

```

1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # 模型设计的代码需要用到上一节数据处理的Python类，定义如下：
5
6  # In[1]:
7
8
9  import random
10 import numpy as np
11 from PIL import Image
12
13 # 数据集总数据数: 1000209
14 # 单条数据集数据: {'usr_info': {'usr_id': 2, 'gender': 0, 'age': 56, 'job': 16},
    'mov_info': {'mov_id': 3654, 'title': [2337, 11, 4926, 26, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], 'category': [8, 7, 14, 0, 0, 0], 'years': 1961}, 'scores': 3.0}
15
16
17 class MovieLen(object):
18     def __init__(self, use_poster):
19         self.use_poster = use_poster
20         # 声明每个数据文件的路径
21         usr_info_path = "./work/ml-1m/users.dat"
22         if use_poster:
23             rating_path = "./work/ml-1m/new_rating.txt"
24         else:
25             rating_path = "./work/ml-1m/ratings.dat"
26
27         movie_info_path = "./work/ml-1m/movies.dat"
28         self.poster_path = "./work/ml-1m/posters/"
29         # 得到电影数据
30         self.movie_info, self.movie_cat, self.movie_title = self.get_movie_info(
            movie_info_path)
31         # 记录电影的最大ID
32         self.max_mov_cat = np.max([self.movie_cat[k] for k in self.movie_cat])
33         self.max_mov_tit = np.max([self.movie_title[k] for k in self.movie_title])
34         self.max_mov_id = np.max(list(map(int, self.movie_info.keys())))
35         # 记录用户数据的最大ID
36         self.max_usr_id = 0
37         self.max_usr_age = 0
38         self.max_usr_job = 0
39         # 得到用户数据
40         self.usr_info = self.get_usr_info(usr_info_path)
41         # 得到评分数据
42         self.rating_info = self.get_rating_info(rating_path)
43         # 构建数据集
44         self.dataset = self.get_dataset(usr_info=self.usr_info,
45                                         rating_info=self.rating_info,
46                                         movie_info=self.movie_info)
47         # 划分数据集，获得数据加载器
48         self.train_dataset = self.dataset[:int(len(self.dataset)*0.9)]
49         self.valid_dataset = self.dataset[int(len(self.dataset)*0.9):]
50         print("##Total dataset instances: ", len(self.dataset))
51         print("##MovieLens dataset information: \nusr num: {}\n"
52               "movies num: {}".format(len(self.usr_info), len(self.movie_info)))
53         # 得到电影数据
54         def get_movie_info(self, path):
55             # 打开文件，编码方式选择ISO-8859-1，读取所有数据到data中
56             with open(path, 'r', encoding="ISO-8859-1") as f:
57                 data = f.readlines()
58             # 建立三个字典，分别用户存放电影所有信息，电影的名字信息、类别信息
59             movie_info, movie_titles, movie_cat = {}, {}, {}
60             # 对电影名字、类别中不同的单词计数
61             t_count, c_count = 1, 1
62
63             count_tit = {}
64             # 按行读取数据并处理
65             for item in data:
66                 item = item.strip().split("::")
67                 v_id = item[0]
68                 v_title = item[1][:7]
69                 cats = item[2].split('|')

```

```

70         v_year = item[1][-5:-1]
71
72         titles = v_title.split()
73         # 统计电影名字的单词，并给每个单词一个序号，放在movie_titles中
74         for t in titles:
75             if t not in movie_titles:
76                 movie_titles[t] = t_count
77                 t_count += 1
78         # 统计电影类别单词，并给每个单词一个序号，放在movie_cat中
79         for cat in cats:
80             if cat not in movie_cat:
81                 movie_cat[cat] = c_count
82                 c_count += 1
83         # 补0使电影名称对应的列表长度为15
84         v_tit = [movie_titles[k] for k in titles]
85         while len(v_tit) < 15:
86             v_tit.append(0)
87         # 补0使电影种类对应的列表长度为6
88         v_cat = [movie_cat[k] for k in cats]
89         while len(v_cat) < 6:
90             v_cat.append(0)
91         # 保存电影数据到movie_info中
92         movie_info[v_id] = {'mov_id': int(v_id),
93                             'title': v_tit,
94                             'category': v_cat,
95                             'years': int(v_year)}
96         return movie_info, movie_cat, movie_titles
97
98     def get_usr_info(self, path):
99         # 性别转换函数, M-0, F-1
100         def gender2num(gender):
101             return 1 if gender == 'F' else 0
102
103         # 打开文件，读取所有行到data中
104         with open(path, 'r') as f:
105             data = f.readlines()
106         # 建立用户信息的字典
107         use_info = {}
108
109         max_usr_id = 0
110         # 按行索引数据
111         for item in data:
112             # 去除每一行中和数据无关的部分
113             item = item.strip().split("::")
114             usr_id = item[0]
115             # 将字符数据转成数字并保存在字典中
116             use_info[usr_id] = {'usr_id': int(usr_id),
117                                 'gender': gender2num(item[1]),
118                                 'age': int(item[2]),
119                                 'job': int(item[3])}
120             self.max_usr_id = max(self.max_usr_id, int(usr_id))
121             self.max_usr_age = max(self.max_usr_age, int(item[2]))
122             self.max_usr_job = max(self.max_usr_job, int(item[3]))
123         return use_info
124     # 得到评分数据
125     def get_rating_info(self, path):
126         # 读取文件里的数据
127         with open(path, 'r') as f:
128             data = f.readlines()
129         # 将数据保存在字典中并返回
130         rating_info = {}
131         for item in data:
132             item = item.strip().split("::")
133             usr_id, movie_id, score = item[0], item[1], item[2]
134             if usr_id not in rating_info.keys():
135                 rating_info[usr_id] = {movie_id: float(score)}
136             else:
137                 rating_info[usr_id][movie_id] = float(score)
138         return rating_info
139     # 构建数据集
140     def get_dataset(self, usr_info, rating_info, movie_info):
141         trainset = []

```

```

142     for usr_id in rating_info.keys():
143         usr_ratings = rating_info[usr_id]
144         for movie_id in usr_ratings:
145             trainset.append({'usr_info': usr_info[usr_id],
146                             'mov_info': movie_info[movie_id],
147                             'scores': usr_ratings[movie_id]})
148     return trainset
149
150 def load_data(self, dataset=None, mode='train'):
151     use_poster = False
152
153     # 定义数据迭代Batch大小
154     BATCHSIZE = 256
155
156     data_length = len(dataset)
157     index_list = list(range(data_length))
158     # 定义数据迭代加载器
159     def data_generator():
160         # 训练模式下，打乱训练数据
161         if mode == 'train':
162             random.shuffle(index_list)
163         # 声明每个特征的列表
164         usr_id_list,usr_gender_list,usr_age_list,usr_job_list = [], [], [], []
165         mov_id_list,mov_tit_list,mov_cat_list,mov_poster_list = [], [], [], []
166         score_list = []
167         # 索引遍历输入数据集
168         for idx, i in enumerate(index_list):
169             # 获得特征数据保存到对应特征列表中
170             usr_id_list.append(dataset[i]['usr_info']['usr_id'])
171             usr_gender_list.append(dataset[i]['usr_info']['gender'])
172             usr_age_list.append(dataset[i]['usr_info']['age'])
173             usr_job_list.append(dataset[i]['usr_info']['job'])
174
175             mov_id_list.append(dataset[i]['mov_info']['mov_id'])
176             mov_tit_list.append(dataset[i]['mov_info']['title'])
177             mov_cat_list.append(dataset[i]['mov_info']['category'])
178             mov_id = dataset[i]['mov_info']['mov_id']
179
180         if use_poster:
181             # 不使用图像特征时，不读取图像数据，加快数据读取速度
182             poster = Image.open(self.poster_path+'mov_id{}.jpg'.format(str(
183                 mov_id[0])))
184             poster = poster.resize([64, 64])
185             if len(poster.size) <= 2:
186                 poster = poster.convert("RGB")
187
188             mov_poster_list.append(np.array(poster))
189
190         score_list.append(int(dataset[i]['scores']))
191         # 如果读取的数据量达到当前的batch大小，就返回当前批次
192         if len(usr_id_list)==BATCHSIZE:
193             # 转换列表数据为数组形式，reshape到固定形状
194             usr_id_arr = np.array(usr_id_list)
195             usr_gender_arr = np.array(usr_gender_list)
196             usr_age_arr = np.array(usr_age_list)
197             usr_job_arr = np.array(usr_job_list)
198
199             mov_id_arr = np.array(mov_id_list)
200             mov_cat_arr = np.reshape(np.array(mov_cat_list), [BATCHSIZE, 6]).
201                 astype(np.int64)
202             mov_tit_arr = np.reshape(np.array(mov_tit_list), [BATCHSIZE, 1,
203                 15]).astype(np.int64)
204
205             if use_poster:
206                 mov_poster_arr = np.reshape(np.array(mov_poster_list)/127.5 -
207                     1, [BATCHSIZE, 3, 64, 64]).astype(np.float32)
208             else:
209                 mov_poster_arr = np.array([0.])
210
211             scores_arr = np.reshape(np.array(score_list), [-1, 1]).astype(np.
212                 float32)

```

```

209         # 放回当前批次数据
210         yield [usr_id_arr, usr_gender_arr, usr_age_arr, usr_job_arr],
                [mov_id_arr, mov_cat_arr, mov_tit_arr
                 , mov_poster_arr], scores_arr

211
212         # 清空数据
213         usr_id_list, usr_gender_list, usr_age_list, usr_job_list = [],
                [], [], []
214         mov_id_list, mov_tit_list, mov_cat_list, score_list = [], [], [],
                []
215         mov_poster_list = []
216     return data_generator
217
218
219 # In[1]:
220
221
222 # 解压数据集
223 get_ipython().system('unzip -o -q -d ~/work/ ~/data/data19736/ml-1m.zip')
224
225
226 # # 模型设计介绍
227
228 #
神经网络模型设计是电影推荐任务中重要的一环。它的作用是提取图像、文本或者语音的特征，利用这些特征完成分类、检测、文本分析等任务。在电影推荐任务中，我们将设计一个神经网络模型，提取用户数据、电影数据的特征向量，然后计算这些向量的相似度，利用相似度的大小去完成推荐。
229 #
230 # 根据第一章中对建模思路的分析，神经网络模型的设计包含如下步骤：
231 # 1. 分别将用户、电影的多个特征数据转换成特征向量。
232 # 2. 对这些特征向量，使用全连接层或者卷积层进一步提取特征。
233 # 3. 将用户、电影多个数据的特征向量融合成一个向量表示，方便进行相似度计算。
234 # 4. 计算特征之间的相似度。
235 #
236 # 依据这个思路，我们设计一个简单的电影推荐神经网络模型：
237 #
238 # <center></center>
239 #
240 # <center><br>图1：网络结构的设计 </br></center>
241 # <br></br>
242 #
243
244 # 该网络结构包含如下内容：
245 #
246 # 1. 提取用户特征和电影特征作为神经网络的输入，其中：
247 # * 用户特征包含四个属性信息，分别是用户ID、性别、职业和年龄。
248 # * 电影特征包含三个属性信息，分别是电影ID、电影类型和电影名称。
249 #
250 # 2.
提取用户特征。使用Embedding层将用户ID映射为向量表示，输入全连接层，并对其他三个属性也做类似的处理。然后将四个属性的特征分别全连接并相加。
251 #
252 # 3.
提取电影特征。将电影ID和电影类型映射为向量表示，输入全连接层，电影名字用文本卷积神经网络得到其定长向量表示。然后将三个属性的特征表示分别全连接并相加。
253 #
254 # 4.
得到用户和电影的向量表示后，计算二者的余弦相似度。最后，用该相似度和用户真实评分的均方差作为该回归模型的损失函数。
255 # ><font
size=2>衡量相似度的计算有多种方式，比如计算余弦相似度、皮尔森相关系数、Jaccard相似系数等等，或者通过计算欧几里得距离、曼哈顿距离、明可夫斯基距离等方式计算相似度。余弦相似度是一种简单好用的向量相似度计算方式，通过计算向量之间的夹角余弦值来评估他们的相似度，本节我们使用余弦相似度计算特征之间的相似度。</font>
256
257 # ### 为何如此设计网络呢？
258 #
259 # 网络的主体框架已经在第一章中做出了分析，但还有一些细节点没有确定。
260 #

```

```

261 # 1. 如何将“数字”转变成“向量”？
262 #
263 # 如NLP章节的介绍，使用词嵌入（Embedding）的方式可将数字转变成向量。
264 #
265 # 2.
如何合并多个向量的信息？例如：如何将用户四个特征（ID、性别、年龄、职业）的向量合并成一个向量？
266 #
267 # 最简单的方式是先将不同特征向量（ID 32维、性别 16维、年龄 16维、职业
16维）通过4个全连接层映射到4个等长的向量（200维度），再将4个等长的向量按位相加即可得到
1个包含全部信息的向量。
268 #
269 #
电影类型的特征是将多个数字（代表出现的单词）转变成的多个向量（6个），可以通过相同的方式
合并成1个向量。
270 #
271 # 3. 如何处理文本信息？
272 #
273 #
如NLP章节的介绍，使用卷积神经网络(CNN)和长短记忆神经网络（LSTM）处理文本信息会有较好的
效果。因为电影标题是相对简单的短文本，所以我们使用卷积网络结构来处理电影标题。
274 #
275 # 4. 尺寸大小应该如何设计？
276 #
这涉及到信息熵的理念：越丰富的信息，维度越高。所以，信息量较少的原始特征可以用更短的向
量表示，例如性别、年龄和职业这三个特征向量均设置成16维，而用户ID和电影ID这样较多信息量
的特征设置成32维。综合了4个原始用户特征的向量和综合了3个电影特征的向量均设计成200维度
，使得它们可以蕴含更丰富的信息。当然，尺寸大小并没有一贯的最优规律，需要我们根据问题的
复杂程度，训练样本量，特征的信息量等多方面信息探索出最有效的设计。
277 #
278 # 第一章的设计思想结合上面几个细节方案，即可得出上图展示的网络结构。
279 #
280 #
接下来我们进入代码实现环节，首先看看如何将数据映射为向量。在自然语言处理中，我们常使用
词嵌入（Embedding）的方式完成向量变换。
281 # # Embedding介绍
282 #
283 #
Embedding是一个嵌入层，将输入的非负整数矩阵中的每个数值，转换为具有固定长度的向量。
284 # Embedding是一个嵌入层，将输入的非负整数矩阵中的每个数值，转换为具有固定长度的向量。
285 #
286 #
在NLP任务中，一般把输入文本映射成向量表示，以便神经网络的处理。在数据处理章节，我们已
经将用户和电影的特征用数字表示。嵌入层Embedding可以完成数字到向量的映射。
287 #
288 #
289 # 飞桨支持[Embedding API](
https://www.paddlepaddle.org.cn/documentation/docs/en/develop/api/paddle/nn/layer/comm
on/Embedding\_en.html
)，该接口根据输入从Embedding矩阵中查询对应Embedding信息，并会根据输入参数num_embeddings
和embedding_dim自动构造一个二维Embedding矩阵。
290 #
291 # > *class* paddle.nn.Embedding *(num_embeddings, embedding_dim, padding_idx=None,
sparse=False, weight_attr=None, name=None)*
292 #
293 # 常用参数含义如下：
294 #
295 # * num_embeddings (int)：表示嵌入字典的大小。
296 # * embedding_dim：表示每个嵌入向量的大小。
297 # * sparse
(bool)：是否使用稀疏更新，在词嵌入权重较大的情况下，使用稀疏更新能够获得更快的训练速度
及更小的内存/显存占用。
298 # * weight_attr (ParamAttr)：指定嵌入向量的配置，包括初始化方法，具体用法请参见
ParamAttr，一般无需设置，默认值为None。
299 #
300 #
301 # 我们需要特别注意，embedding函数在输入Tensor
shape的最后一维后面添加embedding_dim的维度，所以输出的维度数量会比输入多一个。以下面的
代码为例，当输入的Tensor尺寸是[1]、embedding_dim是32时，输出Tensor的尺寸是[1, 32]。
302 #
303 # In[3]:
304
305

