

Topics:

① Recap of MF for RecSys

✓ ② NMF

③ Netflix Priize solution → Paper

✓ ④ Code: RecSys

⑤ MF: feature engineering

→ Text ✓
→ Images ✓
→ Entities ✓
user, webpage,
product, video ...

{ ⑥ Market-basket Analysis
(ARM, Apriori algo) } ✓

Recap

Given

\sim

$$\tilde{A}_{n \times m}$$

Found

2 Matrices

$$B_{n \times d} \cdot C_{d \times m}$$

$$\begin{bmatrix} 1 & 2 & 3 & \dots & d \end{bmatrix} \xrightarrow{\leftarrow u_i^T} \begin{bmatrix} 1 \\ 2 \\ \vdots \\ n \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & \dots & j & \dots & m \end{bmatrix} \xrightarrow{\uparrow i_j \downarrow} \begin{bmatrix} 1 \\ 2 \\ \vdots \\ d \end{bmatrix}$$

$$I_j \in \mathbb{R}^d = c_j$$

$$\tilde{b}_i = u_i \in \mathbb{R}^d$$

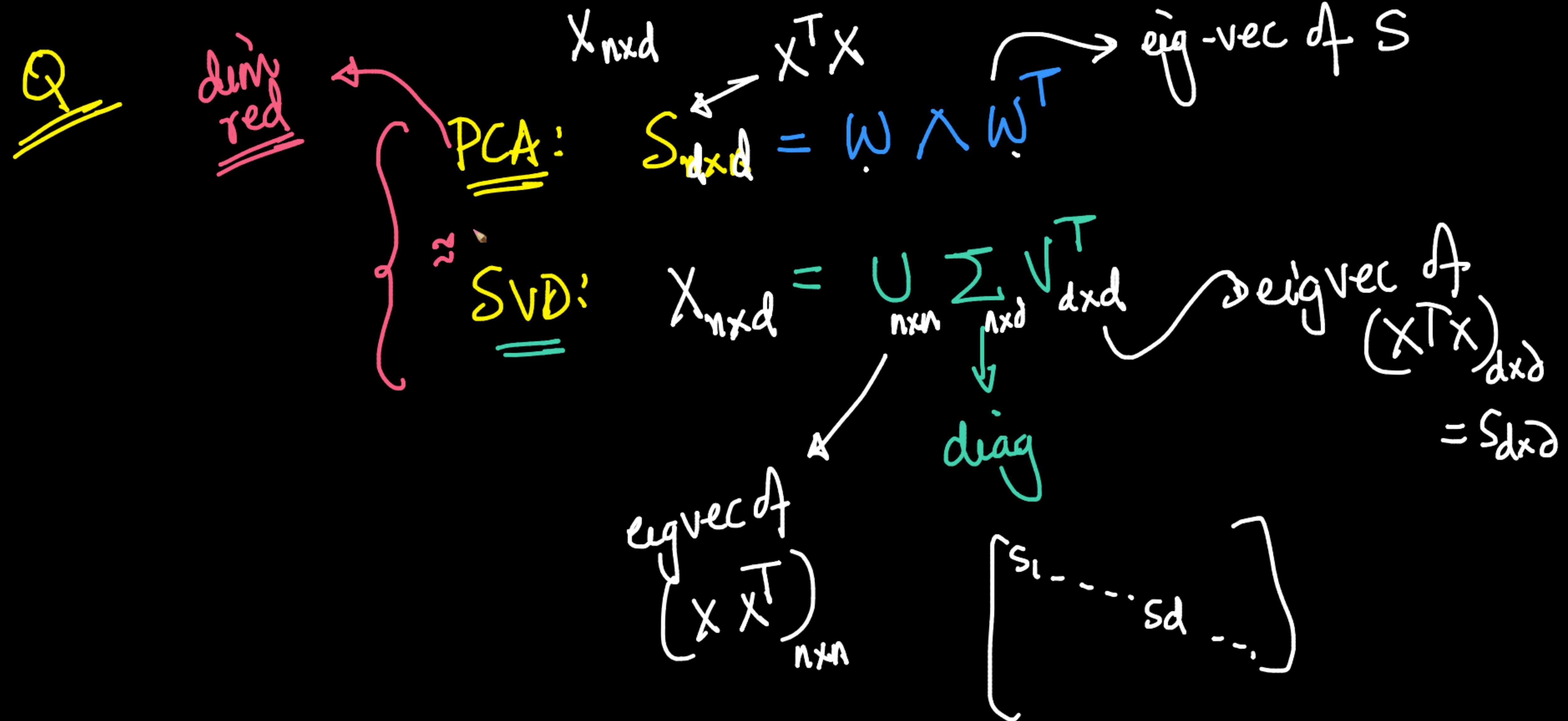
Opt:

$$\min_{B_i, C_j} \sum_{i,j} (A_{ij} - B_i^T C_j)^2$$

A_{ij} ≠ NULL

→ SGD
→ coordinate descent

α → hyper-params
 α → train-test
elbow method



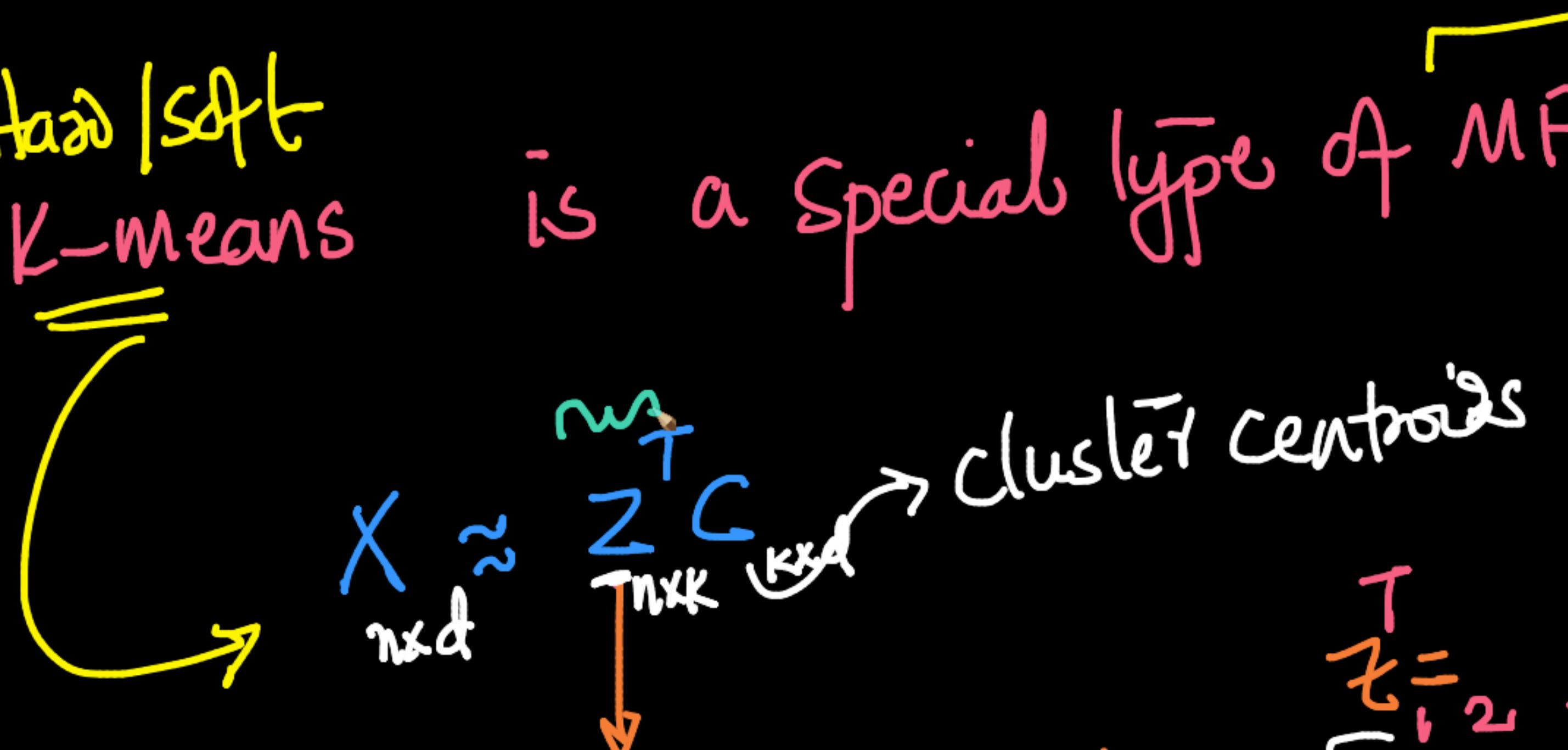
Hard / soft
K-means

is a special type of MF

Hard

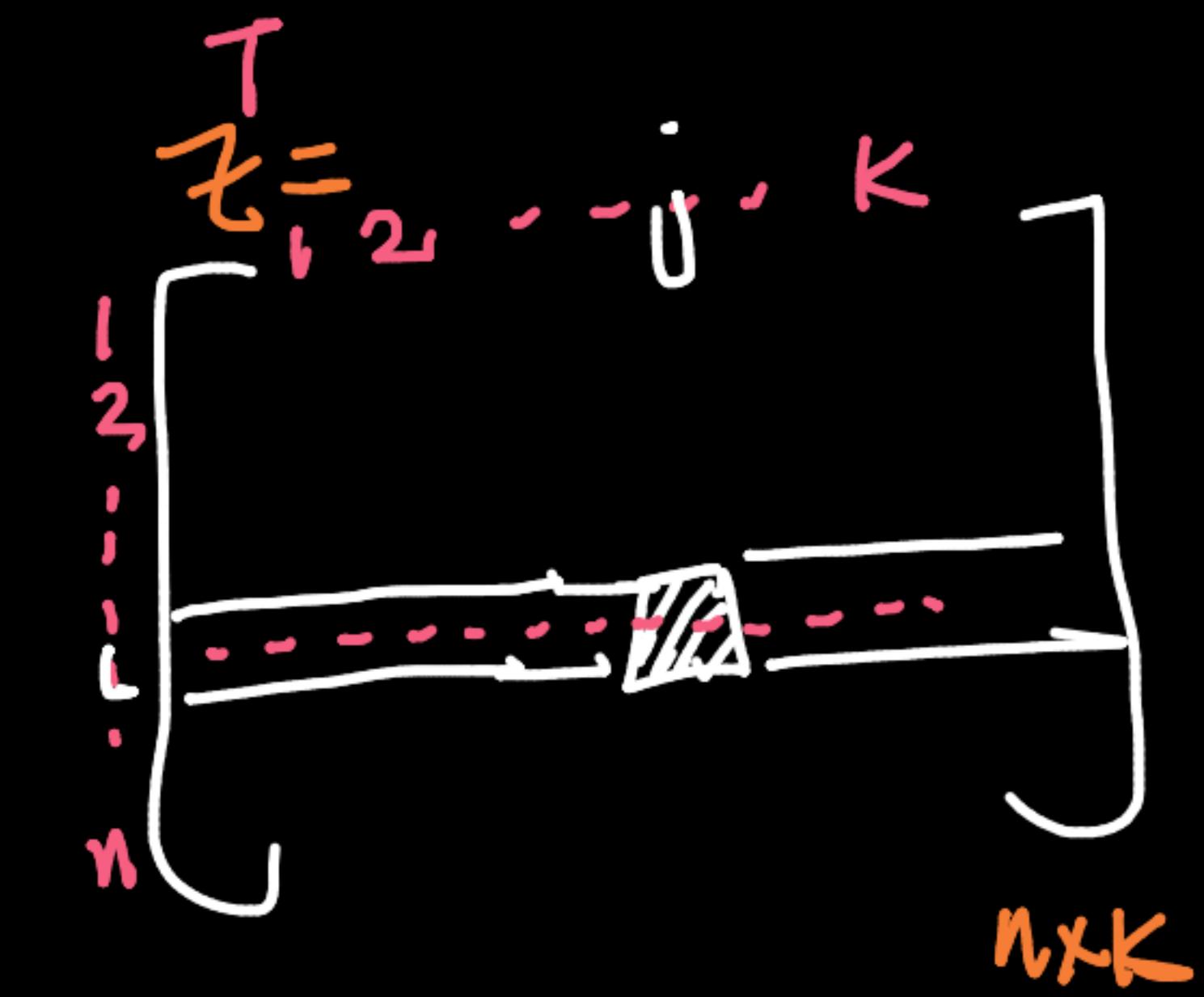
$$z_{ij} = \begin{cases} 1 & \text{if } x_i \in S_j \\ 0 & \text{otherwise} \end{cases}$$

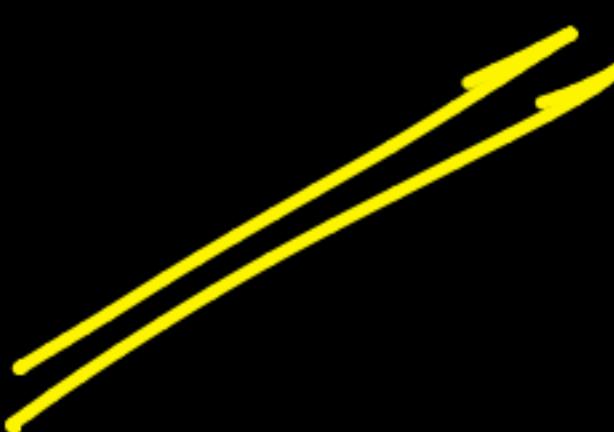
$$\sum_j z_{ij} = 1$$



cluster assignment
matrix

constraints





$$A_{n \times m} = B_{n \times d} \cdot C_{d \times m}$$

[Non-negative
= MF]

