

# 统计信号处理大作业

## 实验报告——应用最小二乘求解协同滤波问题

李 溢 2016011235

2019 年 5 月 23 日

# 目录

1	问题描述	2
2	问题建模	2
3	求解算法	2
3.1	ALS 交替最小二乘	2
3.1.1	优化目标	2
3.1.2	最小二乘闭式解	3
3.1.3	迭代求解	3
3.1.4	实验结果	3
3.2	基于梯度下降的协同滤波	4
3.2.1	前述 ALS 算法的问题	4
3.2.2	基于梯度下降的优化算法	5
3.2.3	梯度下降求解	5
3.2.4	实验结果	5
4	附录：源代码	7
4.1	ALS 算法	7
4.2	梯度下降算法	10

# 1 问题描述

关键词：协同滤波，填充稀疏矩阵

问题的经典背景为：假设有 1000 个用户和 200 部电影，用户  $i$  对电影  $j$  的打分为  $M_{ij}$ ，因此形成了一个  $100 \times 200$  的评分矩阵  $M$ ，由于每个用户不一定对所有电影都打过分，因此矩阵  $M$  不会被完全填充。所要解决的问题就是填充矩阵。

# 2 问题建模

实验中已知稀疏矩阵  $M \in \mathbb{R}^{N \times S}$ ，假设满足  $M = UV^T$ ，其中  $M_{ij} = U(i,:)V(j,:)^T$ ， $U \in \mathbb{R}^{N \times rank}$ ， $V \in \mathbb{R}^{S \times rank}$ ， $rank \ll N, S$  为一个常数，可近似看作  $U$ 、 $V$  的秩。假设矩阵  $M$  低秩，由于矩阵的秩为矩阵奇异值向量的 0 范数，可将其放缩为 1 范数，即矩阵的核范数。矩阵的核范数满足

$$\|X\|_* = \min_{U, V | X=UV^T} \frac{1}{2} (\|U\|_F^2 + \|V\|_F^2)$$

因此将上述问题转化为优化以下目标函数进行求解

$$\min_{U, V} \|W * (M - UV^T)\|_F^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2)$$

其中  $\lambda$  是控制矩阵低秩程度的超参数， $W$  是标志矩阵， $W(i, j) = 1$  表示  $M_{ij} \neq 0$ ， $*$  表示矩阵对应元素相乘。

# 3 求解算法

## 3.1 ALS 交替最小二乘

### 3.1.1 优化目标

考虑优化如下目标函数

$$L(U, V) = \sum_{M_{xi} \neq 0} (M_{xi} - U_{x:} V_i^T)^2 + \lambda (|U_x|^2 + |V_i|^2)$$

其中  $U_{x:}$  表示  $U$  的第  $x$  行， $V_i$  表示  $V$  的第  $i$  行， $M_{xi}$  表示用户  $x$  对用户  $i$  的评分。最后一项为正则化项。

### 3.1.2 最小二乘闭式解

先对  $U_{x:}$  应用最小二乘

$$\frac{\partial L}{\partial U_{x:}} = 2V^T (VU_x^T - M_x^T) + 2\lambda U_x^T \quad (1)$$

令其等于 0, 求出

$$U_x^T = (V^T V + \lambda I)^{-1} V^T M_{x:}^T \quad (2)$$

同理, 对  $V_{i:}$  应用最小二乘, 得到

$$\frac{\partial L}{\partial V_i} = 2U^T (UV_i^T - M_i) + 2\lambda V_{i:}^T = 0 \quad (3)$$

$$V_i^T = (U^T U + \lambda I)^{-1} U^T M_{:i} \quad (4)$$

### 3.1.3 迭代求解

每次先固定  $V$  求解  $U$ , 再固定  $U$  求解  $V$ , 这样交替进行优化, 算法伪代码如下

---

**Algorithm 1** ALS 协同滤波

---

**Input:**  $M \in \mathbb{R}^{N \times S}$ ,  $\lambda$ ,

**Output:**  $X = UV^T$ ,  $U \in \mathbb{R}^{N \times rank}$ ,  $V \in \mathbb{R}^{S \times rank}$ ,  $rank \ll N, S$

random initialization of  $U$  and  $V$

**for** step < max step **do**

$$U^T = (V^T V + \lambda I)^{-1} V^T M^T$$

$$V^T = (U^T U + \lambda I)^{-1} U^T M$$

**end for**

---

### 3.1.4 实验结果

算法收敛时间大概为 1 min, 超参数主要是特征的数量  $rank$  和  $\lambda$ 。

首先我固定  $\lambda = 10$ , 尝试了不同的  $rank$ , 计算在测试集上的 mse

$rank(\lambda=10)$	10	30	100
train mse	88.2968	71.5975	40.3695
test mse	98.635	105.2446	142.0029