```



```

306 import paddle
307 from paddle.nn import Linear, Embedding, Conv2D
308 import numpy as np
309 import paddle.nn.functional as F
310 import paddle.nn as nn
311
312 # 声明用户的最大ID，在此基础上加1（算上数字0）
313 USR_ID_NUM = 6040 + 1
314 # 声明Embedding层，将ID映射为32长度的向量
315 usr_emb = Embedding(num_embeddings=USR_ID_NUM,
316                     embedding_dim=32,
317                     sparse=False)
318 # 声明输入数据，将其转成tensor
319 arr_1 = np.array([1], dtype="int64").reshape((-1))
320 print(arr_1)
321 arr_pdl = paddle.to_tensor(arr_1)
322 print(arr_pdl)
323 # 计算结果
324 emb_res = usr_emb(arr_pdl)
325 # 打印结果
326 print("数字 1 的embedding结果是: ", emb_res.numpy(), "\n形状是: ", emb_res.shape)
327
328 output:
329 [1]
330 Tensor(shape=[1], dtype=int64, place=Place(gpu:0), stop_gradient=True,
331        [1])
332 数字 1 的embedding结果是:  [[ 0.02368815  0.01219996 -0.00823128 -0.02978373
333    0.015901  -0.01567403
334    -0.02949063  0.01960909  0.00287736  0.02580381 -0.01716401  0.02730818
335    -0.00820427  0.01684101 -0.02887885  0.00482129  0.00490872  0.01330269
336    -0.02448237 -0.00270003 -0.01551332 -0.0038403   0.01186426  0.00623586
337     0.01695438 -0.02498322  0.02353216  0.02606978 -0.003106   0.00167086
338    -0.00091827 -0.00629074]]
339 形状是:  [1, 32]
340
341 #
342 # 使用Embedding时，需要注意`num_embeddings`和`embedding_dim`这两个参数。`num_embeddings`表示词表大小；`embedding_dim`表示Embedding层维度。
343 #
344 # 使用的ml-1m数据集的用户ID最大为6040，考虑到0号ID的存在，因此这里我们需要将num_embeddings设置为6041(=6040+1)。embedding_dim表示将数据映射为embedding_dim维度的向量。这里将用户ID数据1转换成了维度为32的向量表示。32是设置的超参数，读者可以自行调整大小。
345 #
346 # 通过上面的代码，我们简单了解了Embedding的工作方式，但是Embedding层是如何将数字映射为高维度的向量的呢？
347 #
348 # 实际上，Embedding层和Conv2D，Linear层一样，Embedding层也有可学习的权重，通过矩阵相乘的方法对输入数据进行映射。Embedding中将输入映射成向量的实际步骤是：
349 #
350 # 1. 将输入数据转换成one-hot格式的向量；
351 #
352 # 2. one-hot向量和Embedding层的权重进行矩阵相乘得到Embedding的结果。
353 #
354 # 下面展示了另一个使用Embedding函数的案例。该案例从0到9的10个ID数字中随机取出了3个，查看使用默认初始化方式的Embedding结果，再查看使用KaimingNormal（0均值的正态分布）初始化方式的Embedding结果。实际上，无论使用哪种参数初始化的方式，这些参数都是要在后续的训练过程中优化的，只是更符合任务场景的初始化方式可以使训练更快收敛，部分场景可以取得略好的模型精度。
355 # In[4]:
356
357 # 声明用户的最大ID，在此基础上加1（算上数字0）
358 USR_ID_NUM = 10
359 # 声明Embedding层，将ID映射为16长度的向量
360 usr_emb = Embedding(num_embeddings=USR_ID_NUM,
361                     embedding_dim=16,
362                     sparse=False)

```

```

363 # 定义输入数据，输入数据为不超过10的整数，将其转成tensor
364 arr = np.random.randint(0, 10, (3)).reshape((-1)).astype('int64')
365 print("输入数据是: ", arr)
366 arr_pd = paddle.to_tensor(arr)
367 emb_res = usr_emb(arr_pd)
368 print("默认权重初始化embedding层的映射结果是: ", emb_res.numpy())
369
370 # 观察Embedding层的权重
371 emb_weights = usr_emb.state_dict()
372 print(emb_weights.keys())
373
374 print("\n查看embedding层的权重形状: ", emb_weights['weight'].shape)
375
376 # 声明Embedding层，将ID映射为16长度的向量，自定义权重初始化方式
377 # 定义KaimingNormal初始化方式
378 init = nn.initializer.KaimingNormal()
379 param_attr = paddle.ParamAttr(initializer=init)
380
381 usr_emb2 = Embedding(num_embeddings=USR_ID_NUM,
382                     embedding_dim=16,
383                     weight_attr=param_attr)
384 emb_res = usr_emb2(arr_pd)
385 print("\nKaimingNormal初始化权重embedding层的映射结果是: ", emb_res.numpy())
386
387 # 上面代码中，我们在[0,
388 10]范围内随机产生了3个整数，因此数据的最大值为整数9，最小为0。因此，输入数据映射为每个
389 one-hot向量的维度是10，定义Embedding权重的第一个维度USR_ID_NUM为10。
390 #
391 # 这里输入的数据shape是[3, 1]，Embedding层的权重形状则是[10,
392 16]，Embedding在计算时，首先将输入数据转换成one-hot向量，one-hot向量的长度和Embedding
393 层的输入参数size的第一个维度有关。比如这里我们设置的是10，所以输入数据将被转换成维度为
394 [3, 10]的one-hot向量，参数size决定了Embedding层的权重形状。最终维度为[3,
395 10]的one-hot向量与维度为[10, 16]Embedding权重相乘，得到最终维度为[3, 16]的映射向量。
396 #
397 #
398 我们也可以对Embedding层的权重进行初始化，如果不设置初始化方式，则采用默认的初始化方式。
399 #
400 #
401 神经网络处理文本数据时，需要用数字代替文本，Embedding层则是将输入数字数据映射成了高维
402 向量，然后就可以使用卷积、全连接、LSTM等网络层处理数据了，接下来我们开始设计用户和电影
403 数据的特征提取网络。
404 #
405 #
406 # 理解Embedding后，我们就可以开始构建提取用户特征的神经网络了。
407 #
408 #
409 # <center></center>
412 # <center><br>图2：提取用户特征网络示意 </br></center>
413 #
414 # 用户特征网络主要包括：
415 # 1. 将用户ID数据映射为向量表示，通过全连接层得到ID特征。
416 # 2. 将用户性别数据映射为向量表示，通过全连接层得到性别特征。
417 # 3. 将用户职业数据映射为向量表示，通过全连接层得到职业特征。
418 # 4. 将用户年龄数据映射为向量表示，通过全连接层得到年龄特征。
419 # 5. 融合ID、性别、职业、年龄特征，得到用户的特征表示。
420 #
421 #
422 在用户特征计算网络中，我们对每个用户数据做embedding处理，然后经过一个全连接层，激活函
423 数使用ReLU，得到用户所有特征后，将特征整合，经过一个全连接层得到最终的用户数据特征，该
424 特征的维度是200维，用于和电影特征计算相似度。
425
426 ## 1. 提取用户ID特征
427 #
428 #
429 开始构建用户ID的特征提取网络，ID特征提取包括两个部分。首先，使用Embedding将用户ID映射
430 为向量；然后，使用一层全连接层和ReLU激活函数进一步提取用户ID特征。
431 #
432 相比较电影类别和电影名称，用户ID只包含一个数字，数据更为简单。这里需要考虑将用户ID映射
433 为多少维度的向量合适，使用维度过大的向量表示用户ID容易造成信息冗余，维度过低又不足以表
434 示该用户的特征。理论上来说，如果使用二进制表示用户ID，用户最大ID是6040，小于2的13次方

```

，因此，理论上使用13维度的向量已经足够了，为了让不同ID的向量更具区分性，我们选择将用户ID映射为维度为32维的向量。

下面是用户ID特征提取代码实现：

In[5]:

自定义一个用户ID数据

usr_id_data = np.random.randint(0, 6040, (2)).reshape((-1)).astype('int64')

print("输入的用户ID是:", usr_id_data)

USR_ID_NUM = 6040 + 1

定义用户ID的embedding层和fc层

usr_emb = Embedding(num_embeddings=USR_ID_NUM,
embedding_dim=32,
sparse=False)

usr_fc = Linear(in_features=32, out_features=32)

usr_id_var = paddle.to_tensor(usr_id_data)

usr_id_feat = usr_fc(usr_emb(usr_id_var))

usr_id_feat = F.relu(usr_id_feat)

print("用户ID的特征是:", usr_id_feat.numpy(), "\n其形状是:", usr_id_feat.shape)

注意到，将用户ID映射为one-hot向量时，Embedding层参数size的第一个参数是，在用户的最大ID基础上加上1。原因很简单，从上一节数据处理已经发现，用户ID是从1开始计数的，最大的用户ID是6040。并且已经知道通过Embedding映射输入数据时，是先把输入数据转换成one-hot向量。向量中只有一个1的向量才被称为one-hot向量，比如，0用四维的one-hot向量表示是[1, 0, 0, 0]，同时，4维的one-hot向量最大只能表示3。所以，要把数字6040用one-hot向量表示，至少需要用6041维度的向量。

接下来我们会看到，类似的Embedding层也适用于处理用户性别、年龄和职业，以及电影ID等特征，实现代码均是类似的。

2. 提取用户性别特征

接下来构建用户性别的特征提取网络，同用户ID特征提取步骤，使用Embedding层和全连接层提取用户性别特征。用户性别不像用户ID数据那样有数千数万种不同数据，性别只有两种可能，不需要使用高维度的向量表示其特征，这里我们将用户性别用为16维的向量表示。

下面是用户性别特征提取实现：

In[6]:

自定义一个用户性别数据

usr_gender_data = np.array((0, 1)).reshape(-1).astype('int64')

print("输入的用户性别是:", usr_gender_data)

用户的性别用0, 1表示

性别最大ID是1，所以Embedding层size的第一个参数设置为1 + 1 = 2

USR_ID_NUM = 2

对用户性别信息做映射，并紧接着一个FC层

USR_GENDER_DICT_SIZE = 2

usr_gender_emb = Embedding(num_embeddings=USR_GENDER_DICT_SIZE,
embedding_dim=16)

usr_gender_fc = Linear(in_features=16, out_features=16)

usr_gender_var = paddle.to_tensor(usr_gender_data)

usr_gender_feat = usr_gender_fc(usr_gender_emb(usr_gender_var))

usr_gender_feat = F.relu(usr_gender_feat)

print("用户性别特征的数据特征是:", usr_gender_feat.numpy(), "\n其形状是:",
usr_gender_feat.shape)


```

473 print("\n性别 0 对应的特征是: ", usr_gender_feat.numpy()[0, :])
474 print("性别 1 对应的特征是: ", usr_gender_feat.numpy()[1, :])
475
476
477 # ## 3. 提取用户年龄特征
478 # 然后构建用户年龄的特征提取网络，同样采用Embedding层和全连接层的方式提取特征。
479 #
480 # 前面我们了解到年龄数据分布是：
481 # * 1: "Under 18"
482 # * 18: "18-24"
483 # * 25: "25-34"
484 # * 35: "35-44"
485 # * 45: "45-49"
486 # * 50: "50-55"
487 # * 56: "56+"
488 #
489 # 得知用户年龄最大值为56，这里仍将用户年龄用16维的向量表示。
490
491 # In[7]:
492
493
494 # 自定义一个用户年龄数据
495 usr_age_data = np.array((1, 18)).reshape(-1).astype('int64')
496 print("输入的用户年龄是:", usr_age_data)
497
498 # 对用户年龄信息做映射，并紧接着一个Linear层
499 # 年龄的最大ID是56，所以Embedding层size的第一个参数设置为56 + 1 = 57
500 USR_AGE_DICT_SIZE = 56 + 1
501
502 usr_age_emb = Embedding(num_embeddings=USR_AGE_DICT_SIZE,
503                          embedding_dim=16)
504 usr_age_fc = Linear(in_features=16, out_features=16)
505
506 usr_age = paddle.to_tensor(usr_age_data)
507 usr_age_feat = usr_age_emb(usr_age)
508 usr_age_feat = usr_age_fc(usr_age_feat)
509 usr_age_feat = F.relu(usr_age_feat)
510
511 print("用户年龄特征的数据特征是: ", usr_age_feat.numpy(), "\n其形状是: ",
512       usr_age_feat.shape)
513 print("\n年龄 1 对应的特征是: ", usr_age_feat.numpy()[0, :])
514 print("年龄 18 对应的特征是: ", usr_age_feat.numpy()[1, :])
515
516 # ## 4. 提取用户职业特征
517 #
518 #
519 # 参考用户年龄的处理方式实现用户职业的特征提取，同样采用Embedding层和全连接层的方式提取
520 # 特征。由上一节信息可以得知用户职业的最大数字表示是20。
521
522 # In[8]:
523
524
525 # 自定义一个用户职业数据
526 usr_job_data = np.array((0, 20)).reshape(-1).astype('int64')
527 print("输入的用户职业是:", usr_job_data)
528
529 # 对用户职业信息做映射，并紧接着一个Linear层
530 # 用户职业的最大ID是20，所以Embedding层size的第一个参数设置为20 + 1 = 21
531 USR_JOB_DICT_SIZE = 20 + 1
532
533 usr_job_emb = Embedding(num_embeddings=USR_JOB_DICT_SIZE, embedding_dim=16)
534 usr_job_fc = Linear(in_features=16, out_features=16)
535
536 usr_job = paddle.to_tensor(usr_job_data)
537 usr_job_feat = usr_job_emb(usr_job)
538 usr_job_feat = usr_job_fc(usr_job_feat)
539 usr_job_feat = F.relu(usr_job_feat)
540
541 print("用户年龄特征的数据特征是: ", usr_job_feat.numpy(), "\n其形状是: ",
542       usr_job_feat.shape)
543 print("\n职业 0 对应的特征是: ", usr_job_feat.numpy()[0, :])
544 print("职业 20 对应的特征是: ", usr_job_feat.numpy()[1, :])

```

```

541
542
543 # ## 5. 融合用户特征
544 #
545 #
特征融合是一种常用的特征增强手段，通过结合不同特征的长处，达到取长补短的目的。简单的融合方法有：特征（加权）相加、特征级联、特征正交等等。此处使用特征融合是为了将用户的多个特征融合到一起，用单个向量表示每个用户，更方便计算用户与电影的相似度。上文使用Embedding加全连接的方法，分别得到了用户ID、年龄、性别、职业的特征向量，可以使用全连接层将每个特征映射到固定长度，然后进行相加，得到融合特征。
546
547 # In[9]:
548
549
550 FC_ID = Linear(in_features=32, out_features=200)
551 FC_JOB = Linear(in_features=16, out_features=200)
552 FC_AGE = Linear(in_features=16, out_features=200)
553 FC_GENDER = Linear(in_features=16, out_features=200)
554
555 # 收集所有的用户特征
556 _features = [usr_id_feat, usr_job_feat, usr_age_feat, usr_gender_feat]
557 _features = [k.numpy() for k in _features]
558 _features = [paddle.to_tensor(k) for k in _features]
559
560 id_feat = F.tanh(FC_ID(_features[0]))
561 job_feat = F.tanh(FC_JOB(_features[1]))
562 age_feat = F.tanh(FC_AGE(_features[2]))
563 genger_feat = F.tanh(FC_GENDER(_features[-1]))
564
565 # 对特征求和
566 usr_feat = id_feat + job_feat + age_feat + genger_feat
567 print("用户融合后特征的维度是：", usr_feat.shape)
568
569
570 # 这里使用全连接层进一步提取特征，而不是直接相加得到用户特征的原因有两点：
571 # * 一是用户每个特征数据维度不一致，无法直接相加；
572 # *
二是用户每个特征仅使用了一层全连接层，提取特征不充分，多使用一层全连接层能进一步提取特征。而且，这里用高维度（200维）的向量表示用户特征，能包含更多的信息，每个用户特征之间的区分也更明显。
573 #
574 #
上述实现中需要对每个特征都使用一个全连接层，实现较为复杂，一种简单的替换方式是，先将每个用户特征沿着长度维度进行级联，然后使用一个全连接层获得整个用户特征向量，两种方式的对比见下图：
575 #
576 # <center></center>
577 # <center> 图3：两种特征方式对比示意 </center>
578 # <br>
579 #
580 # 两种方式均可实现向量的合并，虽然两者的数学公式不同，但它们的表达方式是类似的。
581 #
582 #
583 # 下面是方式2的代码实现。
584
585 # In[10]:
586
587
588 usr_combined = Linear(in_features=80, out_features=200)
589
590 # 收集所有的用户特征
591 _features = [usr_id_feat, usr_job_feat, usr_age_feat, usr_gender_feat]
592
593 print("打印每个特征的维度：", [f.shape for f in _features])
594
595 _features = [k.numpy() for k in _features]
596 _features = [paddle.to_tensor(k) for k in _features]
597
598 # 对特征沿着最后一个维度级联
599 usr_feat = paddle.concat(_features, axis=1)

