

Recommender system

I. Introduction

Matching customers to products is an important practical problem that many companies would be interested in being able to do well. Products information provided by users can be used to develop recommendation systems to help recommend to customers new products that they may like.

Ratings prediction is a common application of such models and is explored in this project. Starting with a database of users/objects ratings, a matrix factorization model is used to predict customers ratings about products they had not rated yet depending on:

- how they rated other products,
- how other customers had rated any of the products.

A rating dataset(1) from the MovieLens website is used for this project.

It contains 100,000 ratings applied to 9,000 movies by 700 users.

Ratings scores are given between 0.5 and 5.

An histogram of the global distribution of the known ratings independently of the type of movies rated, can be seen on Figure 1.

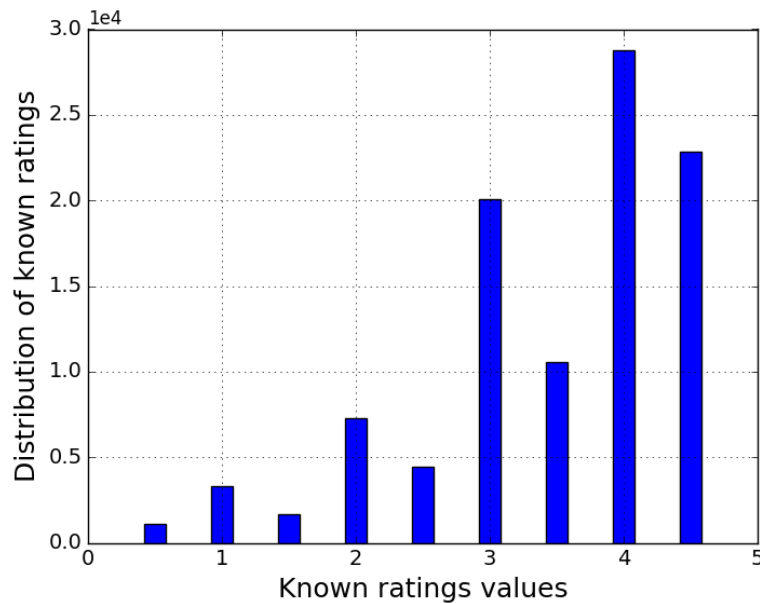


Figure 1: Distribution of known ratings

(1) <https://grouplens.org/datasets/movielens/>

II. Methodology

The ratings matrix characteristic of the recommender system is formed using any possible user/object pair, and the ij^{th} entry in this matrix, m_{ij} , is going to be the rating that the i^{th} user gave to the j^{th} object (Cf. Figure 2).

This matrix will have many missing values, because each user can only rate a fraction of the products. Matrix factorization will try to learn a low-rank factorization of this matrix using only the observed data, while ignoring the data we don't have. From there, it will be possible to fill in all of the missing values that

could be then viewed as predictions for what a user will rate an object and as such, later be used to make recommendations.

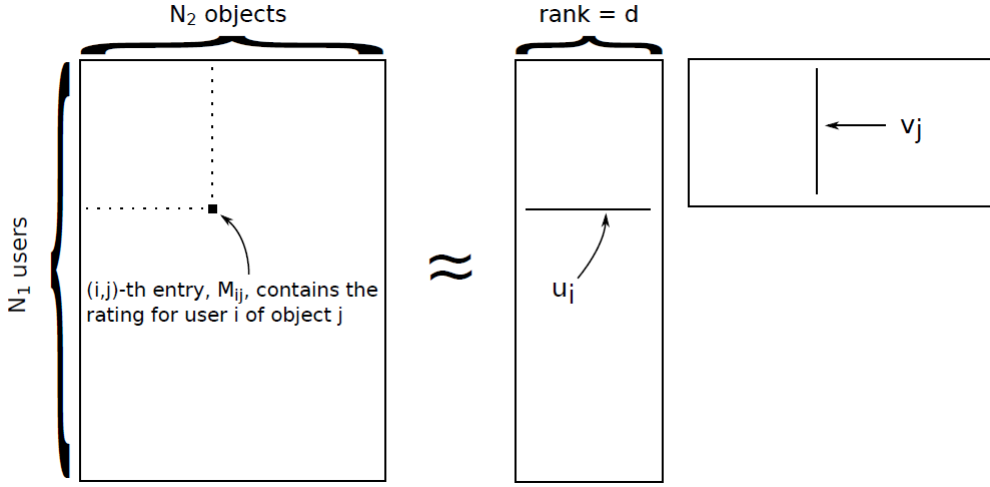


Figure 2: Low-rank matrix factorization

We assume that the ratings matrix M is a N_1 by N_2 matrix of rank d , that can be approximated to the dot product of a N_1 by d U matrix (of objects j rated by users i), times a d by N_2 V matrix (of users i who rated objects j), where $d \ll \min(N_1, N_2)$:

$$M \approx UV^T$$

Probabilistic matrix factorization (PMF) is one particular model for learning a low-rank factorization in this missing data problem, and where we make the following assumptions:

