# Machine Learning Week 9 Section 2

## Recommender Systems

\_\_\_\_\_

When you buy a product online, most websites automatically recommend other products that you may like. Recommender systems look at patterns of activities between different users and different products to produce these recommendations. In this module, we introduce recommender algorithms such as the collaborative filtering algorithm and low-rank matrix factorization.

#### **Predicting Movie Ratings**

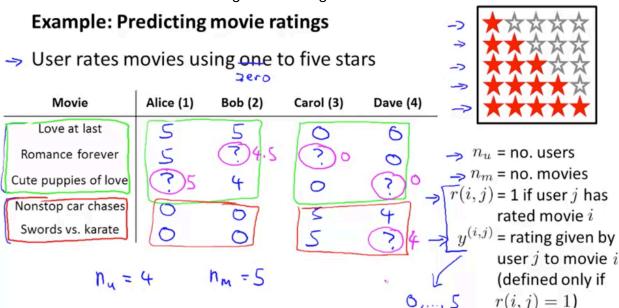
#### **Problem Formulation**

Recommender systems are currently one of the most valuable applications of ML currently used. Ex. Netflix, Amazon, Youtube, etc.

---

Recall, the features we choose to use are very important for ML algorithms. There are some algorithms that are able to learn which features to use. Recommender Systems can do this to some extent.

Problem Formulation: Ex. Predicting movie ratings



Given this dataset, the r(i, j)s and the y(i, j)s we need to predict what the 'question mark' values are.

In our notation, r(i,j)=1 if user j has rated movie i, and  $y^{(i,j)}$  is his rating on that movie. Consider the following example (no. of movies  $n_m=2$ , no. of users  $n_u=3$ ):

	User 1	User 2	User 3
Movie 1	0	1	?
Movie 2	?	5	5

What is r(2, 1)? How about  $y^{(2,1)}$ ?

$$\bigcap r(2,1) = 0, \ y^{(2,1)} = 1$$

$$\bigcap r(2,1)=1,\; y^{(2,1)}=1$$

$$r(2,1) = 0, \ y^{(2,1)} = ext{undefined}$$

Correct

# Content-based Recommendations (Algorithm)

We will now look at an approach to building a recommender system called a content-based approach.

How do we predict what the missing values will be?

Let us say we have a set of features for our movies that measure the degree to which a movie is a particular genre of movie. Then, each movie has a a feature vector x^i from 1 to m with # features n (2 in this example).

We are applying a different instance of linear regression for each user.

Content-base	ed recomr	nende	r systems	Nu = 4	, nm=5	SX.	0.9
Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_1$ (romance)	$\overset{\checkmark}{x_2}$	[0]
Love at last	5	5	0	0	→ 0.9	-> 0	_ \
Romance forever 2	5	?	?	0	-> 1.0	<b>→</b> 0.01	1
Cute puppies of love	74.95	4	0	?	0.99	→ 0	
Nonstop car chases 4	0	0	5	4	→ 0.1	→ 1.0	
Swords vs. karate 5	0	0	5	?	→ 0	→ 0.9	n=5
>> For each use			$\theta^{(j)} \in$	$\mathbb{R}^3$ . Pred	dict user	j as rat	ing

movie i with  $(\theta^{(j)})^T x^{(i)}$  stars.  $\chi^{(3)} = \begin{bmatrix} 0 & qq \\ \hline 0 & 1 \end{bmatrix} \rightleftharpoons \begin{bmatrix} 0 & 0 \\ \hline 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\$ 

We are taking the learned parameters (theta1) for a similar user and applying them to a user with unknown variables (user 3 / x3)

More formally, here is how we can write down the problem:

#### **Problem formulation**

$$\rightarrow r(i,j) = 1$$
 if user j has rated movie i (0 otherwise)

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

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

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

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

$$m^{(j)} = \text{no. of movies rated by user } j$$
  
To learn  $\theta^{(j)}$ :

$$\min_{(i,j)} \frac{1}{2^{m(i)}} \sum_{(i,r(i,j)=1)} \left( (O_{(i)})_{i}(x_{(i)}) - A_{(i,j)} \right)_{j} + \frac{5}{2^{m(i)}} \sum_{k=1}^{k} (O_{(i)}^{k})_{j}$$

Using this approach we get a pretty good estimate of theta j with which to make predictions for user Js movie ratings.