表 1:  $rank$  对 ALS 算法结果的影响

由此看来并不是特征的数量 rank 越大越好, rank 过大时很可能在训练集上产生了过拟合。

然后我固定 rank=10, 取不同的  $\lambda$  值, 计算在测试集上的 mse

$\lambda(\text{rank}=10)$	0.1	1	10	100	500	1000
mse	96.824	96.8548	97.3392	102.6232	131.8353	154.6974

表 2:  $\lambda$  对 ALS 算法结果的影响

因此  $\lambda \leq 100$  时的值对优化实际影响并不大, 增大后优化效果下降。查看后我发现  $\lambda = 100$  时  $U^T U$  与  $V^T V$  对角线元素的平均值大约在 350, 此时式 2 和 4 中  $\lambda I$  一项对于 U、V 的解影响不大。当  $\lambda = 1000$  时, 其对角元素平均值大约在 100。因此可以确认  $\lambda$  影响了 U、V 的低秩程度, 当  $\lambda$  过大时, U、V 的秩偏少, 不足以表征 M 对应的用户和电影的特点, 因此优化效果变差。

注: 实际实验中, 最大迭代次数设为 1000。迭代终止条件除了最大的迭代次数, 我还增加了一个 loss 的变化幅度, 当  $\text{abs}(\text{lossLast} - \text{lossNew})/\text{lossNew} < \epsilon$  时也会停止迭代。其中  $\epsilon$  也是引入的另一个超参数, 不过由于实际实验中迭代时 U 和 V 的更新幅度越到后来越小, 因此这一参数的影响并不很大, 实验中我取  $\epsilon = 1e - 7$

## 3.2 基于梯度下降的协同滤波

### 3.2.1 前述 ALS 算法的问题

根据上面基于 ALS 算法的结果, 可以看到测试集上的 mse 数值在 100 左右, 也就意味着平均预测值与真值的误差约为 10, 而经过计算, 实际测试数据的均值约为 13, 这意味着实际上 ALS 算法得到的结果中, 在测试数据对应的点上的预测值基本为 0。这并不意外, 因为 ALS 的算法推导中存在一定的问题。

我们首先考虑式 1 和 3, 真正的推导过程应为

$$\frac{\partial L}{\partial U_x} = 2 \sum_{i, M_{xi} \neq 0} V_{i:}^T (V_{i:} U_{x:}^T - M_{xi}) + 2\lambda U_x^T \neq 2V^T (V U_x^T - M_x^T) + 2\lambda U_x^T \quad (5)$$

需要注意的是, 这里只取了  $\{i | M_{xi} \neq 0\}$  这一部分的 i 来计算关于  $U_x$  的梯度, 严格来说不能化成式子 1 中的矩阵表达形式, 因为式 1 中实际上相当于把原来没有观测数据的  $M_{xi}$  完全看作了 0, 也即认为这些点有观测数据, 但是观测数据就等于 0, 由此优化出来的 U、V 必然是不对的, 因为优化之后在这些点上的预测时会很接近于 0。

### 3.2.2 基于梯度下降的优化算法

仍旧考虑优化以下目标函数进行求解

$$L(U, V) = \min_{U, V} \|\mathbf{W} * (\mathbf{M} - \mathbf{UV}^T)\|_F^2 + \frac{\lambda}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \quad (6)$$