```

```

600     usr_feat = F.tanh(usr_combined(usr_feat))
601     print("用户融合后特征的维度是: ", usr_feat.shape)
602
603
604     # 上述代码中, 使用了[paddle.concat API](
http://www.paddlepaddle.org.cn/documentation/docs/zh/2.0-rc/api/paddle/tensor/manipulation/concat\_cn.html#concat), 表示沿着第几个维度将输入数据级联到一起。
605     #
606     # > paddle.concat *(x, axis=0, name=None) *
607     #
608     # 常用参数含义如下:
609     #
610     # * x (list|tuple): 待联结的Tensor list或者Tensor tuple
611     # , x中所有Tensor的数据类型应该一致。
612     # * axis (int|Tensor, 可选): 指定对输入x进行运算的轴, 默认值为0。
613     #
614     #
615
616     至此我们已经完成了用户特征提取网络的设计, 包括ID特征提取、性别特征提取、年龄特征提取、
617     职业特征提取和特征融合模块, 下面我们将所有的模块整合到一起, 放到Python类中, 完整代码实
618     现如下:
619
620     # In[11]:
621
622     import random
623     import math
624     class Model(nn.Layer):
625         def __init__(self, use_poster, use_mov_title, use_mov_cat, use_age_job,fc_sizes):
626             super(Model, self).__init__()
627
628             # 将传入的name信息和bool型参数添加到模型类中
629             self.use_mov_poster = use_poster
630             self.use_mov_title = use_mov_title
631             self.use_usr_age_job = use_age_job
632             self.use_mov_cat = use_mov_cat
633             self.fc_sizes = fc_sizes
634
635             # 使用上节定义的数据处理类, 获取数据集的信息, 并构建训练和验证集的数据迭代器
636             Dataset = MovieLen(self.use_mov_poster)
637             self.Dataset = Dataset
638             self.trainset = self.Dataset.train_dataset
639             self.valset = self.Dataset.valid_dataset
640             self.train_loader = self.Dataset.load_data(dataset=self.trainset, mode=
641             'train')
642             self.valid_loader = self.Dataset.load_data(dataset=self.valset, mode='valid')
643
644             """ define network layer for embedding usr info """
645             USR_ID_NUM = Dataset.max_usr_id + 1
646             # 对用户ID做映射, 并紧接着一个FC层
647             self.usr_emb = Embedding(num_embeddings=USR_ID_NUM,embedding_dim=32)
648             self.usr_fc = Linear(32, 32)
649
650             # 对用户性别信息做映射, 并紧接着一个FC层
651             USR_GENDER_DICT_SIZE = 2
652             self.usr_gender_emb = Embedding(num_embeddings=USR_GENDER_DICT_SIZE,
653             embedding_dim=16)
654             self.usr_gender_fc = Linear(16, 16)
655
656             # 对用户年龄信息做映射, 并紧接着一个FC层
657             USR_AGE_DICT_SIZE = Dataset.max_usr_age + 1
658             self.usr_age_emb = Embedding(num_embeddings=USR_AGE_DICT_SIZE,embedding_dim=
659             16)
660             self.usr_age_fc = Linear(16, 16)
661
662             # 对用户职业信息做映射, 并紧接着一个FC层
663             USR_JOB_DICT_SIZE = Dataset.max_usr_job + 1
664             self.usr_job_emb = Embedding(num_embeddings=USR_JOB_DICT_SIZE,embedding_dim=
665             16)
666             self.usr_job_fc = Linear(16, 16)
667
668             # 新建一个FC层, 用于整合用户数据信息

```

```

662 self.usr_combined = Linear(80, 200)
663
664 # 新建一个Linear层，用于整合电影特征
665 self.mov_concat_embed = Linear(in_features=96, out_features=200)
666
667 user_sizes = [200] + self.fc_sizes
668 acts = ["relu" for _ in range(len(self.fc_sizes))]
669 self._user_layers = []
670 for i in range(len(self.fc_sizes)):
671     linear = Linear(
672         in_features=user_sizes[i],
673         out_features=user_sizes[i + 1],
674         weight_attr=paddle.ParamAttr(
675             initializer=nn.initializer.Normal(
676                 std=1.0 / math.sqrt(user_sizes[i])))
677     # 向模型中添加了一个 paddle.nn.Linear 子层
678     self.add_sublayer('linear_user_%d' % i, linear)
679     self._user_layers.append(linear)
680     if acts[i] == 'relu':
681         act = nn.ReLU()
682         # 向模型中添加了一个 paddle.nn.ReLU() 子层
683         self.add_sublayer('user_act_%d' % i, act)
684         self._user_layers.append(act)
685
686
687 # 定义计算用户特征的前向运算过程
688 def get_usr_feat(self, usr_var):
689     """ get usr features"""
690     # 获取到用户数据
691     usr_id, usr_gender, usr_age, usr_job = usr_var
692     # 将用户的ID数据经过embedding和FC计算，得到的特征保存在feats_collect中
693     feats_collect = []
694     usr_id = self.usr_emb(usr_id)
695     usr_id = self.usr_fc(usr_id)
696     usr_id = F.relu(usr_id)
697     feats_collect.append(usr_id)
698
699     # 计算用户的性别特征，并保存在feats_collect中
700     usr_gender = self.usr_gender_emb(usr_gender)
701     usr_gender = self.usr_gender_fc(usr_gender)
702     usr_gender = F.relu(usr_gender)
703
704     feats_collect.append(usr_gender)
705     # 选择是否使用用户的年龄-职业特征
706     if self.use_usr_age_job:
707         # 计算用户的年龄特征，并保存在feats_collect中
708         usr_age = self.usr_age_emb(usr_age)
709         usr_age = self.usr_age_fc(usr_age)
710         usr_age = F.relu(usr_age)
711         feats_collect.append(usr_age)
712         # 计算用户的职业特征，并保存在feats_collect中
713         usr_job = self.usr_job_emb(usr_job)
714         usr_job = self.usr_job_fc(usr_job)
715         usr_job = F.relu(usr_job)
716         feats_collect.append(usr_job)
717
718     # 将用户的特征级联，并通过FC层得到最终的用户特征
719     print([f.shape for f in feats_collect])
720     usr_feat = paddle.concat(feats_collect, axis=1)
721     user_features = F.tanh(self.usr_combined(usr_feat))
722     #通过3层全链接层，获得用于计算相似度的用户特征和电影特征
723     for n_layer in self._user_layers:
724         user_features = n_layer(user_features)
725     return user_features
726
727 #下面使用定义好的数据读取器，实现从用户数据读取到用户特征计算的流程：
728 ## 测试用户特征提取网络
729 fc_sizes=[128, 64, 32]
730 model = Model(use_poster=False, use_mov_title=True, use_mov_cat=True, use_age_job=
True,fc_sizes=fc_sizes)
731 model.eval()
732

```

```

733 data_loader = model.train_loader
734
735 for idx, data in enumerate(data_loader()):
736     # 获得数据，并转为动态图格式，
737     usr, mov, score = data
738     #     print(usr.shape)
739     # 只使用每个Batch的第一条数据
740     usr_v = [[var[0]] for var in usr]
741
742
743     print("输入的用户ID数据: {} \n性别数据: {} \n年龄数据: {} \n职业数据 {}".format(*
744         usr_v))
745
746     usr_v = [paddle.to_tensor(np.array(var)) for var in usr_v]
747     usr_feat = model.get_usr_feat(usr_v)
748     print("计算得到的用户特征维度是: ", usr_feat.shape)
749     break
750
751
752 #
753 上面使用了向量级联+全连接的方式实现了四个用户特征向量的合并，为了捕获特征向量的深层次
754 语义信息，合并后的向量还加入了3层全链接结构。在下面处理电影特征的部分我们会看到使用另
755 外一种向量合并的方式（向量相加）处理电影类型的特征(6个向量合并成1个向量)，然后再加上全
756 连接。
757
758 #
759 # # 电影特征提取网络
760 #
761 #
762 #
763 接下来我们构建提取电影特征的神经网络，与用户特征网络结构不同的是，电影的名称和类别均有
764 多个数字信息，我们构建网络时，对这两类特征的处理方式也不同。
765
766 #
767 # <center></center>
770 #
771 #
772 # 电影特征网络主要包括：
773 # 1. 将电影ID数据映射为向量表示，通过全连接层得到ID特征。
774 # 2. 将电影类别数据映射为向量表示，对电影类别的向量求和得到类别特征。
775 # 3. 将电影名称数据映射为向量表示，通过卷积层计算得到名称特征。
776 #
777
778 # ## 1. 提取电影ID特征
779 #
780 与计算用户ID特征的方式类似，我们通过如下方式实现电影ID特性提取。根据上一节信息得知电影
781 ID的最大值是3952。
782
783 #
784 # In[12]:
785
786 # 自定义一个电影ID数据
787 mov_id_data = np.array((1, 2)).reshape(-1).astype('int64')
788 # 对电影ID信息做映射，并紧接着一个FC层
789 MOV_DICT_SIZE = 3952 + 1
790 mov_emb = Embedding(num_embeddings=MOV_DICT_SIZE, embedding_dim=32)
791 mov_fc = Linear(32, 32)
792
793
794 print("输入的电影ID是:", mov_id_data)
795 mov_id_data = paddle.to_tensor(mov_id_data)
796 mov_id_feat = mov_fc(mov_emb(mov_id_data))
797 mov_id_feat = F.relu(mov_id_feat)
798 print("计算的电影ID的特征是", mov_id_feat.numpy(), "\n其形状是: ", mov_id_feat.shape)
799 print("\n电影ID为 {} 计算得到的特征是: {}".format(mov_id_data.numpy()[0], mov_id_feat.
800     .numpy()[0]))
801 print("电影ID为 {} 计算得到的特征是: {}".format(mov_id_data.numpy()[1], mov_id_feat.
802     numpy()[1]))
803
804

```



```

792
793 # ## 2. 提取电影类别特征
794 #
795 #
与电影ID数据不同的是，每个电影有多个类别，提取类别特征时，如果对每个类别数据都使用一个
全连接层，电影最多的类别数是6，会导致类别特征提取网络参数过多而不利于学习。我们对于电
影类别特征提取的处理方式是：
796 # 1. 通过Embedding网络层将电影类别数字映射为特征向量；
797 # 2. 对Embedding后的向量沿着类别数量维度进行求和，得到一个类别映射向量；
798 # 3. 通过一个全连接层计算类别特征向量。
799 #
800 #
数据处理章节已经介绍到，每个电影的类别数量是不固定的，且一个电影最大的类别数量是6，类
别数量不足6的通过补0到6维。因此，每个类别的数据维度是6，每个电影类别有6个Embedding向量
。我们希望用一个向量就可以表示电影类别，可以对电影类别数量维度降维，
801 # 这里对6个Embedding向量通过求和的方式降维，得到电影类别的向量表示。
802 #
803 # 下面是电影类别特征提取的实现方法：
804
805 # In[13]:
806
807
808 # 自定义一个电影类别数据
809 mov_cat_data = np.array(((1, 2, 3, 0, 0, 0), (2, 3, 4, 0, 0, 0))).reshape(2, -1).
astype('int64')
810 # 对电影ID信息做映射，并紧接着一个Linear层
811 MOV_DICT_SIZE = 6 + 1
812 mov_emb = Embedding(num_embeddings=MOV_DICT_SIZE, embedding_dim=32)
813 mov_fc = Linear(in_features=32, out_features=32)
814
815 print("输入的电影类别是:", mov_cat_data[:, :])
816 mov_cat_data = paddle.to_tensor(mov_cat_data)
817 # 1. 通过Embedding映射电影类别数据；
818 mov_cat_feat = mov_emb(mov_cat_data)
819 # 2. 对Embedding后的向量沿着类别数量维度进行求和，得到一个类别映射向量；
820 mov_cat_feat = paddle.sum(mov_cat_feat, axis=1, keepdim=False)
821
822 # 3. 通过一个全连接层计算类别特征向量。
823 mov_cat_feat = mov_fc(mov_cat_feat)
824 mov_cat_feat = F.relu(mov_cat_feat)
825 print("计算的 movie 类别的特征是", mov_cat_feat.numpy(), "\n其形状是:", mov_cat_feat.
shape)
826 print("\n电影类别为 {} 计算得到的特征是: {}".format(mov_cat_data.numpy()[0, :],
mov_cat_feat.numpy()[0]))
827 print("\n电影类别为 {} 计算得到的特征是: {}".format(mov_cat_data.numpy()[1, :],
mov_cat_feat.numpy()[1]))
828
829
830 #
待合并的6个向量具有相同的维度，直接按位相加即可得到综合的向量表示。当然，我们也可以采
用向量级联的方式，将6个32维的向量级联成192维的向量，再通过全连接层压缩成32维度，代码实
现上要臃肿一些。
831
832 # ## 3. 提取电影名称特征
833 #
834 # 与电影类别数据一样，每个电影名称具有多个单词。我们对于电影名称特征提取的处理方式是：
835 #
836 # 1. 通过Embedding映射电影名称数据，得到对应的特征向量；
837 # 2. 对Embedding后的向量使用卷积层+全连接层进一步提取特征；
838 # 3. 对特征进行降采样，降低数据维度。
839 #
840 #
提取电影名称特征时，使用了卷积层加全连接层的方式提取特征。这是因为电影名称单词较多，最
大单词数量是15，如果采用和电影类别同样的处理方式，即沿着数量维度求和，显然会损失很多信
息。考虑到15这个维度较高，可以使用卷积层进一步提取特征，同时通过控制卷积层的步长，降低
电影名称特征的维度。
841 #
842 #
如果只是简单的经过一层或二层卷积后，特征的维度依然很大，为了得到更低维度的特征向量，有
两种方式，一种是利用求和降采样的方式，另一种是继续使用神经网络层进行特征提取并逐渐降低
特征维度。这里，我们采用“简单求和”的降采样方式压缩电影名称特征的维度，通过飞桨的[reduc
e_sum] (

