

# SVD For matrix Completion Pg 1

The problem:

We have an incomplete Matrix (R)

$$R = \begin{bmatrix} 2 & 1 & \boxed{?} & 5 & 3 \\ 1 & 4 & 3 & 2 & \boxed{?} \\ 2 & \boxed{?} & 4 & 3 & 3 \end{bmatrix}$$

Can we predict the incomplete entries?

⇒ We cannot develop any system that can provide an exact answer

⇒ But we can develop a protocol that can give a very good estimates for the missing values

⇒ But how is that possible?

⇒ ~~Be~~ we presume that the entries of the matrix is an outcome of a physical process & not ~~beac~~ random numbers.

⇒ Therefore the idea is to find ~~the~~ the latent space of the physical process by the available entries. Then use ~~these~~ these latent variables to estimate the missing entries.

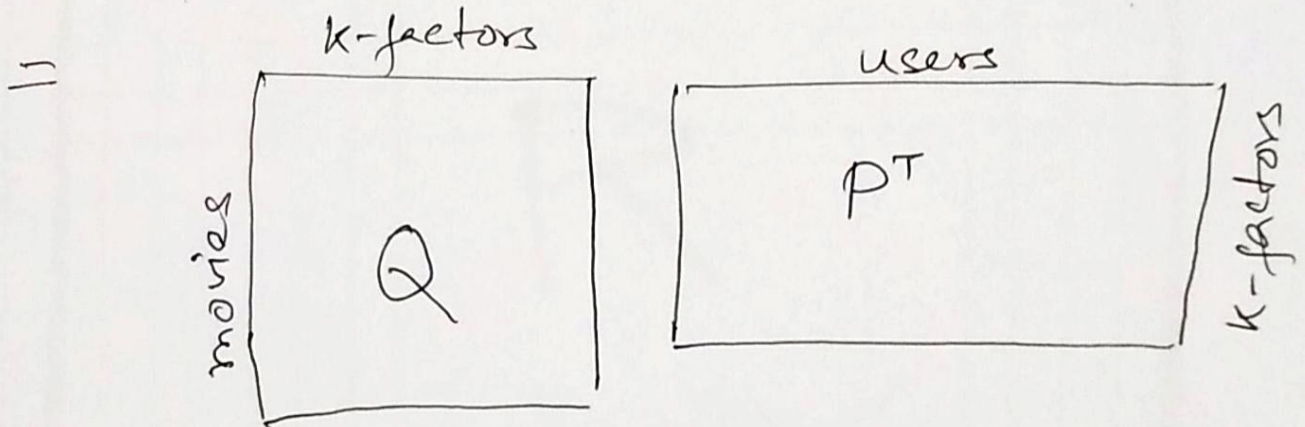
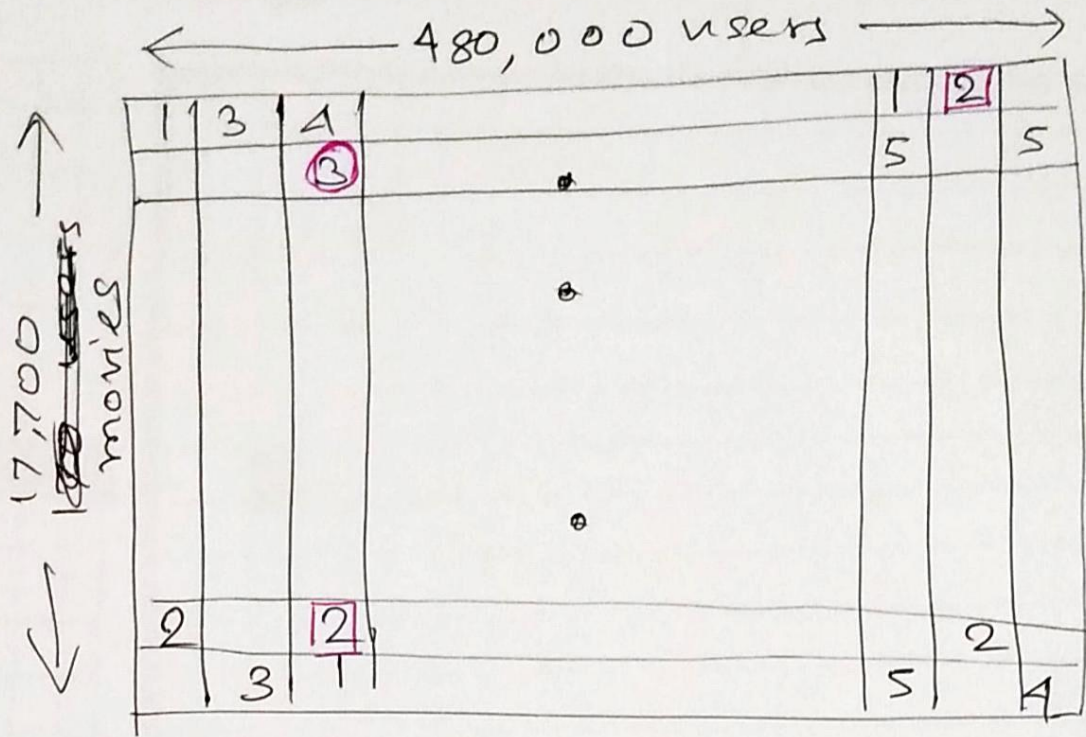