而前述 ALS 的问题实际上就在于计算梯度的过程中忽略了标志矩阵  $\mathbf{W}$  及对应的点乘操作。如果保留矩阵  $\mathbf{W}$ ，将式 6 分别对  $\mathbf{U}$  和  $\mathbf{V}$  求导，得到

$$\frac{\partial L}{\partial \mathbf{U}} = 2(\mathbf{W} * (\mathbf{UV}^T - \mathbf{M})) \mathbf{V} + \lambda \mathbf{U} \quad (7)$$

$$\frac{\partial L}{\partial \mathbf{V}} = 2\mathbf{U}^T (\mathbf{W} * (\mathbf{UV}^T - \mathbf{M})) + \lambda \mathbf{V} \quad (8)$$

很容易验证，如果将上式中的  $\mathbf{W}$  忽略后，就与前面推导的式 1、3 等价，但是实际上  $\mathbf{W}$  并不能忽略，也不能直接分配到括号里边。上面两个式子并不能求出关于  $\mathbf{U}$ 、 $\mathbf{V}$  的闭式解，因此不能直接应用最小二乘法，可以考虑应用梯度下降法优化目标函数。

### 3.2.3 梯度下降求解

每次基于  $L$  对  $\mathbf{U}$  和  $\mathbf{V}$  的梯度分别更新  $\mathbf{U}$  和  $\mathbf{V}$ ，交替进行优化，算法伪代码如下

---

**Algorithm 2** 基于梯度下降的协同滤波

---

**Input:**  $\mathbf{M} \in \mathbb{R}^{N \times S}$ ,  $\lambda$ ,  $lr$ (learning rate)

**Output:**  $\mathbf{X} = \mathbf{UV}^T$ ,  $\mathbf{U} \in \mathbb{R}^{N \times rank}$ ,  $\mathbf{V} \in \mathbb{R}^{S \times rank}$ ,  $rank \ll N, S$

random initialization of  $\mathbf{U}$  and  $\mathbf{V}$

**for** step < max step **do**

$$\frac{\partial L}{\partial \mathbf{U}} = 2(\mathbf{W} * (\mathbf{UV}^T - \mathbf{M})) \mathbf{V} + \lambda \mathbf{U}$$

$$\mathbf{U} = \mathbf{U} - lr * \frac{\partial L}{\partial \mathbf{U}}$$

$$\frac{\partial L}{\partial \mathbf{V}} = 2\mathbf{U}^T (\mathbf{W} * (\mathbf{UV}^T - \mathbf{M})) + \lambda \mathbf{V}$$

$$\mathbf{V} = \mathbf{V} - lr * \frac{\partial L}{\partial \mathbf{V}}$$

**end for**

---

### 3.2.4 实验结果

与 ALS 相似，特征数目  $rank$  过大时会出现严重的过拟合，因此最终取  $rank=10$ ,  $\lambda=10$ ，梯度下降法中主要影响因素为学习率  $lr$  的调整。总的来说，实验中  $lr>0.0002$  时基本很快就会爆炸，即不收敛； $lr$  过小时优化速度太慢。最终我采用了阶梯衰减的学习率调整方案