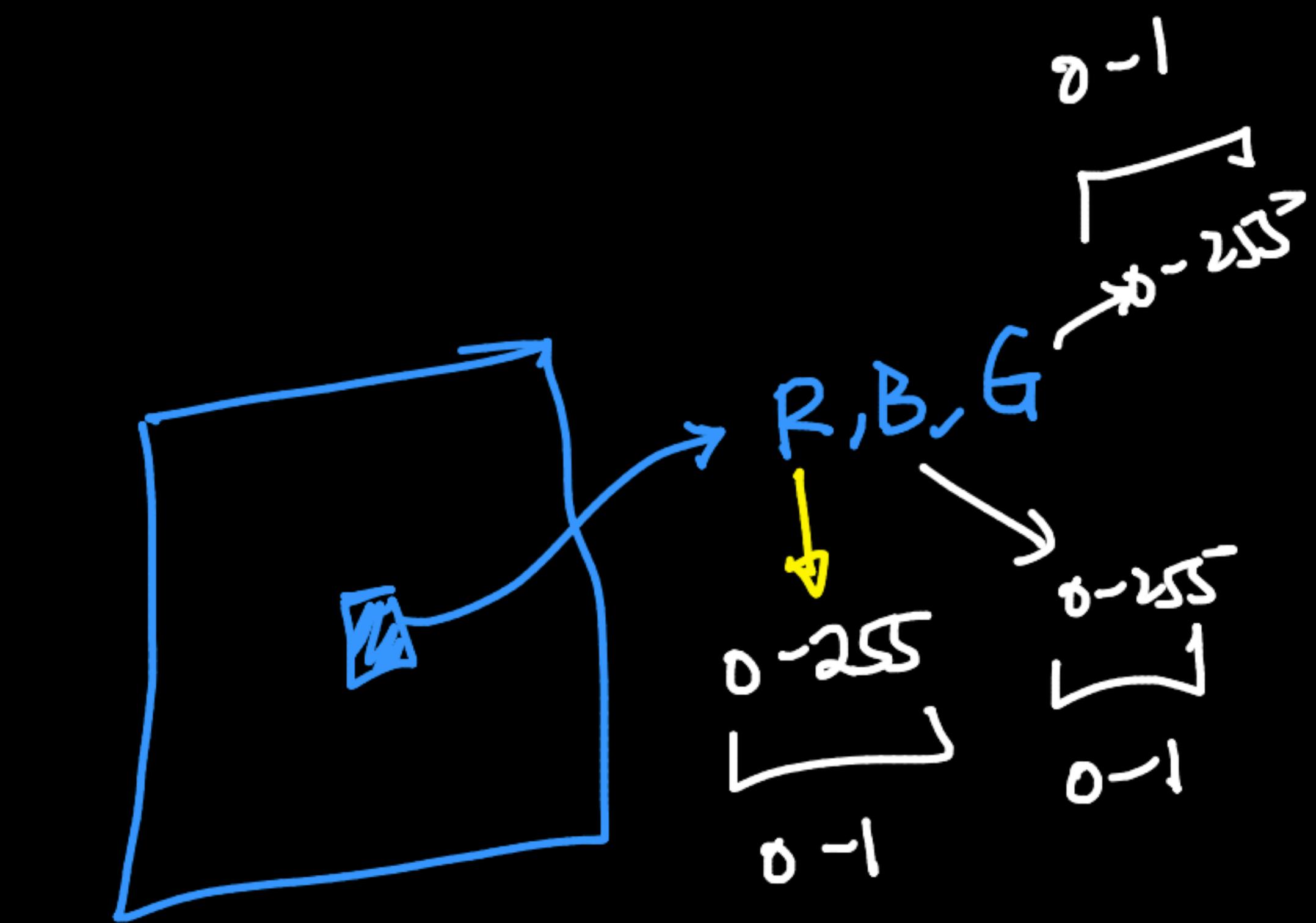
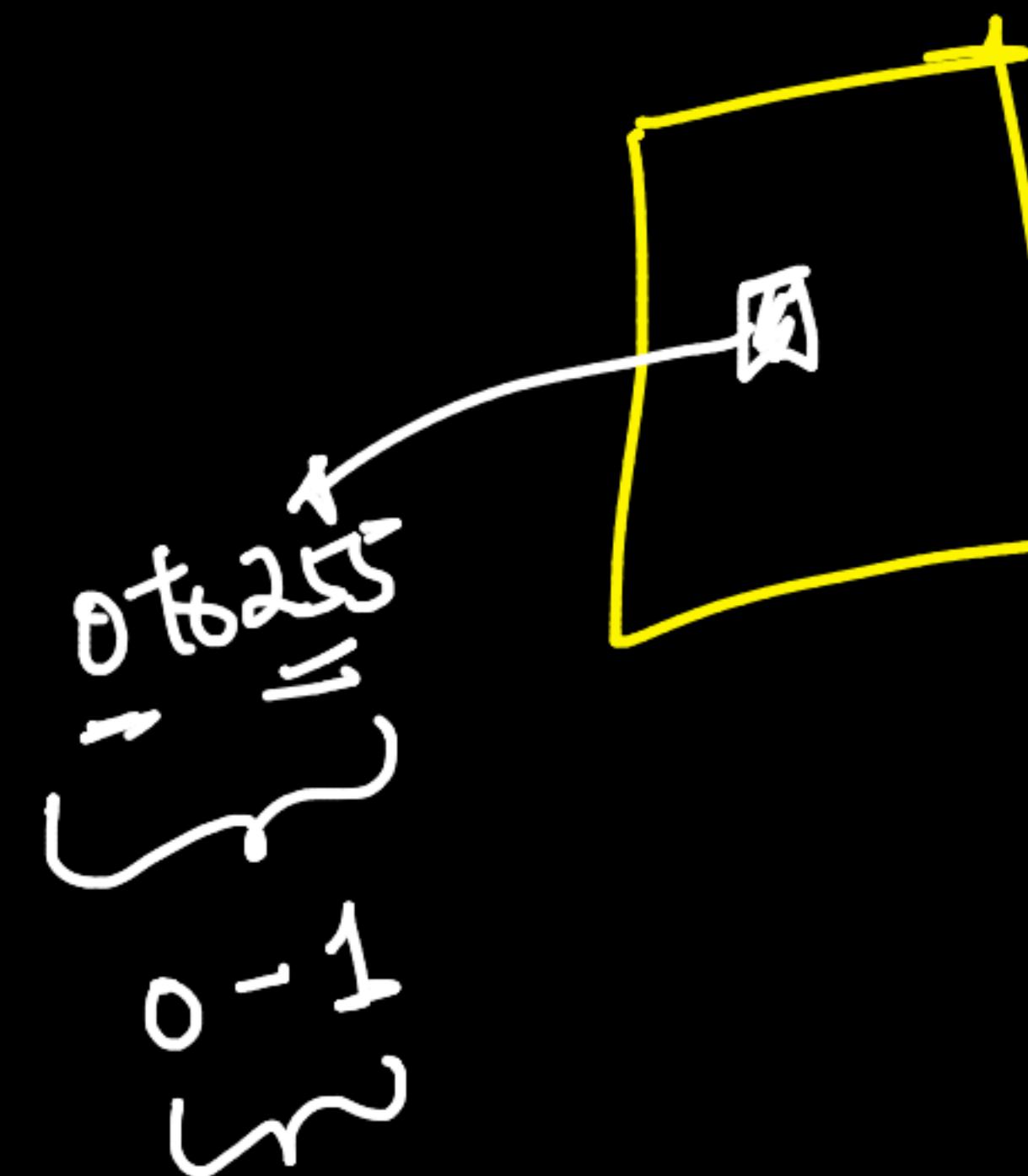
pos

B_{ij} } elements
 C_{ij} } of $B \& C$

↳ non-negative
 ≥ 0

[why?]
→ NMF: generalization of clustering
→ Images → visualize
→ interpretability

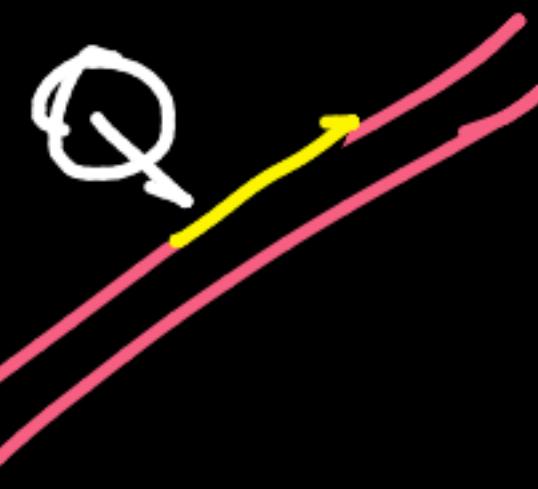
MF \supseteq NMF \supseteq K-means
hard/soft



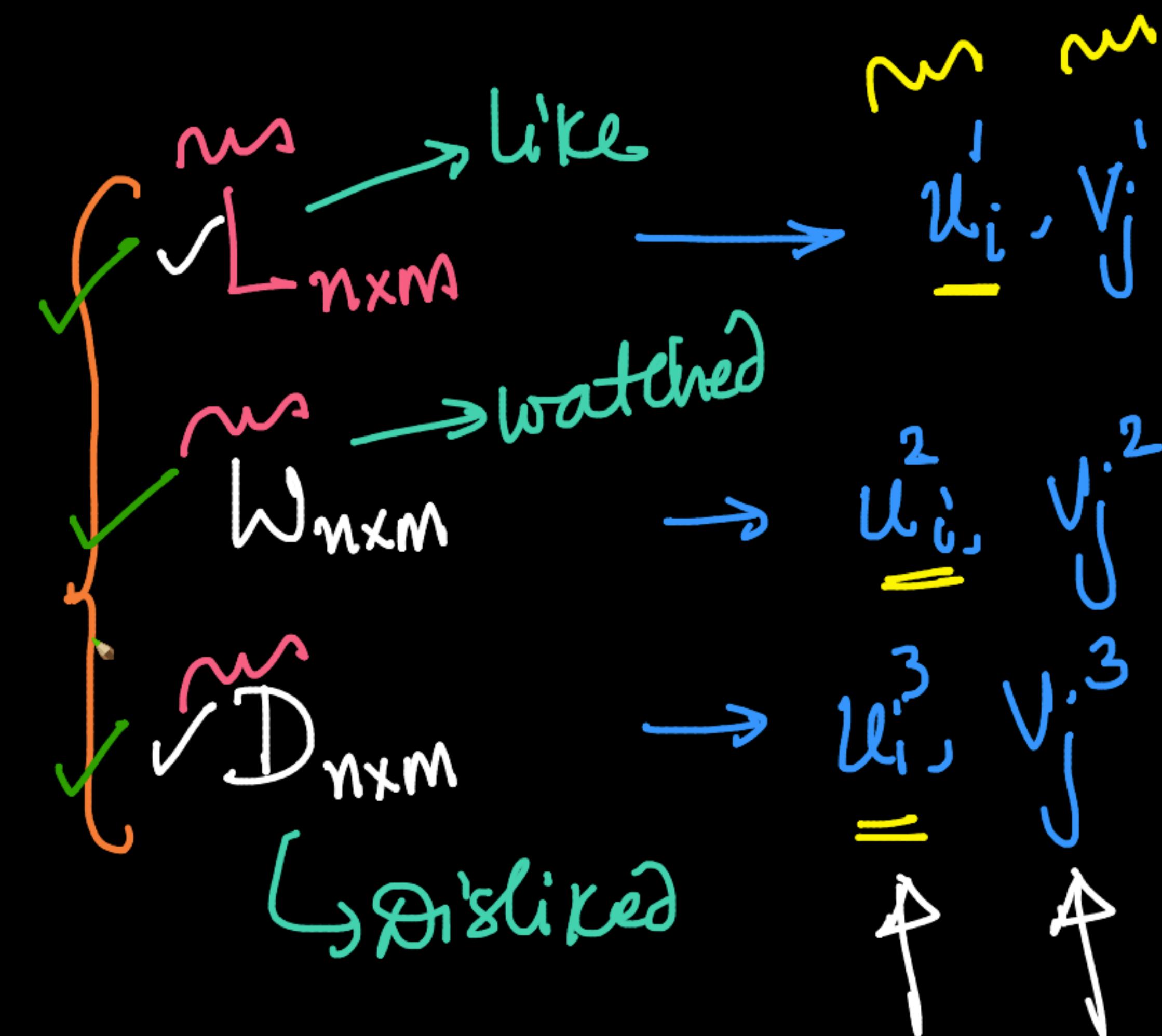
Eigen-faces \leftarrow MF: feature engineering

$$\min_{B, C} \|A - \tilde{B} \cdot C\|_F^2 \equiv \sum_{ij} (A_{ij} - \tilde{B}_i^T \tilde{C})^2$$

s.t. $\tilde{B}_{ij} \geq 0$ \tilde{A}_{ij} }
 $\tilde{C}_{ij} \geq 0$ \tilde{B}_{ij} }

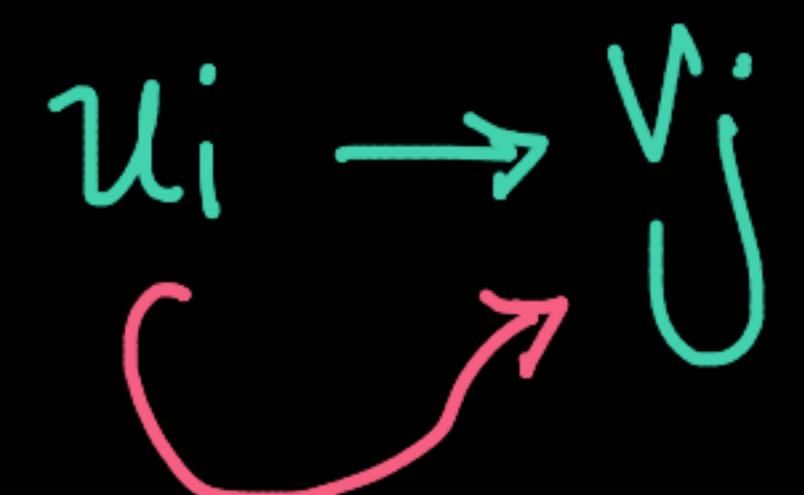


indep
factory



Youtube:

RecSys



option 1:

$$\tilde{L} + \alpha_1 \tilde{W} + \alpha_2 \tilde{D} \rightarrow A_{n \times m}$$

$$\alpha_1 L + \alpha_2 W + \alpha_3 D \quad n \times m$$

option 2:

Independent factorization

$$\left\{ \begin{array}{l} U_i^L, U_i^W, U_i^D \\ V_i^L, V_i^W, V_i^D \end{array} \right.$$

