



Recommender Systems

Predicting movie ratings, collaborative filtering, and low rank matrix factorization.

- 1. Predicting Movie Ratings
 - 1a. Problem Formulation
 - 1b. Content Based Recommendations
- 2. Collaborative Filtering
 - 2a. Introduction
 - 2b. Collaborative Filtering Algorithm
- 3. Low Rank Matrix Factorization

- o 3a. Vectorization: Low Rank Matrix Factorization
- o 3b. Implementation Detail: Mean Normalization

This is an important practical application of machine learning. Netflix, Spotify, Youtube, Amazon and other companies try to recommend things to you every time you use their services.

1. Predicting Movie Ratings

I would like to give full credits to the respective authors as these are my personal python notebooks taken from deep learning courses from Andrew Ng, Data School and Udemy :) This is a simple python notebook hosted generously through Github Pages that is on my main personal notes repository on <https://github.com/ritchieng/ritchieng.github.io>. They are meant for my personal review but I have open-sourced my repository of personal notes as a lot of people found it useful.

Example: Predicting movie ratings

→ User rates movies using ~~one~~
zero to five stars

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	6
Romance forever	5	?	0	0
Cute puppies of love	?	5	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?

$$n_u = 4 \quad n_m = 5$$



→
→
→
→
→

→ n_u = no. users

→ n_m = no. movies

→ $r(i, j) = 1$ if user j has rated movie i
 $y^{(i,j)}$ = rating given by user j to movie i
 (defined only if $r(i, j) = 1$)

1b. Content Based Recommendations

- How do we predict the missing values?

Content-based recommender systems

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$n_u = 4, n_m = 5$	$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$
Love at last 1	5	5	0	0	x_1 (romance) $\rightarrow 0.9 \rightarrow 0$	
Romance forever 2	5	?	?	0	x_2 (action) $\rightarrow 1.0 \rightarrow 0.01$	
Cute puppies of love 3	?	4	0	?	x_3 (drama) $\rightarrow 0.99 \rightarrow 0$	
Nonstop car chases 4	0	0	5	4	x_4 (thriller) $\rightarrow 0.1 \rightarrow 1.0$	
Swords vs. karate 5	0	0	5	?	x_5 (sci-fi) $\rightarrow 0 \rightarrow 0.9$	$n=2$

→ 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.

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

Problem formulation

- If we minimize the following function, we get the parameters to predict

→ $r(i, j) = 1$ if user j has rated movie i (0 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: $\underline{(\theta^{(j)})^T(x^{(i)})}$

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

→ $m^{(j)}$ = no. of movies rated by user j

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

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

We can use other minimization algorithms (other than gradient descent)

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

↙
 $J(\theta^{(1)}, \dots, \theta^{(n_u)})$

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

$\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \dots, \theta^{(n_u)})$

↙ ↘ ↗

2. Collaborative Filtering

2a. Introduction

- Here we will be learning about "Feature Learning"
 - Feature Learning: learning what features to use
- Problem motivation
 - It is inefficient and difficult to ask someone to watch each movie and inform us how romantic or action-packed the movie is

Movie	Alice (1) $\theta^{(1)}$	Bob (2) $\theta^{(2)}$	Carol (3) $\theta^{(3)}$	Dave (4) $\theta^{(4)}$	x_1 (romance)	x_2 (action)
Love at last	5	5	0	0	1.0	0.0
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

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

→ $\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}, \theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$

$\theta^{(j)}$

$(\theta^{(1)})^T x^{(1)} \approx 5$
 $(\theta^{(2)})^T x^{(1)} \approx 5$
 $(\theta^{(3)})^T x^{(1)} \approx 0$
 $(\theta^{(4)})^T x^{(1)} \approx 0$

- Now we've no idea how each movie is romantic (x1) or action-packed (x2)
 - Let's say that
 - Alice (θ_1): likes romance
 - Bob (θ_2): likes romance
 - Carol (θ_3): likes action
 - Dave (θ_4): likes action
 - We can discover x_1 by making sure the following happens
 - $\theta_1 \text{ transpose} * x_1 = 5$
 - $\theta_2 \text{ transpose} * x_1 = 5$
 - $\theta_3 \text{ transpose} * x_1 = 0$
 - $\theta_4 \text{ transpose} * x_1 = 0$
- Optimization algorithm
 - This tries to choose features X_i so that for all the users J that have rated that movie, the algorithm also predicts a value for how that user would have rated that movie that is not too far, in the squared error sense, from the actual value that the user had rated that movie
 - Regularization term to prevent features from becoming too big

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(i)}$:

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

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(1)}, \dots, x^{(n_m)}$:

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

We can guess θ , solve for x , then solve for θ and continue

- There is a more efficient method to do this and we will be discussing this shortly

[Given $x^{(1)}, \dots, x^{(n_m)}$ (and movie ratings),
 can estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$]

[Given $\theta^{(1)}, \dots, \theta^{(n_u)}$,
 can estimate $x^{(1)}, \dots, x^{(n_m)}$]