step	0-12	13-30	31-240	241-1000
lr	0.0002	0.00015	0.00013	0.00011

表 3: 学习率衰减

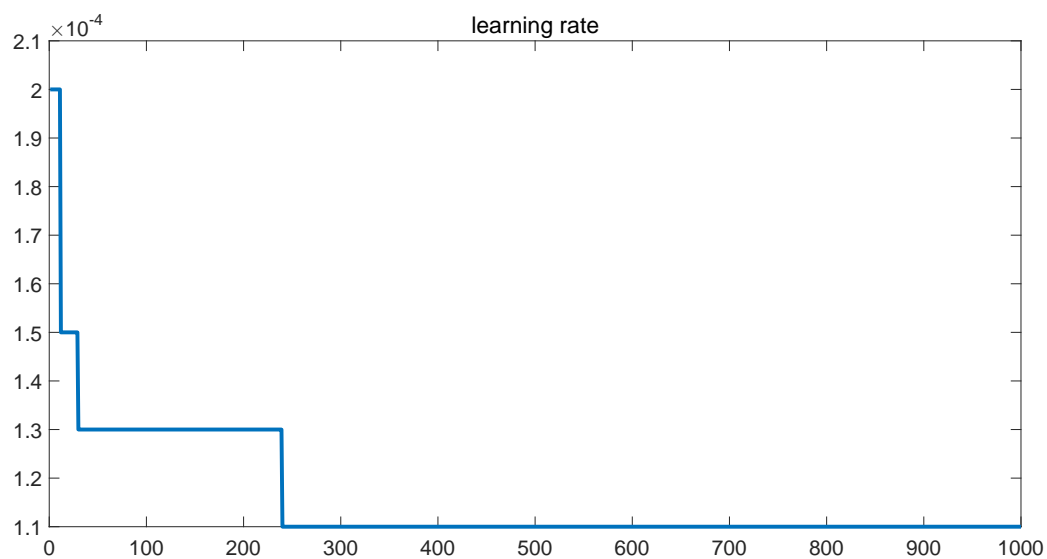


图 1: 学习率衰减方案

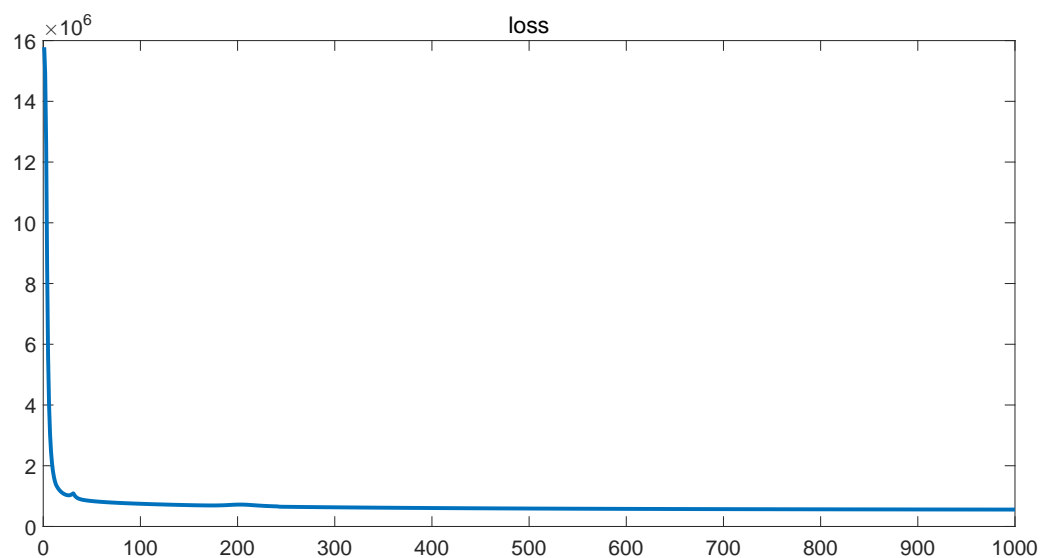


图 2: 梯度下降

实验中任何一个超参数的改变都需要对学习率衰减进行调整，当  $\lambda$  或 rank 较大时就需要调小 lr，因为此时梯度值较大，如果 lr 太大就会不收敛，衰减的区间也需要对应调整，最终的结果基本上是我针对  $\lambda$ 、rank、lr 等参数精细调整过的，也是实验中取得的最好的效果。最终效果如下，可以看到 mse 只有 8，平均每个观测数据的误差为 3 左右，而观测数据的真实平均值约为 13 左右，因此相对较准确，相比于之前的 ALS 算法则有很大的提升。

rank	$\lambda$	train mse	test mse	收敛时间
10	10	4.4583	8.8111	$\sim 1$ min

表 4: 梯度下降算法结果

除此之外，与上面相同，迭代终止条件除了最大迭代次数，我还增加了一个 loss 减小的幅度  $\epsilon = 1e - 5$ 。

## 4 附录：源代码

注：由于编码问题，部分中文注释存在乱码

### 4.1 ALS 算法

```

1 % » ALS  $\mu$  燵
2 clear all; close all; clc;
3 load data_train.mat;
4 %% hyperparameters
5 ranks = [10,30,100];
6 lambdas = [0.1,1,10,100,500,1000]; %  $\lambda^2 \gg \mu$   $\lambda^2 \gg 100\mu^{-1}$ 
7 epsilon = 1e-7;%0.00001; % limit of iteration
8 %% randomly select train data and test data
9 %U = randn(943, rank);
10 %V = randn(1682, rank);
11 idx = randperm(90000,80000);
12 train = data_train(idx,:);
13 temp = 1:90000;
14 temp(idx) = 0;
15 temp(temp==0) = [];
16 test = data_train(temp,:);