Cosine-sim
User-i $\rightarrow V_D$

$$L > W > D$$

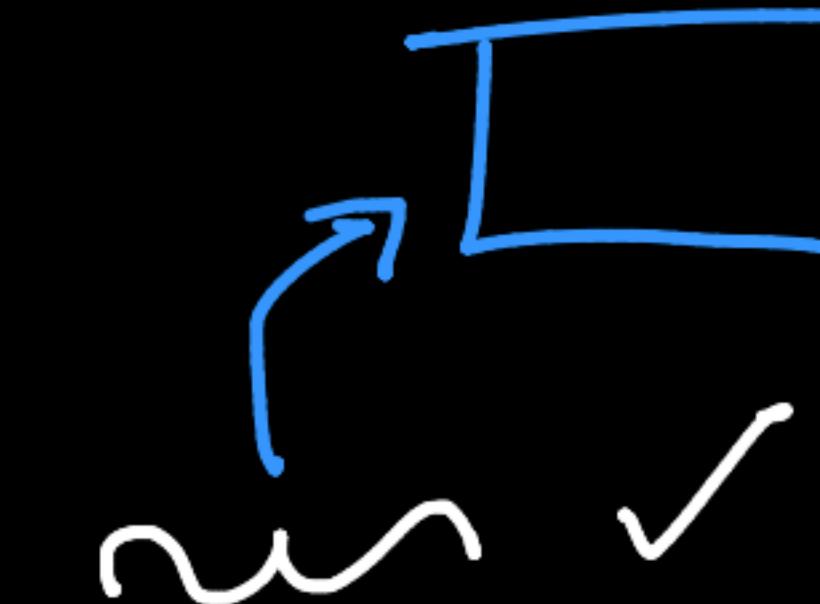
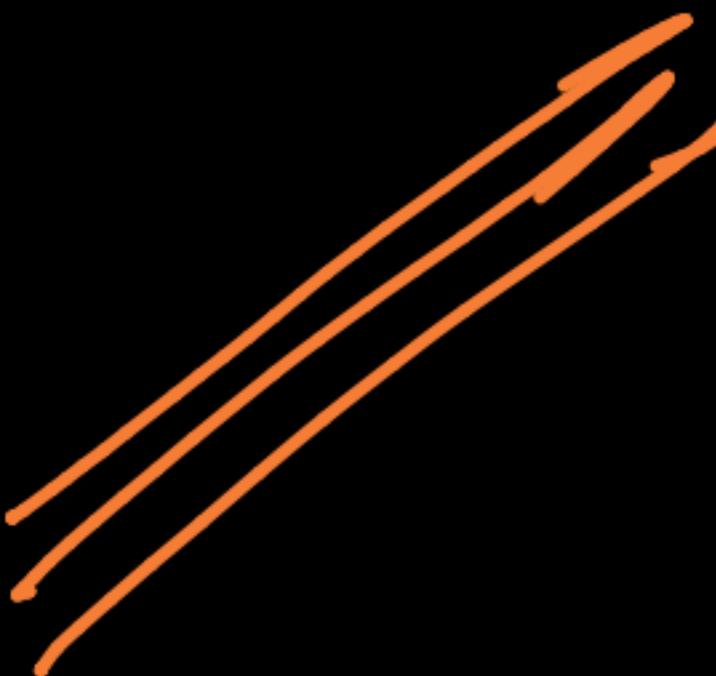
$$\mathcal{U}_{10, \downarrow}^T = \mathcal{L}_{10, \uparrow}$$

$$u_i^L \in \mathbb{R}^d$$

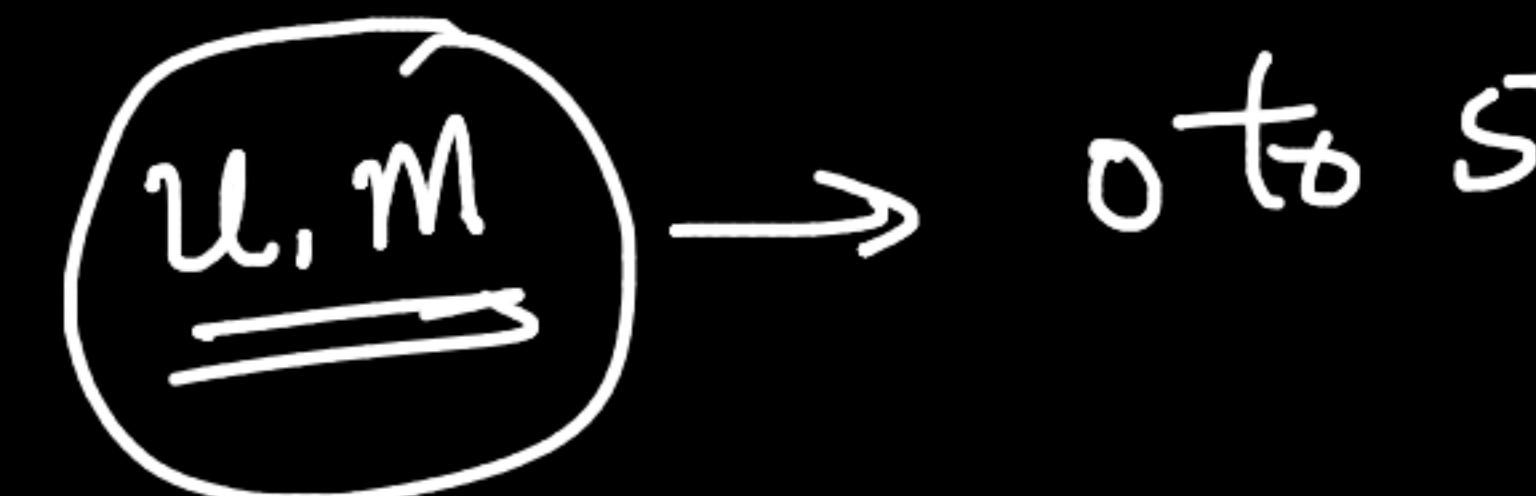
$$u_i^W \in \mathbb{R}^d$$

→ Clustering } ✓
t-SNE/UMAP }

$$\tilde{L} \xrightarrow{n \times n} 0 \text{ or } 1$$



Metric: \overline{RMSE}



never
produced

final-solution → lots of
smaller - models

10% improvement
on \overline{RMSE}

+ bagging (boosting)

Netflix
Prize
2008-09
\$1MM

Paper's Notation

P_u : $\rightarrow n$ users

u_i : user- i

m_j : movie- j

v_i : $\leftarrow m_j$: movie- j

r_{ij} : rating

Dir

Re

$u_i, m_j \rightarrow ?$

Simplest:

$$\min_{v_i, P_u} \sum_{u,i} (r_{ui} - v_i^T P_u)^2$$

$$\beta_u \in \mathbb{R}^d$$

$$q_i \in \mathbb{R}^d$$

d-dim

$$\min_{q_i, \beta_u} \sum_{u,i} \left(q_{ui} - q_i^\top \beta_u \right)^2 + \lambda \left(\sum_i \|q_{il}\|^2 + \sum_u \|\beta_{ul}\|^2 \right)$$

example, Netflix collects star ratings for movies, and TiVo users indicate their preferences for TV shows by pressing thumbs-up and thumbs-down buttons. We refer to explicit user feedback as *ratings*. Usually, explicit feedback comprises a sparse matrix, since any single user is likely to have rated only a small percentage of possible items.

One strength of matrix factorization is that it allows incorporation of additional information. When explicit feedback is not available, recommender systems can infer user preferences using *implicit feedback*, which indirectly reflects opinion by observing user behavior including purchase history, browsing history, search patterns, or even mouse movements. Implicit feedback usually denotes the presence or absence of an event, so it is typically represented by a densely filled matrix.

A BASIC MATRIX FACTORIZATION MODEL

Matrix factorization models map both users and items to a joint latent factor space of dimensionality f , such that user-item interactions are modeled as inner products in that space. Accordingly, each item i is associated with a

Earlier systems relied on imputation to fill in missing ratings and make the rating matrix dense.² However, imputation can be very expensive as it significantly increases the amount of data. In addition, inaccurate imputation might distort the data considerably. Hence, more recent works³⁻⁶ suggested modeling directly the observed ratings only, while avoiding overfitting through a regularized model. To learn the factor vectors (p_u and q_i), the system minimizes the regularized squared error on the set of known ratings:

$$\min_{q^*, p^*} \sum_{(u,i) \in \kappa} (r_{ui} - q_i^T p_u)^2 + \lambda (\|q_i\|^2 + \|p_u\|^2) \quad (2)$$

Here, κ is the set of the (u,i) pairs for which r_{ui} is known (the training set).

The system learns the model by fitting the previously observed ratings. However, the goal is to generalize those previous ratings in a way that predicts future, unknown ratings. Thus, the system should avoid overfitting the observed data by regularizing the learned parameters, whose magnitudes are penalized. The constant λ controls

associated prediction error

$$\underset{def}{e_{ui}} = r_{ui} - q_i^T p_u.$$

Then it modifies the parameters by a magnitude proportional to γ in the opposite direction of the gradient, yielding:

- $q_i \leftarrow q_i + \gamma \cdot (e_{ui} \cdot p_u - \lambda \cdot q_i)$
 - $p_u \leftarrow p_u + \gamma \cdot (e_{ui} \cdot q_i - \lambda \cdot p_u)$
- reg-coeff*
- learning rate*

This popular approach⁴⁻⁶ combines implementation ease with a relatively fast running time. Yet, in some cases, it is beneficial to use ALS optimization.

Alternating least squares

Because both q_i and p_u are unknowns, Equation 2 is not convex. However, if we fix one of the unknowns, the optimization problem becomes quadratic and can be solved optimally. Thus, ALS techniques rotate between fixing the q_i 's and fixing the p_u 's. When all p_u 's are fixed, the system recomputes the q_i 's by solving a least-squares problem, and vice versa. This ensures that each step

item biases can explain, subjecting only the true interaction portion of the data to factor modeling. A first-order approximation of the bias involved in rating r_{ui} is as follows:

$$b_{ui} = \mu + b_i + b_u \quad (3)$$

The bias involved in rating r_{ui} is denoted by b_{ui} and accounts for the user and item effects. The overall average rating is denoted by μ ; the parameters b_u and b_i indicate the observed deviations of user u and item i , respectively, from the average. For example, suppose that you want a first-order estimate for user Joe's rating of the movie *Titanic*. Now, say that the average rating over all movies, μ , is 3.7 stars. Furthermore, *Titanic* is better than an average movie, so it tends to be rated 0.5 stars above the average. On the other hand, Joe is a critical user, who tends to rate 0.3 stars lower than the average. Thus, the estimate for *Titanic*'s rating by Joe would be 3.9 stars ($3.7 + 0.5 - 0.3$). Biases extend Equation 1 as follows:

$$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T p_u \quad (4)$$

have rated only a small percentage of possible items.

One strength of matrix factorization is that it allows incorporation of additional information. When explicit feedback is not available, recommender systems can infer user preferences using *implicit feedback*, which indirectly reflects opinion by observing user behavior including purchase history, browsing history, search patterns, or even mouse movements. Implicit feedback usually denotes the presence or absence of an event, so it is typically represented by a densely filled matrix.

A BASIC MATRIX FACTORIZATION MODEL

Matrix factorization models map both users and items to a joint latent factor space of dimensionality f , such that user-item interactions are modeled as inner products in that space. Accordingly, each item i is associated with a

works¹⁰ suggested modeling directly the observed ratings only, while avoiding overfitting through a regularized model. To learn the factor vectors (p_u and q_i), the system minimizes the regularized squared error on the set of known ratings:

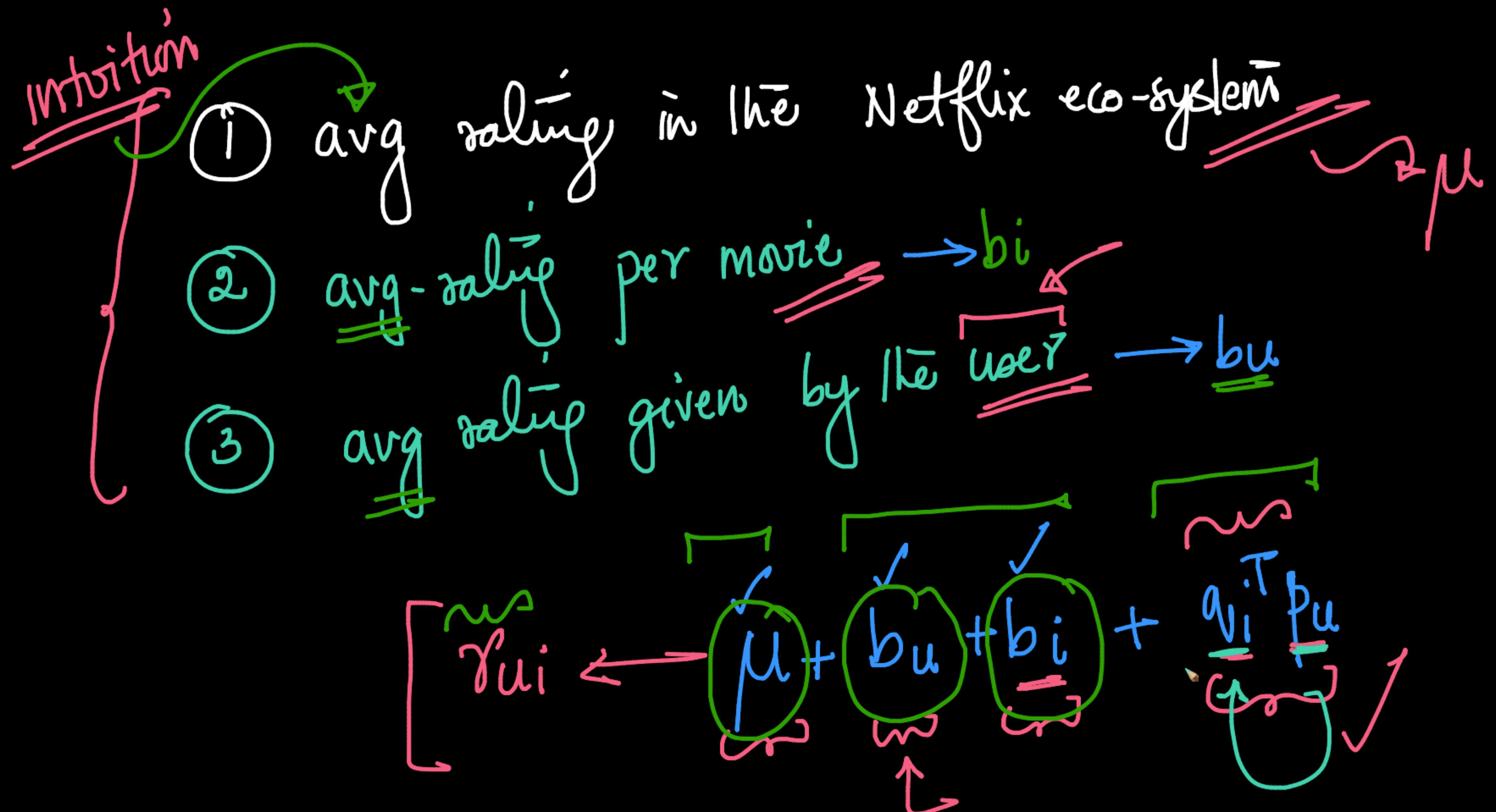
$$\min_{q^*, p^*} \sum_{(u,i) \in \kappa} (r_{ui} - q_i^T p_u)^2 + \lambda (\|q_i\|^2 + \|p_u\|^2)$$

*V·high &
↓*

Here, κ is the set of the (u,i) pairs for which r_{ui} is known (the training set).

