

- ⇒ Recommender Systems
- Recommendation is currently very popular application of ML. eg: Netflix, Amazon, Spotify, iTunes Genius
- Really popular in industry, not so much in academia (so).
- There is an interesting idea that for some specific applications problems, there are certain algorithms that can automatically "learn" a good set of features for you!!
- So, instead of trying to manually pick some features, there might be some settings where you might be able to have an algorithm, just to learn what feature to use, & the recommender system is just one example of that sort of setting.

→ Example motivation: Ranking movies

\* Supp. we're trying to recommend movies to customers.

\* We have a dataset as follows:

⇒ Rating scale  $\rightarrow [0, 5]$

⇒ ?  $\rightarrow$  No rating found

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

\* Let's introduce a bit of notation:

$$n_u = \text{no. of users}$$

$$n_m = \text{no. of movies}$$

$$r(i, j) = 1 \quad \text{if user } j \text{ has rated movie } i$$

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

\* For above example:  $n_u = 3$ ;  $n_m = 2$ ;  $r(2, 1) = 0$ ;  $y^{(5, 2)} = 1$

\* Our objective is that given  $r(i, j)$  &  $y^{(i, j)}$ , ~~fit~~ try to come up with a learning alg. that ~~fit/predict~~ predict the missing values (?)

⇒ Content based recommendations

→ Using a similar example, how do we predict the missing entries?

→ One approach is to introduce 2 features  $x_1$  &  $x_2$  which represent how much romance or how much action a movie may have (on a scale of 0-1).

e.g.:

Movie	Alice(1) $\alpha^1$	Bob(2) $\alpha^2$	Carol(3) $\alpha^3$	$x_1$ (Romance)	$x_2$ (Action)
Fault in our Stars (1)	5	5	?	0.99	0
DDLJ <del>2</del> (2)	5	?	1	1.0	0.01
The Notebook (3)	?	4	0	1.0	0
Die Hard <del>4</del> (4)	0	0	5	0	1
Star Wars <del>5</del> (5)	0	0	5	0.1	1

→ Formally, for convenience sake we also add a  $x_0 = 1$  for every movie. ∴ for DDLJ  $\vec{x}^2 = \begin{bmatrix} x_0^2 \\ x_1^2 \\ x_2^2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1.0 \\ 0.01 \end{bmatrix}$

$$\vec{x}^2 = \begin{bmatrix} x_0^2 \\ x_1^2 \\ x_2^2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1.0 \\ 0.01 \end{bmatrix}$$

→ For filling the missing values, one approach is to do linear regression for every single user:

\* For each user  $j$ , learn a parameter  $\theta^j \in \mathbb{R}^{m+1}$

$m \rightarrow$  no. of categories for movies. For our purpose.

example,  $\theta^j \in \mathbb{R}^{2+1} \therefore \theta^j \in \mathbb{R}^3$

\* Predict user  $j$ 's rating for  $i^{\text{th}}$  movie as

$$\textcircled{O} \quad (\theta^j)^T x^i \text{ stores.}$$

e.g.: Suppose in our prev. example.

For "The Notebook"  $\rightarrow x^3 = \begin{bmatrix} 1 \\ 1.0 \\ 0 \end{bmatrix}$

For Alice, we say we've learned  $\theta^1 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$

$\therefore$  Alice will rate The Notebook :  $(\theta^1)^T x^3 = 1 \times 5 = 5$

$\Rightarrow$  Problem formulation

$\rightarrow$  Some conventional notation :

\*  $r(i, j) = \begin{cases} 1 & \text{if user } j \text{ has rated movie } i \\ 0 & \text{otherwise} \end{cases}$

\*  $y^{(i, j)}$  = 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  $(\theta^j)^T x^i$

\*  $m^j$  = No. of movies rated by user  $j$

→ Objective : Learn  $\theta^s$

→ To "learn"  $\theta^s$  find values which minimize  
the following cost function?  $J(\theta^s)$

$$\min_{\theta^s} \frac{1}{2m^s} \sum_{i:r(i,s)=1} [(\theta^s)^T x^i - y^{(i,s)}]^2$$

\* The cost function is the same as the one used  
in linear regression : Root mean square error

\*  $\sum_{i:r(i,s)=1} (\cdot) \rightarrow$  mean sum over all values  
of  $i$  (all the available movies) when  $r(i,s) = 1$   
(all the films users have rated)

→ We can also add a regularization term to make our eq<sup>n</sup> look as follows:

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

\* The regularization term goes from  $k=1$  through  $\hat{n}$ ,

so  $\theta^*$  ends up being a  $(m+1)$  feature

vector. Don't regularize over the bias terms ( $\theta_0$ ).

