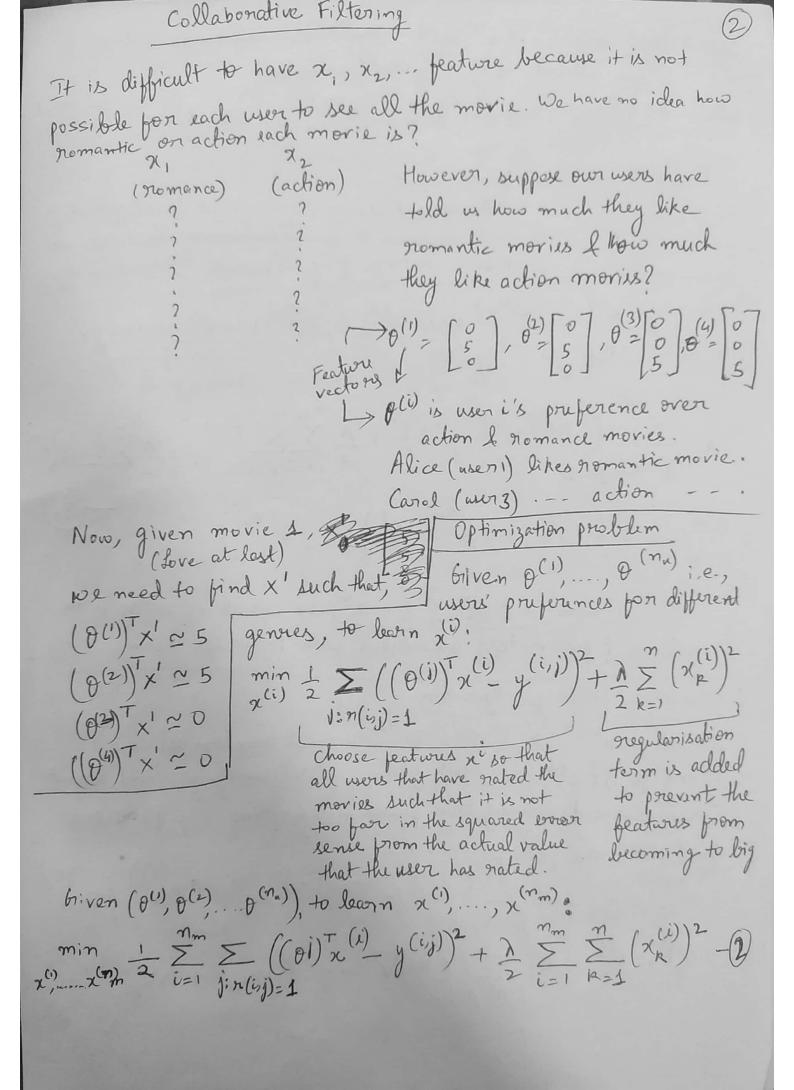
Movie Alice Bob Carol Dave (x) x2 Forestores Foresto	Content Based Recommendations
Romence Forever 5 ?? ? ? 0 1.0 0.01 Cute Papies ?? 4 0 ?? 0 99 0 Non-stop Car closes 0 0 5 4 0.1 1.0 Swords ys 0 0 5 ?? 0 0.9 kanates Our objective is to populate the ? with ratings that they did not wath wath watch. In a number of mories The number of mories The member of movies rated movie i. Y(i,i) = noting given by user j to movie i (defined only if rise) The member of movies rated by user j. The is the feature corresponding to genne i. No is the feature addociated with intercept term. Jet x be the feature vector associated with movie i. Thus, X' = [0.9] / X = [0.9] / X = [0.1] / X = [0.9] Regression Broblem is: For each user j, learn a parameter of belonge with the intercept term, predict user j as rating movie i with (D) Tx'. Suppose 0' is [5]. Since X = [0.93], Alice would	Morrie Alice Bob Carol Dave X1 X2
Cate Popies of fore Non-stop can closes O 0 5 4 0.1 Non-stop can closes O 0 5 7 0.1 O 0.9 Remates O 0 9 Remates O 0 9 Remates O 0 1 1.0 O 0.9 Remates O 0 9 Remates O 0 1 1.0 O 0.9 Remates O 0 9 Remates O 0 1 1.0 O 0 9 Remates O 1 1.0 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 9 O 0 9 Remates O 0 1 1.0 O 0 1 1	Fore at 5 5 0 0
Non-stop Car does Non-stop Car does O 5 4 O 1 1,0 Swords y's 0 0 5 7? 0 0.9 karates O w objective is to populate the ? with ratings that they did not wath wath watch. In = number of wers nm = number of mories ni = 1 if wer i has rated morie i. y(ini) = nating given by wer i to movie i (defined only if risi=1) m; = member of movies rated by user j. xi = is the feature corresponding to genre i. No is the feature addociated with intercept term. Jet xi be the feature vector associated with movie i. Thus, X' = [0,9] , x²=[1,0] , x³=[0,99], x'=[0,1] , x5=[0,9] Regression brobben is: For each user j, learn a parameter of belonge other the Rn+1, where nix the no of features except the feature associated with the intercept term, predict user j as rating movie i with (0) Tx' . Suppose of is [5]. Since x²= [0.93], Alice would	Romance 5 [?] [?] 0 1.0 0.01
