需要读取已经存储为 pkl 文件的训练集数据, 同时若选择加载权重, 则需要读取磁盘存储的 PQ 矩阵, 读入后训练过程内存大约占用 3GB。测试过程占用内存较少。

7.2 梯度爆炸

一开始做 SVD 时，将原始分数输入模型，进行训练集上的训练、梯度，但由于计算 score 值介于 (0-100)，导致在开始训练的 epoch 中 loss 过大导致梯度爆炸，直接程序崩溃。

解决方法

- 调整学习率大小，继续减少其值，但是训练出效果差，在测试集上的 loss 只能减低到原来的 85%。
- 将分数映射，一开始我们直接映射到 (0, 1)，但随着训练的进行，由于每个 epoch 要计算平均 loss 时要做大整数的除法，使得其精度出现问题，loss 下降并不明显，且难以观察是否收敛。
- 最终将 score 映射至 (0, 10)，lr 调整至 $1e-3$ 级别，基本保证 loss 曲线正常。

7.3 参数估计

- **factor** 使用 SVD 分解，旨在于使用 P、Q 矩阵近似原本的 [item, user] 的打分矩阵，factor 过大 (上万)，导致内存无法承载；过小又会导致难以对原矩阵进行预测。但多大算大以及多小算小成为问题。
- **正则惩罚项系数** 惩罚项旨在于防止参数在训练集上拟合过度，但惩罚项又不能过大，否则使得训练的梯度过多来自于惩罚项而非真实预测值和真实值的偏差；过小则导致过拟合现象严重，使其在验证集上表现极差。而我们在后期的实验中就发生很多这样的现象，验证集上的 loss 在下降到一个度后开始反弹、上升。

思考

对于 factor 的估计，我们最终使用 50 的区间，过多会使得训练时间延长，但训练效果并没有明显的提升 (这里效果观察训练集 loss 降低的快慢)，惩罚项系数定位 0.01，我们优先防止其过大导致干扰训练。

7.4 收敛判断

实验过程中，发现训练集上的 loss 曲线一直达不到收敛，而测试集的 loss 曲线同样有这样的问題，要么出现一直降低，但降低速度很慢却没有下降的趋势或者出现反弹。我们认为主要原因在于 factor 的选择不足以很好的拟合原矩阵，其次是我们参数的选取也存在偏差。

7.5 结果估计

下面选一个由 $\text{factor}=50$, $\text{lr}=0.005$, 惩罚项系数 0.05 的结果:

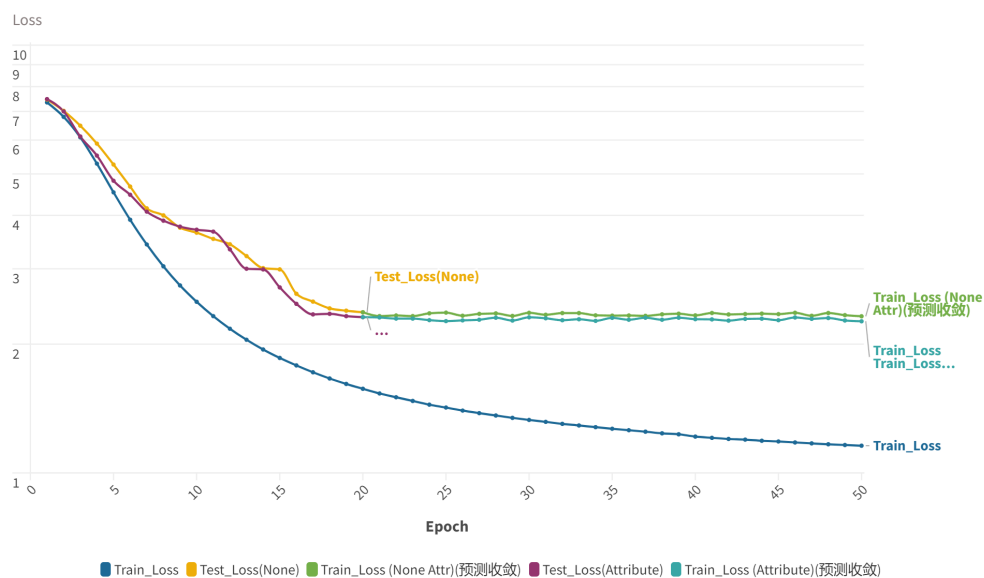


图 7.1: 结果展示

注: 绿色后缀的是根据曲线最后几个结果拟合出来的收敛的猜测值, 因为在实际实验过程中是否收敛难以判断。**分析:** 可以看到, 含有 Attribute 的 loss 曲线下降还是比没有 Attribute 开始下降快一点, 最后预测出来的收敛值也更低, 因此可以认为我们使用 Attribute 标签, 一定程度上可以优化我们原有的 SVD 模型。

致 谢

感谢杨老师耐心细致的讲解，您深入浅出的分析让我们对推荐算法有了清晰的认知并对大数据这个领域产生了浓厚的兴趣。

感谢助教李学长的耐心协调和答疑，为我们排忧解难，更加严谨地从事学术活动。

作为大二的学生凭借着兴趣的驱使在本学期选了这门大三的选修课，感谢老师们的耐心指导，收获颇丰！谨以此报告代表我对大数据计算与应用这门课的喜爱以及对学问的敬重！