

# MatrixFactorization

October 6, 2018

## 0.1 1. Giới thiệu về phương pháp Recommendation system

### 0.1.1 1.1. Tìm quan trọng của recommendation system

Trong cuộc sống hàng ngày chúng ta thường thấy những tình huống khá tình cờ khi các hệ thống lớn có khả năng đề xuất và hiểu sở thích của người dùng và hiển thị những thông tin mà người dùng quan tâm rất chính xác. Chúng ta hãy xem:

- Facebook có khả năng hiển thị trên newfeed những trạng thái của những người mà bạn quan tâm.
- Youtube có thể đề nghị sang những video mà bạn có khả năng yêu thích dựa trên những gì mà bạn đang xem.
- Amazon có thể đề xuất ra những cuốn sách cùng loại vì những cuốn sách mà bạn đã mua hoặc rating cao.
- Google có thể đề xuất quảng cáo và nội dung phù hợp vì những gì bạn tìm kiếm gần đây.

Nếu không có các thuật toán recommendation, trải nghiệm người dùng sẽ kém hơn vì thông tin mà họ thực sự quan tâm không có thể ra mắt đúng thời điểm trong khi thông tin không cần thiết có thể ra mắt nhiều hơn. Hơn nữa là người dùng cảm thấy bị làm phiền và nhiễu loạn thông tin. Trong lĩnh vực marketing hệ thống recommendation liên tục trở nên quan trọng. Mỗi sản phẩm có thể ra mắt đúng người tiêu dùng có nhu cầu sẽ làm tăng doanh thu, giảm chi phí thời gian, chi phí quảng cáo và giúp người tiêu dùng sở hữu được những thứ mà họ cần. Trong lĩnh vực giải trí như video, game, truyền online,... người dùng sẽ hài lòng cao khi tìm được nội dung phù hợp với sở thích của họ. Dù chỉ mới trong 10 năm trở lại đây, song hành cùng thời kỳ bùng nổ internet những công nghệ recommendation system là một lĩnh vực nghiên cứu sôi nổi. Nó đã tạo ra một cuộc cách mạng thay đổi hành vi mua sắm, hành vi giải trí, chiến lược kinh doanh,... trên toàn cầu. Hàng triệu các doanh nghiệp đang hướng tới nó thông qua khai thác nguồn khách hàng vô tận từ tài nguyên mạng và biến kênh bán hàng này thay thế các kênh truyền thống. Lĩnh vực này ngày nay là chìa khóa mở ra giúp các công ty công nghệ Google, Facebook, Amazon, Microsoft,... trở thành những tập đoàn hàng đầu thế giới. Chính vì thế recommendation system luôn có các công ty kinh doanh trên nền tảng online đầu tư nghiên cứu và phát triển để tạo ra một hệ thống thông minh nhằm nâng cao trải nghiệm khách hàng và tối ưu hóa nguồn tài nguyên.

### 0.1.2 1.2. Phương pháp recommendation

Bên trên chúng ta đã biết vai trò của recommendation system trong việc phát triển các lĩnh vực internet và kinh doanh online. Tuy nhiên thực sự bài toán recommendation system là gì? Các phương pháp recommendation system ra sao chúng ta vẫn chưa thực sự hiểu rõ. Theo định nghĩa từ wikipedia thì recommendation system là một nhánh nhỏ của lĩnh vực *hệ thống lọc thông tin* (information filtering system) có nhiệm vụ đề xuất các nội dung yêu thích thông qua rating của một người dùng (user) cho một sản phẩm (item). Dựa trên các sản phẩm phù hợp nhất với người dùng để hiển thị các hệ thống phi dựa trên thông tin về rating của sản phẩm, thông tin người dùng, thông tin về sản phẩm xây dựng thuật toán tối ưu. Dựa trên hàm loss function tính ra sai số để đánh giá các thuật toán và tìm ra một phương pháp có mức độ chính xác nhất. Có rất nhiều các thuật toán

khác nhau c s dng trong recommendation system nhng v c bn chúng bao gm 2 phng pháp chính: Collaborative filtering và Content based filtering. im khác bit c bn gia 2 phng pháp này là:

- Collaborative filtering: Da trên m quan h tng quan v mt hành vi tiêu dùng hoc c trng sn phm tìm ra các users hoc items có chung c tính, s thích. T ó da trên nhng thông tin mà nhóm ngi dùng hoc sn phm liên quan gn nht ã rating ánh giá sn phm mà mt ngi dùng c th cha rating. Tuy nhiên nhc im ca thut toán này là a ra d báo v rating mà không hoàn toàn hieu v user, item mà hoàn toàn da trên quan sát v mc tng ng gia các nhóm users, items d báo. Thut toán c s dng trong các bài toán này ch yu là k-nearest neighbor tìm ra nhóm tng ng và ma trn h s tng quan c s dng o lng mc gn gi v mt hành vi hay c tính phân nhóm.
- Content based filtering: Da trên nhng thông tin và ni dung liên quan n sn phm nh nhà sn xut, th loi, nm sn xut, công dng, c tính,... hoc da trên thông tin ca ngi dùng nh gii tính, tui, ngành ngh,... a ra d báo v rating ca ngi i vi sn phm ó. Thut toán này ch n thun là các phng trình hi qui gia các chiu c tính ca sn phm hoc ngi dùng i vi im rating mà không tn dng c tng quan v mt hành vi gia nhng nhóm ngi dùng hay c trng sn phm nh Collaborative filtering. Trong thc t hành vi ca ngi dùng li cho thy rt gng nhau nu thuc cùng mt nhóm chng hn nh các nhóm nhc thin, nhc vàng, nhc tr, nhc thiu nhi s phù hp vi ngi già, ngi trung niên, ngi tr, thiu nhi. Không xem xét c các yu t tng quan theo nhóm là mt hn ch ln ca content based filtering.

