

---

# Task 1 (visit prediction)

(Kaggle) user name: yfzhou13      display name: Vidal

## Method

We decide to apply One Class Recommendation model to this prediction.

The reason is that in this task, what need to predict is only visit or non-visit. It is quite like matrix factorization problem where all the elements are zero or one instead of other numerical feedback. So we define positive feedback and negative feedback as follows:

- Positive feedback:  $P_u = \{i\}$  is defined as the set of items  $i$  towards which user  $u$  shows explicit positive feedback.
- Negative feedback:  $N_u = \{j\}$  is a set of items  $j$  towards which user does not shown positive feedback until now.

In this particular dataset, each record of user visiting business is seen as positive feedback. Then We assume that users prefer positive feedback business to those negative ones, which is captured by such an objective function:

$$\max_{u,v} \sum_{u,i,j} \ln \sigma(u_u v_i - u_u v_j)$$

Here  $u_u v_i$  and  $u_u v_j$  stands for user  $u$ 's preferences towards business  $i$  and  $j$ , built by a latent factor model.

When conducting optimization, the objective actually changes by adding the regularization component.

$$\max_{u,v} \sum_{u,i,j} \ln \sigma(u_u v_i - u_u v_j) - \sum_u ||u_u||_2^2 - \sum_i ||v_i||_2^2$$

Since the number of entire possible pairs of all positive feedback and negative feedback for every user are extremely large, we optimize it based on stochastic gradient ascent. Specifically, for one user, some of its positive and negative feedback pairs are chosen, which contribute to gradients updating in one iteration.

## Implementation

Particularly, here are some equations for gradient ascent updating.

$$\begin{aligned} z &= \sigma(u_u v_i - u_u v_j) \\ \frac{\partial L}{\partial u_u} &= \sum_{(i,j)} (1 - z)(v_i - v_j) - \lambda u_u \\ \frac{\partial L}{\partial v_i} &= \sum_{(u)} (1 - z)(u_u) - \lambda v_i \\ \frac{\partial L}{\partial v_j} &= \sum_{(u)} (1 - z)(-u_u) - \lambda v_j \end{aligned}$$

To monitor what's going on during the training, we output the value of original objective function, regularized objective function, and also performance on validation set after several iterations (Figure 1).

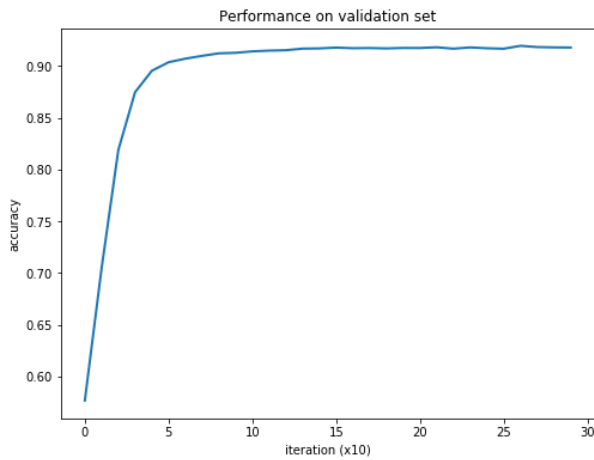


Figure 1

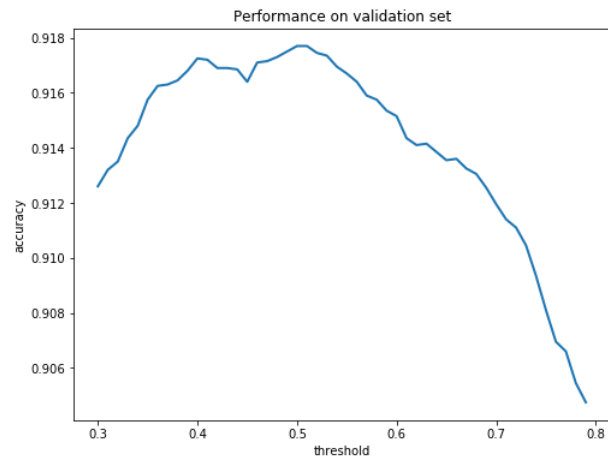


Figure 2

*Observation:* Training is more efficient if users are covered as many as possible even at the cost of reducing sampling pairs of one user.