For recommender systems we need to change some of the notation a bit:

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

$$\min_{(i,j)} \frac{1}{2^{\log 2}} \sum_{(i,i,j)=1}^{\log 2} \left( (Q_{(i)})_{i}(X_{(i)}) - A_{(i,i)} \right)_{5} + \frac{5}{\sqrt{2}} \sum_{k=1}^{\log 2} (Q_{(i)}^{k})_{5}$$

We drop the m(j) but with the minimizer should still get the same value of theta j as before.

More clearly to learn all the parameter vectors theta:

#### Optimization objective:

The whole picture:

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

# Gradient descent update:

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

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

$$\frac{\partial}{\partial \theta_{k}^{(j)}} \ \zeta(\theta^{(j)}, \dots, \theta^{(n_{N})})$$

Note this looks identical to the linear regression we performed before, just without the 1/m term in front of the sum.

<sup>\*\*</sup>Note: This algorithm assumes we have the features of different movies available to us and that the features capture the content of those movies. However in practice we don't typically have such a dataset. Not just for movies but for whatever we are trying to recommend. In the next video we will discuss how to get around this with a different approach.

## Collaborative Filtering

In this video we will discuss another approach to building a recommender system called collaborative filtering. This algorithm has a interesting property called featuring learning in which it can start to learn for itself which features to use.

Here was our previous dataset without the values for the features. Where do we get these features from now?

We only have user ratings from some of our users for some of our movies.

Lets say we can go to (some) of our users and ask them what the value of theta j is for them. (Ask them which movies and which types of movies they like)

With this, we can begin to infer some of the unknowns.

Problem motivation					1	T	X0=
Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_1$ (romance)	$x_2$ (action)	
X(1) Love at last	75	75	20	> 0	#1.0	¥ 0.	0
Romance forever	5	?	?	0	?	?	XW= [1.0]
Cute puppies of love	?	4	0	?		?	[0.0]
Nonstop car chases	0	0	5	4	?	?	<u>_(i)</u>
Swords vs. karate	0	0	5	?	? .	?	~ (L)
$\Rightarrow \boxed{\theta^{(1)} =}$	$\theta^{(2)}$	$ \begin{array}{c} 0 \\ 0 \\ 0 \end{array} $	$\theta^{(3)} = 0$	$\theta^{(4)} =$		(1	(6'0) <sup>T</sup> x'0 <sup>2</sup> 25 (6'0) <sup>T</sup> x'0 <sup>2</sup> 20 (5'0) <sup>T</sup> x'0 <sup>2</sup> 20

Let us formalize this problem:

Given our all thetas for a users that have rated a particular movie, we need to predict x^i.

And then for all our thetas for all movies, we can try to predict X.

# Optimization algorithm

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

Given 
$$\underline{\theta^{(1)}, \dots, \theta^{(n_u)}}$$
, to learn  $\underline{x^{(i)}}$ :
$$\Rightarrow \quad \min_{x^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} (\underline{(\theta^{(j)})^T x^{(i)}} - \underline{y^{(i,j)}})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

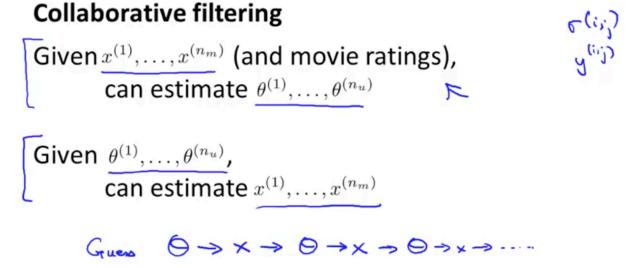
Given 
$$\theta^{(1)}, \dots, \theta^{(n_u)}$$
, to learn  $x^{(1)}, \dots, x^{(n_m)}$ :

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

Putting this all together.

