

Predicting movie ratings

↗ users - n_u
 ↗ items - n_m
 ↗ ratings - $r(i, j) = 1$ if j rates i

Predicting movie ratings

User rates movies using one to five stars

Ratings				
★				
★	★			
★	★	★		
★	★	★	★	
★	★	★	★	★

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)
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	?

$$n_u = 4$$

$$n_m = 5$$

$r(1,1) = 1$ → rated movie
 $r(3,1) = 0$ → not rated

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

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

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

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

+ features		
romance	action	
0.9	0	
1.0	0.01	
0.99	0	
0.1	1.0	
0	0.9	

$$\begin{aligned}
 x^{(1)} &= \begin{bmatrix} 0.9 \\ 0 \end{bmatrix} \\
 x^{(3)} &= \begin{bmatrix} 0.99 \\ 0 \end{bmatrix}
 \end{aligned}$$

For User 1 - Alice movie i

$$w^{(1)} x^{(1)} + b^{(1)} \leftarrow \text{linear regression}$$

$$\begin{aligned}
 w^{(1)} &= \begin{bmatrix} 5 \\ 0 \end{bmatrix} & b^{(1)} &= 0 & x^{(3)} &= \begin{bmatrix} 0.99 \\ 0 \end{bmatrix}
 \end{aligned}$$

$$w^{(1)} \cdot x^{(3)} + b^{(1)} = 4.95$$

for user j ; predict user j 's rating for movie i as

$$w^{(j)} \cdot x^{(i)} + b^{(j)}$$

Cost Function

$$\underset{\substack{\text{parameters} \\ \text{movies user } j \text{ has rated}}}{\underset{i: r(i,j)=1}{\frac{1}{2m^{(j)}} \sum \left(w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)} \right)^2}} + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (w_k^{(j)})^2$$

for regularization
recommendation systems

Cost function

To learn parameters $w^{(j)}, b^{(j)}$ for user j :

$$J(w^{(j)}, b^{(j)}) = \frac{1}{2} \sum_{i: r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (w_k^{(j)})^2$$

To learn parameters $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$ for all users :

$$J\left(\begin{matrix} w^{(1)}, \dots, w^{(n_u)} \\ b^{(1)}, \dots, b^{(n_u)} \end{matrix}\right) = \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i,j)=1} \underbrace{(w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2}_{f(x)} + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Collaborative Filtering :-

Problem motivation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	x_1 (romance)	x_2 (action)
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	?	?	?

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

$$b^{(1)} = 0, b^{(2)} = 0, b^{(3)} = 0, b^{(4)} = 0$$

4:39

using $w^{(j)} \cdot x^{(i)} + b^{(j)}$

$$\left. \begin{array}{l} w^{(1)} \cdot x^{(1)} \approx 5 \\ w^{(2)} \cdot x^{(1)} \approx 5 \\ w^{(3)} \cdot x^{(1)} \approx 0 \\ w^{(4)} \cdot x^{(1)} \approx 0 \end{array} \right\} \rightarrow x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

↑
guess

$$J(x^{(i)}) = \frac{1}{2} \sum_{j: r(i,j) \neq 1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2$$

$$+ \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

To Learn $x^{(1)}, x^{(2)}, \dots, x^{(nm)}$ features,

$$J = \frac{1}{2} \sum_{i=1}^{nm} \sum_{j: r(i,j) \neq 1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{nm} \sum_{k=1}^n (x_k^{(i)})^2$$

Full Cost Function;

→ (+) ←

Minimize \rightarrow Gradient Descent;

Linear regression but

$$w_i = w_i - \alpha \frac{\partial}{\partial w_i} J(w, b, x)$$

$$b = b^{(i)} - \alpha \frac{\partial}{\partial b^{(i)}} J(w, b, x)$$

$$x_k^{(i)} = x_k^{(i)} - \alpha \frac{\partial}{\partial x_k^{(i)}} J(w, b, x)$$

Binary Labels

Binary labels

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

Did user like item? 1 & 0
Spend at least 30s? 1 & 0
Clicked?