Swords v/s 0 0 5 [?] 0 0.9 kanates Our objective is to populate the ? with natings that they did not wath watch. In = number of users nm = number of movies n(i) = 1 if user j has nated movie i. not in it is in the pature associated with movie i (defined only if n's=1) m: = number of movies nated by user j. xi = is the peature corresponding to genre i. Xo is the peature associated with intercept term. Jet x' be the feature vector associated with movie i. Thus, x' = [0.9] , x' = [1.0] , x = [0.9] , x' = [0.1] , x = [0.9] Regnession brobbem is: For each user j, learns a parameter of belong with the Rn+1, where n is the no. of peatures except the peature associated with the intercept term, predict user j as nating movie i with (0.1) x' = [0.9], Alice would	
our objective is to populate the ? with natings that they did not wath watch. In = number of wers nm = number of movies nm' = number of movies notified only if rises y(i,i) = nating given by user j to movie i (defined only if rises) m = member of movies nated by user j. xi = is the feature corresponding to genne i. Xo is the feature addociated with intercept term. Jet xi be the feature vector associated with movie i. Thus, X' = [0.9] , X^2 = [0.9] , X^4 = [0.1] , X^5 = [0.9] Regression Problem is: For each user j, learn a perameter of belonge with the intercept term, predict user j as rating movie i with 10 I Ini Suppose 0 is [5]. Since X^2 = [0.93], Alice would	
our objective is to populate the? with natings that they did not wath watch. In a mumber of users The mumber of movies The pating given by user j to movie i (defined only if rise) m; = mumber of movies noted by user j. m; = mumber of movies noted by user j. xi = is the feature corresponding to genne i. Xo is the feature addociated with intercept term. Jet xi be the feature vector associated with movie i. Thus, X' = [0.9] / X^2 = [0.0] / X = [0.9] , X' = [0.1] / X = [0.9] Regression Problem is: For each user j, learn a parameter of belonge with the intercept term, predict user j as grating movie i with 18 1) The suppose of is [5]. Since X = [0.93], Alice would	
nu = number of users nm = number of mories ni = 1 if wer j has rated morie i. y(i,j) = nating given by user j to movie i (defined only if rides) m; = member of movies rated by user j. xi = is the feature corresponding to genre i. Xo is the feature addociated with intercept term. Jet xi be the feature vector associated with morie i. Thus, X' = [1,9], X^2 = [0,9], X' = [0,1], X^5 = [0,9] Regression Problem is: For each user j, learn a parameter of belonge the R ⁿ⁺¹ , where nix the no. of feature except the feature associated with the intercept term, predict user j as nating movie i with 10) This suppose of is [5]. Since X ³ = [0,93], Alice would	Our objective is to populate the? with natings that they did not wath watch.