Guess $\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \dots$

2b. Collaborative Filtering Algorithm

- Collaborative filtering optimization objective
 - We will take both of these optimization objectives and put them together
 - Now we get to minimize with respect to x and θ simultaneously

$$\begin{aligned} & \rightarrow \text{Given } x^{(1)}, \dots, x^{(n_m)}, \text{ estimate } \theta^{(1)}, \dots, \theta^{(n_u)}: \\ & \quad \min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2 \quad \leftarrow \\ & \rightarrow \text{Given } \theta^{(1)}, \dots, \theta^{(n_u)}, \text{ estimate } x^{(1)}, \dots, x^{(n_m)}: \\ & \quad \min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 \quad \leftarrow \\ & \text{Minimizing } x^{(1)}, \dots, x^{(n_m)} \text{ and } \theta^{(1)}, \dots, \theta^{(n_u)} \text{ simultaneously:} \\ & \quad J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2 \end{aligned}$$

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

Collaborative filtering algorithm

Collaborative filtering algorithm

~~$x_0 = 1$~~ $x \in \mathbb{R}^n$, $\theta \in \mathbb{R}^n$

~~θ_0~~
 ~~θ_1~~
 ~~\vdots~~
 ~~θ_n~~

- 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). E.g. for every $j = 1, \dots, n_u, i = 1, \dots, n_m$:

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

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

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

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

3. Low Rank Matrix Factorization

3a. Vectorization: Low Rank Matrix Factorization

- Low rank matrix factorization
 - Y matrix: all the predicted ratings
 - We conduct matrix factorization by decomposing the matrix into a product of matrices
 - X matrix: features of each movie stacked in rows
 - H matrix: parameters of each user stacked in rows
 - Product of X and H_transpose equals Y matrix

Collaborative filtering

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

$X \Theta^\top$
 $(\Theta^{(i)})^\top (\underline{x}^{(i)})$

Predicted ratings: $(i, j) \rightarrow$

$(\theta^{(1)})^\top (\underline{x}^{(1)})$	$(\theta^{(2)})^\top (\underline{x}^{(1)})$	\dots	$(\theta^{(n_u)})^\top (\underline{x}^{(1)})$
$(\theta^{(1)})^\top (\underline{x}^{(2)})$	$(\theta^{(2)})^\top (\underline{x}^{(2)})$	\dots	$(\theta^{(n_u)})^\top (\underline{x}^{(2)})$
\vdots	\vdots	\vdots	\vdots
$(\theta^{(1)})^\top (\underline{x}^{(n_m)})$	$(\theta^{(2)})^\top (\underline{x}^{(n_m)})$	\dots	$(\theta^{(n_u)})^\top (\underline{x}^{(n_m)})$

$X = \begin{bmatrix} -(\underline{x}^{(1)})^\top \\ -(\underline{x}^{(2)})^\top \\ \vdots \\ -(\underline{x}^{(n_m)})^\top \end{bmatrix}$
 $\Theta = \begin{bmatrix} -(\Theta^{(1)})^\top \\ -(\Theta^{(2)})^\top \\ \vdots \\ -(\Theta^{(n_u)})^\top \end{bmatrix}$

Low rank matrix factorization

Finding related movies

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

$\rightarrow \underline{x}_1 = \text{romance}, \underline{x}_2 = \text{action}, \underline{x}_3 = \text{comedy}, \underline{x}_4 = \dots$

How to find $\underline{x}^{(j)}$ related to $\underline{x}^{(i)}$?

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

5 most similar movies to movie i :

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

3b. Implementation Detail: Mean Normalization

- Issue without mean normalization
 - And for even the Swords vs. Karate, someone rated it 5 stars
 - So some people do like some movies
 - It seems not useful to just predict that Eve is going to rate everything 0 stars
 - If we're predicting that Eve is going to rate everything 0 stars, we also don't have any good way of recommending any movies to her, because you know all of these movies are getting exactly the same predicted rating for Eve so there's no one movie with a higher predicted rating that we could recommend to her

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)
Love at last	5	5	0	0	?
Romance forever	5	?	?	0	?
Cute puppies of love	?	4	0	?	?
Nonstop car chases	0	0	5	4	?
Swords vs. karate	0	0	5	?	?

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

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

\uparrow

$n=2 \quad \underline{\Theta}^{(s)} \in \mathbb{R}^2 \quad \underline{\Theta}^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$(\underline{\Theta}^{(s)})^T \underline{x}^{(i)} = 0$

$\frac{\lambda}{2} [\underline{\Theta}_1^{(s)}]^2 + [\underline{\Theta}_2^{(s)}]^2 \leftarrow$

We can conduct mean normalization to solve this issue

- Our prediction would be the average of each movie

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix} \rightarrow \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user j , on movie i predict:

$$\rightarrow (\underline{\Theta}^{(s)})^T (\underline{x}^{(i)}) + \mu_i$$

learn $\underline{\Theta}^{(s)}, \underline{x}^{(i)}$

User 5 (Eve):

$$\underline{\Theta}^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\underbrace{(\underline{\Theta}^{(s)})^T (\underline{x}^{(i)})}_{=0} + \boxed{\mu_i}$$