It is not difficult to guess that the user-business matrix is quite sparse, in which there are only a few positive label for almost all users. This means if we choose too many positive/negative pairs for one user, it is highly possible that lots of positive records will be repeated more than once. So during gradient updating, it seems more necessary to satisfy 'breath' requirement instead of 'depth'.

## Parameters

There are also some parameters we can control, like reduction dimension and regularization parameter. In addition to regularization parameter. For example, in evaluating, how to choose the threshold which determines visit or non-visit.

### ***Reduction dimension $K$ and regularization $\lambda$***

Technically, as  $K$  becomes larger, the performance should be better in training, validation and test set. Practically, we find that, for example when  $K$  is chosen from 100, 200, 400, 800, the training and validation performance do not show huge difference. However, on test dataset, larger  $K$  reveals stronger ability of generalization obviously. Generalization here points to a model trying to reduce overfitting on training and validation dataset, and keep the prediction ability as accurate as possible. Finally we maintain  $K = 800$ .

Following the pipeline of choosing optimal  $\lambda$ , final  $\lambda = 0.6$

### ***Threshold***

(a) Threshold as a number.

It's hard to estimate a threshold only with prior knowledge. Following the similar idea of choosing parameter lambda, we try a bunch of numbers and find the best.

---

(b) Threshold as a function. (Figure 2)

Different users have distinct preference profiles, thus their threshold might be different naturally. Here we model the threshold as a function, to user activity.

Intuition behind it is that, more active a user is, more business it will prefer. So lower threshold should be modified on highly active users. For example, such function is defined:

$$threshold(u) = upbound - \left( count(u)^\alpha \frac{1}{scale} \right)$$

$count$  is the number of business a user visit before (in training set), representing users' activity. And it is a constant number for one user which can be computed easily.

$upbound$ ,  $scale$ ,  $\alpha$  is hyper parameter we need to tune.

The result shows that threshold function below improves the performance (accuracy) by about 0.3% compared to 91.8% under optimal fixed threshold.

$$threshold(u) = 0.56 - \left( \frac{count(u)^{0.9}}{171} \right)$$

## Results

We divide data into training set and validation set, with quantity of 99% and 1%.

### ***On public leaderboard (seen portion of test data)***

Accuracy	K=100	K=200	K=400	K=800
Fixed threshold	0.91665	0.91705	0.91820	0.91825
Dynamic threshold	0.91790	0.91855	0.92050	0.92170

### ***On private leaderboard (unseen portion of test data)***

Finally, the best accuracy is obtained with 0.92095, under dynamic threshold case of  $K = 800$ ,  $\lambda = 0.6$ . What's more, there is a little difference between seen and unseen data when  $K$  gets larger, while performance dramatically drops if  $K$  is small. I think the reason may be that smaller  $K$  leads to a not bad accuracy on training and validation set but it over fits the data.

---

## Task 2 (rating prediction)

(Kaggle) user name: yfzhou13      display name: Vidal

### Method

We apply latent-factor model to this rating prediction task. It is an idea of dimensional reduction.

$$R = UV^T$$

Based on the simple model in homework 3, which only includes  $\alpha$  and  $\beta$ , the entire latent-factor model is describe as follow.

$$f(u, i) = \alpha + \beta_u + \beta_i + \gamma_u \gamma_i$$

In this model,  $\alpha$  is a constant.  $\beta_u$  and  $\beta_i$  represents bias on different users and items respectively.  $\gamma_u \gamma_i$  is a multiplication of 1 by  $K$  vector and  $K$  by 1 vector, showing preference of this user towards this item.