- The location $u_i \in \mathbb{R}^d$ of the N_1 users and the location $v_j \in \mathbb{R}^d$ of the N_2 objects are going to be generated from zero mean Gaussians with spherical covariances: $u_i \approx N(0, \lambda^{-1}I)$ and $v_j \approx N(0, \lambda^{-1}I)$.
- The distribution on the data given those two locations is going to be: $M_{ij} \approx N(u_i^T v_j, \sigma^2)$, where M_{ij} is an observed value if (i, j) is a measured pair of user/object.

Solving this problem would mean to find the matrix $\hat{M} = UV^T$ that minimizes the sum squared distance to the target matrix M . It is equivalent to maximizing the log of the joint likelihood over U and V , which can be written as:

$$\mathcal{L} = - \sum_{(i,j) \in \Omega} \frac{1}{2\sigma^2} \|M_{ij} - u_i^T v_j\|^2 - \sum_{i=1}^{N_1} \frac{\lambda}{2} \|u_i\|^2 - \sum_{j=1}^{N_2} \frac{\lambda}{2} \|v_j\|^2 + \text{Cst}$$

By taking the derivatives of \mathcal{L} over u_i and v_i independently and setting them to zero, it is then possible to solve for each u_i and v_i individually.

1. Data preprocessing and initialization

A very important step before doing any analysis is to check and reindex the user/object database in order to not have any missing index between any consecutive users/objects.

Because it is difficult to make accurate predictions for users who made very few ratings, any users with less than a certain number of ratings were removed from the database before analysis. Similarly, objects with less than a given number of ratings were removed from the database.

The filtered dataset was split into two sets, a training set used to build the prediction model, and a

testing set containing 10% of the original dataset and that will be used to assess performance. The objects matrix V is initialized as a zero mean Gaussian $v_j \approx N(0, \lambda^{-1}I)$.

2. Iteration process

For each iteration, matrices U then V are respectively calculated and updated in two consecutive steps until convergence of the loss function \mathcal{L} :

- Step 1: Update users locations for $i = 1, \dots, N_1$

$$u_i = \left(\lambda \sigma^2 I + \sum_{j \in \Omega_{u_i}} v_j v_j^\top \right)^{-1} \left(\sum_{j \in \Omega_{u_i}} M_{ij} v_j \right)$$

- Step 2: Update objects locations for $j = 1, \dots, N_2$

$$v_j = \left(\lambda \sigma^2 I + \sum_{i \in \Omega_{v_j}} u_i u_i^\top \right)^{-1} \left(\sum_{i \in \Omega_{v_j}} M_{ij} u_i \right)$$

III. Results

Many different parameters were tested during analysis, but the results presented in this document used the following parameters:

- Number of iterations: 50
- Shape parameter: $\lambda = 2$
- Variance: $\sigma^2 = 0.1$
- The algorithm was used to learn 5 dimensions: $d = 5$.
- Minimum number of rated objects per user: 10.
- Minimum number of ratings per object: 3.

As seen on Figure 3, the objective function was shown to converge towards a maximum.

The ratings matrix M was calculated as the dot product between matrices U and V . The global distribution of the ratings can be seen on Figure 4.

$R < -2$	$-2 \leq R < 0$	$0 \leq R \leq 5$	$5 < R \leq 7$	$7 < R$
1.66%	8.28%	86.93%	2.58%	0.54%

Table 1: Frequencies of predicted ratings R

The first observation we can make from those results is that 13.06% of the ratings are predicted outside the range of valid ratings values.

Furthermore, the root mean square error was found to be equal to 1.03.

Improvements need to be brought to the model in order to restrain predictions to the valid ratings range, as well as decrease the RMSE by introducing for instance regularizing parameters.

