

16. Recommender Systems

Example: Predict movie ratings

User rates movies using zero to five stars.

Movie	Alice	Bob	Carol	Dave
Love at last	5	5	0	0
Romance forever	5	? 4.5	? 0	0
Cute puppies of love	? 5	4	0	? 0
Nonstop car chases	0	0	5	4
Swords vs. Karate	0	0	5	? 4

$n_u = 4$

$n_m = 5$

$n_u = \# \text{ users}$

$n_m = \# \text{ movies}$

$r(i, j) = 1$ if user j has rated movie i

$y(i, j) = \text{rating given by user } j \text{ to movie } i$
(defined only if $r(i, j) = 1$)

Stars 0, ..., 5

Content based recommendation system:

	Alice $\theta^{(1)}$	Bob $\theta^{(2)}$	Carol $\theta^{(3)}$	Dave $\theta^{(4)}$	x_1 (romance) $\theta^{(1)}$	x_2 (action)
$x^{(1)}$ Love at last	5	5	\emptyset	\emptyset	0.9	\emptyset
$x^{(2)}$ Romance forever	5	?	?	\emptyset	1.0	0.01
$x^{(3)}$ Cute puppies of love	?	4	\emptyset	?	0.99	\emptyset
$x^{(4)}$ Nowtop car chases	\emptyset	\emptyset	5	4	0.1	1.0
$x^{(5)}$ Swords vs. karate	\emptyset	\emptyset	5	?	\emptyset	0.9

For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$.

Predict user j as rating movie i with $(\theta^{(j)})^T x^{(i)}$ stars.

Example for Alice \rightarrow Cute puppies of love:

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \quad \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad (\theta^{(1)})^T \cdot x^{(3)} = 5 \cdot 0.99 = 4.95$$

Problem Formulation:

$r(i, j) = 1$ if user j has rated movie i (\emptyset otherwise)

$y^{(i,j)}$ = rating by user j on movie i (if defined)

$\theta^{(j)}$ = parameter vector for user j

$x^{(i)}$ = feature vector for movie i

For user j , movie i , predicted rating: $(\theta^{(j)})^T x^{(i)}$

$m^{(j)}$ = # of movies rated by user j

Linear Regr.
Problem

To learn $\theta^{(j)}$:

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} \underbrace{\left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2}_{\text{prediction}} + \frac{\lambda}{2m^{(j)}} \sum_{i=1}^n (\theta^{(j)}_i)^2$$

\downarrow Sum over movies which have been rated by user.
 \downarrow Squared error

$\theta^{(j)} \in \mathbb{R}^{n+1}$

$n=2$
of features

Optimization objective:To learn $\theta^{(j)}$ (parameter for user j)

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^{(j)})^2$$

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$:

$$\min_{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}} \underbrace{\frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2}_{J(\theta^{(1)}, \dots, \theta^{(n_u)})}$$

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k \neq \emptyset)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \underbrace{\left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \cdot \theta_k^{(j)} \right)}_{\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \dots, \theta^{(n_u)})} \quad (\text{for } k \neq \emptyset)$$

Collaborative Filtering

9

Problem motivation:

movie	Alice $\theta^{(1)}$	Bob $\theta^{(2)}$	Carol $\theta^{(3)}$	Dave $\theta^{(4)}$	assumed initially unknown $x_0=1$	
					x1 Romance	x2 action
$x^{(1)}$ Love at last	5	5	\emptyset	\emptyset	? 1.0	? 0.0
Romance forever	5	?	?	\emptyset	?	?
Cute puppies of love	?	4	\emptyset	?	?	?
Nonstop car chases	\emptyset	\emptyset	5	4	?	?
swords vs. karate	\emptyset	\emptyset	5	?	?	?

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$$

$$\theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$$

$$\theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

$$\theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

User informs about $\theta^{(ij)}$ preferences

multiplier associated to x_2

for $x^{(1)}$

$$\begin{aligned} (\theta^{(1)})^T x^{(1)} &\approx 5 \\ (\theta^{(2)})^T x^{(2)} &\approx 5 \\ (\theta^{(3)})^T x^{(3)} &= 0 \\ (\theta^{(4)})^T x^{(4)} &= 0 \end{aligned}$$

Optimization algorithm: \Rightarrow Given $\theta^{(1)}, \dots, \theta^{(nm)}$, to learn $x^{(i)}$:

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

only for movies which have a rating

not too far in squared error sense from actual value

Given $\theta^{(1)}, \dots, \theta^{(nm)}$, to learn $x^{(1)}, \dots, x^{(nm)}$

$$\min_{x^{(1)}, \dots, x^{(nm)}} \frac{1}{2} \sum_{i=1}^{nm} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{nm} \sum_{k=1}^n (x_k^{(i)})^2$$

Collaborative Filtering:

Given $x^{(1)}, \dots, x^{(nm)}$ (and movie ratings), can estimate $\theta^{(1)}, \dots, \theta^{(nm)}$

Given $\theta^{(1)}, \dots, \theta^{(nm)}$

can estimate $x^{(1)}, \dots, x^{(nm)}$

randomly guess $\theta \xrightarrow{\text{optimize}} x \xrightarrow{\text{optimize}} \theta \rightarrow x \rightarrow \dots$ will converge to a reasonable value of $x^{(i)}$ and $\theta^{(i)}$.