The system learns the model by fitting the previously observed ratings. However, the goal is to generalize those previous ratings in a way that predicts future, unknown ratings. Thus, the system should avoid overfitting the observed data by regularizing the learned parameters, whose magnitudes are penalized. The constant λ controls





$$5 \underbrace{\gamma_{ui}}_{=} = \dots + b_i + a_i \cdot bu$$

$$u=10$$
$$r=10$$



user's biases

$$\gamma_{ui} = 2$$

DTR

min

parameters $\{P_u, q_i\}$
 $\{\mu, b_u, b_i\}$

$$\sum_{u,i} \left(\tilde{\omega}_{ui} - \underbrace{\mu - b_u - b_i}_{\rightarrow P_u^T q_i} \right)^2 + \lambda \sum \left\{ \|\omega_i\|^2 + \|b_i\|^2 \right\} + \mu^2 + b_u^2 + b_i^2$$

Salars

$\omega_{ui} \rightarrow$ more recent ; more weightage

Surprise · A Python scikit for re | RecSys_MF.ipynb - Colaborato | Recommender-Systems-[Netfil | Microsoft Word - BellKorSolutio | +

datajobs.com/data-science-repo/Recommender-Systems-%5BNetflix%5D.pdf

4 / 8 | - 250% + | ☰

Recommender-Systems-[Netflix].pdf

optimally. Thus, ALS techniques rotate between fixing the q_i 's and fixing the p_u 's. When all p_u 's are fixed, the system recomputes the q_i 's by solving a least-squares problem, and vice versa. This ensures that each step decreases Equation 2 until convergence.⁸

While in general stochastic gradient descent is easier and faster than ALS, ALS is favorable in at least two cases. The first is when the system can use parallelization. In ALS, the system computes each q_i independently of the other item factors and computes each p_u independently of the other user factors. This gives rise to potentially massive parallelization of the algorithm.⁹ The second case is for systems centered on implicit data. Because the training set cannot be considered sparse, looping over each single training case—as gradient descent does—would not be practical. ALS can efficiently handle such cases.¹⁰

ADDING BIASES

One benefit of the matrix factorization approach to collaborative filtering is its flexibility in dealing with various

$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T p_u$ (4)

Here, the observed rating is broken down into its four components: global average, item bias, user bias, and user-item interaction. This allows each component to explain only the part of a signal relevant to it. The system learns by minimizing the squared error function:^{4,5}

$$\min_{p^*, q^*, b^*} \sum_{(u,i) \in K} (r_{ui} - \mu - b_u - b_i - p_u^T q_i)^2 + \lambda (\|p_u\|^2 + \|q_i\|^2 + b_u^2 + b_i^2)$$
 (5)

Since biases tend to capture much of the observed signal, their accurate modeling is vital. Hence, other works offer more elaborate bias models.¹¹

ADDITIONAL INPUT SOURCES

Often a system must deal with the cold start problem, wherein many users supply very few ratings, making it

23 / 23

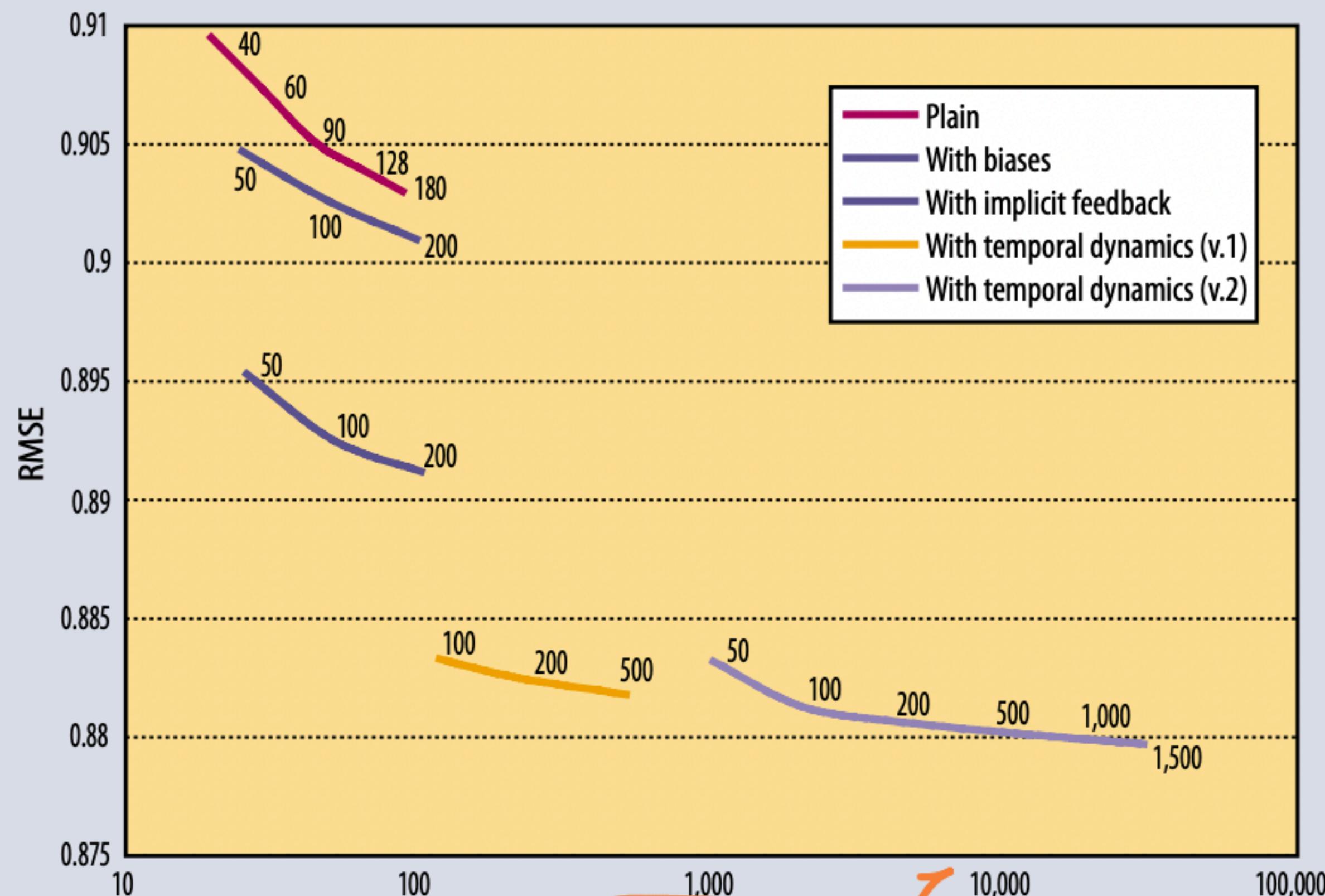


Figure 4. Matrix factorization models' accuracy. The plots show the root-mean-square error of each of four individual factor models (lower is better). Accuracy improves when the factor model's dimensionality (denoted by numbers on the charts) increases. In addition, the more refined factor models, whose descriptions involve more distinct sets of parameters, are more accurate. For comparison, the Netflix system achieves RMSE = 0.9514 on the same data set.

Matrix factorization techniques have become a dominant methodology within collaborative filtering recommenders. Experience with datasets such as the Netflix Prize data has shown that they deliver accuracy superior to classical nearest-neighbor techniques. At the same time, they offer a compact memory-efficient model that systems can learn relatively easily. What makes these techniques even more convenient is that models can integrate naturally many crucial aspects of the data, such as multiple forms of feedback, temporal dynamics, and confidence levels. □

References

1. D. Goldberg et al., "Using Collective Filtering to Weave an Information Tapestry,"

Surprise · A Python scikit for re | RecSys_MF.ipynb - Colaborato | Recommender-Systems-[Netfil | Microsoft Word - BellKorSolutio | +

datajobs.com/data-science-repo/Recommender-Systems-%5BNetflix%5D.pdf

7 / 8 | - 200% + | ☰

Recommender-Systems-[Netflix].pdf

COVER FEATURE

gi ERF
Put ERF

The figure consists of three vertically stacked line plots showing the relationship between the number of parameters (x-axis, logarithmic scale from 10 to 100,000) and the root-mean-square error (RMSE, y-axis, linear scale from 0.875 to 0.91).
The top plot shows the 'Plain' model, with points at dimensionality values 40, 60, 90, 128, 180, and 200. The RMSE decreases as dimensionality increases.
The middle plot shows the 'With biases' model, with points at dimensionality values 50, 100, and 200. The RMSE decreases as dimensionality increases.
The bottom plot shows the 'With temporal dynamics (v.1)' model, with points at dimensionality values 100, 200, 500, 1,000, and 1,500. The RMSE decreases as dimensionality increases.
A legend on the right identifies the five models: Plain (red), With biases (blue), With implicit feedback (green), With temporal dynamics (v.1) (orange), and With temporal dynamics (v.2) (purple).

Model	Dimensionality	RMSE
Plain	40	~0.908
Plain	60	~0.906
Plain	90	~0.905
Plain	128	~0.904
Plain	180	~0.903
Plain	200	~0.902
With biases	50	~0.895
With biases	100	~0.893
With biases	200	~0.891
With temporal dynamics (v.1)	100	~0.884
With temporal dynamics (v.1)	200	~0.883
With temporal dynamics (v.1)	500	~0.882
With temporal dynamics (v.1)	1,000	~0.881
With temporal dynamics (v.1)	1,500	~0.880
With temporal dynamics (v.2)	50	~0.883
With temporal dynamics (v.2)	100	~0.882
With temporal dynamics (v.2)	200	~0.881
With temporal dynamics (v.2)	500	~0.880
With temporal dynamics (v.2)	1,000	~0.879
With temporal dynamics (v.2)	1,500	~0.878

Figure 4. Matrix factorization models' accuracy. The plots show the root-mean-square error of each of four individual factor models (lower is better). Accuracy improves when the factor model's dimensionality (denoted by numbers on the charts) increases. In addition, the more refined factor models, whose descriptions involve more distinct sets of parameters, are more accurate. For example, the Netflix Prize submission with RMSE = 0.9514 on the same data used to train the plain MF model had RMSE = 0.8563.

Matrix factorization techniques have become a dominant methodology within collaborative filtering recommenders. Experience with datasets such as the Netflix Prize data has shown that they deliver accuracy superior to classical nearest-neighbor techniques. At the same time, they offer a compact memory-efficient model that systems can learn relatively easily. What makes these techniques even more convenient is that models can integrate naturally many crucial aspects of the data, such as multiple forms of feedback, temporal dynamics, and confidence levels. □

References

1. D. Goldberg et al., "Using Collaborative Filtering to Weave an Information Tapestry,"

Surprise · A Python scikit for re | RecSys_MF.ipynb - Colaborato | Recommender-Systems-[Netfil | Microsoft Word - BellKorSolutio | +

datajobs.com/data-science-repo/Recommender-Systems-%5BNetflix%5D.pdf

7 / 8 | - 200% + | ☰

Recommender-Systems-[Netflix].pdf

COVER FEATURE

RMSE

Plain

With biases

With implicit feedback

With temporal dynamics (v.1)

With temporal dynamics (v.2)

10 100 1,000 10,000 100,000

Millions of parameters

Model Type	Number of Parameters (Millions)	RMSE
Plain	40	~0.910
Plain	60	~0.908
Plain	90	~0.906
Plain	128	~0.904
Plain	180	~0.903
With biases	50	~0.895
With biases	100	~0.893
With biases	200	~0.891
With implicit feedback	100	~0.884
With implicit feedback	200	~0.883
With implicit feedback	500	~0.882
With implicit feedback	1,000	~0.881
With implicit feedback	1,500	~0.880
With temporal dynamics (v.1)	100	~0.884
With temporal dynamics (v.1)	200	~0.883
With temporal dynamics (v.1)	500	~0.882
With temporal dynamics (v.2)	50	~0.883
With temporal dynamics (v.2)	100	~0.882
With temporal dynamics (v.2)	200	~0.881
With temporal dynamics (v.2)	500	~0.880
With temporal dynamics (v.2)	1,000	~0.879
With temporal dynamics (v.2)	1,500	~0.878

Figure 4. Matrix factorization models' accuracy. The plots show the root-mean-square error of each of four individual factor models (lower is better). Accuracy improves when the factor model's dimensionality (denoted by numbers on the charts) increases. In addition, the more refined factor models, whose descriptions involve more distinct sets of parameters, are more accurate. For example, the Netflix Prize dataset has an RMSE = 0.9514 on the same data set, while the RMSE = 0.8563.