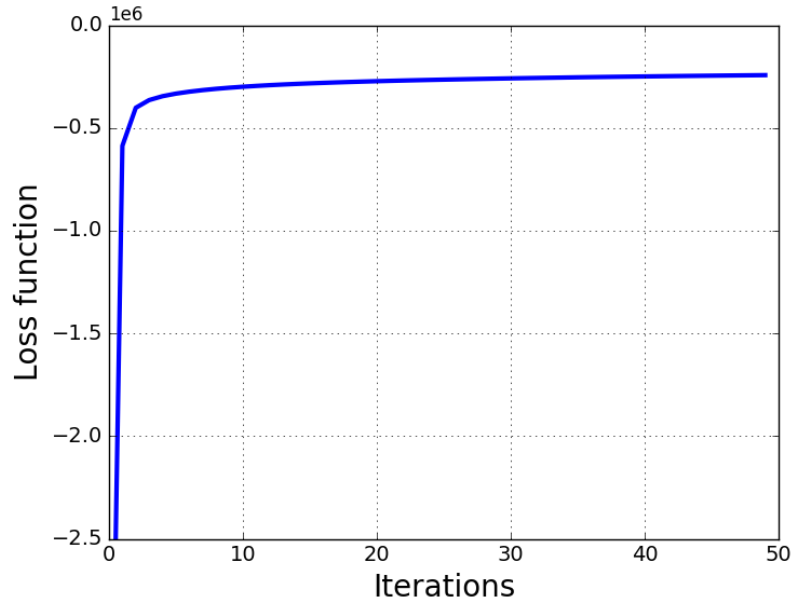


Figure 3: Objective function \mathcal{L} over iteration time

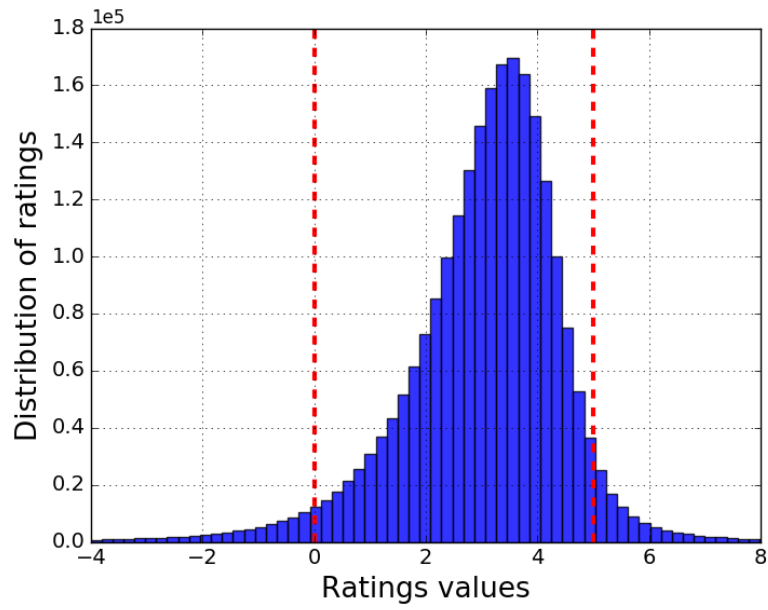


Figure 4: Ratings distribution

IV. Python Code

```

1  from __future__ import division
2  import numpy as np
3  import sys
4  import pandas as pd
5  from numpy.random import multivariate_normal as MvN
6  from matplotlib import pyplot as plt
7
8  """
9  This file has to be run directly from the terminal using the following command li
10     > python PMF.py "ratings.csv"

```

```

11 where 'hw4_PMF.py' is the name of python file and "ratings.csv" contains the ratings
12 """
13
14 try:
15     assert len(sys.argv) > 1, "missing input argument"
16     train_data = np.genfromtxt(sys.argv[1], delimiter=",")
17 except Exception as e:
18     print(e)
19     exit()
20
21 print('data_loaded')
22
23 # Parameters of the PMF objective function to maximize
24 lam = 2          # lambda shape parameter
25 sigma2 = 0.1     # variance
26 d = 5            # rank
27 iterations = 50  # number of iterations
28
29
30 LD = len(train_data)    # nb of data points
31
32 # Load data as dataframe
33 data = pd.DataFrame(train_data)
34
35 # Get index set of objects omega_ui rated by user i
36 omega_ui = {}
37 set_ui = set()
38 for user, obj in zip(data[0], data[1]):
39     if user not in set_ui:
40         set_ui.add(int(user))
41         omega_ui[int(user)] = [int(obj)]
42     else:
43         omega_ui[int(user)].append(int(obj))
44 N1 = len(omega_ui.keys())
45
46 # Get index set of users omega_vj who rated object j
47 omega_vj = {}
48 set_vj = set()
49 for user, obj in zip(data[0], data[1]):
50     if obj not in set_vj:
51         set_vj.add(int(obj))
52         omega_vj[int(obj)] = [int(user)]
53     else:
54         omega_vj[int(obj)].append(int(user))
55 N2 = len(omega_vj.keys())
56
57 # create dictionary with known data
58 Mij = {}
59 set_Mij = set()
60 for user, obj, mij in zip(data[0], data[1], data[2]):
61     Mij[(int(user), int(obj))] = mij
62
63 # Initialisation of vj