Collaborative Filtering algorithm

[5]

Collaborative Filtering optimization objective

~~x_0~~ ~~θ_0~~ not used

use the algorithm \star^1 and \star^2 and use them together

$$x \in \mathbb{R}^n$$

$$\theta \in \mathbb{R}^n$$

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) =$$

$$= \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(i)})^T x^{(j)} - y^{(i,j)})^2 + \underbrace{\frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2}_{\text{optimization objective}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2}_{\text{optimization objective}}$$

to combined \star^1 and \star^2

optimization objective

Simultaneously minimize

$$\min_{\substack{x^{(1)}, \dots, x^{(n_m)} \\ \theta^{(1)}, \dots, \theta^{(n_u)}}} J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$$

Collaborative Filtering algorithm:

- 1) Initialize $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values
- 2) minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm) for every $j=1, \dots, n_u, i=1, \dots, n_m$

$$x \in \mathbb{R}^n$$

$$\theta \in \mathbb{R}^n$$

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

$$\frac{\partial}{\partial x_k^{(i)}} J(\cdot)$$

$$\theta_1$$

$$\theta_n$$

- 3) For a user with parameters θ and a movie with (learned) features x , predict a star rating of $\theta^T x$.

$$\Rightarrow (\theta^{(i)})^T x^{(j)}$$

Vectorization: Low Rank Matrix Factorization

C

Collaborative Filtering:

Movie	Alice	Bob	Carol	Dave
Love at last	5	5	0	0
Romance forever	5	?	?	0
Cute puppies of I.	?	4	0	?
Nonstop car chases	0	0	5	4
Swords vs. Karate	0	0	5	?

$$y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & ? \end{bmatrix}$$

$y(i,j)$

Predicted Ratings:

$$\begin{bmatrix} (\theta^{(1)})^T x^{(1)} & (\theta^{(2)})^T x^{(2)} & \dots & (\theta^{(nw)})^T x^{(1)} \\ (\theta^{(1)})^T x^{(2)} & (\theta^{(2)})^T x^{(2)} & \dots & (\theta^{(nw)})^T x^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ (\theta^{(1)})^T x^{(nw)} & (\theta^{(2)})^T x^{(nw)} & \dots & (\theta^{(nw)})^T x^{(nw)} \end{bmatrix}$$

$$\Rightarrow (\theta^{(j)})^T (x^{(i)})$$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(nw)})^T \end{bmatrix}$$

$$\Theta = \begin{bmatrix} -(\theta^{(1)})^T \\ -(\theta^{(2)})^T \\ \vdots \\ -(\theta^{(nw)})^T \end{bmatrix} \Rightarrow \underline{X\Theta^T}$$

Low rank matrix factorization

Find related movies!

For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$

$x_1 = \text{romance}$, $x_2 = \text{action}$, $x_3 = \text{comedy}$, $x_4 = \dots$

How to find movies j related movie i ?

Small $\|x^{(i)} - x^{(j)}\| \rightarrow$ movie j and i are "similar"

5 most similar movies to movie i :

Find 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$.

Implementation Detail: Mean normalization

3

Users who have not rated any movies:

	Alice	Bob	Carol	Dave	Eve	
Love at last	5	5	0	0	?	0
Romance Forever	5	?	?	0	?	0
Cute puppies of I.	?	4	0	?	?	0
Nonstop car chases	0	0	5	4	?	0
Swords vs. Karate	0	0	5	?	?	0

$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & ? & ? \end{bmatrix}$

0 ← define initially 0 not useful because *

$$\min_{x^{(i)}, \dots, x^{(m)}, \theta^{(1)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(i)})^T x^{(j)} - y_{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_u} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

User has not rated any movie, therefore $r(i,j)$ is always 0

only relevant

$n=2 \quad \theta^{(5)} \in \mathbb{R}^2 \Rightarrow \theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$\frac{\lambda}{2} [(\theta_1^{(5)})^2 + (\theta_2^{(5)})^2]$

* $(\theta^{(5)})^T x^{(i)} = 0$

Mean normalization:

average rating $(5+5+0+0)/4 = 2.5$

subtract average rating from Rating

$5 - 2.5 = 2.5$ $0 - 2.5 = -2.5$

$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & ? & ? \end{bmatrix}$

$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix}$

$\Rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$

$0 - 2.25 = -2.25$

learn $\theta^{(j)}, x^{(i)}$

pretended to be ratings of users

For user j , on movie i predict:

$\Rightarrow (\theta^{(j)})^T (x^{(i)}) + \mu_i$

User 5 (Eve)

$\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \underbrace{(\theta^{(5)})^T (x^{(i)})}_{0} + \mu_i = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix}$ which is the average rating of the movies