```



```

17 % for train
18 M = zeros(943, 1682);
19 W = zeros(943, 1682);
20 M(sub2ind(size(M), train(:,1), train(:,2))) = train(:,3);
21 W(sub2ind(size(W), train(:,1), train(:,2))) = 1;
22 W0 = W;
23 W0(sub2ind(size(W0), test(:,1), test(:,2))) = 1; % all data
24 % for test
25 Mtest = zeros(943, 1682);
26 Wtest = zeros(943, 1682);
27 Mtest(sub2ind(size(M), test(:,1), test(:,2))) = test(:,3);
28 Wtest(sub2ind(size(W), test(:,1), test(:,2))) = 1;
29 s = size(test,1);
30 bestmse=1e10;
31 %% train
32 for i=1:length(ranks)
33     rank = ranks(i);
34     lambda = 10;%lambdas(i);
35     step = 0;
36     U = randn(943, rank);
37     V = randn(1682, rank);
38     I = eye(rank);
39     lossLast = 0;
40     lossNew = norm(W.*(M-U*V'), 'fro').^2 + ...
        lambda/2*(norm(U, 'fro').^2 + norm(V, 'fro').^2);
41     disp(lossNew);
42     while step<=1000 && abs(lossLast-lossNew)/lossNew>epsilon
43         lossLast=lossNew;
44         lastU = U;
45         U = M*V/(V'*V+lambda*I);
46         lastV = V;
47         V = M'*U/(U'*U+lambda*I);
48         lossNew = norm(W.*(M-U*V'), 'fro').^2 + ...
            lambda/2*(norm(U, 'fro').^2 + norm(V, 'fro').^2);
49         diffUmax = max(max(abs(U-lastU))); % ...
            }' ' ¼ U£-V±秒μ ¾ °
50         diffUmean = mean(mean(abs(U-lastU)));

```

```

51     diffVmax = max(max(abs(V-lastV)));
52     diffVmean = mean(mean(abs(V-lastV)));
53     diff = W0.*(M-U*V');
54     diffmax = max(max(abs(diff)));
55     diffmean = mean(mean(abs(diff(W0~=0))));
56     if ~mod(step, 500)
57         disp([num2str(step), '    loss: ', num2str(lossNew), ...
58             '    Δ: ', num2str(abs(lossNew-lossLast)/lossNew), ...
59             '    diffmax: ', num2str(diffmax), ...
60             '    diffmean: ', num2str(diffmean), ...
61             '    diffUmax: ', num2str(diffUmax), ...
62             '    diffUmean: ', num2str(diffUmean), ...
63             '    diffVmean: ', num2str(diffVmean), ...
64             ']);
65
66         testmean = mean(mean(Mtest(Wtest~=0)));    % mean value ...
67             of test data
68         diffptest = Wtest.*(Mtest-U*V');
69         diffptestmax = max(max(abs(diffptest)));    %  $\frac{1}{4} \tilde{\sigma}^2 \mu$  ...
70             UV  $\frac{3}{4} \sigma^2 \frac{3}{4}$ 
71         diffptestmean = mean(mean(abs(diffptest(Wtest~=0))));
72         mse = 1/s * norm(Wtest.*(Mtest-U*V'), 'fro').^2;
73         disp(['    test ', num2str(mse), ...
74             '    mean: ', num2str(testmean), ...
75             '    diffptestmax: ', num2str(diffptestmax), ...
76             '    diffptestmean: ', num2str(diffptestmean)]);
77         disp(' ');
78     end
79     step = step+1;
80 end
81 % train mse
82 s = size(train,1);
83 mse = 1/s * norm(W.*(M-U*V'), 'fro').^2;
84 disp(['train ', num2str(mse)]);
85 % test
86 Mtest = zeros(943, 1682);
87 Wtest = zeros(943, 1682);