SVD & Matrix Completion: Pg 2

This is a billion \$ business:

## Netflix movie rating predictions:

R



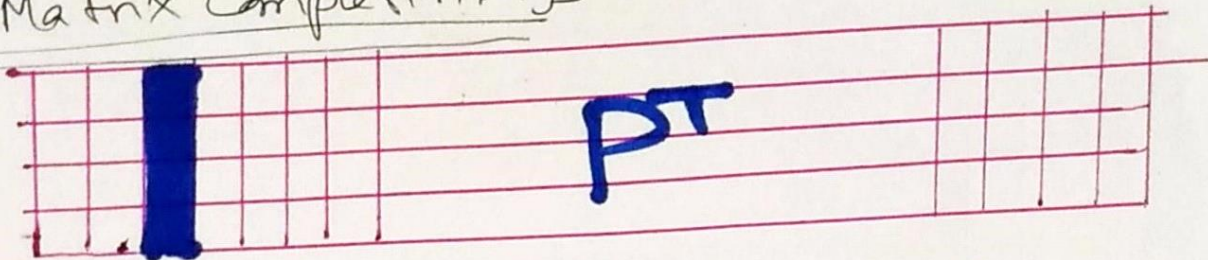
Thin  $k$ -factors

SVD:  $A = U \Sigma V^*$  or  $U \Sigma V^T$

Problem: we have missing values in  $R$



SVD Matrix completion:  $P_{53}$  '  $P$  &  $Q$  are thin matrices



$R$

$$R_{ij} = \sum_k Q_{ik} P_{kj}^T$$

$$R_{73} = \sum_{k=1}^4 Q_{7k} P_{k3}^T$$

$$= \cancel{Q_{71} P_{13}^T} + \cancel{Q_{72} P_{23}^T} + \cancel{Q_{73} P_{33}^T} + \cancel{Q_{74} P_{43}^T}$$