Mi thut toán u có u, nhc im khác nhau và mc hieu qu trong d báo mc yêu thích các cp (user, item) (*ngi dùng, sn phm*) cng khác nhau tùy thuc vào tp d liu. Nhng các thut toán u có im chung ó là s dng d liu mà ngi dùng ã rating i vi các sn phm làm c s d báo rating cho các sn phm cha c ánh giá. Vic này cng gng nh chúng ta chỉ trò chỉ in s vào *ma trn tin ích* (utility matrix). Mt chiu ca ma trn ng vi users và chiu còn li ng vi items. Các ô trên ma trn th hin giá tr rating ca user tng ng lên item. Nh vy s có nhng ô ã c rating bi ngi dùng và các ô còn li cha c rating. Quá trình gii bài toán cng gng nh vic chúng ta i gii ma trn ti nhng ô còn thiu sao cho sai s cui cùng gia d báo và thc t là nh nht.

### 0.1.3 1.3. Gii thiu thut toán matrix factorization

Trong thut toán matrix factorization chúng ta gi nh c trng ca item c th hin qua ma trn  $\mathbf{I}$  và hành vi ca ngi dùng c th hin qua ma trn  $\mathbf{U}$ . Vi mi dòng ca ma trn  $\mathbf{I}$  là mt c trng n (*latent feature*) ca sn phm và mi ct ca  $\mathbf{U}$  là mc yêu thích ca mt ngi dùng i vi c trng n tng ng. Các c trng n này có th coi nh nhng nhân t chính c tng hp t nhieu thông tin liên quan n sn phm tng t nh thành phn chính trong phép phân tích thành phn chính PCA. c trng ca sn phm th m c th hin qua vector dòng  $\mathbf{i}_m$  và hành vi ca ngi dùng th n c th hin qua vector ct  $\mathbf{u}_n$ . Khi ó giá tr d báo mc yêu thích ca mt ngi dùng n lên mt sn phm m s là tích ca 2 vector  $\mathbf{i}_m$  và  $\mathbf{u}_n$ :

$$y_{mn} = \mathbf{i}_m \mathbf{u}_n$$

c lng ca ma trn tin ích  $\hat{\mathbf{Y}}$  s c biu din theo các ma trn hành vi  $\mathbf{I}$  và ma trn ngi dùng  $\mathbf{U}$  nh sau:

$$\hat{\mathbf{Y}} \approx \begin{bmatrix} \mathbf{i}_1 \mathbf{u}_1 & \mathbf{i}_1 \mathbf{u}_2 & \dots & \mathbf{i}_1 \mathbf{u}_N \\ \mathbf{i}_2 \mathbf{u}_1 & \mathbf{i}_2 \mathbf{u}_2 & \dots & \mathbf{i}_2 \mathbf{u}_N \\ \dots & \dots & \ddots & \dots \\ \mathbf{i}_M \mathbf{u}_1 & \mathbf{i}_M \mathbf{u}_2 & \dots & \mathbf{i}_M \mathbf{u}_N \end{bmatrix} = \begin{bmatrix} \mathbf{i}_1 \\ \mathbf{i}_2 \\ \dots \\ \mathbf{i}_M \end{bmatrix} [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \dots \quad \mathbf{u}_N] = \mathbf{I} \mathbf{U}$$

#### 0.1.4 1.4. Thuật toán gradient descent

Giả sử rằng chúng ta đã có thông tin về các ma trận  $\mathbf{U}$  và  $\mathbf{I}$  tức là chúng ta cần thực hiện bây giờ là coi các dòng của ma trận  $\mathbf{I}$  là một item profile và mỗi cột của  $\mathbf{U}$  là một user profile. Giả sử  $\mathbf{I} \in \mathbb{R}^{M \times K}$ ,  $\mathbf{U} \in \mathbb{R}^{K \times N}$ ,  $\mathbf{Y} \in \mathbb{R}^{M \times N}$ . Thông thường ta sẽ chọn số hạng nhỏ hơn số hạng sản phẩm và ngừng. Khi đó  $\mathbf{Y}$  sẽ biểu diễn dạng tích của 2 ma trận có rank nhỏ hơn (*Low-rank Matrix factorization*):

$$\hat{\mathbf{Y}} = \mathbf{I}\mathbf{U}$$

Hàm loss function của thuật toán chính là chuẩn **Frobenius norm** và lệch giữa  $\mathbf{Y}$  và  $\hat{\mathbf{Y}}$  như sau:

$$\mathcal{L}(\mathbf{I}, \mathbf{U}) = \frac{1}{2s} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2$$

tránh hiện tượng overfitting bằng **hệ số điều chỉnh bậc 2** (*l2 - regularization*) của thêm vào:

$$\mathcal{L}(\mathbf{I}, \mathbf{U}) = \frac{1}{2s} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2 + \frac{\lambda_1}{2} \|\mathbf{I}\|_F^2 + \frac{\lambda_2}{2} \|\mathbf{U}\|_F^2 \quad (1.4.1)$$

Nếu coi  $\mathbf{U}$  cố định và cần tối ưu  $\mathbf{I}$ . Bài toán Matrix factorization sẽ trở thành bài toán tối ưu hàm loss function:

$$\mathcal{L}(\mathbf{I}) = \frac{1}{2s} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2 + \frac{\lambda_1}{2} \|\mathbf{I}\|_F^2$$

Nếu coi  $\mathbf{I}$  cố định và cần tối ưu  $\mathbf{U}$ . Hàm loss function sẽ có dạng:

$$\mathcal{L}(\mathbf{U}) = \frac{1}{2s} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2 + \frac{\lambda_2}{2} \|\mathbf{U}\|_F^2$$

Ta nhận thấy hàm loss function là những hàm lồi. Việc tìm nghiệm tối ưu có thể dựa trên bài toán tối ưu lồi bậc 2 (*Quadratic Programming*) hoặc cách đơn giản hơn là thông qua thuật toán gradient descent. Chúng ta sẽ dùng thuật toán Stochastic gradient descent cập nhật lần lượt từng item dựa trên toàn bộ dữ liệu, sau đó lặp lại quá trình này. Vì trong tập hợp ma trận sản phẩm ( $\mathbf{I}$ ) cố định, ta sẽ cần cập nhật ma trận ngừng ( $\mathbf{U}$ ) theo phương pháp gradient descent. Mỗi một lần cập nhật, một ngừng được cập nhật lần nữa. Dựa trên thông tin về những sản phẩm mà ngừng được đánh giá. Vector gradient descent để tính toán cập nhật giá trị của vector  $\mathbf{u}$  ngừng. Quá trình này tiếp tục cho đến khi toàn bộ các vector users được cập nhật. Từng lần như vậy vì trong tập hợp cố định ma trận ngừng cố định và cập nhật ma trận sản phẩm.