M

atrix factorization techniques have become a dominant methodology within collaborative filtering recommenders. Experience with datasets such as the Netflix Prize data has shown that they deliver accuracy superior to classical nearest-neighbor techniques. At the same time, they offer a compact memory-efficient model that systems can learn relatively easily. What makes these techniques even more convenient is that models can integrate naturally many crucial aspects of the data, such as multiple forms of feedback, temporal dynamics, and confidence levels. C

References

- D. Goldberg et al., "Using Collaborative Filtering to Weave an Information Tapestry,"

COVER FEATURE

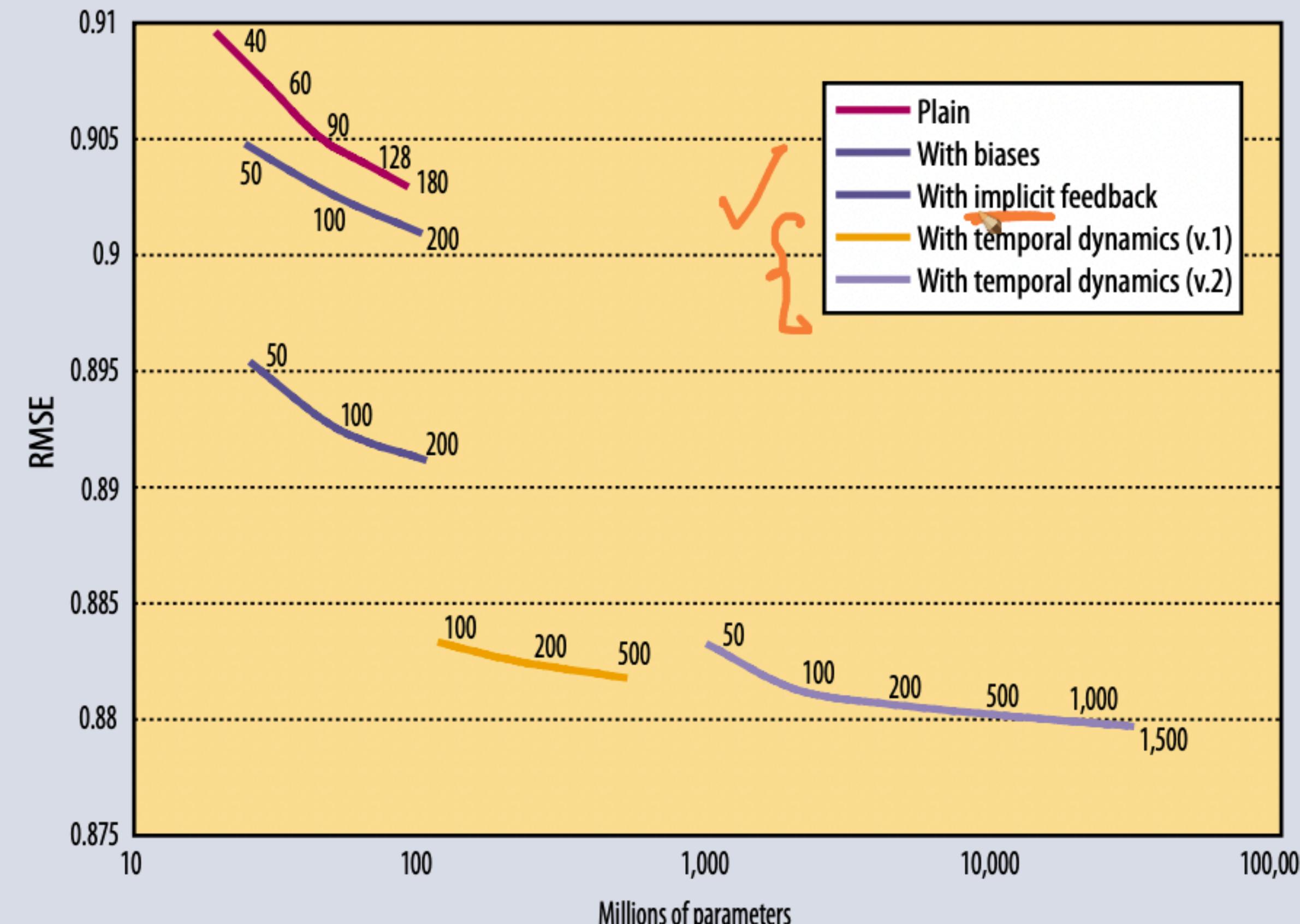


Figure 4. Matrix factorization models' accuracy. The plots show the root-mean-square error of each of four individual factor models (lower is better). Accuracy improves when the factor model's dimensionality (denoted by numbers on the charts) increases. In addition, the more refined factor models, whose descriptions involve more distinct sets of parameters, are more accurate. For example, the Netflix Prize winning system had an RMSE = 0.9514 on the same data set, while the plain MF model had an RMSE = 0.8563.

Matrix factorization techniques have become a dominant methodology within collaborative filtering recommenders. Experience with datasets such as the Netflix Prize data has shown that they deliver accuracy superior to classical nearest-neighbor techniques. At the same time, they offer a compact memory-efficient model that systems can learn relatively easily. What makes these techniques even more convenient is that models can integrate naturally many crucial aspects of the data, such as multiple forms of feedback, temporal dynamics, and confidence levels. □

References

1. D. Goldberg et al., "Using Collaborative Filtering to Weave an Information Tapestry,"

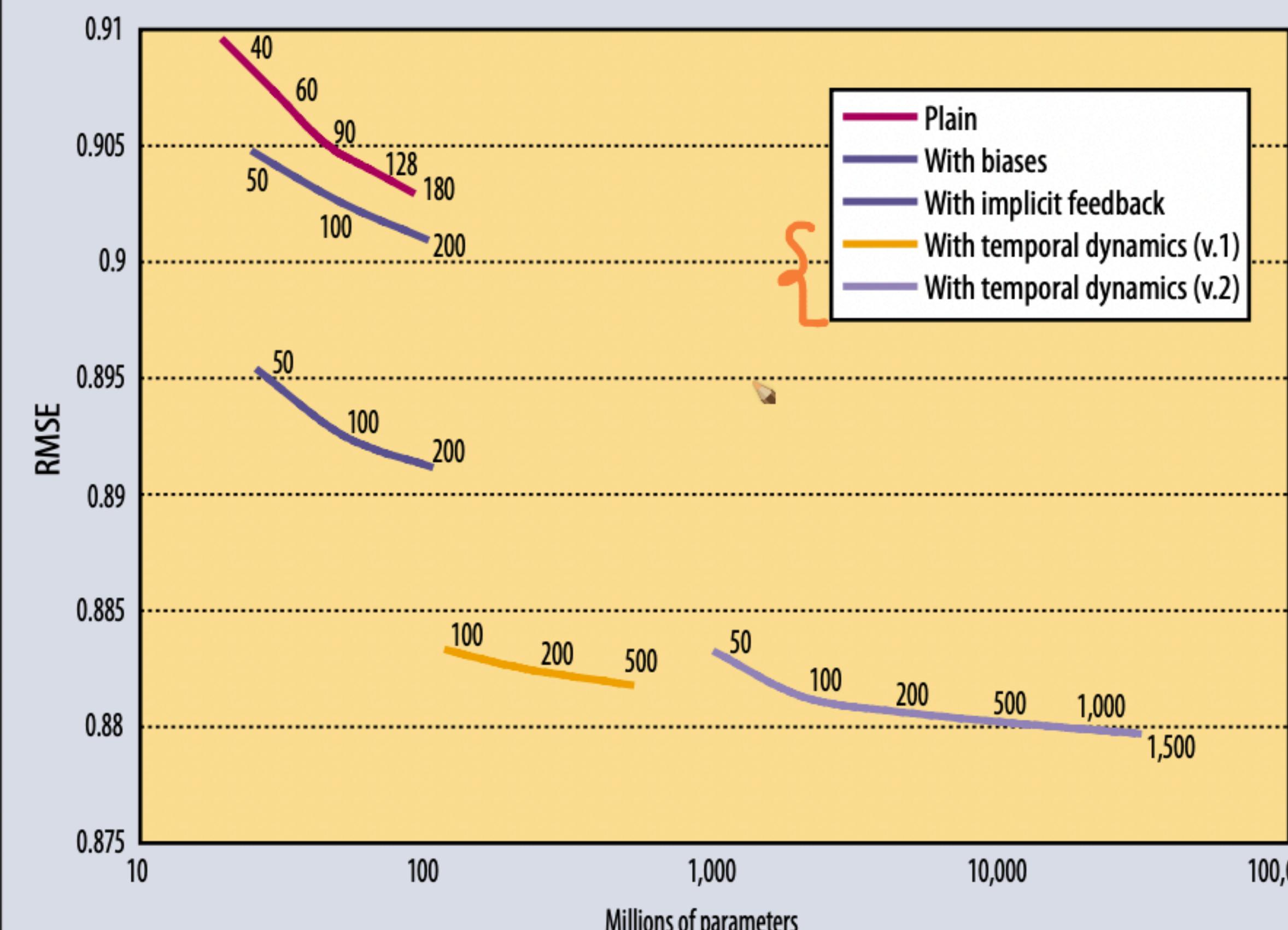


Figure 4. Matrix factorization models' accuracy. The plots show the root-mean-square error of each of four individual factor models (lower is better). Accuracy improves when the factor model's dimensionality (denoted by numbers on the charts) increases. In addition, the more refined factor models, whose descriptions involve more distinct sets of parameters, are more accurate. For example, the $\text{NMF}^{\text{refined}}$ model has an RMSE = 0.9514 on the same data as the NMF model, which has an RMSE = 0.8563.

Matrix factorization techniques have become a dominant methodology within collaborative filtering recommenders. Experience with datasets such as the Netflix Prize data has shown that they deliver accuracy superior to classical nearest-neighbor techniques. At the same time, they offer a compact memory-efficient model that systems can learn relatively easily. What makes these techniques even more convenient is that models can integrate naturally many crucial aspects of the data, such as multiple forms of feedback, temporal dynamics, and confidence levels. □

References

- 1 D. Goldberg et al., "Using Collective Filtering to Weave an Information Tapestry,"

Surprise · A Python scikit for re RecSys_MF.ipynb - Colaborato Recommender-Systems-[Netfil Microsoft Word - BellKorSolutio +

datajobs.com/data-science-repo/Recommender-Systems-%5BNetflix%5D.pdf

7 / 8 | - 200% + | ☰

Temporal-dynamics

Model	Dimensions	Approx. RMSE
Plain	40, 60, 90, 128, 180, 200	0.905 - 0.915
With biases	50, 100, 200	0.895 - 0.905
With implicit feedback	100, 200, 500, 1,000, 1,500	0.882 - 0.892
With temporal dynamics (v.1)	100, 200, 500	0.882 - 0.885
With temporal dynamics (v.2)	50, 100, 200, 500, 1,000	0.882 - 0.885

Figure 4. Matrix factorization models' accuracy. The plots show the root-mean-square error of each of four individual factor models (lower is better). Accuracy improves when the factor model's dimensionality (denoted by numbers on the charts) increases. In addition, the more refined factor models, whose descriptions involve more distinct sets of parameters, are more accurate. For comparison, the Netflix system achieves RMSE = 0.9514 on the same data. The RMSE of the plain MF model is 0.914.

References

1. D. Goldberg et al., "Using Collaborative Filtering to Weave a Formation Tapestry," Comm ACM, vol. 35, 1992, pp.

$$|N(u)|^{-0.5} \sum_{i \in N(u)} x_i^{4.5}$$

Another information source is known user attributes, for example, demographics. Again, for simplicity consider Boolean attributes where user u corresponds to the set of attributes $A(u)$, which can describe gender, age group, Zip code, income level, and so on. A distinct factor vector $y_a \in \mathbb{R}^f$ corresponds to each attribute to describe a user through the set of user-associated attributes:

$$\sum_{a \in A(u)} y_a$$

The matrix factorization model should integrate all signal sources, with enhanced user representation:

$$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T [p_u + |N(u)|^{-0.5} \sum_{i \in N(u)} x_i + \sum_{a \in A(u)} y_a] \quad (6)$$

While the previous examples deal with enhancing user representation—where lack of data is more common—items can get a similar treatment when necessary.

TEMPORAL DYNAMICS

So far, the presented models have been static. In reality, product perception and popularity constantly change as new selections emerge. Similarly, customers' inclinations evolve, leading them to redefine their taste. Thus, the system should account for the dynamic, time-varying nature of user preferences.

might become a fan of crime dramas a year later. Similarly, humans change their perception of certain actors and directors. The model accounts for this effect by taking the user factors (the vector p_u) as a function of time. On the other hand, it specifies static item characteristics, q_i , because, unlike humans, items are static in nature.

Exact parameterizations of time-varying parameters¹¹ lead to replacing Equation 4 with the dynamic prediction rule for a rating at time t :

$$\hat{r}_{ui}(t) = \mu + b_i(t) + b_u(t) + q_i^T p_u(t) \quad (7)$$

INPUTS WITH VARYING CONFIDENCE LEVELS

In several setups, not all observed ratings deserve the same weight or confidence. For example, massive advertising might influence votes for certain items, which do not aptly reflect longer-term characteristics. Similarly, a system might face adversarial users that try to tilt the ratings of certain items.

Another example is systems built around implicit feedback. In such systems, which interpret ongoing user behavior, a user's exact preference level is hard to quantify. Thus, the system works with a cruder binary representation, stating either "probably likes the product" or "probably not interested in the product." In such cases, it is valuable to attach confidence scores with the estimated preferences. Confidence can stem from available

$$|N(u)|^{-0.5} \sum_{i \in N(u)} x_i^{4.5}$$

Another information source is known user attributes, for example, demographics. Again, for simplicity consider Boolean attributes where user u corresponds to the set of attributes $A(u)$, which can describe gender, age group, Zip code, income level, and so on. A distinct factor vector $y_a \in \mathbb{R}^f$ corresponds to each attribute to describe a user through the set of user-associated attributes:

$$\sum_{a \in A(u)} y_a$$

The matrix factorization model should integrate all signal sources, with enhanced user representation:

$$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T [p_u + |N(u)|^{-0.5} \sum_{i \in N(u)} x_i + \sum_{a \in A(u)} y_a] \quad (6)$$

While the previous examples deal with enhancing user representation—where lack of data is more common—items can get a similar treatment when necessary.

TEMPORAL DYNAMICS

So far, the presented models have been static. In reality, product perception and popularity constantly change as new selections emerge. Similarly, customers' inclinations evolve, leading them to redefine their taste. Thus, the system should account for the dynamic, time-varying nature of user preferences.

For example, a user might become a fan of crime dramas a year later. Similarly, humans change their perception of certain actors and directors. The model accounts for this effect by taking the user factors (the vector p_u) as a function of time. On the other hand, it specifies static item characteristics, q_i , because, unlike humans, items are static in nature.

Exact parameterizations of time-varying parameters¹¹ lead to replacing Equation 4 with the dynamic prediction rule for a rating at time t :

$$\hat{r}_{ui}(t) = \mu + b_i(t) + b_u(t) + q_i^T p_u(t) \quad (7)$$

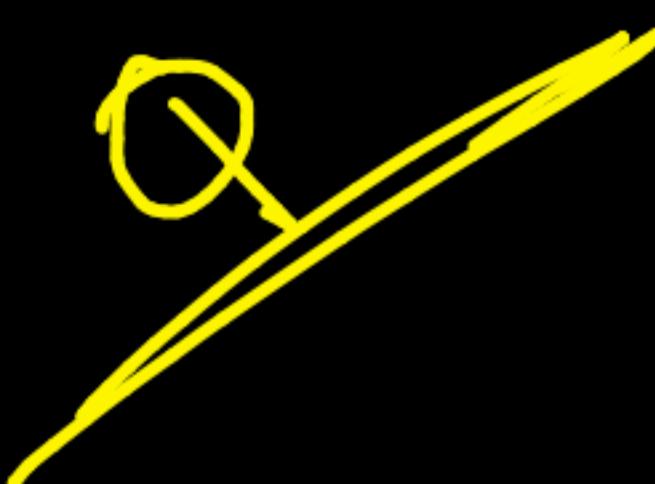
INPUTS WITH VARYING CONFIDENCE LEVELS

In several setups, not all observed ratings deserve the same weight or confidence. For example, massive advertising might influence votes for certain items, which do not aptly reflect longer-term characteristics. Similarly, a system might face adversarial users that try to tilt the ratings of certain items.