For binary j

$$g(w^{(j)} \cdot x^{(i)} + b^{(j)})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

Cost function for binary application

Loss for binary labels = $y^{(i,j)} f_{(w,b,x)}(x) = g(w^{(j)} \cdot x^{(i)} + b^{(j)})$

Cost function for binary application

Previous cost function:

$$\frac{1}{2} \sum_{(i,j): r(i,j)=1} \underbrace{(w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2}_{f(x)} + \frac{\lambda}{2} \sum_{l=1}^{n_m} \sum_{k=1}^n (x_k^{(l)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Loss for binary labels $y^{(i,j)} : f_{(w,b,x)}(x) = g(w^{(j)} \cdot x^{(i)} + b^{(j)})$

$$L(f_{(w,b,x)}(x), y^{(i,j)}) = -y^{(i,j)} \log(f_{(w,b,x)}(x)) - (1 - y^{(i,j)}) \log(1 - f_{(w,b,x)}(x)) \quad \leftarrow \text{Loss for single example}$$

(cost function): $J(w, b, x) = \sum_{(i,j): r(i,j)=1} L(f_{(w,b,x)}(x), y^{(i,j)}) \quad \sum g(w^{(j)} \cdot x^{(i)} + b)$

Mean Normalization

For new user not rated;

$$\left[\begin{array}{cccc} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 6 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{array} \right] \mu = \left[\begin{array}{c} \text{avg.} \\ 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{array} \right]$$

↓ subtract

Normalize
Matrix rows
(can also
do columns)

$$\left[\begin{array}{ccccc} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{array} \right]$$

For user j on movie i predict;

$$w^{(j)} x^{(i)} + b^{(j)} + \mu_i$$

$$\text{now if } w^s \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \times b^{(s)} = 0$$

no matter b^s

$$\underbrace{w^s}_{0} \times \underbrace{b^{(s)}}_{0} + \mu_i$$

$$= 2.5$$

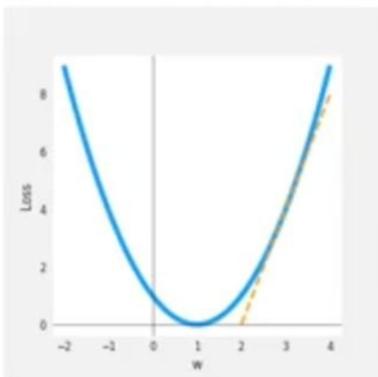
Collaborative Filtering Algorithm in Tensorflow
→ (No loss layer N.N etc) → Without fit

$$J = (wx - 1)^2$$

Gradient descent algorithm
Repeat until convergence

$$w = w - \alpha \frac{d}{dw} J(w, b)$$

Fix b = 0 for this example



Custom Training Loop

```
w = tf.Variable(3.0)
x = 1.0
y = 1.0 # target value
alpha = 0.01
```

```
iterations = 30
for iter in range(iterations):
    # Use TensorFlow's GradientTape to record the steps
    # used to compute the cost J, to enable auto differentiation.
    with tf.GradientTape() as tape:
        fwb = w*x
        costJ = (fwb - y)**2

    # Use the gradient tape to calculate the gradients
    # of the cost with respect to the parameter w.
    [dJdw] = tape.gradient(costJ, [w])

    # Run one step of gradient descent by updating
    # the value of w to reduce the cost.
    w.assign_add(-alpha * dJdw)
```

tf.variables require special function to modify

Auto Diff
software pkg
→ Autograd

Implementation in TensorFlow

Adam
optimization
(not grd. desc)

```
# Instantiate an optimizer.
optimizer = keras.optimizers.Adam(learning_rate=1e-1)
```

iterations = 200

for iter in range(iterations):

Use TensorFlow's GradientTape

to record the operations used to compute the cost

with tf.GradientTape() as tape:

$J(w, b, x)$ # Compute the cost (forward pass is included in cost)
cost_value = cofiCostFuncV(X, W, b, Ynorm, R,
num_users, num_movies, lambda)

normalized rating

Use the gradient tape to automatically retrieve
the gradients of the trainable variables with respect to
the loss

grads = tape.gradient(cost_value, [X, W, b])

derivative of cost func

Run one step of gradient descent by updating
the value of the variables to minimize the loss.
optimizer.apply_gradients(zip(grads, [X, W, b]))

Adam with gradients

Dataset credit: Harper and Konstan, 2019. The MovieLens Datasets: History and Context

$$w = w - \alpha \frac{\partial J}{\partial w}$$

$$b = b - \alpha \frac{\partial J}{\partial b}$$

$$X = X - \alpha \frac{\partial J}{\partial X}$$

Find related items;