→ To make this a little clearer, get rid of the  $m^*$  term.

\*  $m^* \rightarrow$  constant; so won't affect our minimization routine

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

$J(\theta^*)$

→ To get the parameters for all our users, we do the following :

$$\min_{\theta^1, \dots, \theta^{mn}} \left( \frac{1}{2} \sum_{j=1}^{m_u} \sum_{i:n(i,j)=1} [(o_j^i)^T x^i - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{k=1}^m \sum_{n=1}^n \theta_k^2 \right)$$

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

\* In our recommendation system, we want to learn parameters for all the users. So, we add one extra summation term to this which means we determine the  $\theta^j$  value for every user s.t. it minimizes the overall optimization objective / cost function.

We can also write it as :

$$\min_{\theta^1, \dots, \theta^m} J(\theta^1, \dots, \theta^m)$$

$\Rightarrow$  Gradient descent updates

$$\theta_k^{(t+1)} := \theta_k^{(t)} - \alpha \frac{\partial}{\partial \theta_k^{(t)}} J(\theta^{(t)}; \theta_1^{(t)}, \dots, \theta_m^{(t)})$$

$$\rightarrow \frac{\partial}{\partial \theta_k^{(t)}} = \left( \sum_{i: r(i,y)=1} [(\theta^{(t)})^T x^{(i)} - y^{(i,t)}] x_k^{(i)} \right) + \lambda \theta_k^{(t)}$$

$\therefore$  The updates become :

$$\theta_k^{(t+1)} := \theta_k^{(t)} - \alpha \sum_{i: r(i,y)=1} [(\theta^{(t)})^T x^{(i)} - y^{(i,t)}] x_k^{(i)}$$

(for  $k=0$ )

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

(for  $k \neq 0$ )

- The only diff. from the linear regression update is that we've eliminated the  $\frac{1}{m}$  term.
  - This approach is called content-based approach as we assume that we have features regarding the content (e.g. how romantic or action heavy some movie is) which helps us identify things that make them appealing to a user.
- However such features are usually not available & are really hard to find. So, this approach is not that practical!

⇒ Collaborative filtering

→ It can be very difficult to find features such as "amount of romance" or "amount of action" in a movie.

The collaborative filtering alg has a very interesting property - feature learning i.e. it can learn for itself what features it needs to learn!

→ Now, let's make a diff. assumption: ~~these~~

Say, ~~we~~ before training we poll each user & find out how much each user likes romantic & action films & capture that into vectors.

→ eg: Supr. we have a dataset:

Movie	Alice (1)	Bob(2)	Carol (3)	$x_1$ (Romance)	$x_2$ (Action)	$x_3$ ...
<del>Movie</del>	<del>Alice (1)</del>	<del>Bob(2)</del>	<del>Carol (3)</del>	<del><math>x_1</math> (Romance)</del>	<del><math>x_2</math> (Action)</del>	<del><math>x_3</math> ...</del>
Fault in Our Stars	5	5	0	?	?	?
DDLJ	5	?	?	?	?	?
The Notebook	?	4	0	?	?	?
Die Hard	0	0	5	?	?	?
Star Wars	0	0	5	?	?	?

\* We've marked  $x_1$ ,  $x_2$  as ? because unlike content based recommendations, we don't have any knowledge about the "contents" of a film.