et of	nm = number of movies g(i,i) = 1 if user j has nated movie i. y(i,i) = noting given by user j to movie i (defined only if rid=1) m; = member of movies nated by user j. xi = is the feature corresponding to genre i. to is the feature addociated with intercept term. Jet xi be the feature vector associated with movie i. Thus, X' = [1.9], X²=[1.0], X³=[0.99], X'=[0.1], X̄=[0.9] Regression Problem is: For each user j, learn a parameter of belonge then the Rn+1, where nix the no. of peatures except the feature associated with the intercept term, predict user j as nating movie i with (0) Txi Suppose 0' is [5]. Since X³=[0.93], Alice would

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

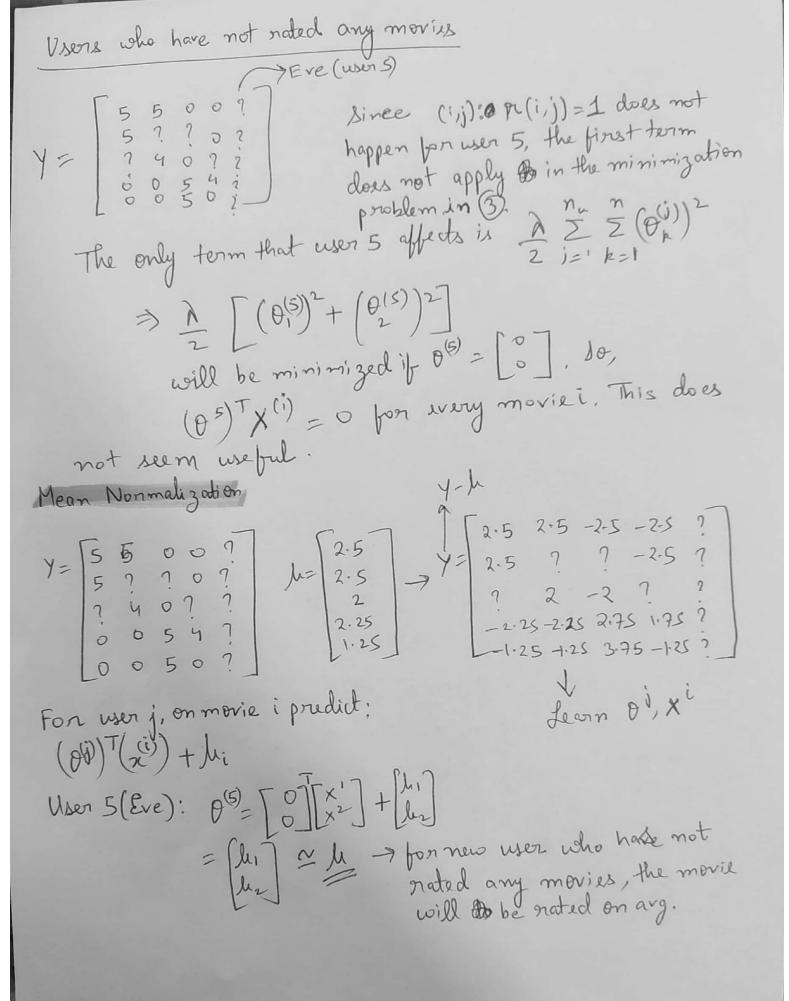
 $\theta^{(j)} = 1$
 $\theta^{(j)} = 1$



can estimate o(1), x(2), ..., p(nu) (parameters) Collaborative Filtering -> Griven O(1), ..., O(nu), can estimate x(1), ..., x(nm) Gruss & -> fearm X -> & -> X..... Collaborative Filtering () Content based filtering will be used to simultaneously learn about movil natings of parameters. It used simultanes Minimizing x(1), ..., x(nm) and o(1), ..., o(nu) simultaneously: insted of going back of forth. $J(\chi^{(i)},\ldots,\chi^{(n_m)},\theta^{(i)},\ldots,\theta^{(n_u)})$ $=\frac{1}{2}\sum_{(i,j): r(i,j)=1}^{n} ((\theta^{(i)})^{T} \chi^{(i)} - y^{(i,j)})^{2} + \frac{\lambda}{2}\sum_{i=1}^{n} \sum_{k=1}^{n} (\chi^{(i)}_{k})^{2} + \frac{\lambda}{2}\sum_{j=1}^{n} \sum_{k=1}^{n} (\theta^{(j)}_{k})^{2}$ min $\mathcal{T}\left(\chi^{(1)},\ldots,\chi^{(n_m)},\Theta^{(1)},\ldots,\Phi^{(n_u)}\right)$ * Here x ER" & DER". We are doing away with peatures of perameters associated with intercept term. Collaborative filtering algorithm