```

https://www.paddlepaddle.org.cn/documentation/docs/zh/2.0-beta/api/paddle/fluid/layers/reduce_sum_cn.html) API实现。

```
843 #
844 # 下面是提取电影名称特征的代码实现：
845 #
846
847 # In[14]:
848
849
850 # 自定义两个电影名称数据
851 mov_title_data = np.array(((1, 2, 3, 4, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0),
852                             (2, 3, 4, 5, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0))).reshape(2,
853                                             , 1, 15).astype('int64')
854
855 # 对电影名称做映射，紧接着FC和pool层
856 MOV_TITLE_DICT_SIZE = 1000 + 1
857 mov_title_emb = Embedding(num_embeddings=MOV_TITLE_DICT_SIZE, embedding_dim=32)
858 mov_title_conv = Conv2D(in_channels=1, out_channels=1, kernel_size=(3, 1), stride=(2,
859                                     1), padding=0)
860 # 使用 3 * 3卷积层代替全连接层
861 mov_title_conv2 = Conv2D(in_channels=1, out_channels=1, kernel_size=(3, 1), stride=1,
862                             padding=0)
863
864 mov_title_data = paddle.to_tensor(mov_title_data)
865 print("电影名称数据的输入形状：", mov_title_data.shape)
866 # 1. 通过Embedding映射电影名称数据；
867 mov_title_feat = mov_title_emb(mov_title_data)
868 print("输入通过Embedding层的输出形状：", mov_title_feat.shape)
869 # 2. 对Embedding后的向量使用卷积层进一步提取特征；
870 mov_title_feat = F.relu(mov_title_conv(mov_title_feat))
871 print("第一次卷积之后的特征输出形状：", mov_title_feat.shape)
872 mov_title_feat = F.relu(mov_title_conv2(mov_title_feat))
873 print("第二次卷积之后的特征输出形状：", mov_title_feat.shape)
874
875 batch_size = mov_title_data.shape[0]
876 # 3.
877 最后对特征进行降采样，keepdim=False会让输出的维度减少，而不是用[2,1,1,32]的形式占位；
878 mov_title_feat = paddle.sum(mov_title_feat, axis=2, keepdim=False)
879 print("reduce_sum降采样后的特征输出形状：", mov_title_feat.shape)
880
881 mov_title_feat = F.relu(mov_title_feat)
882 mov_title_feat = paddle.reshape(mov_title_feat, [batch_size, -1])
883 print("电影名称特征的特征输出形状：", mov_title_feat.shape)
884
885 print("\n计算的 movie 名称的特征是", mov_title_feat.numpy(), "\n其形状是：",
886       mov_title_feat.shape)
887 print("\n电影名称为 {} 计算得到的特征是：{}".format(mov_title_data.numpy()[0, :, 0],
888       mov_title_feat.numpy()[0]))
889 print("\n电影名称为 {} 计算得到的特征是：{}".format(mov_title_data.numpy()[1, :, 0],
890       mov_title_feat.numpy()[1]))
891
892 # 上述代码中，通过Embedding层已经获得了维度是[batch_size, 1, 15,
893 32]电影名称特征向量，因此，该特征可以视为是通道数量为1的特征图，很适合使用卷积层进一步
894 提取特征。这里我们使用两个3*3大小的卷积核的卷积层提取特征，输出通道保持不变，仍然
895 是1。特征维度中15是电影名称中单词的数量（最大数量），使用3*3的卷积核，由于卷
896 积感受野的原因，进行卷积时会综合多个单词的特征，同时设置卷积的步长参数stride为(2,
897 1)，即可对电影名称的维度降维，同时保持每个名称的向量长度不变，以防过度压缩每个名称特征
898 的信息。
899 #
900 #
901 从输出结果来看，第一个卷积层之后的输出特征维度依然较大，可以使用第二个卷积层进一步提取
902 特征。获得第二个卷积的特征后，特征的维度已经从7*32，降低到了5*32，因此可
903 以直接使用求和（向量按位相加）的方式沿着电影名称维度进行降采样（5*32 ->
904 1*32），得到最终的 movie 名称特征向量。
905 #
906 # 需要注意的是，降采样后的数据尺寸依然比下一层要求的输入向量多出一维 [2, 1,
907 32]，所以最终输出前需调整下形状。
908
909 ## 4. 融合电影特征
910 #
911 与用户特征融合方式相同，电影特征融合采用特征级联加全连接层的方式，将电影特征用一个200
912 维的向量表示。
```

```

893
894 # In[15]:
895
896
897 mov_combined = Linear(in_features=96, out_features=200)
898 # 收集所有的电影特征
899 _features = [mov_id_feat, mov_cat_feat, mov_title_feat]
900 _features = [k.numpy() for k in _features]
901 _features = [paddle.to_tensor(k) for k in _features]
902
903 # 对特征沿着最后一个维度级联
904 mov_feat = paddle.concat(_features, axis=1)
905 mov_feat = mov_combined(mov_feat)
906 mov_feat = F.tanh(mov_feat)
907 print("融合后的电影特征维度是: ", mov_feat.shape)
908
909
910 #
911 至此已经完成了电影特征提取的网络设计，包括电影ID特征提取、电影类别特征提取和电影名称特
912 征提取。
913
914 # 下面将这些模块整合到一个Python类中，完整代码如下：
915
916 # In[16]:
917
918 class MovModel(nn.Layer):
919     def __init__(self, use_poster, use_mov_title, use_mov_cat, use_age_job, fc_sizes):
920         super(MovModel, self).__init__()
921
922         # 将传入的name信息和bool型参数添加到模型类中
923         self.use_mov_poster = use_poster
924         self.use_mov_title = use_mov_title
925         self.use_usr_age_job = use_age_job
926         self.use_mov_cat = use_mov_cat
927         self.fc_sizes = fc_sizes
928
929         # 获取数据集的信息，并构建训练和验证集的数据迭代器
930         Dataset = MovieLen(self.use_mov_poster)
931         self.Dataset = Dataset
932         self.trainset = self.Dataset.train_dataset
933         self.valset = self.Dataset.valid_dataset
934         self.train_loader = self.Dataset.load_data(dataset=self.trainset, mode=
935             'train')
936         self.valid_loader = self.Dataset.load_data(dataset=self.valset, mode='valid')
937
938         """ define network layer for embedding usr info """
939         # 对电影ID信息做映射，并紧接着一个Linear层
940         MOV_DICT_SIZE = Dataset.max_mov_id + 1
941         self.mov_emb = Embedding(num_embeddings=MOV_DICT_SIZE, embedding_dim=32)
942         self.mov_fc = Linear(32, 32)
943
944         # 对电影类别做映射
945         CATEGORY_DICT_SIZE = len(Dataset.movie_cat) + 1
946         self.mov_cat_emb = Embedding(num_embeddings=CATEGORY_DICT_SIZE, embedding_dim
947             =32)
948         self.mov_cat_fc = Linear(32, 32)
949
950         # 对电影名称做映射
951         MOV_TITLE_DICT_SIZE = len(Dataset.movie_title) + 1
952         self.mov_title_emb = Embedding(num_embeddings=MOV_TITLE_DICT_SIZE,
953             embedding_dim=32)
954         self.mov_title_conv = Conv2D(in_channels=1, out_channels=1, kernel_size=(3, 1
955             ), stride=(2, 1), padding=0)
956         self.mov_title_conv2 = Conv2D(in_channels=1, out_channels=1, kernel_size=(3,
957             1), stride=1, padding=0)
958
959         # 新建一个Linear层，用于整合电影特征
960         self.mov_concat_embed = Linear(in_features=96, out_features=200)
961
962         # 电影特征和用户特征使用了不同的全连接层，不共享参数
963         movie_sizes = [200] + self.fc_sizes

```

```

958     acts = ["relu" for _ in range(len(self.fc_sizes))]
959     self._movie_layers = []
960     for i in range(len(self.fc_sizes)):
961         linear = Linear(
962             in_features=movie_sizes[i],
963             out_features=movie_sizes[i + 1],
964             weight_attr=paddle.ParamAttr(
965                 initializer=nn.initializer.Normal(
966                     std=1.0 / math.sqrt(movie_sizes[i])))
967         self.add_sublayer('linear_movie_%d' % i, linear)
968         self._movie_layers.append(linear)
969         if acts[i] == 'relu':
970             act = nn.ReLU()
971             self.add_sublayer('movie_act_%d' % i, act)
972             self._movie_layers.append(act)
973
974     # 定义电影特征的前向计算过程
975     def get_mov_feat(self, mov_var):
976         """ get movie features"""
977         # 获得电影数据
978         mov_id, mov_cat, mov_title, mov_poster = mov_var
979         feats_collect = []
980         # 获得batchsize的大小
981         batch_size = mov_id.shape[0]
982         # 计算电影ID的特征, 并存在feats_collect中
983         mov_id = self.mov_emb(mov_id)
984         mov_id = self.mov_fc(mov_id)
985         mov_id = F.relu(mov_id)
986         feats_collect.append(mov_id)
987
988         # 如果使用电影的种类数据, 计算电影种类特征的映射
989         if self.use_mov_cat:
990             # 计算电影种类的特征映射, 对多个种类的特征求和得到最终特征
991             mov_cat = self.mov_cat_emb(mov_cat)
992             print(mov_cat.shape)
993             mov_cat = paddle.sum(mov_cat, axis=1, keepdim=False)
994
995             mov_cat = self.mov_cat_fc(mov_cat)
996             feats_collect.append(mov_cat)
997
998         if self.use_mov_title:
999             # 计算电影名字的特征映射, 对特征映射使用卷积计算最终的特征
1000             mov_title = self.mov_title_emb(mov_title)
1001             mov_title = F.relu(self.mov_title_conv2(F.relu(self.mov_title_conv(
1002                 mov_title))))
1003
1004             mov_title = paddle.sum(mov_title, axis=2, keepdim=False)
1005             mov_title = F.relu(mov_title)
1006             mov_title = paddle.reshape(mov_title, [batch_size, -1])
1007             feats_collect.append(mov_title)
1008
1009         # 使用一个全连接层, 整合所有电影特征, 映射为一个200维的特征向量
1010         mov_feat = paddle.concat(feats_collect, axis=1)
1011         mov_features = F.tanh(self.mov_concat_embed(mov_feat))
1012         for n_layer in self._movie_layers:
1013             mov_features = n_layer(mov_features)
1014         return mov_features
1015
1016     # 由上述电影特征处理的代码可以观察到:
1017     # * 电影ID特征的计算方式和用户ID的计算方式相同。
1018     # *
1019     # 对于包含多个元素的电影类别数据, 采用将所有元素的映射向量求和的结果, 然后加上全连接结构
1020     # 作为最终的电影类别特征表示。考虑到电影类别的数量有限, 这里采用简单的求和特征融合方式。
1021     # *
1022     # 对于电影的名称数据, 其包含的元素数量多于电影种类元素数量, 则采用卷积计算的方式, 之后再
1023     # 将计算的特征沿着数据维度进行求和。读者也可自行设计这部分特征的计算网络, 并观察最终训练
1024     # 结果。
1025
1026     # 下面使用定义好的数据读取器, 实现从电影数据中提取电影特征。
1027
1028     # In[17]:

```

```

1024
1025
1026 ## 测试电影特征提取网络
1027 fc_sizes=[128, 64, 32]
1028 model = MovModel(use_poster=False, use_mov_title=True, use_mov_cat=True, use_age_job=
True,fc_sizes=fc_sizes)
1029 model.eval()
1030
1031 data_loader = model.train_loader
1032
1033 for idx, data in enumerate(data_loader()):
1034     # 获得数据，并转为动态图格式
1035     usr, mov, score = data
1036     # 只使用每个Batch的第一条数据
1037     mov_v = [var[0:1] for var in mov]
1038
1039     _mov_v = [np.squeeze(var[0:1]) for var in mov]
1040     print("输入的电影ID数据: {} \n类别数据: {} \n名称数据: {}".format(*_mov_v))
1041     mov_v = [paddle.to_tensor(var) for var in mov_v]
1042     mov_feat = model.get_mov_feat(mov_v)
1043     print("计算得到的电影特征维度是: ", mov_feat.shape)
1044     break
1045
1046
1047
1048 # # 相似度计算
1049 #
1050 #
1051 # 计算得到用户特征和电影特征后，我们还需要计算特征之间的相似度。如果一个用户对某个电影很
1052 # 感兴趣，并给了五分评价，那么该用户和电影特征之间的相似度是很高的。
1053 #
1054 # 衡量向量距离（相似度）有多种方案：欧式距离、曼哈顿距离、切比雪夫距离、余弦相似度等，本
1055 # 节我们使用忽略尺度信息的余弦相似度构建相似度矩阵。余弦相似度又称为余弦相似性，是通过计
1056 # 算两个向量的夹角余弦值来评估他们的相似度，如下图，两条红色的直线表示两个向量，之间的夹
1057 # 角可以用来表示相似度大小，角度为0时，余弦值为1，表示完全相似。
1058 #
1059 # 
1063 #
1064 #
1065 # 余弦相似度的公式为：
1066 #
1067 # 
$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

1068 #
1069 #
1070 # 下面是计算相似度的实现方法，输入用户特征和电影特征，计算出两者之间的相似度。另外，我们
1071 # 将用户对电影的评分作为相似度衡量的标准，由于相似度的数据范围是[0,
1072 # 1]，还需要把计算的相似度扩大到评分数据范围，评分分为1-5共5个档次，所以需要将相似度扩大
1073 # 5倍。使用飞桨[scale API](
1074 # http://www.paddlepaddle.org.cn/documentation/docs/zh/2.0-rc/api/paddle/fluid/layers/scale\_cn.html#scale), 可以对输入数据进行缩放。计算余弦相似度可以使用[cosine_similarity
1075 # API](
1076 # https://www.paddlepaddle.org.cn/documentation/docs/zh/2.0-rc/api/paddle/nn/functional/common/cosine\_similarity\_cn.html#cosine-similarity) 完成。
1077 #
1078 # In[18]:
1079 #
1080 #
1081 def similarty(usr_feature, mov_feature):
1082     res = F.cosine_similarity(usr_feature, mov_feature)
1083     res = paddle.scale(res, scale=5)
1084     return usr_feat, mov_feat, res
1085
1086 # 使用上文计算得到的用户特征和电影特征计算相似度
1087 usr_feat, mov_feat, _sim = similarty(usr_feat, mov_feat)
1088 print("相似度得分是: ", np.squeeze(_sim.numpy()))
1089
1090

```



```

1077
1078 # 从结果中我们发现相似度很小，主要有以下原因：
1079 # 1. 神经网络并没有训练，模型参数都是随机初始化的，提取出的特征没有规律性。
1080 # 2. 计算相似度的用户数据和电影数据相关性很小。
1081 #
1082 # 下一节我们就开始训练，让这个网络能够输出有效的用户特征向量和电影特征向量。
1083
1084 # ## 总结
1085 #
1086 #
本节中，我们介绍了个性化推荐的模型设计，包括用户特征网络、电影特征网络和特征相似度计算
三部分。
1087 #
1088 #
其中，用户特征网络将用户数据映射为固定长度的特征向量，电影特征网络将电影数据映射为固定
长度的特征向量，最终利用余弦相似度计算出用户特征和电影特征的相似度。相似度越大，表示用
户对该电影越喜欢。
1089 #
1090 # 以下为模型设计的完整代码：
1091
1092 # In[19]:
1093
1094
1095 class Model(nn.Layer):
1096     def __init__(self, use_poster, use_mov_title, use_mov_cat, use_age_job):
1097         super(Model, self).__init__()
1098
1099         # 将传入的name信息和bool型参数添加到模型类中
1100         self.use_mov_poster = use_poster
1101         self.use_mov_title = use_mov_title
1102         self.use_usr_age_job = use_age_job
1103         self.use_mov_cat = use_mov_cat
1104
1105         # 获取数据集的信息，并构建训练和验证集的数据迭代器
1106         Dataset = MovieLen(self.use_mov_poster)
1107         self.Dataset = Dataset
1108         self.trainset = self.Dataset.train_dataset
1109         self.valset = self.Dataset.valid_dataset
1110         self.train_loader = self.Dataset.load_data(dataset=self.trainset, mode=
'train')
1111         self.valid_loader = self.Dataset.load_data(dataset=self.valset, mode='valid')
1112
1113         """ define network layer for embedding usr info """
1114         USR_ID_NUM = Dataset.max_usr_id + 1
1115         # 对用户ID做映射，并紧接着一个Linear层
1116         self.usr_emb = Embedding(num_embeddings=USR_ID_NUM, embedding_dim=32, sparse=
False)
1117         self.usr_fc = Linear(in_features=32, out_features=32)
1118
1119         # 对用户性别信息做映射，并紧接着一个Linear层
1120         USR_GENDER_DICT_SIZE = 2
1121         self.usr_gender_emb = Embedding(num_embeddings=USR_GENDER_DICT_SIZE,
embedding_dim=16)
1122         self.usr_gender_fc = Linear(in_features=16, out_features=16)
1123
1124         # 对用户年龄信息做映射，并紧接着一个Linear层
1125         USR_AGE_DICT_SIZE = Dataset.max_usr_age + 1
1126         self.usr_age_emb = Embedding(num_embeddings=USR_AGE_DICT_SIZE, embedding_dim=
16)
1127         self.usr_age_fc = Linear(in_features=16, out_features=16)
1128
1129         # 对用户职业信息做映射，并紧接着一个Linear层
1130         USR_JOB_DICT_SIZE = Dataset.max_usr_job + 1
1131         self.usr_job_emb = Embedding(num_embeddings=USR_JOB_DICT_SIZE, embedding_dim=
16)
1132         self.usr_job_fc = Linear(in_features=16, out_features=16)
1133
1134         # 新建一个Linear层，用于整合用户数据信息
1135         self.usr_combined = Linear(in_features=80, out_features=200)
1136
1137         """ define network layer for embedding mov info """
1138         # 对电影ID信息做映射，并紧接着一个Linear层