\* Supr. we will ask user to find out how much each user likes romantic & action films. This results in the following parameter set:

$$\theta^1 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \quad \theta^2 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \quad \theta^3 = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

\* Each column of these  $\theta$  vectors can represent anything, but these values will be DIFFERENT for each user acc. to their likes & choices. For convenience sake, assume  $\theta_2 \rightarrow$  liking for romantic film &  $\theta_3 \rightarrow$  liking for film

→ If we can get those parameters from the user  
we can infer the missing values from our table!!

e.g.: For the movie "Fault in our stars",

We know from the o vectors that Alice & Bob  
love romantic films & Coral hates them.  
(by using this)  
Based on the fact that

Alice & Bob loved the movie ~~Fault in our stars~~ (Fault in our stars)

& Coral hated it, so we can conclude that

"Fault in our stars" is a romantic film.

→ This is a bit simplified by what we're really asking is: "What feature vector should  $x'$  be  
so that"

$$*(\theta^1)^T x' \approx 5$$

$$*(\theta^2)^T x' \approx 5$$

$$*(\theta^3)^T x' \approx 0$$

From this, we can guess  $x'$  may be  $\begin{bmatrix} 1 \\ ? \\ 0 \end{bmatrix} \rightarrow \text{Romantic}$

⇒ Formalizing the collaborative filtering algorithm

⇒ Optimization algorithm

Given  $\theta^1, \dots, \theta^{m_u}$  [i.e. given the parameter vectors to learn  $x^i$ : for each user's preferences]

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

\* Basically means, that out of all the ~~possible~~ <sup>possible</sup> given by users,  <sup>$(\theta^1, \dots, \theta^{m_u})$</sup>  select the feature vector  $x^i$  which minimizes the mean square error.

→ Since, we want to learn <sup>the</sup> ~~all~~ features for all the films, so add all the square errors to make a combined optimization objective:

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

⇒ How does collaborative filtering fare with content-based recommendation systems?

→ Collaborative filtering : Given  $\theta^1, \dots, \theta^m$  (movie ratings) estimate  $x^1, \dots, x^m$

→ Content based recommendation : Given  $x^1, \dots, x^m$  (movie ratings) estimate  $\theta^1, \dots, \theta^m$

→ Which one to choose first is kind of like the chicken-egg problem.

→ Ideal way :

- \* Randomly initialize  $\theta$
- \* Generate  $x$  with collaborative filtering
- \* Then use content based recommendation to improve  $\theta$
- \* Then use that ( $\uparrow$ ) to improve  $x$
- \* . . . . so on

This works really well practically!

\* The algo is termed collaborative filtering because in this all the users come together & collaborate by rating the movies & filling the pool, effectively helping the algorithm to churn out better recommendation for everybody!

⇒ Collaborative filtering Algorithm

→ We discussed earlier about the ideal way of doing collaborative filtering with content based model where we go back & forth between the 2 models. Calculate  $\theta \rightarrow x \rightarrow \theta \rightarrow x \dots$

→ There's a more efficient algorithm which can solve for  $\theta$ 's &  $x$  simultaneously:

→ Optimization objective

→ Minimizing  $x^1, \dots, x^{m_m}, \theta^1, \dots, \theta^{m_u}$  simultaneously

$$J(x^1, \dots, x^{m_m}, \theta^1, \dots, \theta^{m_u}) = \frac{1}{2} \sum_{(i,j) : \lambda(i,j)=1} [(x^j)^T \theta^i - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{j=1}^{m_u} \sum_{k=1}^m (\theta_k^j)^2 + \frac{\lambda}{2} \sum_{i=1}^{m_m} \sum_{k=1}^n (x_k^i)^2$$

→ Objective:

$$\min_{x^1, \dots, x^{m_m}, \theta^1, \dots, \theta^{m_u}} J(x^1, \dots, x^{m_m}, \theta^1, \dots, \theta^{m_u})$$

$$\rightarrow \text{The term } \sum_{(i,j) \in \{(i,j) : i \neq j\}} [(O^T x^i - y^{(i,j)})]^2$$

\* Sums over all pairs  $(i,j)$  such that  $i \neq j$

\* This has a really interesting property

- Keeps  $\alpha$  constant: Becomes similar to

$$\left[ \sum_{j=1}^{m_n} \sum_{i: n(i,j)=1} [(O^T x^i - y^{(i,j)})]^2 \right] \quad \begin{array}{l} \text{Given} \\ \text{label} \end{array}$$

- Keeps  $O$  constant: Becomes similar to

$$\left[ \sum_{i=1}^{m_m} \sum_{j: n(i,j)=1} [(O^T x^i - y^{(i,j)})]^2 \right] \quad \begin{array}{l} \text{Given} \\ \text{label} \end{array}$$

\* In order to come up with just one optimization function, we treat  $J$  as a "func" of both  $x$  &  $\alpha$ .