Another example is systems built around implicit feedback. In such systems, which interpret ongoing user behavior, a user's exact preference level is hard to quantify. Thus, the system works with a cruder binary representation, stating either "probably likes the product" or "probably not interested in the product." In such cases, it is valuable to attach confidence scores with the estimated preferences. Confidence can stem from available

μ : const ✓

{ bu, bi } time-varying → gave very high
wi
pu reduction in RMSE



Youtube

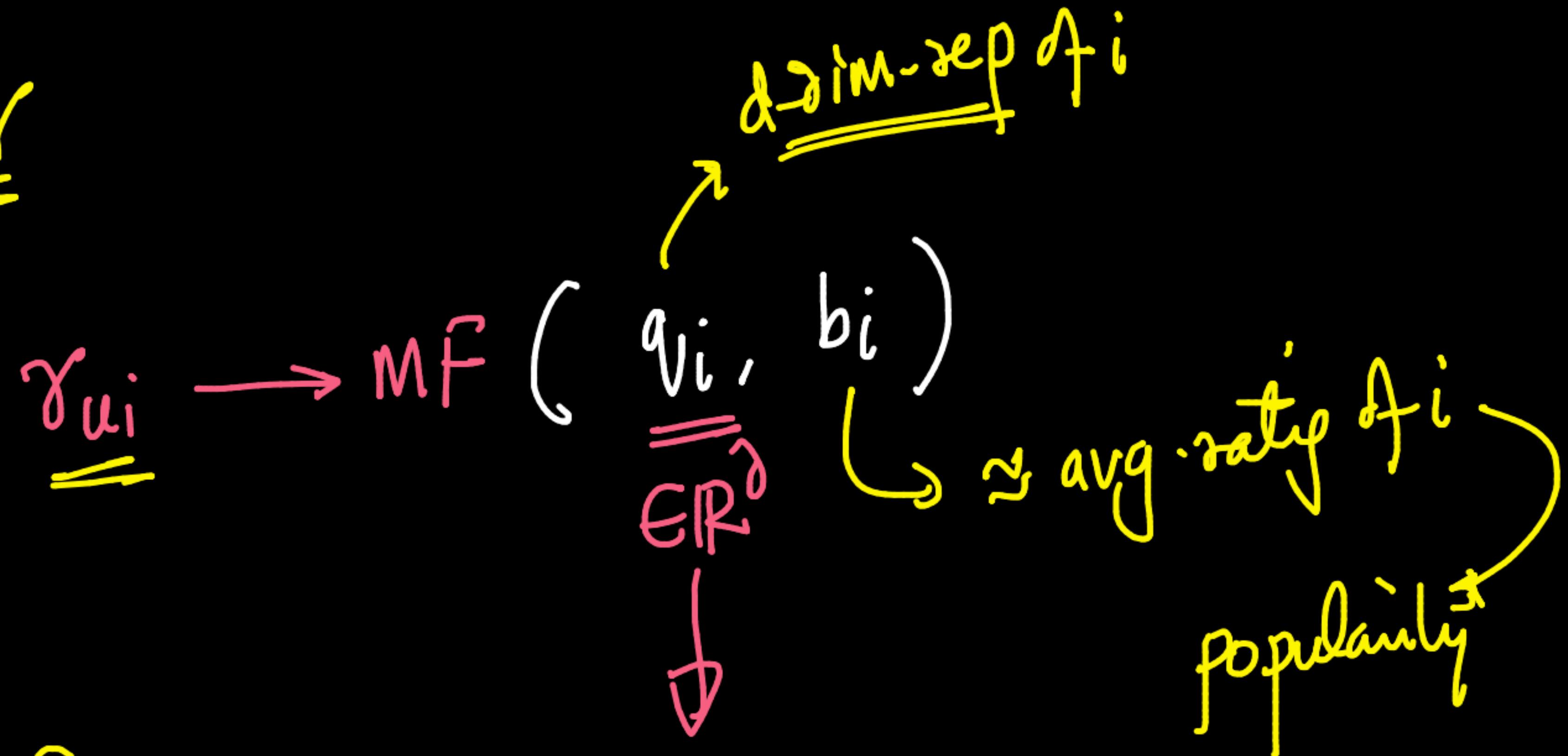
recently aware
time-based rec sys

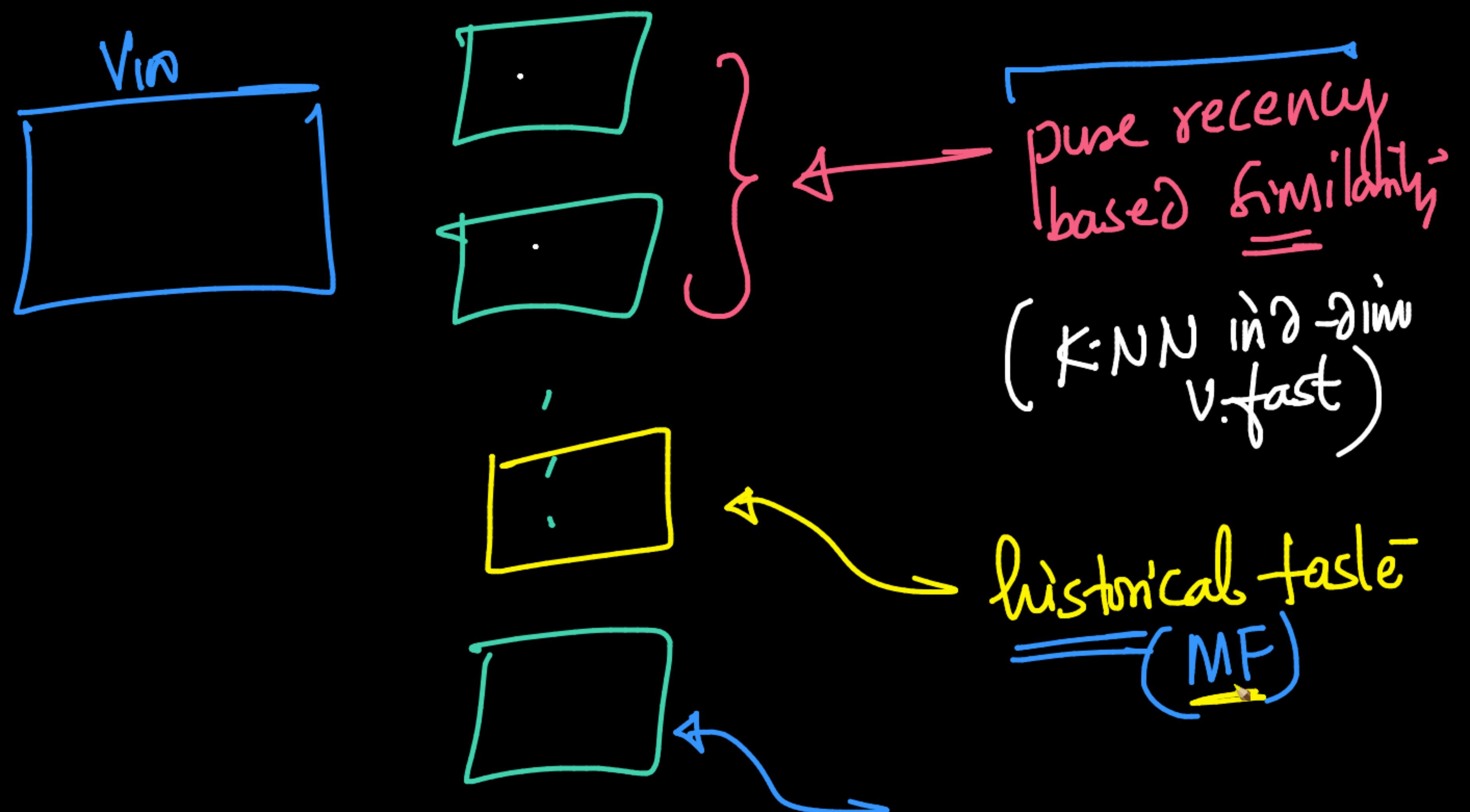
- update models: b_u, b_i, w_{ij}, P_u periodically
- more weightage to more recent data
 - add a w_{ui} based on recently

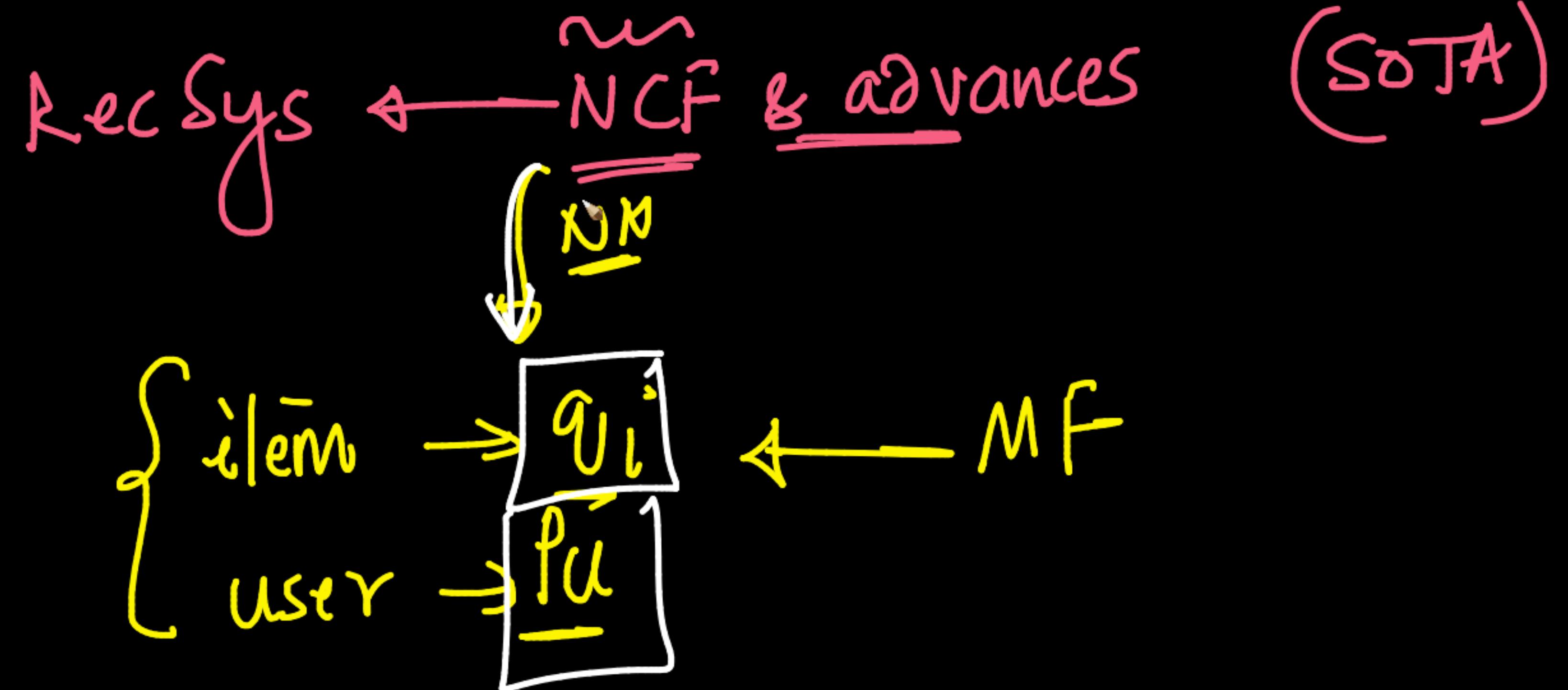
Simpler

Nightly

best \rightarrow {top 100 videos that are good to recommend}









+ Code + Text

Connect ▾

1

BUILDING WHEELS FOR COLLECTED PACKAGES: SCIRL-SURPRISE

A set of small, light-gray navigation icons located at the bottom right of the screen. From left to right, they include: a vertical arrow pointing up; a vertical arrow pointing down; a circular arrow indicating a loop or refresh; a speech bubble icon; a gear icon representing settings; a square icon with an upward-pointing arrow; a trash can icon; and three vertical dots representing more options.

Building wheel for scikit-surprise (setup.py) ... done

```
    Created wheel for scikit-surprise: filename=scikit_surprise-1.1.1-cp37-cp37m-linux_x86_64.whl size  
    Stored in directory: /root/.cache/pip/wheels/76/44/74/b498c42be47b2406bd27994e16c5188e337c657025ab
```

Successfully built scikit-surprise

Installing collected packages: scikit-surprise, surprise

Successfully installed scikit-surprise-1.1.1 surprise-0.

```
[ ] #Source: https://surprise.readthedocs.io/en/stable/getting\_started.html
from surprise import SVD
from surprise import Dataset
from surprise.model_selection import train_test_split
from surprise import accuracy
```

```
[ ] # Load the movielens-100k dataset (download it if needed)
data = Dataset.load_builtin('ml-100k')
```

```
[ ] # sample random trainset and testset  
# test set is made of 25% of the ratings.
```

Welcome to Surprise' document | RecSys_MF.ipynb - Colaboratory | Recommender-Systems-[Netflix] | Microsoft Word - BellKorSolutio | MovieLens 100K Dataset | GroupLens.org

groupLens.org/datasets/movielens/100k/

groupLens about datasets publications blog

MovieLens 100K Dataset

MovieLens 100K movie ratings. Stable benchmark dataset. 100,000 ratings from 1000 users on 1700 movies. Released 4/1998.

- [README.txt](#)
- [ml-100k.zip](#) (size: 5 MB, [checksum](#))
- [Index of unzipped files](#)

Permalink: <https://grouplens.org/datasets/movielens/100k/>

Datasets

[MovieLens](#)

[WikiLens](#)

[Book-Crossing](#)

[Book Genome Dataset](#)

[Jester](#)

[EachMovie](#)

[HetRec 2011](#)

[Serendipity 2018](#)

[Personality 2018](#)

[Learning from Sets of Items 2019](#)

search grouplens.org 

Stay in Touch [Contact Us](#) [Follow us on Twitter](#)

Project Links [MovieLens](#)

Code [GroupLens on GitHub](#)

Cyclopath

Welcome to Surprise' document x | RecSys_MF.ipynb - Colaboratory x | MovieLens 100K Dataset | GroupLens x | Recommender-Systems-[Netflix] x | Microsoft Word - BellKorSolutio x | +

groupLens.org/datasets/movielens/100k/

groupLens about datasets publications blog

MovieLens 100K Dataset

MovieLens 100K movie ratings. Stable benchmark dataset. 100,000 ratings from 1000 users on 1700 movies.
Released 4/1998.