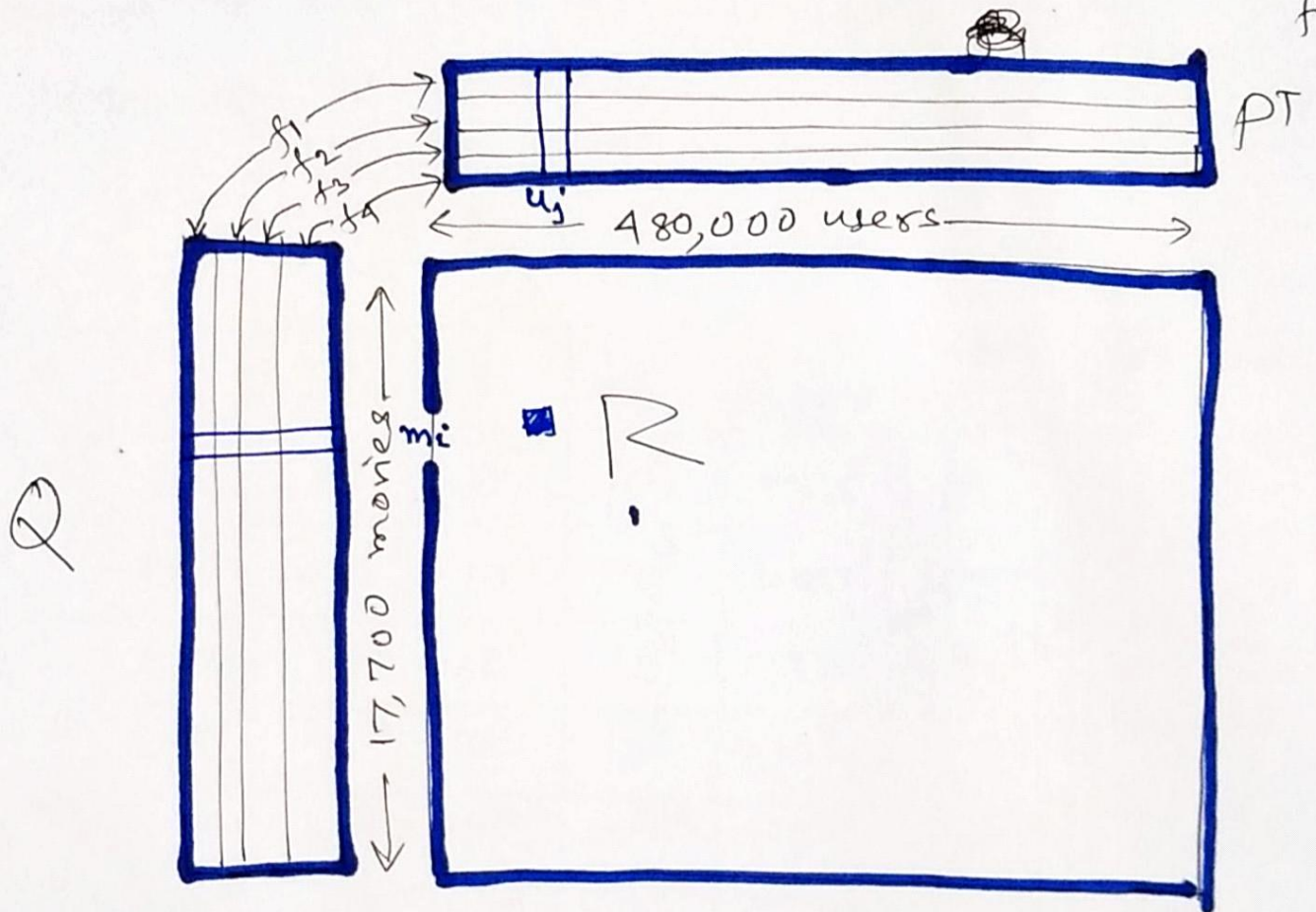
$$= \vec{Q}_7 \cdot \vec{P}_2^T$$

$$= Q_{71} P_{13}^T + Q_{72} P_{23}^T + Q_{73} P_{33}^T + Q_{74} P_{43}^T$$



# SVD & Matrix Completion pg 4

Factors  $k$ :  $f_1$  sci-fi,  $f_2$  romance,  $f_3$  comedy,  $f_4$  crime-thriller etc.



$$R_{ij} = Q_i \cdot \vec{Q}_i \cdot \vec{P}_j^T$$

$$= \sum_k Q_{ik} P_{kj}^T$$

$$= \sum_k Q_{ik} P_{jk}$$

← Note: transpose removed 'coz index exchanged

But what to do with missing values in  $R$ ?

⇒ We can not perform the SVD in the traditional sense: Numerical methods will fall apart!!!



## SVD & Matrix Completion: Pg 5

Let us perform a back-of-the-envelope calculation to understand whether it is a solvable problem or not!