[?]

\* We've dropped the  $\alpha_0 = 1$  term.  $\therefore x, \alpha \in \mathbb{R}^n$

As we are now learning all the features so if the system needs a feature always = 1, it can learn that by itself.

$\Rightarrow$  Full Collaborative Filtering Algorithm

- 1) Initialize  $x^1, \dots, x^{m_m}, \theta^1, \dots, \theta^{m_u}$  to small random values.
- \* This serves to break symmetry & ensures that the algorithm learns features that are different from each other.
- 2) Minimize  $J(x^1, \dots, x^{m_m}, \theta^1, \dots, \theta^{m_u})$  using gradient descent (or an advanced optimization algo).
 
$$x_k^i := x_k^i - \alpha \left( \sum_{j: r(i,j)=1} [(\theta^j)^T x^i - y^{(i,j)}] \frac{\partial J}{\partial x_k^i} + \lambda x_k^i \right)$$

$$\theta_k^j := \theta_k^j - \alpha \left( \sum_{i: r(i,j)=1} [(\theta^j)^T x^i - y^{(i,j)}] x_k^i + \lambda \theta_k^j \right) \frac{\partial J}{\partial \theta_k^j}$$
- 3) For a user with parameters  $\theta^j$  & a movie with (binary) features  $x$ , predict the rating by  $\theta^j x$ .

⇒ Vectorization: Low Rank matrix factorization

→ Suppose we have a dataset:

Movie	User1	User2	User3	User4
M1	5	5	0	0
M2	5	?	?	0
M3	?	4	0	?
M4	0	0	5	4
M5	0	0	5	?

→ We can also vectorize the above dataset in  $\mathbf{Y}$ , s.t.

$$\mathbf{Y} = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & ? \end{bmatrix} = y^{(i,j)} \text{ for each user } j \text{ & movie } i$$

→ Vectorizing the predicted rating:

$$\begin{bmatrix} (\theta^1)^T x^1 & (\theta^2)^T x^1 & \dots & (\theta^m)^T x^1 \\ (\theta^1)^T x^2 & (\theta^2)^T x^2 & \dots & (\theta^m)^T x^2 \\ (\theta^1)^T x^{m+1} & (\theta^2)^T x^{m+1} & \dots & (\theta^m)^T x^{m+1} \end{bmatrix}$$

→ We can define 2 more matrices:

$$X = \begin{bmatrix} (x^1)^T \\ (x^2)^T \\ \vdots \\ (x^{m_m})^T \end{bmatrix} \rightarrow \text{Each row containing features of a particular movie}$$

$$\theta \cdot \theta = \begin{bmatrix} (\theta^1)^T \\ (\theta^2)^T \\ \vdots \\ (\theta^{m_w})^T \end{bmatrix} \rightarrow \text{Each row containing weights for the features of a given user}$$

\* Given  $X \cdot \theta$ , we can get the full matrix  $\hat{Y}$  of all predicted ratings of all movies of all users by:  $\hat{Y} = X \cdot \theta^T$

\* ~~This~~ This algorithm is called low rank matrix factorization.

This comes from the property that  $X \cdot \theta^T$  creates a low rank matrix

\* Now we can apply gradient descent to:

$$\min(Y - \hat{Y}) = \min(Y - X \cdot \theta^T)$$

- ⇒ Finding related movies
- For each product  $i$ , we learn a feature vector  $x^i \in \mathbb{R}^m$ : eg:  $x_1 = \text{romance}$ ,  $x_2 = \text{action}$ ,  $x_3 = \text{comedy}$ , . . .
- \* After learning, it might difficult to characterize feature into human understandable form.
- How to find ~~movies~~ of related to movie  $i$ ?