##### Tối ưu ma trận ngừng:

Đơn giản hóa quá trình tính toán, ta có thể biểu diễn hàm loss function theo từng loss function của từng user như sau:

$$\mathcal{L}(\mathbf{U}) = \frac{1}{2s} \sum_{n=1}^N \sum_{m:r_{mn}=1} (y_{mn} - i_m \cdot u_n)^2 + \frac{\lambda_1}{2} \|\mathbf{U}\|_F^2$$

Trong đó  $r_{mn}$  là phần tử thực của ma trận rating  $R \in \mathbb{R}^{M \times N}$  có giá trị 0 hoặc 1.  $r_{mn} = 1$  ám chỉ sản phẩm  $m$  đã được rating bởi user  $n$  và bằng 0 trong tập hợp chưa có rating. Nếu  $r_{mn} = 1$  trong hàm loss function là để ra những sản phẩm đã được rating bởi user  $n$ . Khi đó nếu coi  $\hat{\mathbf{I}}_n$  là ma trận các sản phẩm đã được rating của user  $n$  và  $\hat{y}_n$  là vector kết quả rating ngừng thì hàm loss function vì user  $n$  có thể viết gọn như sau:

$$\mathcal{L}(\mathbf{U} | \nabla f | \nabla = \setminus) = \frac{1}{2s} \sum_{m:r_{mn}=1} (y_{mn} - \mathbf{i}_m \cdot \mathbf{u}_n)^2 + \frac{\lambda_1}{2} \|\mathbf{u}_n\|^2 = \frac{1}{2s} \|\hat{y}_n - \hat{\mathbf{I}}_n \cdot \mathbf{u}_n\|^2 + \frac{\lambda_1}{2} \|\mathbf{u}_n\|^2$$

o hàm của nó ngừng:

$$\frac{\partial \mathcal{L}(\mathbf{U} | \mathcal{I})}{\partial \mathbf{u}_n} = -\frac{1}{s} \hat{\mathbf{I}}_n^T (\hat{y}_n - \hat{\mathbf{I}}_n \cdot \mathbf{u}_n) + \lambda_1 \mathbf{u}_n$$

Công thức cập nhật nghiệm cho mỗi cột của ma trận người dùng:

$$\mathbf{u}'_n = \mathbf{u}_n - \theta \left( -\frac{1}{s} \hat{\mathbf{I}}_n^T (\hat{y}_n - \hat{\mathbf{I}}_n \cdot \mathbf{u}_n) + \lambda_1 \mathbf{u}_n \right)$$

**Tích nhân ma trận:**

Hoàn toàn tương tự ta cũng có tích nhân ma trận sản phẩm, hàm loss function tích nhân item = m:

$$\mathcal{L}(\mathbf{I} | \mathcal{I}) = \frac{1}{2s} \sum_{n:r_{mn}=1} (y_{mn} - \mathbf{i}_m \cdot \mathbf{u}_n)^2 + \frac{\lambda_1}{2} \|\mathbf{i}_m\|^2 = \frac{1}{2s} \|\hat{y}_m - \mathbf{i}_m \cdot \hat{\mathbf{U}}_m\|^2 + \frac{\lambda_1}{2} \|\mathbf{i}_m\|^2$$

o hàm tương ứng tích nhân item s là:

$$\frac{\partial \mathcal{L}(\mathbf{I} | \mathcal{I})}{\partial \mathbf{i}_m} = -\frac{1}{s} (\hat{y}_m - \mathbf{i}_m \cdot \hat{\mathbf{U}}_m) \hat{\mathbf{U}}_m^T + \lambda_1 \mathbf{i}_m$$

Công thức cập nhật nghiệm cho mỗi dòng của ma trận sản phẩm:

$$\mathbf{i}'_m = \mathbf{i}_m - \theta \left( -\frac{1}{s} (\hat{y}_m - \mathbf{i}_m \cdot \hat{\mathbf{U}}_m) \hat{\mathbf{U}}_m^T + \lambda_1 \mathbf{i}_m \right)$$

## 0.2 Xây dựng code thuật toán

### 0.2.1 Thực hành trên bộ dữ liệu movie length 1M

Chúng ta sẽ thực hiện phương pháp matrix factorization trên bộ dữ liệu [Movie length 1M](#) gồm 1 triệu các lượt ratings cho khoảng 4000 bộ phim có thu thập 6000 người dùng. Rất phù hợp vì thuật toán đã xây dựng, các xử lý dữ liệu sẽ thực hiện trên ma trận. Load dữ liệu vào như sau:

```
In [1]: import pandas as pd
import numpy as np
columns = ['user_id', 'item_id', 'rating', 'timestamp']
movie_length = pd.read_csv('ml-1m/ratings.dat', header = 0, \
                           names = columns, sep = '::', engine = 'python')
movie_length = movie_length.sort_values(['user_id', 'item_id'])
movie_length.head()
```

```
Out[1]:
```

	user_id	item_id	rating	timestamp
39	1	1	5	978824268
24	1	48	5	978824351
38	1	150	5	978301777
43	1	260	4	978300760
22	1	527	5	978824195

Kích thước các dữ liệu và số lượng users, items

```
In [2]: print('Data movie length shape: %s'%str(movie_length.shape))
print('No customers: %s'%str(np.unique(movie_length.iloc[:, 0]).shape[0]))
print('No movies: %s'%str(np.unique(movie_length.iloc[:, 1]).shape[0]))
```

Data movie length shape: (1000208, 4)  
No customers: 6040  
No movies: 3706

Thng kê mô t s tn sut rating ca các users:

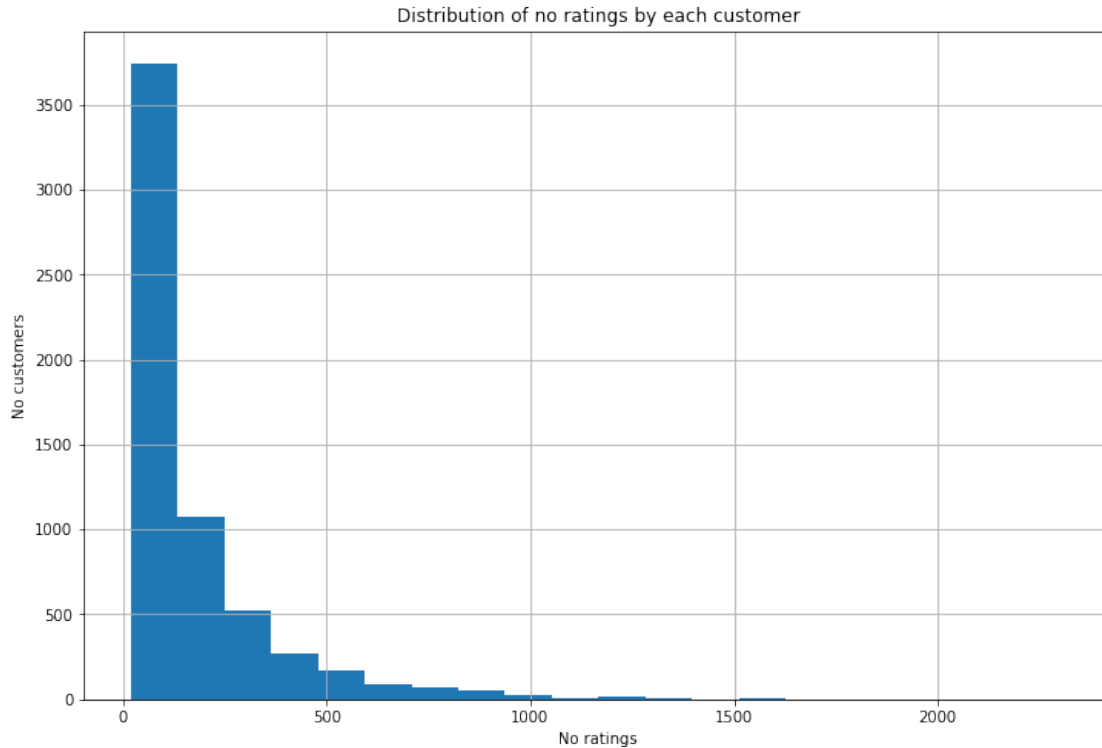
```
In [3]: movie_length['user_id'].value_counts().describe()
```

```
Out[3]: count      6040.000000  
        mean       165.597351  
        std        192.747126  
        min         20.000000  
        25%         44.000000  
        50%         96.000000  
        75%        208.000000  
        max        2314.000000  
        Name: user_id, dtype: float64
```

- Bình quân mt user rate tng cng 165 b phim.
- User thp nht rate 20 b phim và user nhieu nht rate 2314 b phim.
- Không s lng b phim rating ph bin ca mt user là t 44 b ti 208 b phim (chim 50%).

```
In [4]: import matplotlib.pyplot as plt  
        %matplotlib inline  
        movie_length[['user_id', 'item_id']].groupby(['user_id']).count().\n        hist(bins = 20, figsize = (12, 8))  
        plt.title('Distribution of no ratings by each customer')  
        plt.xlabel('No ratings')  
        plt.ylabel('No customers')
```

```
Out[4]: Text(0,0.5,'No customers')
```



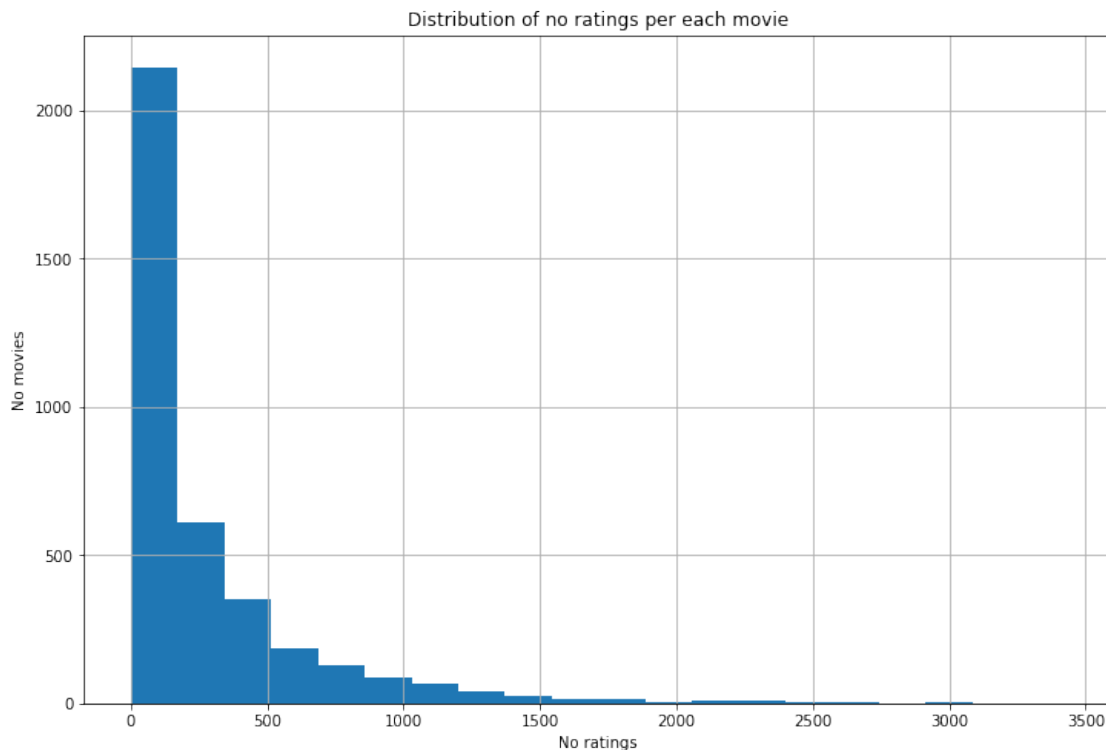
Nhìn vào phân bố số lượng user theo mức rating ta có thể thấy mục đích của chúng ta có hình ảnh không cân bằng khi có nhiều user rating rất ít và nhiều user rating nhiều hơn. Tuy nhiên đây không phải là bài toán classification nên việc mục đích có kích thước mất cân bằng cũng không ảnh hưởng gì đến chính xác của thuật toán. Hiện nay thuật toán matrix factorization xây dựng hàm loss function riêng cho từng user nên việc user này rating bao nhiêu sản phẩm không ảnh hưởng đến kết quả dự báo rating của user khác. Hoàn toàn ngược lại ta cũng thường kê số lượng các user rating vì từng bộ phim.

```
In [5]: movie_length['item_id'].value_counts().describe()
```

```
Out[5]: count    3706.000000
        mean      269.888829
        std       384.046815
        min        1.000000
        25%        33.000000
        50%       123.500000
        75%       350.000000
        max      3428.000000
        Name: item_id, dtype: float64
```

```
In [6]: movie_length[['user_id', 'item_id']].groupby(['item_id']).count().\
        hist(bins = 20, figsize = (12, 8))
        plt.title('Distribution of no ratings per each movie')
        plt.xlabel('No ratings')
        plt.ylabel('No movies')
```