Size of  $R \rightarrow m \times n$   
" "  $Q \rightarrow m \times k$   
" "  $P^T \rightarrow k \times n$

m	n	k	R	P	Q	Total Storage (P+Q)	Compression ratio R/(P+Q)	Comment
$10^3$	$10^3$	100	$10^6$	$10^5$	$10^5$	$2 \times 10^5$	5	
$10^6$	$10^6$	100	$10^{12}$	$10^8$	$10^8$	$2 \times 10^8$	5000	
$10^9$	$10^9$	100	$10^{18}$	$10^{11}$	$10^{11}$	$2 \times 10^{11}$	$5 \times 10^6$	
$10^9$	$10^5$	100	$10^{14}$	$10^7$	$10^{11}$	$\sim 10^{11}$	$10^3$	
$10^{23}$	$10^{14}$	100	$10^{27}$	$10^{16}$	$10^{25}$	$\sim 10^{25}$	$10^2$	
$10^{23}$	$10^{23}$	100	$10^{46}$	$10^{25}$	$10^{25}$	$\sim 10^{25}$	$10^{21}$	

→ Number of unknowns  $\sim 10^{11}$

Even if 1% Entries in R is available

We have  $10^{14} \times 10^{-2} = 10^{12}$  equations available

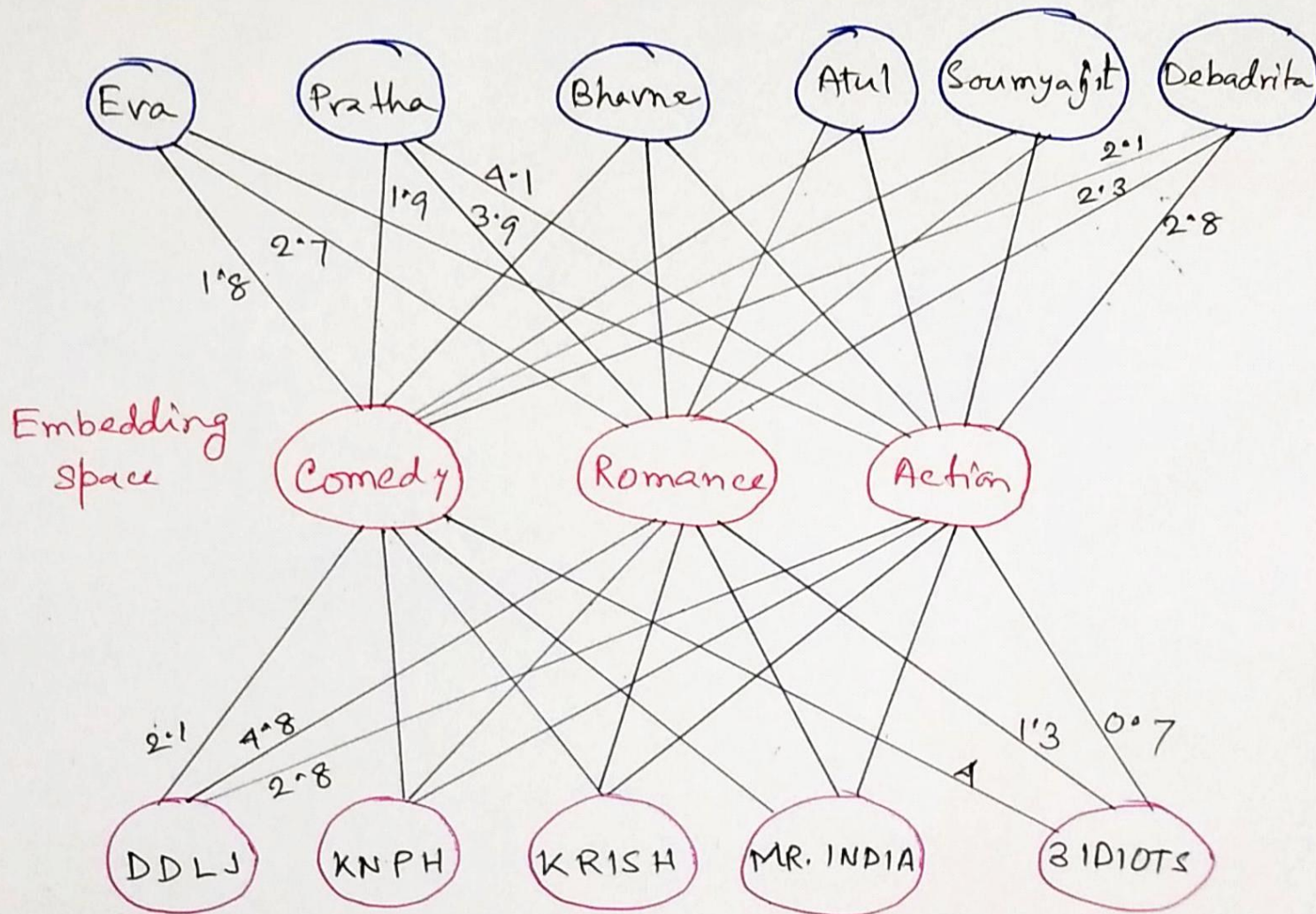
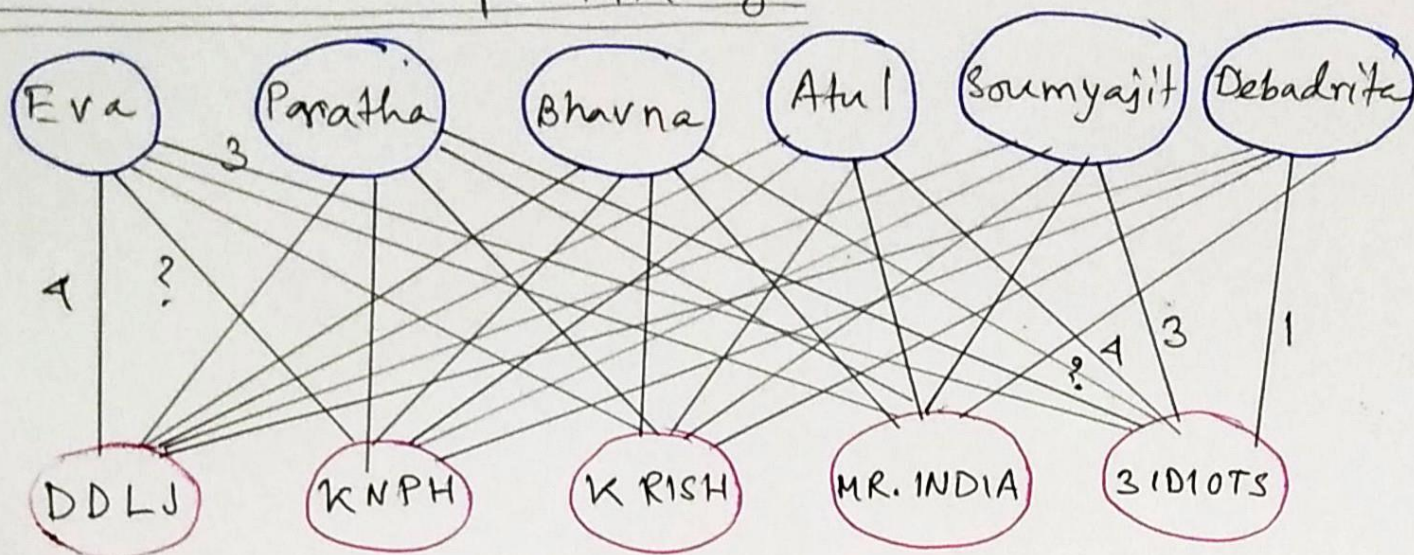
⇒ So ~~these~~ we still have nearly 10 times more equations available than needed

⇒ These systems of equations must be solvable.

⇒ Solution method: → one possible candidate: Gradient descent



# SVD & Matrix completion : Pg 6



Discuss the connection with Auto-encoders.