```

```

64 def init_V():
65     V = MvN(mean=np.zeros((N2)), cov=np.identity(N2) * 1 / lam, size = d)
66     return V
67
68 # update ui
69 def update_ui(V,i):
70     i += 1
71     ui = np.zeros((5,1))
72     VT = V.T
73
74     # first sum
75     sum1 = 0
76     for j in omega_ui[i]:
77         vjT = VT[j-1]
78         vj = vjT.T
79         sum1 += np.outer(vj,vjT)
80     sum1 += lam*sigma2*np.identity(5)
81
82     # second sum
83     sum2 = 0
84     for j in omega_ui[i]:
85         vjT = VT[j-1]
86         vj = vjT.T
87         sum2 += Mij[(i, j)]*vj
88     a = np.zeros((d,d))
89     for k in range(d):
90         a[k][0] = sum2[k]
91     sum2 = a
92
93     # product
94     prod = np.dot(np.linalg.inv(sum1), sum2)
95     for k in range(d):
96         ui[k] = prod[k][0]
97
98     return ui.reshape(-1)
99
100 # update vj
101 def update_vj(U,j):
102     j += 1
103     vj = np.zeros((5,1))
104     # first sum
105     sum1 = 0
106     for i in omega_vj[j]:
107         ui = U[i-1]
108         uiT = ui.T
109         sum1 += np.outer(ui,uiT)
110     sum1 += lam*sigma2*np.identity(5)
111
112     # second sum
113     sum2 = 0
114     for i in omega_vj[j]:
115         ui = U[i-1]
116         sum2 += Mij[(i, j)]*ui

```

```

117     a = np.zeros((d,d))
118     for k in range(d):
119         a[k][0] = sum2[k]
120     sum2 = a
121
122     # product
123     prod = np.dot(np.linalg.inv(sum1), sum2)
124     for k in range(d):
125         vj[k] = prod[k][0]
126
127     return vj.reshape(-1)
128
129 # Calculate loss function
130 def loss(U,V):
131     # first sum
132     sum1 = 0
133     VT = V.T
134     for ij in Mij.keys():
135         i = ij[0]
136         j = ij[1]
137         vjT = VT[j-1]
138         sum0 = Mij[ij] - np.dot(U[i-1], vjT)
139         sum1 += (sum0)**2
140     sum1 *= 1/(2*sigma2)
141
142     # second sum
143     sum2 = 0
144     for i in range(N1):
145         sum2 += np.linalg.norm(U[i], ord=2)**2
146     sum2 *= lam / 2
147
148     # third sum
149     sum3 = 0
150     for j in range(N2):
151         sum3 += np.linalg.norm(VT[j], ord=2)**2
152     sum3 *= lam / 2
153
154     # global sum
155     L = - sum1 - sum2 - sum3
156
157     return L
158
159 def PMF(V):
160     #Initialisation
161     VT = np.zeros((N2,d))
162     U = np.zeros((N1,d))
163
164     # update ui
165     for i in range(N1):
166         U[i] = update_ui(V, i)
167
168
169     # update vj

```

```

170     for j in range(N2):
171         VT[j] = update_vj(U, j)
172     V = VT.T
173
174     # Calculate loss function
175     L = loss(U, V)
176
177     return L, U, V
178
179
180 V = init_V()
181 L_save = []
182 for iteri in range(iterations):
183     # print(iteri)
184     L, U, V = PMF(V)
185     VT=V.T
186     L_save.append(L)
187     if iteri == 9:
188         np.savetxt("U-10.csv", U, delimiter=",")
189         np.savetxt("V-10.csv", VT, delimiter=",")
190     if iteri == 24:
191         np.savetxt("U-25.csv", U, delimiter=",")
192         np.savetxt("V-25.csv", VT, delimiter=",")
193     if iteri == 49:
194         np.savetxt("U-50.csv", U, delimiter=",")
195         np.savetxt("V-50.csv", VT, delimiter=",")
196
197 np.savetxt("objective.csv", L_save, delimiter=",")
198
199
200 # Plot objective function
201 plt.plot(L_save)
202 plt.grid()
203 plt.xlabel('Iterations')
204 plt.ylabel('Loss function')
205 plt.show()

```