1. Initialize $\chi^{(1)}, \dots, \chi^{(n_m)}, 0, \dots, 0$ (nu) to small random values.

2. Minimize $J(\chi^{(1)}, \dots, \chi^{(n_m)}, \eta^{(1)}, \eta^{(n_m)}, \eta^{(1)}, \eta^{(n_m)})$ using gradient descent (or an advanced optimization problem algorithm). E.g. for every



Alternative to Collaborative Filtering Take all the elements of the table and create a matrix Y = \begin{align*} 5 & 5 & 0 & \begin{align*} \int \text{minings} & \text{matrix of the table and create a matrix Y} \\ = \begin{align*} 5 & 5 & 0 & \text{Indicted natings} & \text{minings} &	3. For a user with parameters of and a movie with (learned) 3
Take all the elements of the table and create a matrix y = \[\begin{align*} 5 & 5 & 0 & \cdot & \frac{\text{Predicted}}{\text{Predicted}} & \text{ratings} \; \begin{align*} 5 & 7 & 7 & \cdot & \c	beatures x, predict a star rating o'x.
Take all the elements of the table and create a matrix y = \[\begin{align*} 5 & 5 & 0 & \cdot & \frac{\text{Predicted}}{\text{Predicted}} & \text{ratings} \; \begin{align*} 5 & 7 & 7 & \cdot & \c	Alternative to Collaborative Filtering
Finding related movies For each product i, we learn a feature vector $\chi^{(i)} \subset \mathbb{R}^n$ $\chi_{i} = 1$ $\chi_{i} = 1$ $\chi_{i} = 1$ Finding related movies For each product i, we learn a feature vector $\chi_{i} \subset \mathbb{R}^n$ $\chi_{i} = 1$ Here to find movies $\chi_{i} = 1$ Here to find movies $\chi_{i} = 1$ $\chi_{i} =$	Take all the elements of the table and conste a matrix of
fow rank matrix $7 \times . D$ where $(D^{j})^{T} \times i$ is the rating given by user j to morie i . $ X = \begin{bmatrix} -(x^{(1)})^{T} \\ -(x^{(2)})^{T} \end{bmatrix} \qquad D = \begin{bmatrix} 0 \\ 0 \end{bmatrix} D^{(2)} \qquad D^{(n)} $ Finding related movies For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^{n}$ $x_{i=1} = x_{i} = x$	= 5500 [nedicted natings: 5770 [(01)] T(1) (02) T(2) (01) T(2) (01) T(2) (02) T(2) (02) T(2)
$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \end{bmatrix}$ $= \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(m)})^T \end{bmatrix}$ $= \begin{bmatrix} -(x^{(m)})^T \\ -(x^{(m$	7000) 20 1000
$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \end{bmatrix}$ $= \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(m)})^T \end{bmatrix}$ $= \begin{bmatrix} -(x^{(m)})^T \\ -(x^{(m$	fow rank mothine 7X. By where (D) is the rating given by user j to morie i.
For each product i, we learn a feature vector $\chi^{(i)} \in \mathbb{R}^n$ χ_{12} romance, χ_{2} = action, χ_{3} = comedy, χ_{4} = How to find movies j related to movie i ? $ \chi^{(i)} - \chi^{(j)} $ small \Rightarrow movie j f i are similar.	$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \end{bmatrix}$ $Q = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ $Q(n)$
For each product i, we learn a feature vector $\chi^{(i)} \in \mathbb{R}^n$ χ_{12} romance, χ_{2} = action, χ_{3} = comedy, χ_{4} = How to find movies j related to movie i ? $ \chi^{(i)} - \chi^{(j)} $ small \Rightarrow movie j f i are similar.	Finding related movies
	For each product i, we learn a feature vector $\chi^{(i)} \in \mathbb{R}^n$ χ_{12} romance, χ_{2} = action, χ_{3} = comedy, χ_{4} = How to find movies j related to movie i ? $ \chi^{(i)} - \chi^{(j)} $ small \Rightarrow movie j f i are similar.