```

```

1139 MOV_DICT_SIZE = Dataset.max_mov_id + 1
1140 self.mov_emb = Embedding(num_embeddings=MOV_DICT_SIZE, embedding_dim=32)
1141 self.mov_fc = Linear(in_features=32, out_features=32)
1142
1143 # 对电影类别做映射
1144 CATEGORY_DICT_SIZE = len(Dataset.movie_cat) + 1
1145 self.mov_cat_emb = Embedding(num_embeddings=CATEGORY_DICT_SIZE, embedding_dim
=32, sparse=False)
1146 self.mov_cat_fc = Linear(in_features=32, out_features=32)
1147
1148 # 对电影名称做映射
1149 MOV_TITLE_DICT_SIZE = len(Dataset.movie_title) + 1
1150 self.mov_title_emb = Embedding(num_embeddings=MOV_TITLE_DICT_SIZE,
embedding_dim=32, sparse=False)
1151 self.mov_title_conv = Conv2D(in_channels=1, out_channels=1, kernel_size=(3, 1
), stride=(2,1), padding=0)
1152 self.mov_title_conv2 = Conv2D(in_channels=1, out_channels=1, kernel_size=(3,
1), stride=1, padding=0)
1153
1154 # 新建一个FC层，用于整合电影特征
1155 self.mov_concat_embed = Linear(in_features=96, out_features=200)
1156
1157 user_sizes = [200] + self.fc_sizes
1158 acts = ["relu" for _ in range(len(self.fc_sizes))]
1159 self._user_layers = []
1160 for i in range(len(self.fc_sizes)):
1161     linear = Linear(
1162         in_features=user_sizes[i],
1163         out_features=user_sizes[i + 1],
1164         weight_attr=paddle.ParamAttr(
1165             initializer=nn.initializer.Normal(
1166                 std=1.0 / math.sqrt(user_sizes[i])))
1167     self.add_sublayer('linear_user %d' % i, linear)
1168     self._user_layers.append(linear)
1169     if acts[i] == 'relu':
1170         act = nn.ReLU()
1171         self.add_sublayer('user_act %d' % i, act)
1172         self._user_layers.append(act)
1173
1174 #电影特征和用户特征使用了不同的全连接层，不共享参数
1175 movie_sizes = [200] + self.fc_sizes
1176 acts = ["relu" for _ in range(len(self.fc_sizes))]
1177 self._movie_layers = []
1178 for i in range(len(self.fc_sizes)):
1179     linear = nn.Linear(
1180         in_features=movie_sizes[i],
1181         out_features=movie_sizes[i + 1],
1182         weight_attr=paddle.ParamAttr(
1183             initializer=nn.initializer.Normal(
1184                 std=1.0 / math.sqrt(movie_sizes[i])))
1185     self.add_sublayer('linear_movie %d' % i, linear)
1186     self._movie_layers.append(linear)
1187     if acts[i] == 'relu':
1188         act = nn.ReLU()
1189         self.add_sublayer('movie_act %d' % i, act)
1190         self._movie_layers.append(act)
1191
1192 # 定义计算用户特征的前向运算过程
1193 def get_usr_feat(self, usr_var):
1194     """ get usr features"""
1195     # 获取到用户数据
1196     usr_id, usr_gender, usr_age, usr_job = usr_var
1197     # 将用户的ID数据经过embedding和Linear计算，得到的特征保存在feats_collect中
1198     feats_collect = []
1199     usr_id = self.usr_emb(usr_id)
1200     usr_id = self.usr_fc(usr_id)
1201     usr_id = F.relu(usr_id)
1202     feats_collect.append(usr_id)
1203
1204     # 计算用户的性别特征，并保存在feats_collect中
1205     usr_gender = self.usr_gender_emb(usr_gender)
1206     usr_gender = self.usr_gender_fc(usr_gender)

```

```

1207     usr_gender = F.relu(usr_gender)
1208     feats_collect.append(usr_gender)
1209     # 选择是否使用用户的年龄-职业特征
1210     if self.use_usr_age_job:
1211         # 计算用户的年龄特征，并保存在feats_collect中
1212         usr_age = self.usr_age_emb(usr_age)
1213         usr_age = self.usr_age_fc(usr_age)
1214         usr_age = F.relu(usr_age)
1215         feats_collect.append(usr_age)
1216         # 计算用户的职业特征，并保存在feats_collect中
1217         usr_job = self.usr_job_emb(usr_job)
1218         usr_job = self.usr_job_fc(usr_job)
1219         usr_job = F.relu(usr_job)
1220         feats_collect.append(usr_job)
1221
1222     # 将用户的特征级联，并通过Linear层得到最终的用户特征
1223     usr_feat = paddle.concat(feats_collect, axis=1)
1224     user_features = F.tanh(self.usr_combined(usr_feat))
1225     #通过3层全链接层，获得用于计算相似度的用户特征和电影特征
1226     for n_layer in self._user_layers:
1227         user_features = n_layer(user_features)
1228
1229     return user_features
1230
1231     # 定义电影特征的前向计算过程
1232 def get_mov_feat(self, mov_var):
1233     """ get movie features"""
1234     # 获得电影数据
1235     mov_id, mov_cat, mov_title, mov_poster = mov_var
1236     feats_collect = []
1237     # 获得batchsize的大小
1238     batch_size = mov_id.shape[0]
1239     # 计算电影ID的特征，并存在feats_collect中
1240     mov_id = self.mov_emb(mov_id)
1241     mov_id = self.mov_fc(mov_id)
1242     mov_id = F.relu(mov_id)
1243     feats_collect.append(mov_id)
1244
1245     # 如果使用电影的种类数据，计算电影种类特征的映射
1246     if self.use_mov_cat:
1247         # 计算电影种类的特征映射，对多个种类的特征求和得到最终特征
1248         mov_cat = self.mov_cat_emb(mov_cat)
1249         mov_cat = paddle.sum(mov_cat, axis=1, keepdim=False)
1250
1251         mov_cat = self.mov_cat_fc(mov_cat)
1252         feats_collect.append(mov_cat)
1253
1254     if self.use_mov_title:
1255         # 计算电影名字的特征映射，对特征映射使用卷积计算最终的特征
1256         mov_title = self.mov_title_emb(mov_title)
1257         mov_title = F.relu(self.mov_title_conv2(F.relu(self.mov_title_conv(
1258             mov_title))))
1259         mov_title = paddle.sum(mov_title, axis=2, keepdim=False)
1260         mov_title = F.relu(mov_title)
1261         mov_title = paddle.reshape(mov_title, [batch_size, -1])
1262         feats_collect.append(mov_title)
1263
1264     # 使用一个全连接层，整合所有电影特征，映射为一个200维的特征向量
1265     mov_feat = paddle.concat(feats_collect, axis=1)
1266     mov_features = F.tanh(self.mov_concat_embed(mov_feat))
1267
1268     for n_layer in self._movie_layers:
1269         mov_features = n_layer(mov_features)
1270
1271     return mov_features
1272
1273     # 定义个性化推荐算法的前向计算
1274 def forward(self, usr_var, mov_var):
1275     # 计算用户特征和电影特征
1276     usr_feat = self.get_usr_feat(usr_var)
1277     mov_feat = self.get_mov_feat(mov_var)

```

```

1278         #通过3层全连接层，获得用于计算相似度的用户特征和电影特征
1279         for n_layer in self._user_layers:
1280             user_features = n_layer(user_features)
1281
1282         for n_layer in self._movie_layers:
1283             mov_features = n_layer(mov_features)
1284
1285         # 根据计算的特征计算相似度
1286         res = F.cosine_similarity(user_features, mov_features)
1287         # 将相似度扩大范围到和电影评分相同数据范围
1288         res = paddle.scale(res, scale=5)
1289         return usr_feat, mov_feat, res
1290
1291
1292 # train.py
1293 #!/usr/bin/env python
1294 # coding: utf-8
1295
1296 # 启动训练前，复用前面章节的数据处理和神经网络模型代码，已阅读可直接跳过。
1297 #
1298
1299 # In[ ]:
1300
1301
1302 import random
1303 import numpy as np
1304 from PIL import Image
1305
1306 import paddle
1307 from paddle.nn import Linear, Embedding, Conv2D
1308 import paddle.nn.functional as F
1309 import math
1310
1311 class MovieLen(object):
1312     def __init__(self, use_poster):
1313         self.use_poster = use_poster
1314         # 声明每个数据文件的路径
1315         usr_info_path = "./work/ml-1m/users.dat"
1316         if use_poster:
1317             rating_path = "./work/ml-1m/new_rating.txt"
1318         else:
1319             rating_path = "./work/ml-1m/ratings.dat"
1320
1321         movie_info_path = "./work/ml-1m/movies.dat"
1322         self.poster_path = "./work/ml-1m/posters/"
1323         # 得到电影数据
1324         self.movie_info, self.movie_cat, self.movie_title = self.get_movie_info(
            movie_info_path)
1325         # 记录电影的最大ID
1326         self.max_mov_cat = np.max([self.movie_cat[k] for k in self.movie_cat])
1327         self.max_mov_tit = np.max([self.movie_title[k] for k in self.movie_title])
1328         self.max_mov_id = np.max(list(map(int, self.movie_info.keys())))
1329         # 记录用户数据的最大ID
1330         self.max_usr_id = 0
1331         self.max_usr_age = 0
1332         self.max_usr_job = 0
1333         # 得到用户数据
1334         self.usr_info = self.get_usr_info(usr_info_path)
1335         # 得到评分数据
1336         self.rating_info = self.get_rating_info(rating_path)
1337         # 构建数据集
1338         self.dataset = self.get_dataset(usr_info=self.usr_info,
1339                                         rating_info=self.rating_info,
1340                                         movie_info=self.movie_info)
1341         # 划分数据集，获得数据加载器
1342         self.train_dataset = self.dataset[:int(len(self.dataset)*0.9)]
1343         self.valid_dataset = self.dataset[int(len(self.dataset)*0.9):]
1344         print("##Total dataset instances: ", len(self.dataset))
1345         print("##MovieLens dataset information: \nusr num: {}\n"
1346               "movies num: {}".format(len(self.usr_info), len(self.movie_info)))
1347         # 得到电影数据
1348         def get_movie_info(self, path):

```

```

1349 # 打开文件，编码方式选择ISO-8859-1，读取所有数据到data中
1350 with open(path, 'r', encoding="ISO-8859-1") as f:
1351     data = f.readlines()
1352 # 建立三个字典，分别用户存放电影所有信息，电影的名字信息、类别信息
1353 movie_info, movie_titles, movie_cat = {}, {}, {}
1354 # 对电影名字、类别中不同的单词计数
1355 t_count, c_count = 1, 1
1356
1357 count_tit = {}
1358 # 按行读取数据并处理
1359 for item in data:
1360     item = item.strip().split("::")
1361     v_id = item[0]
1362     v_title = item[1][:7]
1363     cats = item[2].split('|')
1364     v_year = item[1][-5:-1]
1365
1366     titles = v_title.split()
1367     # 统计电影名字的单词，并给每个单词一个序号，放在movie_titles中
1368     for t in titles:
1369         if t not in movie_titles:
1370             movie_titles[t] = t_count
1371             t_count += 1
1372     # 统计电影类别单词，并给每个单词一个序号，放在movie_cat中
1373     for cat in cats:
1374         if cat not in movie_cat:
1375             movie_cat[cat] = c_count
1376             c_count += 1
1377     # 补0使电影名称对应的列表长度为15
1378     v_tit = [movie_titles[k] for k in titles]
1379     while len(v_tit) < 15:
1380         v_tit.append(0)
1381     # 补0使电影种类对应的列表长度为6
1382     v_cat = [movie_cat[k] for k in cats]
1383     while len(v_cat) < 6:
1384         v_cat.append(0)
1385     # 保存电影数据到movie_info中
1386     movie_info[v_id] = {'mov_id': int(v_id),
1387                        'title': v_tit,
1388                        'category': v_cat,
1389                        'years': int(v_year)}
1390 return movie_info, movie_cat, movie_titles
1391
1392 def get_usr_info(self, path):
1393     # 性别转换函数，M-0, F-1
1394     def gender2num(gender):
1395         return 1 if gender == 'F' else 0
1396
1397     # 打开文件，读取所有行到data中
1398     with open(path, 'r') as f:
1399         data = f.readlines()
1400     # 建立用户信息的字典
1401     use_info = {}
1402
1403     max_usr_id = 0
1404     # 按行索引数据
1405     for item in data:
1406         # 去除每一行中和数据无关的部分
1407         item = item.strip().split("::")
1408         usr_id = item[0]
1409         # 将字符数据转成数字并保存在字典中
1410         use_info[usr_id] = {'usr_id': int(usr_id),
1411                            'gender': gender2num(item[1]),
1412                            'age': int(item[2]),
1413                            'job': int(item[3])}
1414         self.max_usr_id = max(self.max_usr_id, int(usr_id))
1415         self.max_usr_age = max(self.max_usr_age, int(item[2]))
1416         self.max_usr_job = max(self.max_usr_job, int(item[3]))
1417     return use_info
1418 # 得到评分数据
1419 def get_rating_info(self, path):
1420     # 读取文件里的数据