In the previous video we showed that if we have a set of movie ratings (the r(i,j)s and the y(i,j)s) we can learn our parameters theta and then make predictions for different users.

In this video, we showed that if your users are willing to give you parameters theta, we can then estimate features for the different movies.



Notice this is a chicken & egg situation. With the ratings (features), we can learn the parameters, but with the features, we can learn the parameters.

What we can do is start with a random theta, estimate our features, and then iterate, improving our theta, our features, etc. as the system matures.

With each user contributing ratings on a subset of the data, the users are <u>collaborating</u> to all improve the system and algorithm.

In the next video we will discuss an even better technique for collaborative filtering.

## Collaborative Filtering Algorithm

From the ideas we discussed in the previous video of how we can use features to find parameters or parameters to find features, we are going to put those ideas together to come up with a collaborative filtering algorithm.

# Collaborative filtering optimization objective

$$\Rightarrow \text{Given } x^{(1)}, \dots, x^{(n_m)}, \text{ estimate } \theta^{(1)}, \dots, \theta^{(n_u)} :$$

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

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

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

It turns out that instead of going back and forth, there is a more efficient algorithm that does not need to go back and forth and can instead solve for theta and x simultaneously and put them into the same objective:

Minimizing 
$$x^{(1)}, \ldots, x^{(n_m)}$$
 and  $\theta^{(1)}, \ldots, \theta^{(n_u)}$  simultaneously:

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

First, looking at the two equations above, notice that the squared error sum in the middle is the same but using i or j, and then summed over all movies that have ratings. As such we can combine those terms into a combined optimization objective + the respective regularizations.

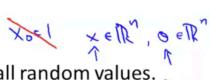
Notice in the combined minimization objective, if you were to hold the x's constant or the theta's constant you would effectively be solving for the above equations.

\*\*Note: Previously we have been using the convention to add a bias term x0 = 1, for this implementation we will not be doing that. As such theta will also be an element of R^n like x. (And we are regularizing all our terms, no 0 term that doesn't get regularized)

This is because the algorithm now has the flexibility to learn it's own features and for example could set it's own bias term x1 = 1 if it needed to.

Putting everything together:

## Collaborative filtering algorithm



- Collaborative filtering algorithm

  1. Initialize  $x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)}$  to small random values.  $\rightarrow$  2. Minimize  $J(x^{(1)},\ldots,x^{(n_m)},\theta^{(1)},\ldots,\theta^{(n_u)})$  using gradient descent (or an advanced optimization algorithm). E.g. for every  $j = 1, ..., n_u, i = 1, ..., n_m$ :

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

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

$$\frac{\partial}{\partial x_{k}^{(i)}}$$

$$\frac{\partial}{\partial x_{k}^{(i)}}$$

For a user with parameters  $\theta$  and a movie with (learned) features  $\underline{x}$ , predict a star rating of  $\underline{\theta^T x}$ .

<sup>\*\*</sup>Notice that we initialize our x's and theta's to small random variables to perform symmetry breaking (similar to the random initialization of a neural network's parameters) and ensures the algorithm learns features that are different from one another.

## **Low Rank Matrix Factorization**

#### Vectorization: Low Rank Matrix Factorization

In this video we will discuss the vectorization implementation of the collaborative filtering algorithm and some other tips and tricks.

---

We want to write out an alternative way of writing out the predictions of the collaborative filtering algorithm.

#### **Collaborative filtering** Movie Alice (1) Bob (2) Carol (3) Dave (4) 0 5 5 0 Love at last $Y = \begin{bmatrix} 5 & 7 & 7 & 0 \\ 5 & 7 & 7 & 0 \\ 7 & 4 & 0 & 7 \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$ Romance forever 5 ? ? 0 Cute puppies of 4 0 love Nonstop car chases Swords vs. karate

The rating user j would give to movie i. is at index (i,j) of the predicted ratings matrix.