- [README.txt](#)
- [ml-100k.zip](#) (size: 5 MB, [checksum](#))
- [Index of unzipped files](#)

Permalink: <https://grouplens.org/datasets/movielens/100k/>

10⁵

$10^3 \quad 1.7 \times 10^3$

1.7×10^6

Datasets

[MovieLens](#)

[WikiLens](#)

[Book-Crossing](#)

[Book Genome Dataset](#)

[Jester](#)

[EachMovie](#)

[HetRec 2011](#)

[Serendipity 2018](#)

[Personality 2018](#)

[Learning from Sets of Items 2019](#)

search grouplens.org 

Stay in Touch

[Contact Us](#)

[Follow us on Twitter](#)

Project Links

[MovieLens](#)

Code

[GroupLens on GitHub](#)

Cyclonpath

40 / 40

+ Code + Text

Connect ▾

↑ ↓ ↻ ⌂ ⚙️ 📁 🗑️ ⋮

```
# Load the movielens-100k dataset (download it if needed),
data = Dataset.load_builtin('ml-100k')

[ ] # sample random trainset and testset
# test set is made of 25% of the ratings.
trainset, testset = train_test_split(data, test_size=.25)

# We'll use the famous SVD algorithm.
algo = SVD()

# Train the algorithm on the trainset, and predict ratings for the testset
algo.fit(trainset)
predictions = algo.test(testset)

# Then compute RMSE
accuracy.rmse(predictions)
```

RMSE: 0.9463

0.9462710000722717

+ Code + Text

Connect ▾



```
# Load the movielens-100k dataset (download it if needed),
data = Dataset.load_builtin('ml-100k')

[ ] # sample random trainset and testset
# test set is made of 25% of the ratings.
trainset, testset = train_test_split(data, test_size=.25)

# We'll use the famous SVD algorithm.
algo = SVD() ✓

# Train the algorithm on the trainset, and predict ratings for the testset
algo.fit(trainset)
predictions = algo.test(testset)

# Then compute RMSE
accuracy.rmse(predictions)
```

RMSE: 0.9463

0.9462710000722717

Getting Started — Surprise 1 dc x RecSys_MF.ipynb - Colaboratory x MovieLens 100K Dataset | Grouped x Recommender-Systems-[Netflix] x Microsoft Word - BellKorSolutio x +

colab.research.google.com/drive/1wrVVng8EFinYUm6MQ8X7X3BOudz2PG6x#scrollTo=V3uInQ_y9q70

+ Code + Text Connect |   | 

test set is made of 25% of the ratings.
trainset, testset = train_test_split(data, test_size=.25)

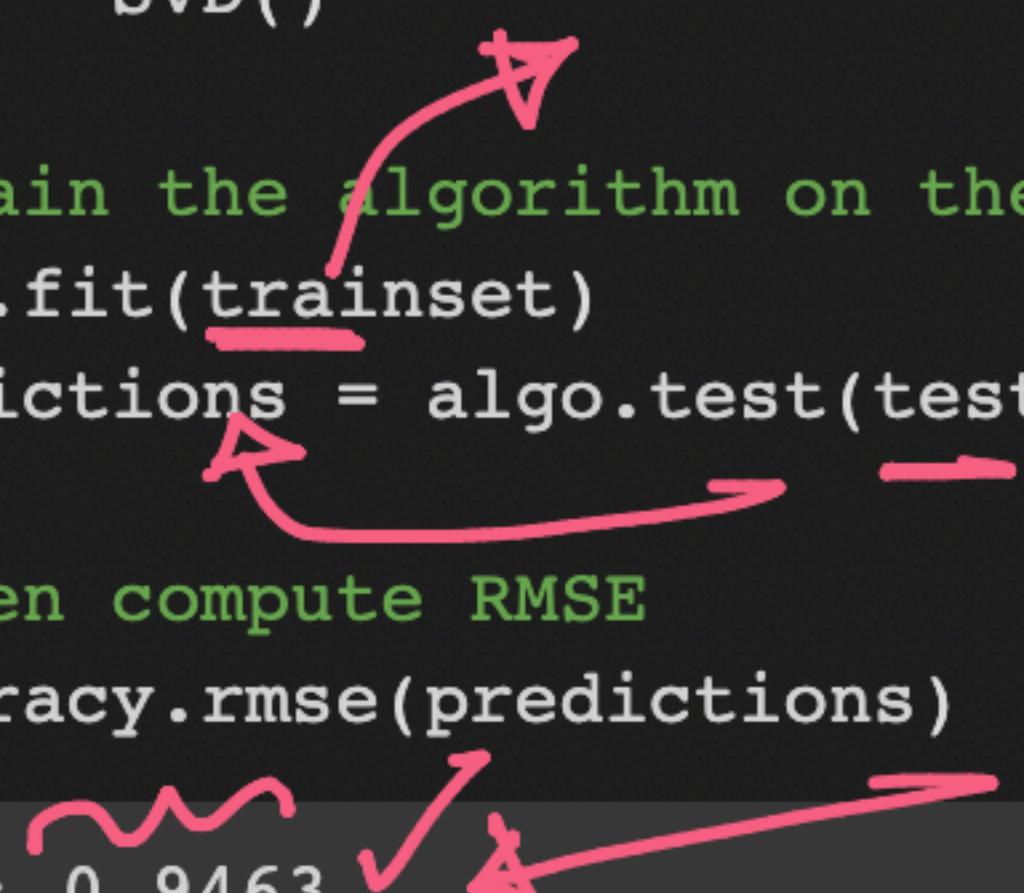
We'll use the famous SVD algorithm.
algo = SVD()

Train the algorithm on the trainset, and predict ratings for the testset
algo.fit(trainset)
predictions = algo.test(testset)

Then compute RMSE
accuracy.rmse(predictions)

RMSE: 0.9463
0.9462710000722717

0 - 5 .



0 - 5 .

[] uid = str(196) # raw user id (as in the ratings file). They are **strings**!
iid = str(302) # raw item id (as in the ratings file). They are **strings**!

get a prediction for specific users and items.
pred = algo.predict(uid, iid)
pred

43 / 43

+ Code + Text

Connect ▾



```
algo.fit(trainset)
predictions = algo.test(testset)
```



```
# Then compute RMSE
```

```
accuracy.rmse(predictions)
```

RMSE: 0.9463

0.9462710000722717

```
[ ] uid = str(196) # raw user id (as in the ratings file). They are **strings**!  
    iid = str(302) # raw item id (as in the ratings file). They are **strings**!
```

```
# get a prediction for specific users and items
```

```
pred = algo.predict(uid, iid)
```

pred

```
Prediction(uid='196', iid='302', r_ui=None, est=4.227320528336027, details={'was_impossible':  
False})
```

surprise.readthedocs.io/en/stable/prediction_algorithms.html

Surprise
stable

Search docs

USER GUIDE

Getting Started

Using prediction algorithms

Baselines estimates configuration

Similarity measure configuration

How to build your own prediction algorithm

Notation standards, References

FAQ

API REFERENCE

prediction_algorithms package

The model_selection package

similarities module

accuracy module

dataset module

Trainset class

Reader class

dump module

Docs » Using prediction algorithms

Edit on GitHub

Using prediction algorithms

Surprise provides a bunch of built-in algorithms. All algorithms derive from the `AlgoBase` base class, where are implemented some key methods (e.g. `predict`, `fit` and `test`). The list and details of the available prediction algorithms can be found in the `prediction_algorithms` package documentation.

Every algorithm is part of the global Surprise namespace, so you only need to import their names from the Surprise package, for example:

```
from surprise import KNNBasic
algo = KNNBasic()
```

Some of these algorithms may use **baseline estimates**, some may use a **similarity measure**. We will here review how to configure the way baselines and similarities are computed.

Baselines estimates configuration

Note

This section only applies to algorithms (or similarity measures) that try to minimize the following regularized squared error (or equivalent):

$$\sum_{r_{ui} \in R_{train}} (r_{ui} - (\mu + b_u + b_i))^2 + \lambda (b_u^2 + b_i^2).$$

For algorithms using

MF → Surprise

Q

✓ { lots of millions of params
↔
deep learning

✓ desktop/laptop → cloud ✓
↔

→ NMF & MF → Spark (MLops)
↔

Using prediction algorithms — x | RecSys_MF.ipynb - Colaboratory x | MovieLens 100K Dataset | Grouped x | Recommender-Systems-[Netflix] x | Microsoft Word - BellKorSolutio x | sklearn.decomposition.NMF — x +

scikit-learn.org/stable/modules/generated/sklearn.decomposition.NMF.html

scikit learn Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.1.2 Other versions

Please cite us if you use the software.

sklearn.decomposition.NMF Examples using sklearn.decomposition.NMF

`class sklearn.decomposition.NMF(n_components=None, *, init='None', solver='cd', beta_loss='frobenius', tol=0.0001, max_iter=200, random_state=None, alpha='deprecated', alpha_W=0.0, alpha_H='same', l1_ratio=0.0, verbose=0, shuffle=False, regularization='deprecated')`

[source]

Non-Negative Matrix Factorization (NMF).

Find two non-negative matrices, i.e. matrices with all non-negative elements, (W , H) whose product approximates the non-negative matrix X . This factorization can be used for example for dimensionality reduction, source separation or topic extraction.

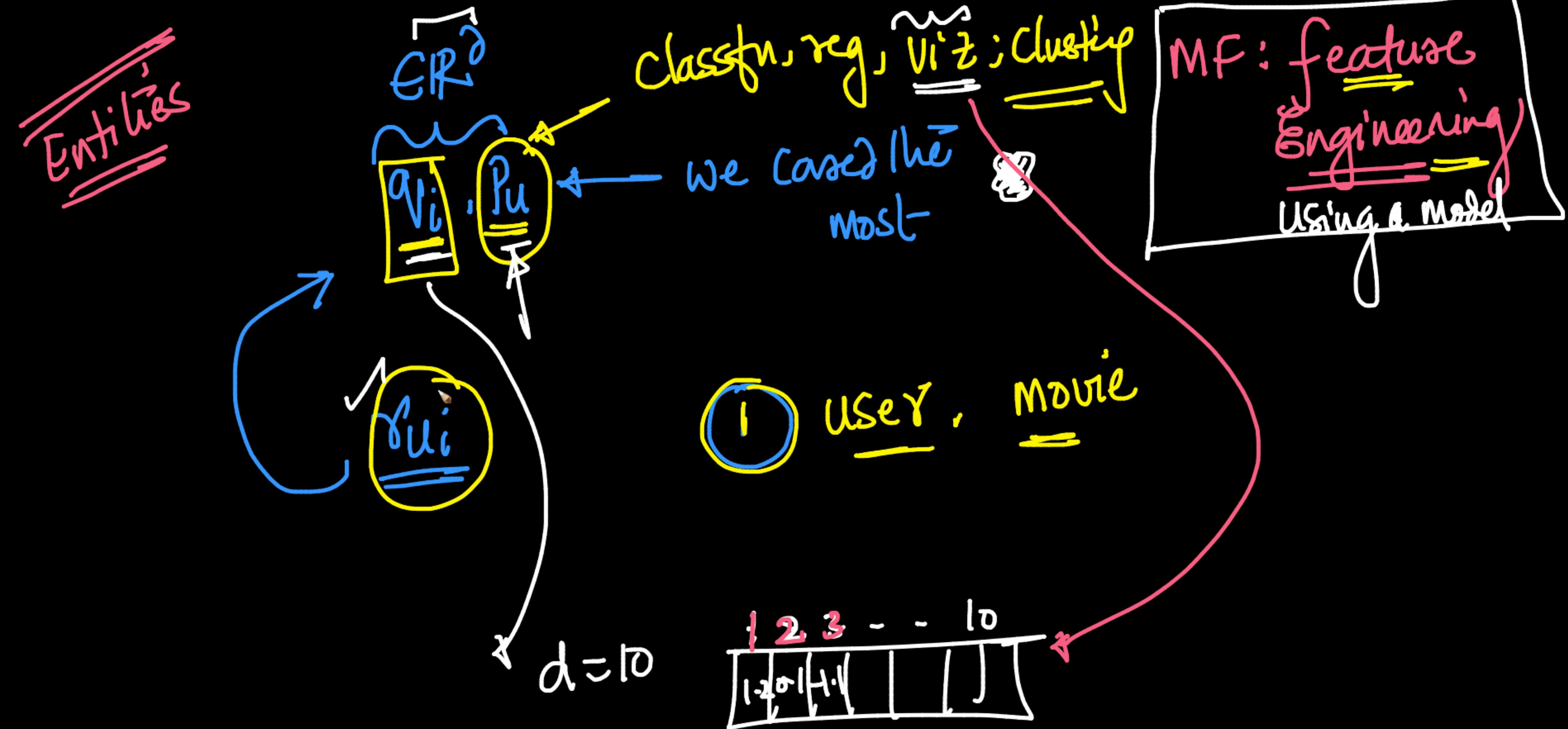
The objective function is:

$$\begin{aligned} L(W, H) = & 0.5 * ||X - WH||_{loss}^2 \\ & + \text{alpha_W} * \text{l1_ratio} * \text{n_features} * ||\text{vec}(W)||_1 \\ & + \text{alpha_H} * \text{l1_ratio} * \text{n_samples} * ||\text{vec}(H)||_1 \\ & + 0.5 * \text{alpha_W} * (1 - \text{l1_ratio}) * \text{n_features} * ||W||_{Fro}^2 \\ & + 0.5 * \text{alpha_H} * (1 - \text{l1_ratio}) * \text{n_samples} * ||H||_{Fro}^2 \end{aligned}$$

Where:

$$||A||_{Fro}^2 = \sum_{i,j} A_{ij}^2 \text{ (Frobenius norm)}$$
$$||\text{vec}(A)||_1 = \sum_{i,j} \text{abs}(A_{ij}) \text{ (Elementwise L1 norm)}$$

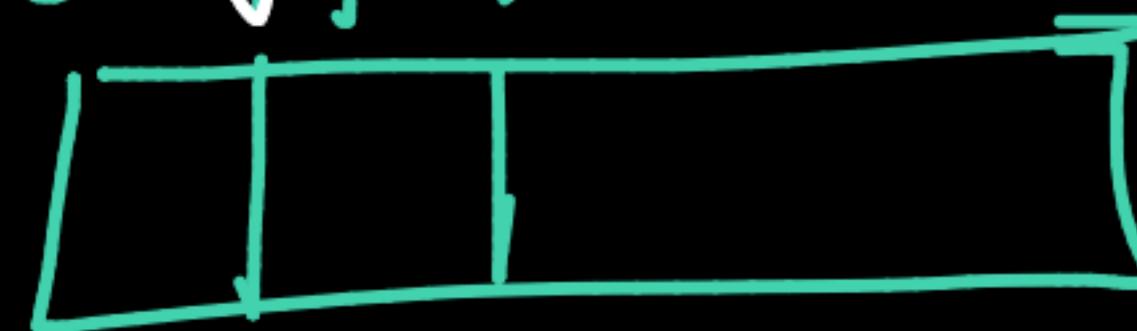
The generic norm supported beta-divergence loss. The choice



Manual FE

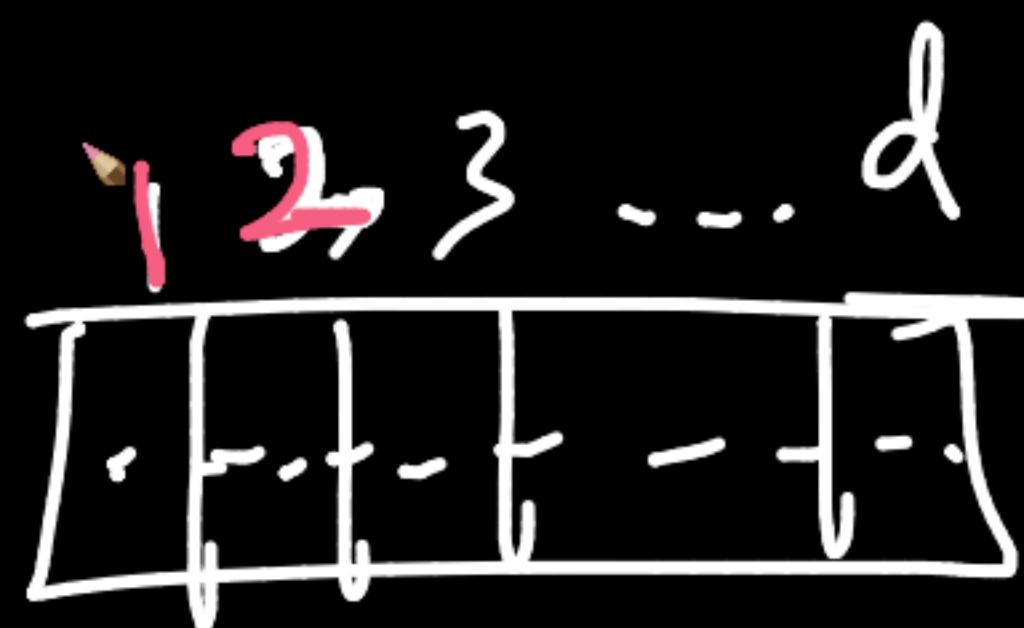
use Y:

Count Pin . . .