Find item k with $x^{(k)}$

↳ similar to $x^{(i)}$

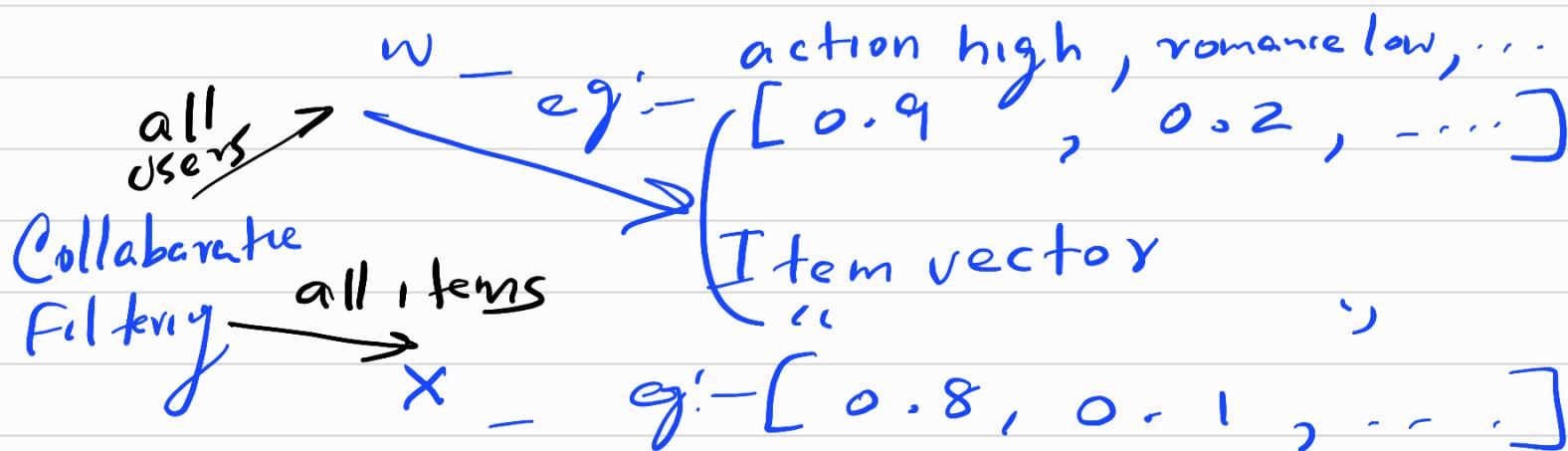
$$\sum_{l=1}^n (x_{li}^{(k)} - x_{li}^{(i)})^2 \Rightarrow \|x^{(k)} - x^{(i)}\|^2$$

closest vector dist.

Limitations of Collaborative filtering;

- Cold - Start Problem \rightarrow new item few users rated
or new user few rated
- Use side information

two vectors \rightarrow User vector



$w \times b$

Collaborative Filtering

Reccomend based on ratings of users giving similar ratings

eg:- $x_u^{(j)}$ for user j

$x_m^{(i)}$ for movie i

feature vectors
diff. in size

Content Based Filtering

recommend based on features of user & item

- age/gender/...
Movies watched

- year/reviews/genre
Avg. rating - - -

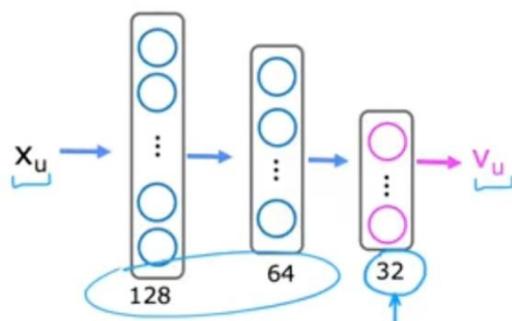
Content Based; match movie i with user j

$$w^{(j)} \cdot x^{(i)} + b$$

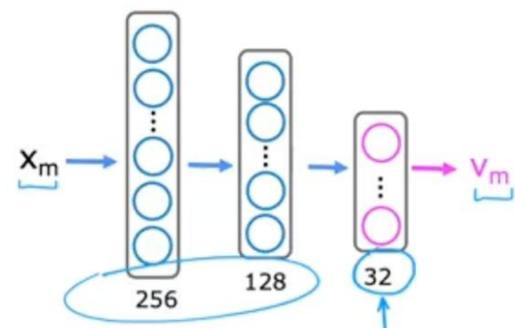
from $x_u^{(i)}$ $\rightarrow v_u^{(j)}$ \quad $v_m^{(i)} \leftarrow$ from $x_m^{(i)}$