```

```

86     Mtest(sub2ind(size(M), test(:,1), test(:,2))) = test(:,3);
87     Wtest(sub2ind(size(W), test(:,1), test(:,2))) = 1;
88     s = size(test,1);
89     mse = 1/s * norm(Wtest.*(Mtest-U*V'), 'fro').^2;
90     if mse < bestmse
91         bestmse = mse;
92         X = U*V';
93         save X_als.mat X;
94     end
95     disp(['test ', num2str(mse), ...
96         ' step ', num2str(step)]);
97 end

```

## 4.2 梯度下降算法

```

1  % »    GD  $\mu$     烽
2  clear all; close all; clc;
3  load data_train.mat;
4  %% hyperparameters
5  rank = 10;
6  bound = [12, 30, 240, 600];
7  lrs = [0.00015, 0.00013, 0.00011, 0.00011];
8  lr = 0.0002;
9  lambda = 10;
10 epsilon = 0.00001;          % limit of iteration
11 %% learn
12 U = randn(943, rank);
13 V = randn(1682, rank);
14 epoch = 5;
15 %% randomly select train data and test data
16 idx = randperm(90000,80000);
17 train = data_train(idx,:);
18 temp = 1:90000;
19 temp(idx) = 0;
20 temp(temp==0) = [];

```

```

21 test = data_train(temp,:);
22 %% train
23 M = zeros(943, 1682);
24 W = zeros(943, 1682);
25 M(sub2ind(size(M), train(:,1), train(:,2))) = train(:,3);
26 W(sub2ind(size(W), train(:,1), train(:,2))) = 1;
27 step = 0;
28 train_loss = zeros(1001);
29 % , a -1/4 * I 1e07, 1e04, 1e05
30 lossLast = norm(W.*(M-U*V'), 'fro').^2 + ...
    lambda/2*(norm(U, 'fro').^2 + norm(V, 'fro').^2);
31 gradU = 2*W.*(U*V'-M)*V + lambda*U; % , a -100 -1/4
32 lastGradU = gradU;
33 U = U - lr * gradU;
34 gradV = 2*(U'*(W.*(U*V'-M)))' + lambda*V;
35 lastGradV = gradV;
36 V = V - lr * gradV;
37 lossNew = norm(W.*(M-U*V'), 'fro').^2 + lambda/2*(norm(U, 'fro').^2 ...
    + norm(V, 'fro').^2);
38 while step<1000 && abs(lossLast-lossNew)/lossNew>epsilon
39     lossLast = lossNew;
40     gradU = 2*W.*(U*V'-M)*V + lambda*U;
41     lastGradU = mean(abs(gradU(:)));
42     U = U - lr * gradU;
43     gradV = 2*(U'*(W.*(U*V'-M)))' + lambda*V;
44     lastGradV = mean(abs(gradV(:)));
45     V = V - lr * gradV;
46     lossNew = norm(W.*(M-U*V'), 'fro').^2 + ...
        lambda/2*(norm(U, 'fro').^2 + norm(V, 'fro').^2);
47     if ~mod(step, 50)
48         disp([num2str(step), ' loss: ', num2str(lossNew), ...
49             ' Δ: ', num2str(abs(lossNew-lossLast)/lossNew)]);
50     end
51     temp = find(bound==step);
52     if any(temp)
53         lr = lrs(temp);
54     end

```

```

55     step = step+1;
56     train_loss(step) = lossNew;
57 end
58 s = size(train,1);
59 mse = 1/s * norm(W.*(M-U*V'), 'fro').^2;
60 disp(['train ', num2str(mse)]);
61 %% test
62 M = zeros(943, 1682);
63 W = zeros(943, 1682);
64 M(sub2ind(size(M), test(:,1), test(:,2))) = test(:,3);
65 W(sub2ind(size(W), test(:,1), test(:,2))) = 1;
66 s = size(test,1);
67 mse = 1/s * norm(W.*(M-X), 'fro').^2;
68 disp(['test ', num2str(mse)]);
69 %% save
70 X = U*V';
71 save X.mat X;
72 plot(train_loss(1:step), 'linewidth', 3);
73 title('loss');
74 set(gca, 'fontsize', 16);

```