Thus optimization objective is square error plus regularization.

$$\min_{\alpha, \beta, \gamma} \sum_{u,i} (f(u, i) - R_{u,i})^2 + \lambda [\sum_u \beta_u^2 + \sum_i \beta_i^2 + \sum_u \|\gamma_u\|_2^2 + \sum_i \|\gamma_i\|_2^2]$$

### Implementation

Basically, we train these unknown vectors or matrices on training dataset, then tune parameter  $\lambda$  on validation set. Because the optimization objective is not convex, we can just solve it approximately. Here are particular steps for training.

- 1. Initialize  $\alpha$ ,  $\beta$ ,  $\gamma$  randomly.
- 2. Find the solution for  $\alpha$ ,  $\beta$ , using derivatives updating until convergence.
- 3. Keep  $\alpha$ ,  $\beta$  unchanged, compute gradient of  $\gamma_u$  and  $\gamma_i$ .
- 4. Update  $\gamma_u$  and  $\gamma_i$  with the gradient computed.
- 5. Repeat step 3 and step 4.

In step 2, let derivative equal to 0.

$$\begin{aligned}\alpha &= \frac{\sum_{u,i} (R_{u,i} - \beta_u - \beta_i)}{N_{train}} \\ \beta_u &= \frac{\sum_{i \in I_u} (R_{u,i} - \alpha - \beta_i)}{\lambda + |I_u|} \\ \beta_i &= \frac{\sum_{u \in U_i} (R_{u,i} - \alpha - \beta_u)}{\lambda + |U_i|}\end{aligned}$$

In step 3, compute gradients as follows.

$$\begin{aligned}\frac{\partial L}{\partial \gamma_u} &= \sum_{i \in I_u} (\alpha + \beta_i + \beta_u + \gamma_u \gamma_i - R_{u,i}) \gamma_i + \lambda \gamma_u \\ \frac{\partial L}{\partial \gamma_i} &= \sum_{u \in U_i} (\alpha + \beta_i + \beta_u + \gamma_u \gamma_i - R_{u,i}) \gamma_u + \lambda \gamma_i\end{aligned}$$

---

## Results

We divide data into training set and validation set, with quantity of 99% and 1%.

### ***On public leaderboard (seen portion of test data)***

Only with  $\alpha$  and  $\beta$ , we obtain accuracy of 0.74583 on test data under optimal  $\lambda = 4.1$ . By adding  $\gamma$ , our results become more accurate.

RMSE	Lambda=5	Lambda=10
K=1	0.74572	0.74546
K=2	0.74566	0.74534
K=5	0.74564	0.74548

### ***On private leaderboard (unseen portion of test data)***

Best performance is obtained at  $\lambda = 10.0$ , with RMSE of 0.75746. This is a considerable drop compared to another part of test data. The reason might be iteration times is not enough to train a model which avoid serious overfitting.

---

## Appendix:

### 1. Visit prediction: one-class

```
161 def OneClass(lam, K, learnRate, max_iter):
162     gammaU = np.random.rand(Lu, K)/1 - 0.5
163     gammaB = np.random.rand(K, Lb)/1 - 0.5
164     accRec = []
165     for it in range(max_iter):
166         objective = 0
167         regularization = 0
168         gu = np.zeros((Lu, K))
169         gb = np.zeros((K, Lb))
170
171         for (user, busi) in trainPair:
172             u = userDict[user]
173             i = busiDict[busi]
174             j = busiDict[findNegBusi(user)]
175             z = sigmoid(np.dot(gammaU[u,:],gammaB[:,i]) - \
176                         np.dot(gammaU[u,:],gammaB[:,j]))
177             gu[u,:] += (1-z)*(gammaB[:,i]-gammaB[:,j])
178             gb[:,i] += (1-z)*(gammaU[u,:])
179             gb[:,j] += (1-z)*(-gammaU[u,:])
180
181             objective += log(z)
182
183         gu -= lam*gammaU
184         gb -= lam*gammaB
185         gammaU += learnRate*gu
186         gammaB += learnRate*gb
187
188         regularization = objective - lam*np.sum(np.square(gammaU)) - \
189                               lam*np.sum(np.square(gammaB))
190         if (it+1)%2==0:
191             print 'iteration: ' + str(it+1) + '\t' + str(regularization) \
192                   + '\t' + str(objective)
193         if (it+1)%10==0:
194             accCur = max([Accuracy(gammaU,gammaB,t) for t in \
195                           [0.3,0.4,0.5,0.6,0.7]])
196             accRec.append(accCur)
197             print accCur
198
199     return gammaU, gammaB, accRec
```