Neural network architecture

$x_u \rightarrow v_u$ User network



$x_m \rightarrow v_m$ Movie network



Prediction: $v_u \cdot v_m$

$g(v_u \cdot v_m)$ to predict the probability that $y^{(i,j)}$ is 1

Cost Function;

$$J = \sum_{(i,j); r(i,j)=1} (V_u^{(i)} \cdot V_m^{(i)} - y^{(i,j)})^2 + \text{NN regularization}$$

similar movie;

$$\| V_m^{(k)} - V_m^{(i)} \|^2 \rightarrow \text{small}$$

* pre compute ahead of time

* feature selection

* computational reg.

for large sets

Retrieval + Ranking

find 10 most similar

3 genres most viewed
top 20 of country

combine (remove dup) & k / already purch. (\leftarrow)

more results performance ↑

Rank
→ finetuned model

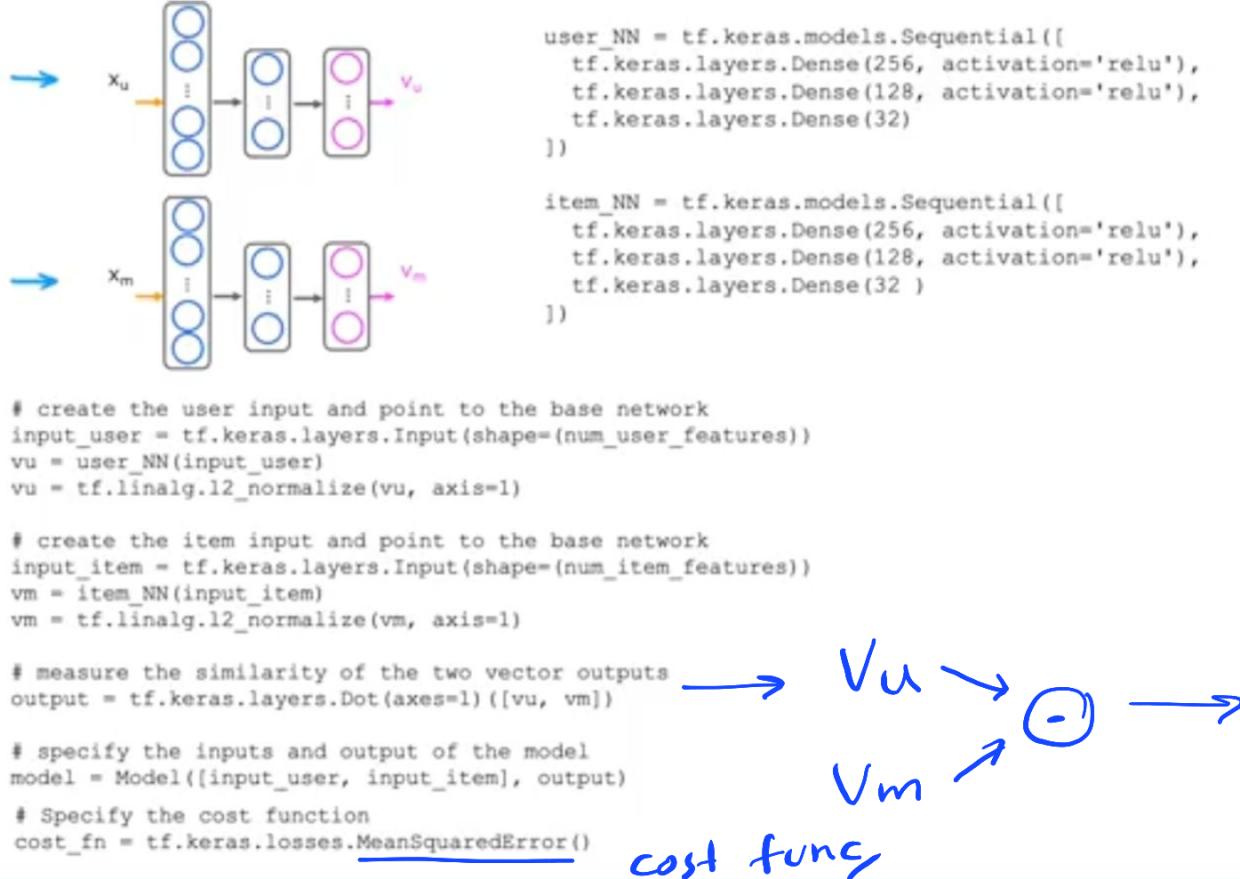
Ethical concerns;

→ Advertising?

Transparency?

Maximize profit?

Amelioration → Do not accept ads
from exploitative businesses



Feature vectors can be any size
as long as $v_u = v_m$