```



```

1421     with open(path, 'r') as f:
1422         data = f.readlines()
1423     # 将数据保存在字典中并返回
1424     rating_info = {}
1425     for item in data:
1426         item = item.strip().split("::")
1427         usr_id, movie_id, score = item[0], item[1], item[2]
1428         if usr_id not in rating_info.keys():
1429             rating_info[usr_id] = {movie_id: float(score)}
1430         else:
1431             rating_info[usr_id][movie_id] = float(score)
1432     return rating_info
1433 # 构建数据集
1434 def get_dataset(self, usr_info, rating_info, movie_info):
1435     trainset = []
1436     for usr_id in rating_info.keys():
1437         usr_ratings = rating_info[usr_id]
1438         for movie_id in usr_ratings:
1439             trainset.append({'usr_info': usr_info[usr_id],
1440                             'mov_info': movie_info[movie_id],
1441                             'scores': usr_ratings[movie_id]})
1442     return trainset
1443
1444 def load_data(self, dataset=None, mode='train'):
1445     use_poster = False
1446
1447     # 定义数据迭代Batch大小
1448     BATCHSIZE = 256
1449
1450     data_length = len(dataset)
1451     index_list = list(range(data_length))
1452     # 定义数据迭代加载器
1453     def data_generator():
1454         # 训练模式下，打乱训练数据
1455         if mode == 'train':
1456             random.shuffle(index_list)
1457         # 声明每个特征的列表
1458         usr_id_list, usr_gender_list, usr_age_list, usr_job_list = [], [], [], []
1459         mov_id_list, mov_tit_list, mov_cat_list, mov_poster_list = [], [], [], []
1460         score_list = []
1461         # 索引遍历输入数据集
1462         for idx, i in enumerate(index_list):
1463             # 获得特征数据保存到对应特征列表中
1464             usr_id_list.append(dataset[i]['usr_info']['usr_id'])
1465             usr_gender_list.append(dataset[i]['usr_info']['gender'])
1466             usr_age_list.append(dataset[i]['usr_info']['age'])
1467             usr_job_list.append(dataset[i]['usr_info']['job'])
1468
1469             mov_id_list.append(dataset[i]['mov_info']['mov_id'])
1470             mov_tit_list.append(dataset[i]['mov_info']['title'])
1471             mov_cat_list.append(dataset[i]['mov_info']['category'])
1472             mov_id = dataset[i]['mov_info']['mov_id']
1473
1474         if use_poster:
1475             # 不使用图像特征时，不读取图像数据，加快数据读取速度
1476             poster = Image.open(self.poster_path+'mov_id{}.jpg'.format(str(
1477                 mov_id[0])))
1478             poster = poster.resize([64, 64])
1479             if len(poster.size) <= 2:
1480                 poster = poster.convert("RGB")
1481
1482             mov_poster_list.append(np.array(poster))
1483
1484         score_list.append(int(dataset[i]['scores']))
1485         # 如果读取的数据量达到当前的batch大小，就返回当前批次
1486         if len(usr_id_list) == BATCHSIZE:
1487             # 转换列表数据为数组形式，reshape到固定形状
1488             usr_id_arr = np.array(usr_id_list)
1489             usr_gender_arr = np.array(usr_gender_list)
1490             usr_age_arr = np.array(usr_age_list)
1491             usr_job_arr = np.array(usr_job_list)

```

```

1492         mov_id_arr = np.array(mov_id_list)
1493         mov_cat_arr = np.reshape(np.array(mov_cat_list), [BATCHSIZE, 6]).
            astype(np.int64)
1494         mov_tit_arr = np.reshape(np.array(mov_tit_list), [BATCHSIZE, 1,
            15]).astype(np.int64)
1495
1496         if use_poster:
1497             mov_poster_arr = np.reshape(np.array(mov_poster_list)/127.5 -
            1, [BATCHSIZE, 3, 64, 64]).astype(np.float32)
1498         else:
1499             mov_poster_arr = np.array([0.])
1500
1501         scores_arr = np.reshape(np.array(score_list), [-1, 1]).astype(np.
            float32)
1502
1503         # 放回当前批次数据
1504         yield [usr_id_arr, usr_gender_arr, usr_age_arr, usr_job_arr],
            [mov_id_arr, mov_cat_arr, mov_tit_arr
            , mov_poster_arr], scores_arr
1505
1506         # 清空数据
1507         usr_id_list, usr_gender_list, usr_age_list, usr_job_list = [],
            [], [], []
1508         mov_id_list, mov_tit_list, mov_cat_list, score_list = [], [], [],
            []
1509         mov_poster_list = []
1510     return data_generator
1511
1512 # In[ ]:
1513
1514
1515
1516 class Model(paddle.nn.Layer):
1517     def __init__(self, use_poster, use_mov_title, use_mov_cat, use_age_job,fc_sizes):
1518         super(Model, self).__init__()
1519
1520         # 将传入的name信息和bool型参数添加到模型类中
1521         self.use_mov_poster = use_poster
1522         self.use_mov_title = use_mov_title
1523         self.use_usr_age_job = use_age_job
1524         self.use_mov_cat = use_mov_cat
1525         self.fc_sizes=fc_sizes
1526
1527         # 获取数据集的信息，并构建训练和验证集的数据迭代器
1528         Dataset = MovieLen(self.use_mov_poster)
1529         self.Dataset = Dataset
1530         self.trainset = self.Dataset.train_dataset
1531         self.valset = self.Dataset.valid_dataset
1532         self.train_loader = self.Dataset.load_data(dataset=self.trainset, mode=
            'train')
1533         self.valid_loader = self.Dataset.load_data(dataset=self.valset, mode='valid')
1534
1535         usr_embedding_dim=32
1536         gender_embedding_dim=16
1537         age_embedding_dim=16
1538
1539         job_embedding_dim=16
1540         mov_embedding_dim=16
1541         category_embedding_dim=16
1542         title_embedding_dim=32
1543
1544         """ define network layer for embedding usr info """
1545         USR_ID_NUM = Dataset.max_usr_id + 1
1546
1547         # 对用户ID做映射，并紧接着一个Linear层
1548         self.usr_emb = Embedding(num_embeddings=USR_ID_NUM, embedding_dim=
            usr_embedding_dim, sparse=False)
1549         self.usr_fc = Linear(in_features=usr_embedding_dim, out_features=32)
1550
1551         # 对用户性别信息做映射，并紧接着一个Linear层
1552         USR_GENDER_DICT_SIZE = 2
1553         self.usr_gender_emb = Embedding(num_embeddings=USR_GENDER_DICT_SIZE,

```

```

embedding_dim=gender_embedding_dim)
1554 self.usr_gender_fc = Linear(in_features=gender_embedding_dim, out_features=16)
1555
1556 # 对用户年龄信息做映射，并紧接着一个Linear层
1557 USR_AGE_DICT_SIZE = Dataset.max_usr_age + 1
1558 self.usr_age_emb = Embedding(num_embeddings=USR_AGE_DICT_SIZE, embedding_dim=
age_embedding_dim)
1559 self.usr_age_fc = Linear(in_features=age_embedding_dim, out_features=16)
1560
1561 # 对用户职业信息做映射，并紧接着一个Linear层
1562 USR_JOB_DICT_SIZE = Dataset.max_usr_job + 1
1563 self.usr_job_emb = Embedding(num_embeddings=USR_JOB_DICT_SIZE, embedding_dim=
job_embedding_dim)
1564 self.usr_job_fc = Linear(in_features=job_embedding_dim, out_features=16)
1565
1566 # 新建一个Linear层，用于整合用户数据信息
1567 self.usr_combined = Linear(in_features=80, out_features=200)
1568
1569 """ define network layer for embedding usr info """
1570 # 对电影ID信息做映射，并紧接着一个Linear层
1571 MOV_DICT_SIZE = Dataset.max_mov_id + 1
1572 self.mov_emb = Embedding(num_embeddings=MOV_DICT_SIZE, embedding_dim=
mov_embedding_dim)
1573 self.mov_fc = Linear(in_features=mov_embedding_dim, out_features=32)
1574
1575 # 对电影类别做映射
1576 CATEGORY_DICT_SIZE = len(Dataset.movie_cat) + 1
1577 self.mov_cat_emb = Embedding(num_embeddings=CATEGORY_DICT_SIZE, embedding_dim
=category_embedding_dim, sparse=False)
1578 self.mov_cat_fc = Linear(in_features=category_embedding_dim, out_features=32)
1579
1580 # 对电影名称做映射
1581 MOV_TITLE_DICT_SIZE = len(Dataset.movie_title) + 1
1582 self.mov_title_emb = Embedding(num_embeddings=MOV_TITLE_DICT_SIZE,
embedding_dim=title_embedding_dim, sparse=False)
1583 self.mov_title_conv = Conv2D(in_channels=1, out_channels=1, kernel_size=(3, 1
), stride=(2,1), padding=0)
1584 self.mov_title_conv2 = Conv2D(in_channels=1, out_channels=1, kernel_size=(3,
1), stride=1, padding=0)
1585
1586 # 新建一个Linear层，用于整合电影特征
1587 self.mov_concat_embed = Linear(in_features=96, out_features=200)
1588
1589 user_sizes = [200] + self.fc_sizes
1590 acts = ["relu" for _ in range(len(self.fc_sizes))]
1591 self._user_layers = []
1592 for i in range(len(self.fc_sizes)):
1593     linear = paddle.nn.Linear(
1594         in_features=user_sizes[i],
1595         out_features=user_sizes[i + 1],
1596         weight_attr=paddle.ParamAttr(
1597             initializer=paddle.nn.initializer.Normal(
1598                 std=1.0 / math.sqrt(user_sizes[i])))
1599     self.add_sublayer('linear_user_%d' % i, linear)
1600     self._user_layers.append(linear)
1601     if acts[i] == 'relu':
1602         act = paddle.nn.ReLU()
1603         self.add_sublayer('user_act_%d' % i, act)
1604         self._user_layers.append(act)
1605
1606 #电影特征和用户特征使用了不同的全连接层，不共享参数
1607 movie_sizes = [200] + self.fc_sizes
1608 acts = ["relu" for _ in range(len(self.fc_sizes))]
1609 self._movie_layers = []
1610 for i in range(len(self.fc_sizes)):
1611     linear = paddle.nn.Linear(
1612         in_features=movie_sizes[i],
1613         out_features=movie_sizes[i + 1],
1614         weight_attr=paddle.ParamAttr(
1615             initializer=paddle.nn.initializer.Normal(
1616                 std=1.0 / math.sqrt(movie_sizes[i])))
1617     self.add_sublayer('linear_movie_%d' % i, linear)

```

```

1618         self._movie_layers.append(linear)
1619         if acts[i] == 'relu':
1620             act = paddle.nn.ReLU()
1621             self.add_sublayer('movie_act_%d' % i, act)
1622             self._movie_layers.append(act)
1623
1624 # 定义计算用户特征的前向运算过程
1625 def get_usr_feat(self, usr_var):
1626     """get user features"""
1627     # 获取到用户数据
1628     usr_id, usr_gender, usr_age, usr_job = usr_var
1629     # 将用户的ID数据经过embedding和Linear计算,得到的特征保存在feats_collect中
1630     feats_collect = []
1631     usr_id = self.usr_emb(usr_id)
1632     usr_id = self.usr_fc(usr_id)
1633     usr_id = F.relu(usr_id)
1634     feats_collect.append(usr_id)
1635
1636     # 计算用户的性别特征,并保存在feats_collect中
1637     usr_gender = self.usr_gender_emb(usr_gender)
1638     usr_gender = self.usr_gender_fc(usr_gender)
1639     usr_gender = F.relu(usr_gender)
1640     feats_collect.append(usr_gender)
1641     # 选择是否使用用户的年龄-职业特征
1642     if self.use_usr_age_job:
1643         # 计算用户的年龄特征,并保存在feats_collect中
1644         usr_age = self.usr_age_emb(usr_age)
1645         usr_age = self.usr_age_fc(usr_age)
1646         usr_age = F.relu(usr_age)
1647         feats_collect.append(usr_age)
1648         # 计算用户的职业特征,并保存在feats_collect中
1649         usr_job = self.usr_job_emb(usr_job)
1650         usr_job = self.usr_job_fc(usr_job)
1651         usr_job = F.relu(usr_job)
1652         feats_collect.append(usr_job)
1653
1654     # 将用户的特征级联,并通过Linear层得到最终的用户特征
1655     usr_feat = paddle.concat(feats_collect, axis=1)
1656     user_features = F.tanh(self.usr_combined(usr_feat))
1657
1658     #通过3层全链接层,获得用于计算相似度的用户特征和电影特征
1659     for n_layer in self._user_layers:
1660         user_features = n_layer(user_features)
1661
1662     return user_features
1663
1664 # 定义电影特征的前向计算过程
1665 def get_mov_feat(self, mov_var):
1666     """get movie features"""
1667     # 获得电影数据
1668     mov_id, mov_cat, mov_title, mov_poster = mov_var
1669     feats_collect = []
1670     # 获得batchsize的大小
1671     batch_size = mov_id.shape[0]
1672     # 计算电影ID的特征,并存在feats_collect中
1673     mov_id = self.mov_emb(mov_id)
1674     mov_id = self.mov_fc(mov_id)
1675     mov_id = F.relu(mov_id)
1676     feats_collect.append(mov_id)
1677
1678     # 如果使用电影的种类数据,计算电影种类特征的映射
1679     if self.use_mov_cat:
1680         # 计算电影种类的特征映射,对多个种类的特征求和得到最终特征
1681         mov_cat = self.mov_cat_emb(mov_cat)
1682         mov_cat = paddle.sum(mov_cat, axis=1, keepdim=False)
1683
1684         mov_cat = self.mov_cat_fc(mov_cat)
1685         feats_collect.append(mov_cat)
1686
1687     if self.use_mov_title:
1688         # 计算电影名字的特征映射,对特征映射使用卷积计算最终的特征
1689         mov_title = self.mov_title_emb(mov_title)

```

```

1690         mov_title = F.relu(self.mov_title_conv2(F.relu(self.mov_title_conv(
1691             mov_title))))
1692         mov_title = paddle.sum(mov_title, axis=2, keepdim=False)
1693         mov_title = F.relu(mov_title)
1694         mov_title = paddle.reshape(mov_title, [batch_size, -1])
1695
1696         feats_collect.append(mov_title)
1697
1698         # 使用一个全连接层，整合所有电影特征，映射为一个200维的特征向量
1699         mov_feat = paddle.concat(feats_collect, axis=1)
1700         mov_features = F.tanh(self.mov_concat_embed(mov_feat))
1701
1702         for n_layer in self._movie_layers:
1703             mov_features = n_layer(mov_features)
1704
1705         return mov_features
1706
1707     # 定义个性化推荐算法的前向计算
1708     def forward(self, usr_var, mov_var):
1709         # 计算用户特征和电影特征
1710         user_features = self.get_usr_feat(usr_var)
1711         mov_features = self.get_mov_feat(mov_var)
1712
1713         # 根据计算的特征计算相似度
1714         sim = F.common.cosine_similarity(user_features, mov_features).reshape([-1, 1])
1715         # 使用余弦相似度算子，计算用户和电影的相似程度
1716         # sim = F.cosine_similarity(user_features, mov_features,
1717         #                             axis=1).reshape([-1, 1])
1718         # 将相似度扩大范围到和电影评分相同数据范围
1719         res = paddle.scale(sim, scale=5)
1720         return user_features, mov_features, res
1721
1722 # In[ ]:
1723
1724 # 解压数据集
1725 get_ipython().system('unzip -o -q -d ~/work/ ~/data/data19736/ml-1m.zip')
1726
1727
1728 # # 模型训练
1729 #
1730 #
1731 # 在模型训练前需要定义好训练的参数，包括是否使用GPU、设置损失函数、选择优化器以及学习率
1732 # 等。
1733 #
1734 # 在本次任务中，由于数据较为简单，我们选择在CPU上训练，优化器使用Adam，学习率设置为0.01
1735 # ，一共训练5个epoch。
1736 #
1737 #
1738 # 然而，针对推荐算法的网络，如何设置损失函数呢？在CV和NLP章节中的案例多是分类问题，采用
1739 # 交叉熵作为损失函数。但在电影推荐中，可以作为标签的只有评分数据，因此，我们用评分数据作
1740 # 为监督信息，神经网络的输出作为预测值，使用均方差（Mean Square
1741 # Error）损失函数去训练网络模型。
1742 #
1743 # <font
1744 # size=2>说明：使用均方差损失函数即使用回归的方法完成模型训练。电影的评分数据只有5个，是
1745 # 否可以使用分类损失函数完成训练呢？事实上，评分数据是一个连续数据，如评分3和评分4是接近
1746 # 的，如果使用分类的方法，评分3和评分4是两个类别，容易割裂评分间的连续性。
1747 #
1748 # 很多互联网产品会以用户的点击或消费数据作为训练数据，这些数据是二分类问题（点或不点，买
1749 # 或不买），可以采用交叉熵等分类任务的损失函数。
1750 # </font>
1751 #
1752 # 整个训练过程和其他的模型训练大同小异，不再赘述。
1753
1754 # In[ ]:
1755
1756 def train(model):