If we have 2 matrices  $x^i$  &  $x^j$ ; we

want to minimize  $\|x^i - x^j\|$  → dist. b/w those 2 movies

eg: 5 most similar movies to movie  $i$

Find the 5 movies  $j$  with smallest  $\|x^i - x^j\|$

$\Rightarrow$  Implementation detail: Mean normalization

→ Say we have a dataset where a user hasn't rated any movie<sup>(us)</sup>:

Movie	$U_1$	$U_2$	$U_3$	$U_4$	$U_5$
$M_1$	5	5	0	0	?
$M_2$	5	?	?	0	?
$M_3$	?	4	0	?	?
$M_4$	0	0	5	4	?
$M_5$	0	0	5	?	?

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

→ Say  $\theta \in \mathbb{R}^2$  as  $m = 2$

want to learn  $\theta^S$

→ Let's analyze the optimization objective

$$\min_{\theta^1, \dots, \theta^{mn}, x^1, \dots, x^{mr}} \left\{ \frac{1}{2} \sum_{(i,j) : r(i,j)=1} [(x^j)^T \theta^i - y^{(i,j)}]^2 + \lambda \sum_{k=1}^m \sum_{i=1}^{m_r} (x_k^i)^2 \right\}$$

Not relevant while minimizing  $\theta^1, \dots, \theta^{mn}$

Not relevant as there are no  $r(i,j)=1$  for  $x^j$  from  $m_r$ .

$$+ \frac{\lambda}{2} \sum_{j=1}^{m_r} \sum_{k=1}^m (\theta_k^j)^2$$

\* So, we're only minimizing  $\boxed{\frac{\lambda}{2} \sum_{j=1}^{m_r} \sum_{k=1}^m (\theta_k^j)^2}$ ,

which for our example can be simplified

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

\* Of course, if the goal is to minimize this term

then  $\theta^5 = \boxed{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}$ . As there's no data, 

the alg will minimize it to the minimum possible value.

\* As, this user will have indirectly explicitly rated all movies to 0, so we won't be able to recommend anything to this user

→ This is where mean normalization comes in!

⇒ Mean normalization

→ We can normalize the data relative to the mean.

→ Firstly, we'll use ~~use~~ a  $Y$  matrix to store the data from previous ratings:

\*  $i^{\text{th}}$  row : Rating for  $i^{\text{th}}$  movie

\*  $j^{\text{th}}$  column : Ratings given by  $j^{\text{th}}$  user

→ We now a  $m \times m$  dimensional column vector  $u$ :

$$u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_{m \times m} \end{bmatrix}$$

where

$$u_i = \frac{\sum_{j: r(i,j)=1} Y_{i,j}}{\sum_j r(i,j)}$$

= Sum of ratings for  $i^{\text{th}}$  movie  
whichever user 'j' has rated it

No. of user 'j' who  
rated the movie  $i$

\*  $u_i$  = Mean of the previous ratings for the  $i^{\text{th}}$  movie.

→ We can normalize the data by subtracting  $\mu$ , the mean rating, from the actual ratings for each user.

→ eg: for our previous example

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ . & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}; \quad \mu = \begin{bmatrix} (5+5+0+0)/4 \\ (5+0)/2 \\ (4+0)/2 \\ (5+4+0+0)/4 \\ (0+0+5+0)/4 \end{bmatrix} = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix}$$

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

→ By subtracting the mean of a movie from its ratings, we've essentially normalized each film to have an average rating of 0!

- \* Now, we can use this new Y matrix to learn  $\theta^*$  &  $\sigma^2$ .
- \* Slight modification to linear regression prediction
  - Include the mean normalization term:
$$(\theta^*)^T x^i + u_i$$
- Now, for a new user the initial predicted value will be equal to the  $u$ -term instead of simply being initialized to 0, which makes for more sense intuitively. (As we're predicting the average rating, which is the <sup>rest</sup> <sub>we can do</sub> we can do)
- eg: For  $\theta_0^*, \theta_1^*, \dots$
- Movies with no ratings: Can apply mean-normalization to rows!
- \* But not relevant, as you shouldn't recommend on unrated movie.