Out[6]: Text(0,0.5,'No movies')



Mt vài ánh giá: \* S ln 1 b phim c rating ít nht là 1 ln. \* S ln 1 b phim c rating nhiều nht là 3428 ln. \* Mc rating ph bin ca mt b phim là t 33 n 123 ln.

## 0.3 2.2. Thuật toán matrix factorization

### 0.3.1 2.2.1. Các hàm trong thuật toán

u tiên ta s tin hành chia mu train và test theo t l sao cho s lng rating trong tp train chim 2/3 s lng các lt rating. Cách chia mu hp lý nht là m bo t l s lng ratings xut hin trong tp train i vì s lng ratings xut hin trong tp test ca cùng mt user là bng nhau. Cách chia này m bo s công bng i vì các user khi không có user nào có quá nhiều d liu train và d liu test ít hoc d liu train quá ít nhng d liu test li quá nhiều. Giá tr c d báo t mô hình mà d liu train quá ít s thng không chun xác và làm sai lch kt qu kim tra sai s trên test.

```
In [7]: #declare split_rate for train/total ratings
split_rate = 2/3
```

```
def split_train_test(dataset):
    gb = dataset.groupby('user_id')
    ls = [gb.get_group(x) for x in gb.groups]
    items = [x for x in gb.groups]
    index_size = [{'i': i, 'index': gb.groups[i], 'size': len(gb.groups[i])} for i in items]
    index_train = pd.Int64Index([])
```

```

index_test = pd.Int64Index([])
for x in index_size:
    np.random.shuffle(x['index'].values)
    le = int(x['size']*split_rate)
    index_train = index_train.append(x['index'][:le])
    index_test = index_test.append(x['index'][le:])
train = dataset.iloc[index_train].values
test = dataset.iloc[index_test].values
#minus id to 1 to index start from 0
train[:, 0] -= 1
train[:, 1] -= 1
test[:, 0] -= 1
test[:, 1] -= 1
return train, test

```

```
train, test = split_train_test(movie_length)
```

### Tin hành xây dựng thuật toán:

Các biến chính cho thuật toán dựa vào bao gồm: \*  $n\_users$ : Số lượng user (chính là  $N$  trong thuật toán). \*  $n\_items$ : Số lượng item (chính là  $M$  trong thuật toán). \*  $K$ : Số lượng nhân tố ẩn  $K$  (giá trị  $K$  trong thuật toán). \*  $\theta$ : Tham số  $\theta$  cập nhật trong thuật toán gradient descent. \*  $split\_rate$ : Tỷ lệ chia dữ liệu train/test. \*  $\lambda$ : Tham số điều chỉnh các thành phần điều chỉnh  $L_2$  - regularization (nghiên cứu thử nghiệm  $\lambda_1 = \lambda_2$ ). \*  $I$ : Ma trận sản phẩm. \*  $U$ : Ma trận người dùng.

Lưu ý rằng các ma trận  $I, U$  xây dựng dựa trên các *nhân tố ẩn* (latent feature) nên bạn cần xác định các nhân tố này và phải khi đó giá trị ngẫu nhiên cho chúng. Số lượng nhân tố ẩn  $K$  là một giá trị tùy ý bạn có thể lựa chọn. Theo [Matrix Factorization For Recommendation System](#) thì khi số lượng nhân tố ẩn càng nhiều thuật toán càng chính xác hơn nhưng cũng làm gia tăng chi phí tính toán. Nếu thì những mô hình xây dựng dựa trên nhân tố ẩn càng tinh luyện có mức khác biệt lớn càng ra kết quả chính xác hơn các nhân tố ngẫu nhiên.

```

In [8]: n_users = np.max(train[:, 0] + 1) #plus one because index start from 0
        n_items = np.max(train[:, 1] + 1)
        n_ratings = train.shape[0]
        print('N user dimension: %s'%n_users)
        print('M item dimension: %s'%n_items)
        print('S Number of rating: %s'%n_ratings)
        K = 2
        theta = 0.75
        lamda = 0.2
        #Initialize random matrix according to Gauss distribution
        I = np.random.randn(n_items, K)
        U = np.random.randn(K, n_users)

```

N user dimension: 6040

M item dimension: 3952

S Number of rating: 664826

```

In [9]: import scipy.sparse as sparse
        #Rating matrix

```



```

Y = np.zeros(shape = (n_items, n_users))
print('Y utility matrix shape: %s'%str(Y.shape))
Y = sparse.coo_matrix((train[:, 2], (train[:, 1], train[:, 0])),\
                      shape = (n_items, n_users), dtype = np.float).toarray()

```

Y utility matrix shape: (3952, 6040)

Không phi hoàn toàn các giá trị trên ma trận  $Y$  về rating. Vì vậy ma trận  $R$  có ra nhím ảnh hưởng các vị trí có rating của  $Y$  bằng giá trị 1 và chưa có rating bằng 0.

```

In [10]: R = sparse.coo_matrix((np.ones((n_ratings,)), (train[:, 1], train[:, 0])),\
                              shape = (n_items, n_users)).toarray()

