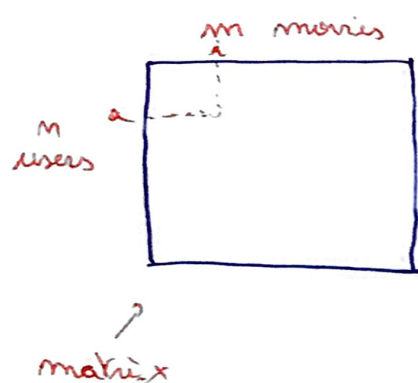


Objectives:

- Understand the problem definition and assumptions of Recommender Systems
- Understand the impact of similarity measures in the K-Nearest Neighbor method
- Understand the need to impose the low rank assumption in collaborative filtering
- iteratively find values of U and V (given $X = UV^T$) in collaborative filtering

Introduction: Netflix

Netflix challenge: 0.5 million users = m
18 000 movies = m

y_{ai}

The interesting part about it is that most of this matrix is actually empty. (only ranked by 5000 users)

We want to predict what will be the opinion of this user for the movie he or she hasn't seen.

We don't want to use the regression here for many reasons:

1. not clear what features are we going to record, and it will be really huge feature vectors.
2. In order for us to do linear regression for a specific user, we actually need to have enough ranking of this user. So if we just have a few rankings, we would be unable to recommend this user any new products.

K-Nearest Neighbor Method:

k means, how big should be your advisory pool on how many neighbors you want to look at

$$\text{predicted} \rightarrow \hat{y}_{ai} = \frac{\sum_{b \in KNN(a,i)} y_{bi}}{k}$$

I'm going to select all the nearest neighbors who did watch this movie and then divide by k

We will identify users similar to me, take their values and take the average.

We assume we have similarity

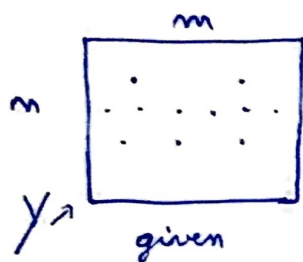
$$\hat{y}_{ai} = \frac{\sum_{b \in KNN(a,i)} \text{sim}(a,b) y_{bi}}{\sum_{b \in KNN(a,i)} \text{sim}(a,b)}$$

↑
similarity

We take the score of KNN and compare the similarity between a and b and use it as a weighting factor.

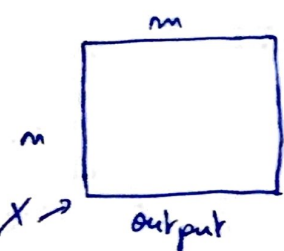
Collaborative Filtering: The Naive approach

23



← Sparse matrix: very few entries

Our goal is to build another matrix X , same size.



← I need every single entry full.

The Algorithm is given Y and will output X , which will contain prediction for every single user and movie.

For each predicted matrix we need to see how good is this matrix

$$J(X) = \sum_{(a,i) \in D} \underbrace{\frac{(Y_{ai} - X_{ai})^2}{2}}_{\text{loss}} + \frac{\lambda}{2} \underbrace{\sum_{(a,i)} X_{ai}^2}_{\text{Regularization} \rightarrow \|X\|^2}$$

objective

$$D = \{(a,i) \mid Y_{ai} \text{ is given}\}$$

1. $(a,i) \in D$

$$\frac{\partial J(X_{ai})}{\partial X_{ai}} = \frac{\partial \left(\frac{(Y_{ai} - X_{ai})^2}{2} + \frac{\lambda}{2} X_{ai}^2 \right)}{\partial X_{ai}} = 0 \implies X_{ai} = \frac{Y_{ai}}{1 + \lambda}$$

2. $(a,i) \notin D$

$$\frac{\partial J(X_{ai})}{\partial X_{ai}} = \frac{\partial \left(\frac{\lambda}{2} X_{ai}^2 \right)}{\partial X_{ai}} = 0 \implies X_{ai} = 0$$

When we did our estimation, we said for all the entries that you didn't know what the value should be, you just put zeros everywhere. And for those that you knew what the value should be you actually corrupt the value.

There is something wrong in this particular way of thinking.

We are taking every single user, compute the derivative and that will cost a lot of time. and the problem with this approach is that there is no connection between assigned values for all different entries of X .

We don't discover any groupings between either users or movies.

Collaborative Filtering with Matrix Factorization:

What we will do now is to somehow tie all these parameters, not to do independent estimation, and we want to decrease the number of parameters.

Assumption: X is low Rank (Rank captures how much dependency do you see between the entries of the matrix)

$$\text{Rank 1: } \overset{m}{\underset{n}{\begin{bmatrix} 1 & 2 & 3 \\ 5 & 10 & 15 \end{bmatrix}}} = \underset{\substack{\mathcal{O}(n+m) \\ \text{parameters}}}{\begin{bmatrix} 1 \\ 5 \end{bmatrix}} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 10 \\ 50 \end{bmatrix} \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix}$$

Why is it interesting for us to factorize? To decrease the number of parameters

$$\underset{u}{\begin{bmatrix} 1 \\ 5 \end{bmatrix}} \underset{v}{\begin{bmatrix} 1 & 2 & 3 \end{bmatrix}} = u \cdot v^T \quad X_{ai} = u_a \cdot v_i$$

The higher you make the Rank \Rightarrow The more multifaceted representation of the user and the movie you can have.