### 2. Rating prediction: latent-factor model

```

81# latent factor model (alpha, beta)
82def LFM1(lam, max_iter):
83    alpha = 0.1
84    betaU = {i:0.0 for i in uniUser}
85    betaB = {i:0.0 for i in uniBusi}
86    for it in range(max_iter):
87        alpha = 0
88        loss = 0
89        for (user, busi, rating) in trainPair:
90            alpha += (rating - betaU[user] - betaB[busi])/len(trainPair)
91        for (user, busi, rating) in trainPair:
92            diff = alpha + betaU[user] + betaB[busi] - rating
93            loss += diff**2
94        squareError = loss
95        for (user, busi_list) in userBusi_train.items():
96            betaU[user] = 0
97            for (busi,rating) in busi_list:
98                betaU[user] += (rating - alpha - betaB[busi])/(lam + \
99                    len(busi_list))
100            loss += lam*(betaU[user]**2)
101        for (busi, user_list) in busiUser_train.items():
102            betaB[busi] = 0
103            for (user,rating) in user_list:
104                betaB[busi] += (rating - alpha - betaU[user])/(lam + \
105                    len(user_list))
106            loss += lam*(betaB[busi]**2)
107        if (it+1)%5==0:
108            print 'iteration: ' + str(it+1)
109            print loss, squareError
110            meanSquaredError(validPair,[alpha,betaU,betaB])
111    return alpha, betaU, betaB

```

```

161# latent factor model (alpha, beta, gamma)
162def LFM3(lam, K, learnRate, max_iter):
163    alpha, betaU, betaB = LFM1(lam = 4.0, max_iter = 30)
164    errList = []
165    gammaU = {i:np.random.rand(K)/10-0.05 for i in uniUser}
166    gammaB = {i:np.random.rand(K)/10-0.05 for i in uniBusi}
167    for it in range(max_iter):
168        loss = 0
169
170        gra_gammaU = {user:lam*gammaU[user] for user in uniUser}
171        gra_gammaB = {busi:lam*gammaB[busi] for busi in uniBusi}
172
173        for (user, busi, rating) in trainPair:
174            diff = alpha + betaU[user] + betaB[busi] + \
175                np.inner(gammaU[user],gammaB[busi]) - rating
176            gra_gammaU[user] += gammaB[busi]*diff
177            gra_gammaB[busi] += gammaU[user]*diff
178            loss += diff**2
179        squareError = loss
180        for (user, busi_list) in userBusi_train.items():
181            gammaU[user] -= learnRate*gra_gammaU[user]
182            loss += lam*(betaU[user]**2+sum(gammaU[user]**2))
183        for (busi, user_list) in busiUser_train.items():
184            gammaB[busi] -= learnRate*gra_gammaB[busi]
185            loss += lam*(betaB[busi]**2+sum(gammaB[busi]**2))
186        if (it+1)%10==0:
187            print 'iteration: ' + str(it+1)
188            print loss, squareError
189            err = meanSquaredError(validPair,[alpha,betaU,betaB,gammaU,gammaB])
190            errList.append(err)
191    return alpha,betaU,betaB,gammaU,gammaB,errList

```