```

thủ thuật gradient descent hiệu quả hơn chúng ta cần chuẩn hóa ma trận  $Y$  về giá trị kỳ vọng bằng 0 bằng cách trừ đi giá trị rating trong vector rating của mỗi user về trung bình của vector rating đó.

```

In [11]: def standardize_Y(Y):
    sum_rating = Y.sum(axis = 0)
    u_rating = np.count_nonzero(Y, axis = 0)
    u_mean = sum_rating/u_rating
    for n in range(n_users):
        for m in range(n_items):
            if Y[m, n] != 0:
                Y[m, n] -= u_mean[n]
    return Y, u_mean

```

```

Y_stad, u_mean = standardize_Y(Y)

```

Sau khi chuẩn hóa ma trận  $Y$  thì phần quan trọng nhất là áp dụng thủ thuật gradient descent tối ưu hóa các hệ số của ma trận  $U, I$ . Dựa trên lý thuyết về thủ thuật đã xây dựng môc 1.4 xây dựng các hàm số update ma trận:

**Thủ thuật gradient descent cho ma trận người dùng:**

```

In [12]: def updateU(U):
    for n in range(n_users):
        # Matrix items include all items is rated by user n
        i Rated = np.where(Y_stad[:, n] != 0)[0] #item's index rated by n
        In = I[i Rated, :]
        if In.shape[0] == 0:
            U[:, n] = 0
        else:
            s = In.shape[0]
            u_n = U[:, n]
            y_n = Y_stad[i Rated, n]
            grad = -1/s * np.dot(In.T, (y_n - np.dot(In, u_n))) + lamda*u_n
            U[:, n] -= theta*grad
    return U

```

**Thủ thuật gradient descent cho ma trận sản phẩm:**

```
In [13]: def updateI(I):
    for m in range(n_items):
        # Matrix users who rated into item m
        i_rated = np.where(Y_stad[m, :] != 0)[0] #user's index rated into m
        Um = U[:, i_rated]
        if Um.shape[1] == 0:
            I[m, :] = 0
        else:
            s = Um.shape[1]
            i_m = I[m, :]
            y_m = Y_stad[m, i_rated]
            grad = -1/s * np.dot(y_m - np.dot(i_m, Um), Um.T) + lamda*i_m
            I[m, :] -= theta*grad
    return I
```

#### Xây dựng hàm dự báo ma trận Y:

Dựa trên ma trận  $U$  và  $I$  ta có thể tính toán ma trận dự báo của  $Y$  là  $\hat{Y}$  theo công thức (1.4.0) và xây dựng hàm `pred()`. Các kết quả dự báo của chuyển hóa ngược lại rating bằng cách cộng thêm trung bình rating của mỗi user vào các giá trị ratings thực cùng 1 user. Một số kết quả vượt quá mức giá trị của rating là  $[1, 5]$  và khi đó sẽ gán lại về 2 hoặc 5. Ma trận thu được sẽ thỏa mãn tính chất của  $\hat{Y}$  là kết quả dự báo rating của người dùng i về sản phẩm tng ng. Do ta chỉ cần ánh giá trị  $Y$  trên những cặp (user,item) đã có rating nên trong hàm `pred_train_test()` ta cần đưa vào ma trận rating  $R$  thay thế những vị trí chưa có rating bằng 0. Lý do hàm này có tên là `pred_train_test()` là chúng ta cũng có thể thực hiện tương tự cho tập test khi thay thế giá trị chưa rating bằng 0.

```
In [14]: def pred(U, I):
    #predict utility matrix base on formula  $\hat{Y} = I.U$ 
    Y_hat = np.dot(I, U)
    #invert to forecast values by plus user's mean ratings
    for n in range(n_users):
        Y_hat[:, n] += u_mean[n]
    #convert to interger values because of rating is integer
    Y_hat = Y_hat.astype(np.int32)
    #replace values > 5 by 5 and values < 1 by 1
    Y_hat[Y_hat > 5] = 5
    Y_hat[Y_hat < 1] = 1
    return Y_hat

def pred_train_test(Y_hat, R):
    #replace values have not yet rated by 0
    Y_pred = Y_hat.copy()
    Y_pred[R == 0] = 0
    return Y_pred
```

#### Xây dựng hàm loss function:

Hàm loss function sẽ xây dựng dựa trên công thức (1.4.1) như sau:

```
In [15]: def loss(Y, Y_hat):
    error = Y - Y_hat
```

```

    loss_value = 1/(2*n_ratings)*np.linalg.norm(error, 'fro')**2 + \
    lamda/2*(np.linalg.norm(I, 'fro')**2 + np.linalg.norm(U, 'fro')**2)
    return loss_value

```

Sử dụng ma trận  $\hat{Y}$  để báo trên tập test

```

In [16]: Y_test = sparse.coo_matrix((test[:, 2], (test[:, 1], test[:, 0])), \
                                     shape = (n_items, n_users), dtype = np.float).toarray()
    R_test = sparse.coo_matrix((np.ones(test.shape[0]), (test[:, 1], test[:, 0])), \
                                shape = (n_items, n_users), dtype = np.float).toarray()

```

### Xây dựng hàm tính RMSE:

Sau khi tính được ma trận dự báo  $\hat{Y}$  trên tập test kết hợp với ma trận tin tức  $Y$  của tập test để tính được RMSE trên tập test như sau:

```

In [17]: import math
    def RMSE(Y_test, Y_pred):
        error = Y_test - Y_pred
        n_ratings = test.shape[0]
        rmse = math.sqrt(np.linalg.norm(error, 'fro')**2/n_ratings)
        return rmse

```

### Xây dựng vòng lặp tối ưu chính:

Sau khi đã thiết kế các hàm tính toán *loss function*, *RMSE* và các hàm tối ưu *gradient descent* ta sẽ tiến hành xây dựng vòng lặp tối ưu cập nhật các ma trận  $U$  và  $I$  và ánh xạ giá trị qua các mô hình thông qua giá trị của *loss function* và *RMSE*.

```

In [19]: def fit(Umatrix, Imatrix, Ytrain, Ytest, n_iter, log_iter):
    for i in range(n_iter):
        #update U and I
        Umatrix = updateU(Umatrix)
        Imatrix = updateI(Imatrix)
        #calculate Y_hat
        Y_hat = pred(Umatrix, Imatrix)
        #calculate Y_hat_train by replace non ratings by 0
        Y_pred_train = pred_train_test(Y_hat, R)
        #calculate loss function
        loss_value = loss(Ytrain, Y_pred_train)
        #calculate Y_pred on test dataset
        Y_pred_test = pred_train_test(Y_hat, R_test)
        #calculate RMSE
        rmse = RMSE(Ytest, Y_pred_test)
        if i % log_iter == 0:
            print('Iteration: {}; RMSE: {}; Loss value: {}'.format(i, rmse, loss_value))
    return Y_hat, Y_pred_test
Y_hat, Y_pred = fit(Umatrix = U, Imatrix = I, Ytrain = Y, Ytest = Y_test, n_iter = 100)