Alternating Minimization:

$$\mathcal{J}(X) = \sum_{(a,i) \in D} (Y_{ai} - X_{ai})^2 / 2 + \frac{\lambda}{2} \sum_{(a,i)} X_{ai}^2 \quad \text{and } X_{ai} = u_a v_i ; \text{rank 1}$$

↓

$$\mathcal{J}(u, v) = \sum_{(a,i) \in D} (Y_{ai} - u_a v_i)^2 + \frac{\lambda}{2} \sum_{a=1}^n u_a^2 + \frac{\lambda}{2} \sum_{i=1}^m v_i^2$$

Assume: $\overset{3 \text{ movies}}{\begin{bmatrix} 5 & ? & 7 \\ 1 & 2 & ? \end{bmatrix}}$

$$\underset{2 \text{ users}}{Y} = \begin{bmatrix} 5 & ? & 7 \\ 1 & 2 & ? \end{bmatrix} \quad u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$$

① Initialization: $v = \begin{bmatrix} 2 \\ 7 \\ 8 \end{bmatrix}$

$$u v^T = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \begin{bmatrix} 2 & 7 & 8 \end{bmatrix} = \begin{bmatrix} 2u_1 & 7u_1 & 8u_1 \\ 2u_2 & 7u_2 & 8u_2 \end{bmatrix}$$

For user 1: u_1 $\mathcal{J}(u, v) = ?$

$$\frac{\partial}{\partial u_1} \left[\frac{(5 - 2u_1)^2}{2} + \frac{(7 - 8u_1)^2}{2} + \frac{\lambda}{2} u_1^2 \right] = -66 + (68 + \lambda)u_1 = 0 \Rightarrow u_1 = \frac{66}{\lambda + 68}$$

$$u_2 = \frac{16}{\lambda + 53}$$

$$\text{if } \lambda = 1 \Rightarrow u_1 = \frac{66}{69} \text{ and } u_2 = \frac{16}{54}$$

② We will take these u_s and recompute the v_s

$$u = \begin{bmatrix} \frac{66}{69} \\ \frac{16}{54} \end{bmatrix} \text{ and } v = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$$

$$u \cdot v^T = \begin{bmatrix} \frac{66}{69} \\ \frac{16}{54} \end{bmatrix} [v_1 \ v_2 \ v_3] = \begin{bmatrix} \frac{66}{69} v_1 & \frac{66}{69} v_2 & \frac{66}{69} v_3 \\ \frac{16}{54} v_1 & \frac{16}{54} v_2 & \frac{16}{54} v_3 \end{bmatrix}$$

Fixing V and Finding U:

2 users, 3 movies

$$Y = \begin{bmatrix} 1 & 8 & ? \\ 2 & ? & 5 \end{bmatrix} \quad \text{The goal is to find } u \text{ and } v \text{ such that } X = u v^T$$

We initialize $v = [4 \ 2 \ 1]^T$

$$X = u v^T = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} [4 \ 2 \ 1] = \begin{bmatrix} 4u_1 & 2u_1 & u_1 \\ 4u_2 & 2u_2 & u_2 \end{bmatrix}$$

missing: We ignore the missing values in the expression.

$$\frac{\partial}{\partial u_1} \left[\frac{(1-4u_1)^2}{2} + \frac{(8-2u_1)^2}{2} + \frac{1}{2} u_1^2 \right] = -4(1-4u_1) - 2(8-2u_1) + 1u_1$$

$$\Rightarrow (1+20)u_1 - 20 = 0 \Rightarrow u_1 = \frac{20}{1+20}$$

$$\frac{\partial}{\partial u_2} \left[\frac{(2-4u_2)^2}{2} + \frac{(5-u_2)^2}{2} + \frac{1}{2} u_2^2 \right] = -4(2-4u_2) - (5-u_2) + 1u_2$$

$$\Rightarrow (1+17)u_2 - 13 = 0 \Rightarrow u_2 = \frac{13}{1+17}$$

kernels:

$$k(x, q) = (x^T q + 1)^2 \quad \text{Polynomial Kernel}$$

$$k(x, q) = \phi(x)^T \phi(q)$$

$$\underline{\phi(x) = ?} \quad \phi(x) = [\phi_1(x_1, x_2), \dots, \phi_N(x_1, x_2)]$$

$$\begin{aligned} k(x, q) &= (x^T q + 1)^2 = \left(1 + \sum_{i=1}^2 x_i q_i \right)^2 = (x_1 q_1 + x_2 q_2 + 1)^2 \\ &= x_1^2 q_1^2 + x_2^2 q_2^2 + 2x_1 x_2 q_1 q_2 + 2x_1 q_1 + 2x_2 q_2 + 1 \end{aligned}$$

$$\Rightarrow \phi(x) = [x_1^2, x_2^2, \sqrt{2}x_1 x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1]$$