```



```

1747 # 配置训练参数
1748 lr = 0.001
1749 Epoches = 10
1750 paddle.set_device('cpu')
1751
1752 # 启动训练
1753 model.train()
1754 # 获得数据读取器
1755 data_loader = model.train_loader
1756 # 使用adam优化器，学习率使用0.01
1757 opt = paddle.optimizer.Adam(learning_rate=lr, parameters=model.parameters())
1758
1759 for epoch in range(0, Epoches):
1760     for idx, data in enumerate(data_loader()):
1761         # 获得数据，并转为tensor格式
1762         usr, mov, score = data
1763         usr_v = [paddle.to_tensor(var) for var in usr]
1764         mov_v = [paddle.to_tensor(var) for var in mov]
1765         scores_label = paddle.to_tensor(score)
1766         # 计算出算法的前向计算结果
1767         _, _, scores_predict = model(usr_v, mov_v)
1768         # 计算loss
1769         loss = F.square_error_cost(scores_predict, scores_label)
1770         avg_loss = paddle.mean(loss)
1771
1772         if idx % 500 == 0:
1773             print("epoch: {}, batch_id: {}, loss is: {}".format(epoch, idx,
1774                 avg_loss.numpy()))
1775
1776         # 损失函数下降，并清除梯度
1777         avg_loss.backward()
1778         opt.step()
1779         opt.clear_grad()
1780
1781         # 每个epoch 保存一次模型
1782         paddle.save(model.state_dict(), './checkpoint/epoch'+str(epoch)+'pdparams')
1783
1784 # In[ ]:
1785
1786 # 启动训练
1787 fc_sizes=[128, 64, 32]
1788 use_poster, use_mov_title, use_mov_cat, use_age_job = False, True, True, True
1789 model = Model(use_poster, use_mov_title, use_mov_cat, use_age_job, fc_sizes)
1790 train(model)
1791
1792
1793
1794 #
1795 从训练结果来看，Loss保持在1以下的范围，主要是因为使用的均方差Loss，计算得到预测评分和
1796 真实评分的均方差，真实评分的数据是1-5之间的整数，评分数据较大导致计算出来的Loss也偏大
1797 。
1798 #
1799 # 不过不用担心，我们只是通过训练神经网络提取特征向量，Loss只要收敛即可。
1800
1801 # 对训练的模型在验证集上做评估，除了训练所使用的Loss之外，还有两个选择：
1802 # 1.
1803 评分预测精度ACC (Accuracy)：将预测的float数字转成整数，计算预测评分和真实评分的匹配度。
1804 评分误差在0.5分以内的算正确，否则算错误。
1805 # 2. 评分预测误差 (Mean Absolut Error) MAE：计算预测评分和真实评分之间的平均绝对误差。
1806 # 3. 均方根误差 (Root Mean Squard Error) RMSE：计算预测评分和真实值之间的平均平方误差
1807 #
1808 # 下面是使用训练集评估这两个指标的代码实现。
1809
1810 # In[ ]:
1811
1812 from math import sqrt
1813 def evaluation(model, params_file_path):
1814     model_state_dict = paddle.load(params_file_path)
1815     model.load_dict(model_state_dict)
1816     model.eval()

```

```

1813
1814     acc_set = []
1815     avg_loss_set = []
1816     squaredError=[]
1817     for idx, data in enumerate(model.valid_loader()):
1818         usr, mov, score_label = data
1819         usr_v = [paddle.to_tensor(var) for var in usr]
1820         mov_v = [paddle.to_tensor(var) for var in mov]
1821
1822         _, _, scores_predict = model(usr_v, mov_v)
1823
1824         pred_scores = scores_predict.numpy()
1825
1826         avg_loss_set.append(np.mean(np.abs(pred_scores - score_label)))
1827         squaredError.extend(np.abs(pred_scores - score_label)**2)
1828
1829         diff = np.abs(pred_scores - score_label)
1830         diff[diff>0.5] = 1
1831         acc = 1 - np.mean(diff)
1832         acc_set.append(acc)
1833     RMSE=sqrt(np.sum(squaredError) / len(squaredError))
1834     # print("RMSE = ", sqrt(np.sum(squaredError) / len(squaredError))) #均方根误差RMSE
1835     return np.mean(acc_set), np.mean(avg_loss_set), RMSE
1836
1837
1838 # In[ ]:
1839
1840
1841 param_path = "./checkpoint/epoch"
1842 for i in range(10):
1843     acc, mae, RMSE = evaluation(model, param_path+str(i)+'_pdparams')
1844     print("ACC:", acc, "MAE:", mae, 'RMSE:', RMSE)
1845
1846
1847 #
1848 # 上述结果中，我们采用了ACC和MAE指标测试在验证集上的评分预测的准确性，其中ACC值越大越好
1849 # ，MAE值越小越好，RMSE越小也越好。
1850 #
1851 # <<font
1852 # size=2>可以看到ACC和MAE的值不是很理想，但是这仅仅是对于评分预测不准确，不能直接衡量推
1853 # 荐结果的准确性。考虑到我们设计的神经网络是为了完成推荐任务而不是评分任务，所以：
1854 # <br>1.
1855 # 只针对预测评分任务来说，我们设计的模型不够合理或者训练数据不足，导致评分预测不理想；
1856 # <br>2.
1857 # 从损失函数的收敛可以知道网络的训练是有效的，但评分预测的好坏不能完全反映推荐结果的好坏
1858 # 。</font>
1859 #
1860 #
1861 # 到这里，我们已经完成了推荐算法的前三步，包括：数据的准备、神经网络的设计和神经网络的训
1862 # 练。
1863 #
1864 # 目前还需要完成剩余的两个步骤：
1865 #
1866 # 1. 提取用户、电影数据的特征并保存到本地；
1867 #
1868 # 2. 利用保存的特征计算相似度矩阵，利用相似度完成推荐。
1869 #
1870 #
1871 # 下面，我们利用训练的神经网络提取数据的特征，进而完成电影推荐，并观察推荐结果是否令人满
1872 # 意。
1873 #
1874 # # 保存特征
1875 #
1876 #
1877 # 训练完模型后，我们得到每个用户、电影对应的特征向量，接下来将这些特征向量保存到本地，这
1878 # 样在进行推荐时，不需要使用神经网络重新提取特征，节省时间成本。
1879 #
1880 # 保存特征的流程是：
1881 # - 加载预训练好的模型参数。
1882 # -
1883 # 输入数据集的数据，提取整个数据集的用户特征和电影特征。注意数据输入到模型前，要先转成内

```

置的tensor类型并保证尺寸正确。

- 分别得到用户特征向量和电影特征向量，使用Pickle库保存字典形式的特征向量。

#

#

使用用户和电影ID为索引，以字典格式存储数据，可以通过用户或者电影的ID索引到用户特征和电影特征。

#

#

下面代码中，我们使用了一个Pickle库。Pickle库为python提供了一个简单的持久化功能，可以很容易的将Python对象保存到本地，但缺点是保存的文件可读性较差。

```
# In[ ]:
```

```
from PIL import Image
```

```
# 加载第三方库Pickle，用来保存Python数据到本地
```

```
import pickle
```

```
# 定义特征保存函数
```

```
def get_usr_mov_features(model, params_file_path, poster_path):
```

```
    paddle.set_device('cpu')
```

```
    usr_pkl = {}
```

```
    mov_pkl = {}
```

```
    # 定义将list中每个元素转成tensor的函数
```

```
    def list2tensor(inputs, shape):
```

```
        inputs = np.reshape(np.array(inputs).astype(np.int64), shape)
```

```
        return paddle.to_tensor(inputs)
```

```
    # 加载模型参数到模型中，设置为验证模式eval()
```

```
    model_state_dict = paddle.load(params_file_path)
```

```
    model.load_dict(model_state_dict)
```

```
    model.eval()
```

```
    # 获得整个数据集的数据
```

```
    dataset = model.Dataset.dataset
```

```
for i in range(len(dataset)):
```

```
    # 获得用户数据，电影数据，评分数据
```

```
    #
```

```
    本案例只转换所有在样本中出现过的user和movie，实际中可以使用业务系统中的全量数据
```

```
    usr_info, mov_info, score = dataset[i]['usr_info'], dataset[i]['mov_info'],
```

```
    dataset[i]['scores']
```

```
    usrid = str(usr_info['usr_id'])
```

```
    movid = str(mov_info['mov_id'])
```

```
    # 获得用户数据，计算得到用户特征，保存在usr_pkl字典中
```

```
    if usrid not in usr_pkl.keys():
```

```
        usr_id_v = list2tensor(usr_info['usr_id'], [1])
```

```
        usr_age_v = list2tensor(usr_info['age'], [1])
```

```
        usr_gender_v = list2tensor(usr_info['gender'], [1])
```

```
        usr_job_v = list2tensor(usr_info['job'], [1])
```

```
        usr_in = [usr_id_v, usr_gender_v, usr_age_v, usr_job_v]
```

```
        usr_feat = model.get_usr_feat(usr_in)
```

```
        usr_pkl[usrid] = usr_feat.numpy()
```

```
    # 获得电影数据，计算得到电影特征，保存在mov_pkl字典中
```

```
    if movid not in mov_pkl.keys():
```

```
        mov_id_v = list2tensor(mov_info['mov_id'], [1])
```

```
        mov_tit_v = list2tensor(mov_info['title'], [1, 1, 15])
```

```
        mov_cat_v = list2tensor(mov_info['category'], [1, 6])
```

```
        mov_in = [mov_id_v, mov_cat_v, mov_tit_v, None]
```

```
        mov_feat = model.get_mov_feat(mov_in)
```

```
        mov_pkl[movid] = mov_feat.numpy()
```

```
print(len(mov_pkl.keys()))
```

```
# 保存特征到本地
```

```

1935     pickle.dump(usr_pkl, open('./usr_feat.pkl', 'wb'))
1936     pickle.dump(mov_pkl, open('./mov_feat.pkl', 'wb'))
1937     print("usr / mov features saved!!!")
1938
1939
1940 param_path = "./checkpoint/epoch9.pdparams"
1941 poster_path = "./work/ml-lm/posters/"
1942 get_usr_mov_features(model, param_path, poster_path)
1943
1944
1945 #
保存好有效代表用户和电影的特征向量后，在下一节我们讨论如何基于这两个向量构建推荐系统。

1946
1947 # ## 作业 10-2
1948 #
1949 # 1.
以上算法使用了用户与电影的所有特征（除Poster外），可以设计对比实验，验证哪些特征是重要的，把最终的特征挑选出来。为了验证哪些特征起到关键作用，读者可以启用或弃用其中某些特征，或者加入电影海报特征，观察是否对模型Loss或评价指标有提升。

1950 # 1.
加入电影海报数据，验证电影海报特征（Poster）对推荐结果的影响，实现并分析推荐结果（有没有效果？为什么？）。

1951
1952
1953 # 综合 py
1954 #!/usr/bin/env python
1955 # coding: utf-8
1956
1957 # 训练并保存好模型，我们可以开始实践电影推荐了，推荐方式可以有多种，比如：
1958 # 1. 根据一个电影推荐其相似的电影。
1959 # 2. 根据用户的喜好，推荐其可能喜欢的电影。
1960 # 3. 给指定用户推荐与其喜好相似的用户喜欢的电影。
1961 #
1962 #
1963 # 这里我们实现第二种推荐方式，另外两种留作实践作业。
1964 #
1965 # # 根据用户喜好推荐电影
1966 #
1967 #
在前面章节，我们已经完成了神经网络的设计，并根据用户对电影的喜好（评分高低）作为训练指标完成训练。神经网络有两个输入，用户数据和电影数据，通过神经网络提取用户特征和电影特征，并计算特征之间的相似度，相似度的大小和用户对该电影的评分存在对应关系。即如果用户对这个电影感兴趣，那么对这个电影的评分也是偏高的，最终神经网络输出的相似度就更大一些。完成训练后，我们就可以开始给用户推荐电影了。

1968 #
1969 #
根据用户喜好推荐电影，是通过计算用户特征和电影特征之间的相似性，并排序选取相似度最大的结果来进行推荐，流程如下：

1970 #
1971 #
1972 # <center></center>
1973 #
1974 #
1975 # 从计算相似度到完成推荐的过程，步骤包括：
1976 #
1977 # 1. 读取保存的特征，根据一个给定的用户ID、电影ID，我们可以索引到对应的特征向量。
1978 # 2. 通过计算用户特征和其他电影特征向量的相似度，构建相似度矩阵。
1979 # 3.
对这些相似度排序后，选取相似度最大的几个特征向量，找到对应的电影ID，即得到推荐清单。

1980 # 4.
加入随机选择因素，从相似度最大的top_k结果中随机选取pick_num个推荐结果，其中pick_num必须小于top_k。

1981 #
1982 #
1983
1984 # ## 1. 读取特征向量
1985
1986 #
上一节我们已经训练好模型，并保存了电影特征，因此可以不用经过计算特征的步骤，直接读取特

```

```

1987     征。
1988     #
1989     特征以字典的形式保存，字典的键值是用户或者电影的ID，字典的元素是该用户或电影的特征向量
1990     。
1991     #
1992     # 下面实现根据指定的用户ID和电影ID，索引到对应的特征向量。
1993
1994     # In[1]:
1995
1996     get_ipython().system(' unzip -o data/data19736/ml-1m.zip -d /home/aistudio/work/')
1997     # ! unzip -o data/data20452/save_feat.zip -d /home/aistudio/
1998     get_ipython().system(' unzip -o data/data20452/save_feature_v1.zip -d
1999     /home/aistudio/')
2000
2001     # In[2]:
2002
2003     import pickle
2004     import numpy as np
2005
2006     mov_feat_dir = 'mov_feat.pkl'
2007     usr_feat_dir = 'usr_feat.pkl'
2008
2009     usr_feats = pickle.load(open(usr_feat_dir, 'rb'))
2010     mov_feats = pickle.load(open(mov_feat_dir, 'rb'))
2011
2012     usr_id = 2
2013     usr_feat = usr_feats[str(usr_id)]
2014
2015     mov_id = 1
2016     # 通过电影ID索引到电影特征
2017     mov_feat = mov_feats[str(mov_id)]
2018
2019     # 电影特征的路径
2020     movie_data_path = "./work/ml-1m/movies.dat"
2021     mov_info = {}
2022     # 打开电影数据文件，根据电影ID索引到电影信息
2023     with open(movie_data_path, 'r', encoding="ISO-8859-1") as f:
2024         data = f.readlines()
2025         for item in data:
2026             item = item.strip().split("::")
2027             mov_info[str(item[0])] = item
2028
2029     usr_file = "./work/ml-1m/users.dat"
2030     usr_info = {}
2031     # 打开文件，读取所有行到data中
2032     with open(usr_file, 'r') as f:
2033         data = f.readlines()
2034         for item in data:
2035             item = item.strip().split("::")
2036             usr_info[str(item[0])] = item
2037
2038     print("当前的用户是: ")
2039     print("usr_id:", usr_id, usr_info[str(usr_id)])
2040     print("对应的特征是: ", usr_feats[str(usr_id)])
2041
2042     print("\n当前电影是: ")
2043     print("mov_id:", mov_id, mov_info[str(mov_id)])
2044     print("对应的特征是: ")
2045     print(mov_feat)
2046
2047     # 以上代码中，我们索引到 usr_id = 2 的用户特征向量，以及 mov_id = 1 的电影特征向量。
2048
2049     # ## 2. 计算用户和所有电影的相似度，构建相似度矩阵
2050     #
2051     # 如下示例均以向 userid = 2
2052     的用户推荐电影为例。与训练一致，以余弦相似度作为相似度衡量。
2053
2054     # In[3]:

```



```

2054
2055
2056 import paddle
2057
2058 # 根据用户ID获得该用户的特征
2059 usr_ID = 2
2060 # 读取保存的用户特征
2061 usr_feat_dir = 'usr_feat.pkl'
2062 usr_feats = pickle.load(open(usr_feat_dir, 'rb'))
2063 # 根据用户ID索引到该用户的特征
2064 usr_ID_feat = usr_feats[str(usr_ID)]
2065
2066 # 记录计算的相似度
2067 cos_sims = []
2068 # 记录下与用户特征计算相似的电影顺序
2069
2070 # 索引电影特征，计算和输入用户ID的特征的相似度
2071 for idx, key in enumerate(mov_feats.keys()):
2072     mov_feat = mov_feats[key]
2073     usr_feat = paddle.to_tensor(usr_ID_feat)
2074     mov_feat = paddle.to_tensor(mov_feat)
2075
2076     # 计算余弦相似度
2077     sim = paddle.nn.functional.common.cosine_similarity(usr_feat, mov_feat)
2078     # 打印特征和相似度的形状
2079     if idx==0:
2080         print("电影特征形状: {}, 用户特征形状: {},
2081               相似度结果形状: {}, 相似度结果: {}".format(mov_feat.shape, usr_feat.shape,
2082                                                             sim.numpy().shape, sim.numpy()))
2081
2082     # 从形状为 (1, 1) 的相似度sim中获得相似度值sim.numpy()[0]，并添加到相似度列表cos_sims中
2083     cos_sims.append(sim.numpy()[0])
2084
2085 # ## 3. 对相似度排序，选出最大相似度
2086 #
2087 #
2088 # 使用np.argsort()函数完成从小到大的排序，注意返回值是原列表位置下标的数组。因为cos_sims
2089 # 和
2090 # mov_feats.keys()的顺序一致，所以都可以用index数组的内容索引，获取最大的相似度值和对应
2091 # 电影。
2092 #
2093 # 处理流程是先计算相似度列表
2094 # cos_sims，将其排序后返回对应的下标列表index，最后从cos_sims和mov_info中取出相似度值和
2095 # 对应的电影信息。
2096 #
2097 #
2098 # 这个处理流程只是展示推荐系统的推荐效果，实际中推荐系统需要采用效率更高的工程化方案，建立“召回+排序”的检索系统。这些检索系统的架构才能应对推荐系统对大量线上需求的实时响应。
2099
2100 # In[4]:
2101
2102 # 对相似度排序，获得最大相似度在cos_sims中的位置
2103 index = np.argsort(cos_sims)
2104 # 打印相似度最大的前topk个位置
2105 topk = 5
2106 print("相似度最大的前{}个索引是{}\n对应的相似度是: {}\n".format(topk, index[-topk:],
2107 [cos_sims[k] for k in index[-topk:]]))
2108
2109 for i in index[-topk:]:
2110     print("对应的电影分别是: movie:{}".format(mov_info[list(mov_feats.keys())[i]]))
2111
2112 # 以上结果可以看出，给用户推荐的电影多是Drama、War、Thriller类型的电影。
2113 #
2114 # 是不是到这里就可以把结果推荐给用户了？还有一个小步骤我们继续往下看。
2115
2116 # ## 4. 加入随机选择因素，使得每次推荐的结果有“新鲜感”
2117 #
2118 #

```

为了确保推荐的多样性，维持用户阅读推荐内容的“新鲜感”，每次推荐的结果需要有所不同，我们随机抽取top_k结果中的一部分，作为给用户的推荐。比如从相似度排序中获取10个结果，每次随机抽取6个结果推荐给用户。

```
2113 #
2114 #
2115 # 使用np.random.choice函数实现随机从top_k中选择一个未被选的电影，不断选择直到选择列表res
2116 # 长度达到pick_num为止，其中pick_num必须小于top_k。
2117 #
2118 # 读者可以反复运行本段代码，观测推荐结果是否有所变化。
2119 #
2120 # 代码实现如下：
2121 #
2122 # In[5]:
2123
2124 top_k, pick_num = 10, 6
2125
2126 # 对相似度排序，获得最大相似度在cos_sims中的位置
2127 index = np.argsort(cos_sims)[-top_k:]
2128
2129 print("当前的用户是：")
2130 # usr_id, usr_info 是前面定义、读取的用户ID、用户信息
2131 print("usr_id:", usr_id, usr_info[str(usr_id)])
2132 print("推荐可能喜欢的电影是：")
2133 res = []
2134
2135 # 加入随机选择因素，确保每次推荐的结果稍有差别
2136 while len(res) < pick_num:
2137     val = np.random.choice(len(index), 1)[0]
2138     idx = index[val]
2139     mov_id = list(mov_feats.keys())[idx]
2140     if mov_id not in res:
2141         res.append(mov_id)
2142
2143 for id in res:
2144     print("mov_id:", id, mov_info[str(id)])
2145
2146 # 最后，我们将根据用户ID推荐电影的实现封装成一个函数，方便直接调用，其函数实现如下。
2147 #
2148 # In[6]:
2149
2150 # 定义根据用户兴趣推荐电影
2151 def recommend_mov_for_usr(usr_id, top_k, pick_num, usr_feat_dir, mov_feat_dir,
2152 mov_info_path):
2153     assert pick_num <= top_k
2154     # 读取电影和用户的特征
2155     usr_feats = pickle.load(open(usr_feat_dir, 'rb'))
2156     mov_feats = pickle.load(open(mov_feat_dir, 'rb'))
2157     usr_feat = usr_feats[str(usr_id)]
2158
2159     cos_sims = []
2160
2161     # with dygraph.guard():
2162     paddle.disable_static()
2163     # 索引电影特征，计算和输入用户ID的特征的相似度
2164     for idx, key in enumerate(mov_feats.keys()):
2165         mov_feat = mov_feats[key]
2166         usr_feat = paddle.to_tensor(usr_feat)
2167         mov_feat = paddle.to_tensor(mov_feat)
2168         # 计算余弦相似度
2169         sim = paddle.nn.functional.common.cosine_similarity(usr_feat, mov_feat)
2170
2171         cos_sims.append(sim.numpy()[0])
2172     # 对相似度排序
2173     index = np.argsort(cos_sims)[-top_k:]
2174
2175     mov_info = {}
2176     # 读取电影文件里的数据，根据电影ID索引到电影信息
2177     with open(mov_info_path, 'r', encoding="ISO-8859-1") as f:
2178         data = f.readlines()
```

```

2179         for item in data:
2180             item = item.strip().split("::")
2181             mov_info[str(item[0])] = item
2182
2183         print("当前的用户是: ")
2184         print("usr_id:", usr_id)
2185         print("推荐可能喜欢的电影是: ")
2186         res = []
2187
2188         # 加入随机选择因素, 确保每次推荐的都不一样
2189         while len(res) < pick_num:
2190             val = np.random.choice(len(index), 1)[0]
2191             idx = index[val]
2192             mov_id = list(mov_feats.keys())[idx]
2193             if mov_id not in res:
2194                 res.append(mov_id)
2195
2196         for id in res:
2197             print("mov_id:", id, mov_info[str(id)])
2198
2199
2200 # In[7]:
2201
2202
2203 movie_data_path = "./work/ml-lm/movies.dat"
2204 top_k, pick_num = 10, 6
2205 usr_id = 2
2206 recommend_mov_for_usr(usr_id, top_k, pick_num, 'usr_feat.pkl', 'mov_feat.pkl',
2207                        movie_data_path)
2208
2209 #
2210 从上面的推荐结果来看, 给ID为2的用户推荐的电影多是Drama、War类型的。我们可以通过用户的ID
2211 从已知的评分数据中找到其评分最高的电影, 观察和推荐结果的区别。
2212
2213 #
2214 下面代码实现给定用户ID, 输出其评分最高的topk个电影信息, 通过对比用户评分最高的电影和当前
2215 推荐的电影结果, 观察推荐是否有效。
2216
2217 # In[8]:
2218
2219 # 给定一个用户ID, 找到评分最高的topk个电影
2220
2221 usr_a = 2
2222 topk = 10
2223
2224 #####
2225 ## 获得ID为usr_a的用户评分过的电影及对应评分 ##
2226 #####
2227 rating_path = "./work/ml-lm/ratings.dat"
2228 # 打开文件, ratings_data
2229 with open(rating_path, 'r') as f:
2230     ratings_data = f.readlines()
2231
2232 usr_rating_info = {}
2233 for item in ratings_data:
2234     item = item.strip().split("::")
2235     # 处理每行数据, 分别得到用户ID, 电影ID, 和评分
2236     usr_id, movie_id, score = item[0], item[1], item[2]
2237     if usr_id == str(usr_a):
2238         usr_rating_info[movie_id] = float(score)
2239
2240 # 获得评分过的电影ID
2241 movie_ids = list(usr_rating_info.keys())
2242 print("ID为 {} 的用户, 评分过的电影数量是: ".format(usr_a), len(movie_ids))
2243
2244 #####
2245 ## 选出ID为usr_a评分最高的前topk个电影 ##
2246 #####
2247 ratings_topk = sorted(usr_rating_info.items(), key=lambda item: item[1])[-topk:]
2248

```

```

2246 movie_info_path = "./work/ml-lm/movies.dat"
2247 # 打开文件，编码方式选择ISO-8859-1，读取所有数据到data中
2248 with open(movie_info_path, 'r', encoding="ISO-8859-1") as f:
2249     data = f.readlines()
2250
2251 movie_info = {}
2252 for item in data:
2253     item = item.strip().split("::")
2254     # 获得电影的ID信息
2255     v_id = item[0]
2256     movie_info[v_id] = item
2257
2258 for k, score in ratings_topk:
2259     print("电影ID: {}, 评分是: {}, 电影信息: {}".format(k, score, movie_info[k]))
2260
2261
2262 #
2263 通过上述代码的输出可以发现，Drama类型的电影是用户喜欢的类型，可见推荐结果和用户喜欢的
2264 电影类型是匹配的。但是推荐结果仍有一些不足的地方，这些可以通过改进神经网络模型等方式来
2265 进一步调优。
2266
2267 # # 从推荐案例的三点思考
2268 #
2269 # 1. Deep Learning is all about "Embedding
2270 Everything"。不难发现，深度学习建模是套路满满的。任何事物均用向量的方式表示，可以直接
2271 基于向量完成“分类”或“回归”任务；也可以计算多个向量之间的关系，无论这种关系是“相似性”还
2272 是“比较排序”。在深度学习兴起不久的2015年，当时AI相关的国际学术会议上，大部分论文均是
2273 将某个事物Embedding后再进行挖掘，火热的程度仿佛即使是路边一块石头，也要Embedding一下看
2274 看是否能挖掘出价值。直到近些年，能够Embedding的事物基本都发表过论文，Embedding的方法也
2275 变得成熟，这方面的论文才逐渐有减少的趋势。
2276
2277 #
2278 # 2.
2279 在深度学习兴起之前，不同领域之间的迁移学习往往要用到很多特殊设计的算法。但深度学习兴起
2280 后，迁移学习变得尤其自然。训练模型和使用模型未必是同样的方式，中间基于Embedding的向量
2281 表示，即可实现不同任务交换信息。例如本章的推荐模型使用用户对电影的评分数据进行监督训练
2282 ，训练好的特征向量可以用于计算用户与用户的相似度，以及电影与电影之间的相似度。对特征向
2283 量的使用可以极其灵活，而不局限于训练时的任务。
2284
2285 #
2286 # 3.
2287 网络调参：神经网络模型并没有一套理论上可推导的最优规则，实际中的网络设计往往是在理论和
2288 经验指导下的“探索”活动。例如推荐模型的每层网络尺寸的设计遵从了信息熵的原则，原始信息量
2289 越大对应表示的向量长度就越长。但具体每一层的向量应该有多长，往往是根据实际训练的效果进
2290 行调整。所以，建模工程师被称为数据处理工程师和调参工程师是有道理的，大量的精力花费在处
2291 理样本数据和模型调参上。
2292
2293 #
2294 # <center></center>
2297
2298 #
2299 # <br>
2300 #
2301 # # 在工业实践中的推荐系统
2302 #
2303 #
2304 本章介绍了比较简单的推荐系统构建方法，在实际应用中，验证一个推荐系统的好坏，除了预测准
2305 确度，还需要考虑多方面的因素，比如多样性、新颖性，甚至商业目标匹配度等。要实践一个好的
2306 推荐系统，值得更深入的探索研究。下面将工业实践推荐系统还需要考虑的主要问题做一个概要性
2307 的介绍。
2308
2309 #
2310 #
2311 1. **推荐来源**：推荐来源会更加多样化，除了使用深度学习模型的方式，还大量使用标签匹配
2312 的个性化推荐方式。此外，推荐热门的内容，具有时效性的内容和一定探索性的内容，都非常关键
2313 。对于新闻类的内容推荐，用户不希望地球人都在谈论的大事自己毫无所知，期望更快更全面的了
2314 解。如果用户经常使用的推荐产品总推荐“老三样”，会使得用户丧失“新鲜感”而流失。因此，除了
2315 推荐一些用户喜欢的内容之外，谨慎的推荐一些用户没表达过喜欢的内容，可探索用户更广泛的兴
2316 趣领域，以便有更多不重复的内容可以向用户推荐。
2317
2318 #
2319 #
2320 2. **检索系统**：将推荐系统构建成“召回+排序”架构的高性能检索系统，以更短的特征向量建
2321 倒排索引。在“召回+排序”的架构下，通常会训练出两种不同长度的特征向量，使用较短的特征向
2322 量做召回系统，从海量候选中筛选出几十个可能候选。使用较短的向量做召回，性能高但不够准确

```

，然后使用较长的特征向量做几十个候选的精细排序，因为待排序的候选很少，所以性能低一些也影响不大。

#

#

3. ****冷启动问题****：现实中推荐系统往往要在产品运营的初期一起上线，但这时候系统尚没有用户行为数据的积累。这时，我们往往建立一套专家经验的规则系统，比如一个在美妆行业工作的店小二对各类女性化妆品偏好是非常了解的。通过规则系统运行一段时间积累数据后，再逐渐转向机器学习的系统。很多推荐系统也会主动向用户收集一些信息，比如大家注册一些资讯类APP时，经常会要求选择一些兴趣标签。

#

#

4. ****推荐系统的评估****：推荐系统的评估不仅是计算模型Loss所能代表的，是使用推荐系统用户的综合体验。除了采用更多代表不同体验的评估指标外（准确率、召回率、覆盖率、多样性等），还会从两个方面收集数据做分析：

#

（1）行为日志：如用户对推荐内容的点击率，阅读市场，发表评论，甚至消费行为等。

#

#

（2）人工评估：选取不同的具有代表性的评估员，从兴趣相关度、内容质量、多样性、时效性等多个维度评估。如果评估员就是用户，通常是以问卷调查的方式下发和收集。

#

#

其中，多样性的指标是针对探索性目标的。而推荐的覆盖度也很重要，代表了所有的内容有多少能够被推荐系统送到用户面前。如果推荐每次只集中在少量的内容，大部分内容无法获得用户流量的话，会影响系统内容生态的健康。比如电商平台如果只推荐少量大商家的产品给用户，多数小商家无法获得购物流量，会导致平台上的商家集中度越来越高，生态不再繁荣稳定。

#

#

从上述几点可见，搭建一套实用的推荐系统，不只是一个有效的推荐模型。要从业务的需求场景出发，构建完整的推荐系统，最后再实现模型的部分。如果技术人员的视野只局限于模型本身，是无法在工业实践中搭建一套有业务价值的推荐系统的。

#

<center></center>

#

<center>图3：推荐系统的全流程</center>

#

#

作业

#

1、设计并完成两个推荐系统，根据相似用户推荐电影（user-based）和根据相似电影推荐电影（item-based），并分析三个推荐系统的推荐结果差异。

#

#

上文中，我们已经将映射后的用户特征和电影特征向量保存在了本地，通过两者的相似度计算结果进行推荐。实际上，我们还可以计算用户之间的相似度矩阵和电影之间的相似度矩阵，实现根据相似用户推荐电影和根据相似电影推荐电影。

#

#

#

2、构建一个【热门】、【新品】和【个性化推荐】三条推荐路径的混合系统。构建更贴近真实场景的推荐系统，而不仅是个性化推荐模型，每次推荐10条，三种各占比例2、3、5条，每次的推荐结果不同。

#

3、推荐系统的案例，实现本地的版本（非AI Studio上实现），进行训练和预测并截图提交。有助于大家掌握脱离AI Studio平台，使用本地机器完成建模的能力。

#