```

Iteration: 0; RMSE: 1.0978996819873132; Loss value: 367.43589641528513

Iteration: 10; RMSE: 1.0972517727989382; Loss value: 371.61999405456453

```

Iteration: 20; RMSE: 1.0967013607569651; Loss value: 376.1721780533387
Iteration: 30; RMSE: 1.096062264477316; Loss value: 380.4612541939793
Iteration: 40; RMSE: 1.0953942145953313; Loss value: 384.2559773463809
Iteration: 50; RMSE: 1.0948442303904773; Loss value: 387.81426251576477
Iteration: 60; RMSE: 1.0945337212732151; Loss value: 390.15217556003586
Iteration: 70; RMSE: 1.0943198545615616; Loss value: 391.6867939243047
Iteration: 80; RMSE: 1.0941863369527691; Loss value: 392.6789203399845
Iteration: 90; RMSE: 1.0940269119371213; Loss value: 393.30221888979895

```

### 0.3.2 2.2.2. Xây dựng class MF

Trên các hàm đã xây dựng (2.2.1) ta sẽ thiết kế class MF có chức năng xử lý dữ liệu, fitting model, đánh giá kết quả model và đưa ra các recommend cho khách hàng về sản phẩm như sau:

**Class Data xử lý dữ liệu:**

```

In [20]: class Data(object):
        """
        This class used to manage data.
        Two arguments:
        dataset: pandas data frame include user_id, item_id and rating
        split_rate: number train ratings/ total ratings
        """
        def __init__(self, dataset, split_rate):
            self.dataset = dataset
            self.split_rate = split_rate
            self.train, self.test = self.split_train_test(self.dataset)
            self.Ytrain, self.Rtrain = self.utility_matrix(self.train)
            self.Ytest, self.Rtest = self.utility_matrix(self.test)
            self.Ystad, self.u_mean = self.standardize_Y(self.Ytrain)
            self.n_users = np.max(self.train[:, 0] + 1) #plus one because index start from 0
            self.n_items = np.max(self.train[:, 1] + 1)
            self.n_ratings = self.train.shape[0]

        def split_train_test(self, dataset):
            "split train and test"
            gb = dataset.groupby('user_id')
            ls = [gb.get_group(x) for x in gb.groups]
            items = [x for x in gb.groups]
            index_size = [{'i': i, 'index': gb.groups[i], 'size': len(gb.groups[i])} for i in range(len(ls))]
            index_train = pd.Int64Index([])
            index_test = pd.Int64Index([])
            for x in index_size:
                np.random.shuffle(x['index'].values)
                le = int(x['size']*split_rate)
                index_train = index_train.append(x['index'][:le])
                index_test = index_test.append(x['index'][le:])
            train = dataset.iloc[index_train].values

```

```

test = dataset.iloc[index_test].values
#minus id to 1 to index start from 0
train[:, 0] -= 1
train[:, 1] -= 1
test[:, 0] -= 1
test[:, 1] -= 1
return train, test

def utility_matrix(self, data_mtx):
    "create Y and R matrix"
    Y = np.zeros(shape = (n_items, n_users))
    Y = sparse.coo_matrix((data_mtx[:, 2], (data_mtx[:, 1], data_mtx[:, 0])), \
                          shape = (n_items, n_users), dtype = np.float).toarray()
    R = sparse.coo_matrix((np.ones((data_mtx.shape[0],)), (data_mtx[:, 1], data_mtx[:, 0])), \
                          shape = (n_items, n_users)).toarray()

    return Y, R

def standardize_Y(self, Y):
    "standard data to mean ratings of each user = 0"
    sum_rating = Y.sum(axis = 0)
    u_rating = np.count_nonzero(Y, axis = 0)
    u_mean = sum_rating/u_rating
    for n in range(n_users):
        for m in range(n_items):
            if Y[m, n] != 0:
                Y[m, n] -= u_mean[n]
    return Y, u_mean

```

Class model xây dựng và đánh giá model:

```

In [21]: class Model():
    """
    This class manage update U and I matrix, predict and evaluate error
    Four arguments:
    data: instance from Data class which supplies the data for model
    theta: learning rate
    lamda: regularization parameter
    K: number of latent factors
    """
    def __init__(self, data, theta, lamda, K):
        self.data = data
        self.theta = theta
        self.lamda = lamda
        self.K = K
        self.I = np.random.randn(data.n_items, K)
        self.U = np.random.randn(K, data.n_users)

```

```

def updateU(self):
    for n in range(self.data.n_users):
        # Matrix items include all items is rated by user n
        i_rated = np.where(self.data.Ystad[:, n] != 0)[0] #item's index rated by user n
        In = self.I[i_rated, :]
        if In.shape[0] == 0:
            self.U[:, n] = 0
        else:
            s = In.shape[0]
            u_n = self.U[:, n]
            y_n = self.data.Ystad[i_rated, n]
            grad = -1/s * np.dot(In.T, (y_n - np.dot(In, u_n))) + self.lamda*u_n
            self.U[:, n] -= self.theta*grad

def updateI(self):
    for m in range(self.data.n_items):
        # Matrix users who rated into item m
        i_rated = np.where(self.data.Ystad[m, :] != 0)[0] #user's index rated into item m
        Um = self.U[:, i_rated]
        if Um.shape[1] == 0:
            self.I[m, :] = 0
        else:
            s = Um.shape[1]
            i_m = self.I[m, :]
            y_m = self.data.Ystad[m, i_rated]
            grad = -1/s * np.dot(y_m - np.dot(i_m, Um), Um.T) + self.lamda*i_m
            self.I[m, :] -= self.theta*grad

def pred(self, I, U):
    #predict utility matrix base on formula Yhat = I.U
    Yhat = np.dot(I, U)
    #invert to forecast values by plus user's mean ratings
    for n in range(self.data.n_users):
        Yhat[:, n] += self.data.u_mean[n]
    #convert to interger values because of rating is integer
    Yhat = Yhat.astype(np.int32)
    #replace values > 5 by 5 and values < 1 by 1
    Yhat[Yhat > 5] = 5
    Yhat[Yhat < 1] = 1
    return Yhat

def pred_train_test(self, Yhat, R):
    #replace values have not yet rated by 0
    Y_pred = Yhat.copy()
    Y_pred[R == 0] = 0
    return Y_pred

def loss(self, Y, Yhat):

