

Week 9

Recommender Systems - Collaborative filtering

Problem motivation

feature vector:	$\theta^{(1)}$	$\theta^{(2)}$	$\theta^{(3)}$	$\theta^{(4)}$	x_1	x_2
movie	Alice (1)	Bob (2)	Card (3)	Dave (4)	(romance)	(action)
Love at first sight	5	5	0	0	0.9	0
Romantic for fools	5	?	?	0	1.0	0.01
~~~~~	?	4	0	?	0.99	0
~~~~~	0	0	5	4	0.1	1.0
~~~~~	0	0	5	4	0	0.9

Suppose we have a dataset where we do not know the values of  $x_1$  &  $x_2$  i.e

	$x_1$	$x_2$	
probability	1.0	0.0	no idea
	?	?	how romantic
	?	?	each movie
	?	?	is. No idea
	?	?	how action
	?	?	picked each
	?	?	movie is

$$\theta^1 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^2 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^3 = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}, \theta^4 = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} \quad x_0 = 1$$

Alice & Bob both really like romance  
 Card & Dave both tell us that they really like action

∴ what value of  $x^{(1)}$  gives!

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

## Week 9

### Recommender Systems - Collaborative filtering

one specific  
movie case:

#### Optimization algorithm

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

Let's say that the users have told us their preferences.

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

regularization term (prevents features from becoming too big)

we want to choose values

squared error term

for the feature vector  $x$  such that the squared error is minimised.

we want to learn the features for all the movies:

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

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

sum over all  $nm$  movies etc

minimise the above will hopefully produce a reasonable set of features for all the movies.



Week 9

# Recommender Systems - Collaborative filtering algorithms

Collaborative filtering optimization objective

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

all pairs of  $i, j$

$$r(i, j) = r(j, i) = 1$$

where  $r(i, j) = 1$

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

sum over all users that have rated that movie.

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

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

for every movie  $i$

minimising  $x^{(1)}, \dots, x^{(nm)}$  and  $\theta^{(1)}, \dots, \theta^{(nu)}$  simultaneously

$$J(x^{(1)}, \dots, x^{(nm)}, \theta^{(1)}, \dots, \theta^{(nu)}) = \frac{1}{2} \sum_{i, j: r(i, j)=1} (x^{(ij)} - \theta^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{nm} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{nu} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\min_{\substack{x^{(1)}, \dots, x^{(nm)} \\ \theta^{(1)}, \dots, \theta^{(nu)}}}$$

instead of  $\theta \rightarrow x \rightarrow \theta \rightarrow$  we are going to minimise with both sets of parameters simultaneously.

$x_0 = 1$  corresponds to an intercept term  
doing it this way however yields:

$$x_0 = 1; \begin{cases} x \in \mathbb{R}^n \\ \theta \in \mathbb{R}^n \end{cases}, x \in \mathbb{R}^{n+1}$$

Week 9

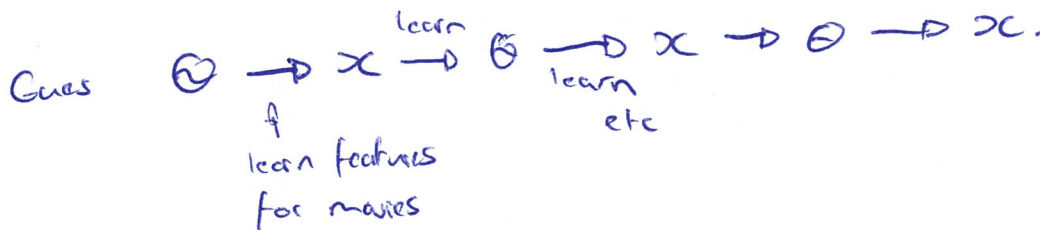
## Recommender Systems - Collaborative Filtering

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

also if we know the features  $x^{(i)}, \dots, x^{(n)}$ , we can learn  
the parameters  $\theta$

also

Given  $\theta^{(i)}, \dots, \theta^{(n)}$ , can estimate  $x^{(i)}, \dots, x^{(n)}$



actually works ✓✓

Week 9

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

## Recommender Systems - Collaborative filtering algorithm

↑  
∴ regularizing everything.

- 1.) Initialize  $x^{(1)}, \dots, x^{(nm)}, \theta^{(1)}, \dots, \theta^{(nm)}$  to small random values
- 2.) Minimise  $J(x^{(1)}, \dots, x^{(nm)}, \theta^{(1)}, \dots, \theta^{(nm)})$  using gradient descent (or some other optimization algorithm e.g. for every  $j=1, \dots, n_u, i=1, \dots, n_m$  :

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

from the partial derivatives of the cost function

- 3) for a user with parameters  $\theta$  and a movie with (learned) features  $x$ , predict a star rating of  $\theta^T x$ .  
 $(\theta^{(j)})^T (x^{(i)})$  is going to rate movie  $i$



# Week 9

## Recommender Systems: Vectorization low rank matrix factorization

### Collaborative filtering

movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at first	5	5	0	0
romance forever	5	?	?	0
~~~~~	?	4	0	?
~~~~~	0	0	5	4
~~~~~	0	0	5	?

$n_m = 5$ (movies)

$n_u = 4$ (users)

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

y_{ij} movies i users j

Collaborative filtering

$$Y = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

predicted ratings:

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

$(\theta^{(1)})^T (x^{(1)})$ predicted rating user 1 on movie 1

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

$$\oplus = \begin{bmatrix} - (\theta^{(1)})^T - \\ - (\theta^{(2)})^T - \\ \vdots \\ - (\theta^{(n_u)})^T - \end{bmatrix}$$

capital θ

Collaborative filtering is also called Low-rank matrix factorization.

Week 9

Recommender systems: vectorization, low rank matrix factorization

finding 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 to movie i ?

movie i has a feature vector $\|x^{(i)} - x^{(j)}\|$ ← different movie j
if the difference between movie i & j is minimised, this is a pretty strong indication that they are similar...

eg. 5 most similar movies to movie i :

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

Week 9

Recommender Systems - Implementational detail: mean normalization

Users who have not rated any movies

lets say we add a 5th user "eve"...

$$n=2 \text{ (romance (action))}, \Theta^{(5)} \in \mathbb{R}^2 \quad \Theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

the first term of the cost function plays no role because there are no cases where $(i,j) \text{ s.t. } r(i,j) \neq 1$

the only term that affects is $= \frac{1}{2} \sum_{j=1}^m \sum_{k=1}^n (\theta_k^{(j)})^2$

so we want to minimize $\frac{\lambda}{2} [(\theta_1^{(5)})^2 + (\theta_2^{(5)})^2]$

$$\therefore (\Theta^{(5)})^T x^{(i)} = 0$$

\therefore eve will rate all movies as 0 stars.

Mean Normalization:

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

eve's unrated movies.

for user j , on movie i predict:

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

user 5 (eve):

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

$$(\Theta^{(5)})^T (x^{(i)}) + \mu_i = 0$$

\therefore eve's parameters will be the average

the average value that each movie obtains

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

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

subtract off the mean rating.

pretend that this was the data that I got from my users.

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

stay the same