```

```

        error = Y-Yhat
        n_ratings = np.sum(Y != 0)
        loss_value = 1/(2*n_ratings)*np.linalg.norm(error, 'fro')**2 + \
            lamda/2*(np.linalg.norm(self.I, 'fro')**2 + \
                np.linalg.norm(self.U, 'fro')**2)
        return loss_value

    def RMSE(self, Y, Yhat):
        error = Y - Yhat
        n_ratings = np.sum(Y != 0)
        rmse = math.sqrt(np.linalg.norm(error, 'fro')**2/n_ratings)
        return rmse

```

**Xây dựng class MF quản lý model và data:**

```

In [22]: class MF():
        """
        This class used to manage model and data
        Two main arguments:
        data: control the data
        model: control the functions which execute model
        """
        def __init__(self, data, model, n_iter, print_log_iter):
            self.data = data
            self.model = model
            self.n_iter = n_iter
            self.print_log_iter = print_log_iter
            self.Y_pred_train = None
            self.Y_pred_test = None
            self.Yhat = None

        def fit(self):
            for i in range(self.n_iter):
                #update U and I
                self.model.updateU()
                self.model.updateI()
                #calculate Y_hat
                self.Yhat = self.model.pred(self.model.I, self.model.U)
                #calculate Y_pred_train by replace non ratings by 0
                self.Y_pred_train = self.model.pred_train_test(self.Yhat, self.data.Rtrain)
                self.Y_pred_test = self.model.pred_train_test(self.Yhat, self.data.Rtest)
                if i % self.print_log_iter == 0:
                    print('Iteration: {}; RMSE: {}; Loss value: {}'.\
                        format(i, self.model.RMSE(self.data.Ytest, self.Y_pred_test), \
                            self.model.loss(self.data.Ytrain, self.Y_pred_train)))

        def recommend_for_user(self, user_id, k_neighbors):
            recm = np.concatenate((np.arange(1, self.Y_pred_test.shape[0]+1).reshape(-1, 1),

```

```

                                self.Y_pred_test[:, user_id - 1].reshape(-1, 1)), axis
recm.sort(axis = 0)
print('Top %s item_id recommended to user_id %s: %s'%\
      (k_neighbors, user_id, str(recm[-k_neighbors:, 0])))

In [23]: data = Data(dataset = movie_length, split_rate = 2/3)
         model = Model(data = data, theta = 0.75, lamda = 0.1, K = 2)
         mf = MF(data = data, model = model, n_iter = 100, print_log_iter = 10)
         mf.fit()

Iteration: 0; RMSE: 1.2230236685162748; Loss value: 681.8471565701848
Iteration: 10; RMSE: 1.1780718373263668; Loss value: 286.40240426148296
Iteration: 20; RMSE: 1.087117979096145; Loss value: 518.2976746384784
Iteration: 30; RMSE: 1.08441027204107; Loss value: 534.4106319526331
Iteration: 40; RMSE: 1.08385608972465; Loss value: 546.1164972768984
Iteration: 50; RMSE: 1.0828157171073196; Loss value: 559.6753058301773
Iteration: 60; RMSE: 1.079754860855491; Loss value: 575.3332305811103
Iteration: 70; RMSE: 1.0750987045005322; Loss value: 590.821120508004
Iteration: 80; RMSE: 1.0709110406988624; Loss value: 602.5417655852891
Iteration: 90; RMSE: 1.0685236770672915; Loss value: 608.9641965725858

```

Ta nhìn thấy kết quả cả hàm loss function và RMSE giảm dần sau các vòng lặp. Điều này cho thấy thuật toán gradient descent đã phát huy tác dụng trong việc làm giảm sai số. Tuy nhiên đôi khi chúng ta sẽ gặp tình huống loss function và RMSE tăng dần. Có nhiều nguyên nhân dẫn tới điều này chẳng hạn như học suất learning rate và regularization coefficient quá cao làm cho thuật toán nhay ra khi các trọng số học cũng có thể các loss function chỉ tăng tạm thời và giảm sau đó. Khi đi chuyển qua các bài tập tiếp theo. Vì khi nhìn 1 chúng ta cần điều chỉnh lại các học suất learning rate và regularization thuật toán đi chuyển ứng dụng thì nhìn thấy điều này. Khi nhìn thấy 2 không quá nghiêm trọng bởi thuật toán có thể đi chuyển thì nhìn thấy ngay sau đó.

Recommend 10 sản phẩm tìm kiếm cho user\_id = 200:

```
In [24]: mf.recommend_for_user(user_id = 200, k_neighbors = 10)
```

```
Top 10 item_id recommended to user_id 200: [3943 3944 3945 3946 3947 3948 3949 3950 3951 3952]
```

## 0.4 2.3. Tài liệu tham khảo

1. [Recommendation System - Stanford](#)
2. [Collaborative Filtering Youtube - Stanford](#)
3. [Recommendation System - Machine Learning - Andrew Ng](#)
4. [Matrix Factorization - Machine Learning C Bn - Tiep Huu Vu](#)
5. [Matrix Factorization techniques for recommender systems](#)
6. [Learning from Incomplete Ratings Using Non-negative Matrix Factorization](#)
7. [Matrix Factorization - Albert Au Yeung](#)
8. [Algorithms for Non-negative Matrix Factorization - Daniel D.Lee and H. Sebastian Seung](#)