MF · FE

User:



no interpretability

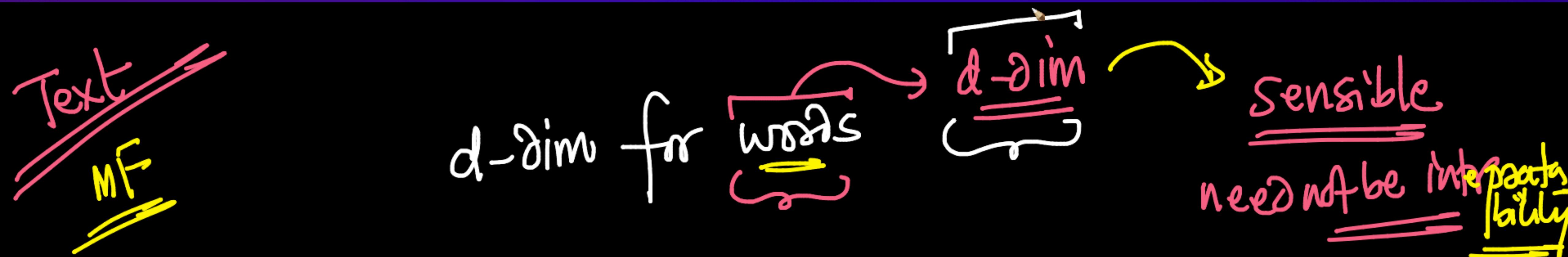
} → Viz(UMAP
tSNE)

Text

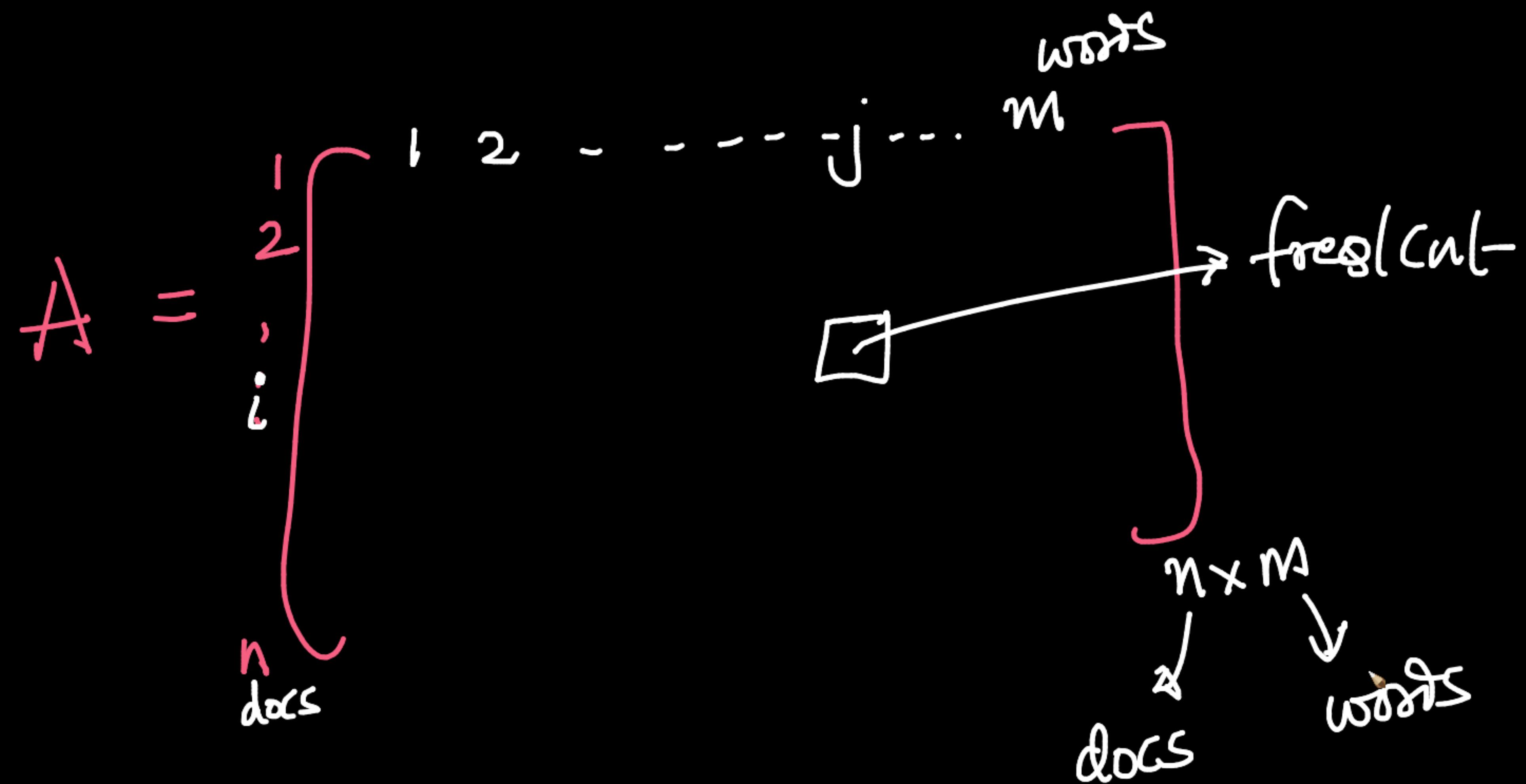
d.

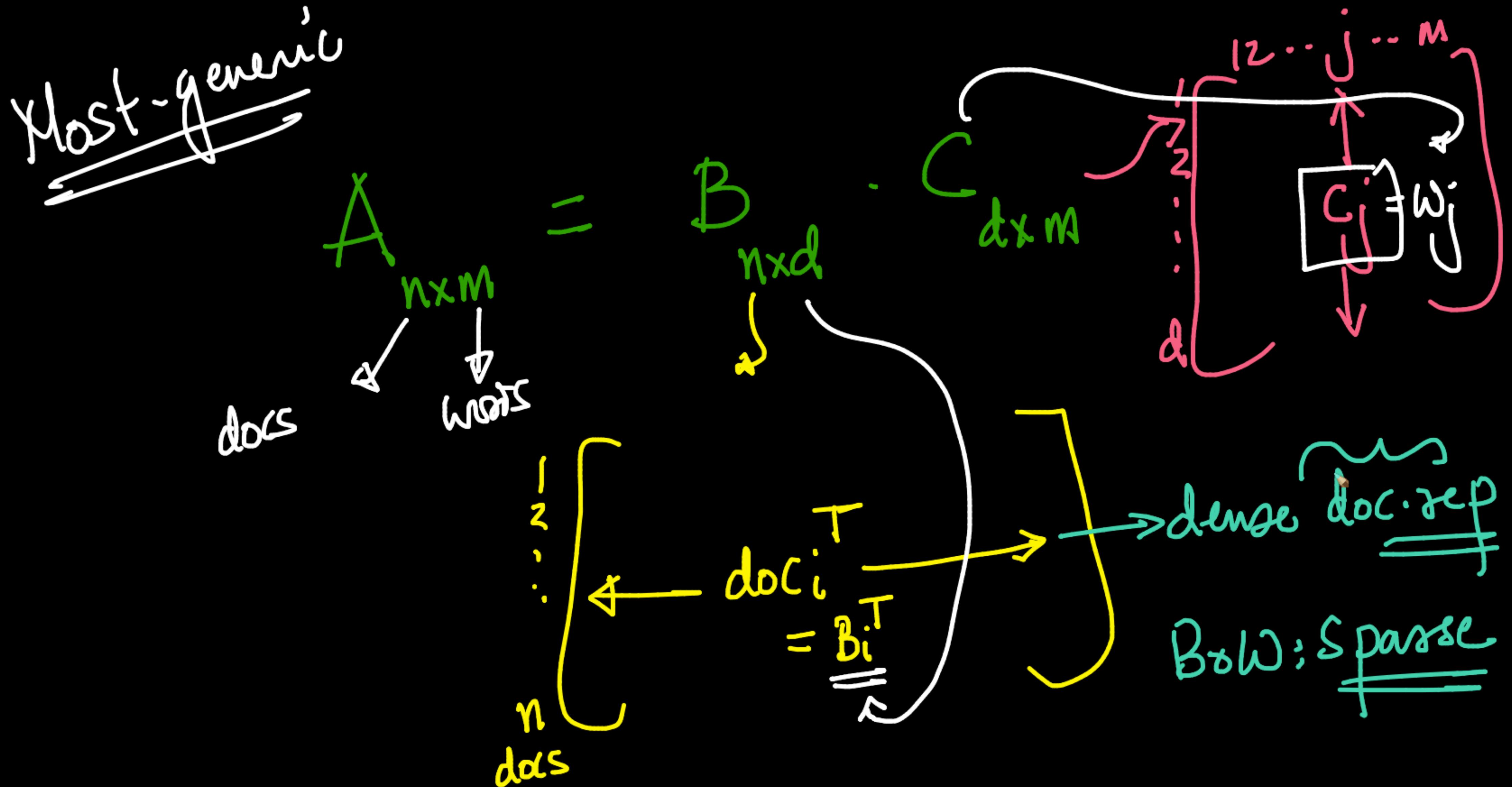
Text

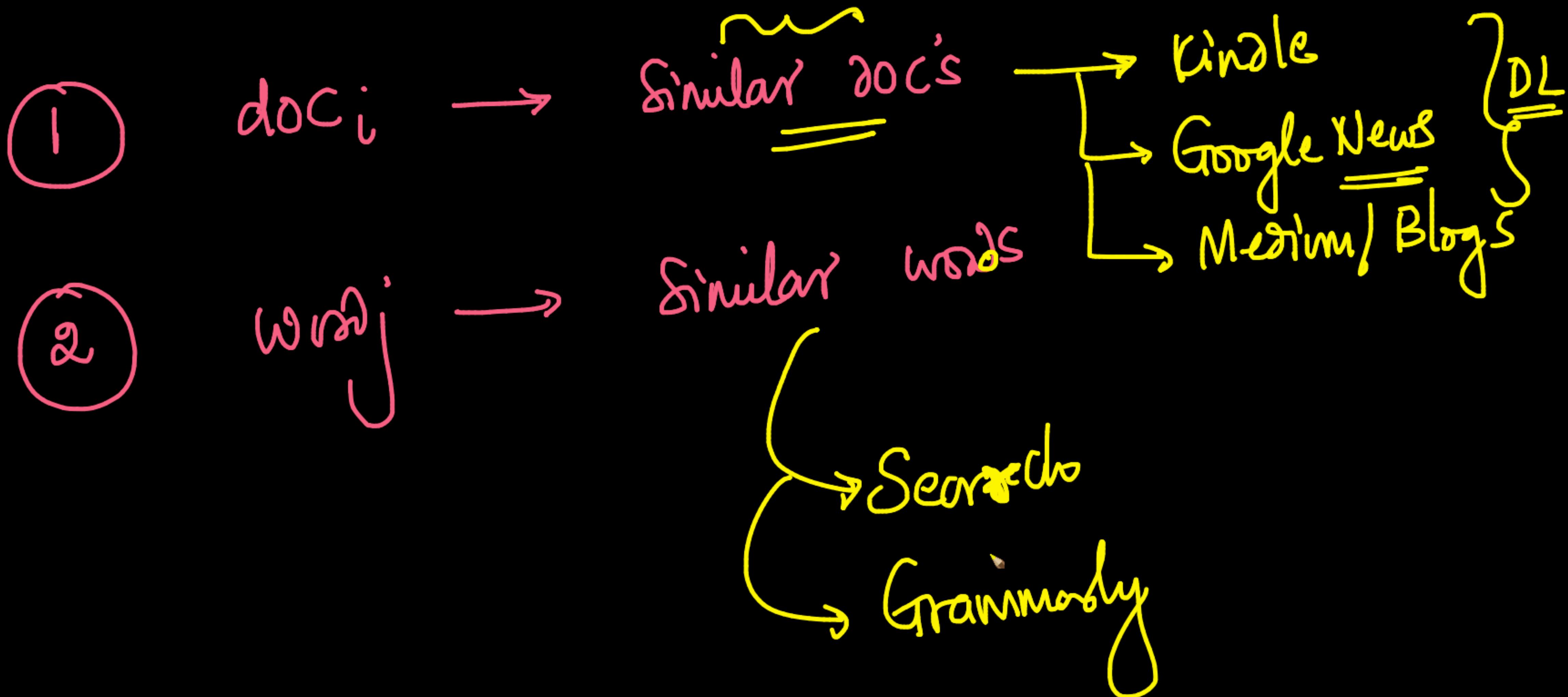
d.



~~Large~~ Corpus of text \rightarrow ~~wikipedia~~
~~Dataset~~

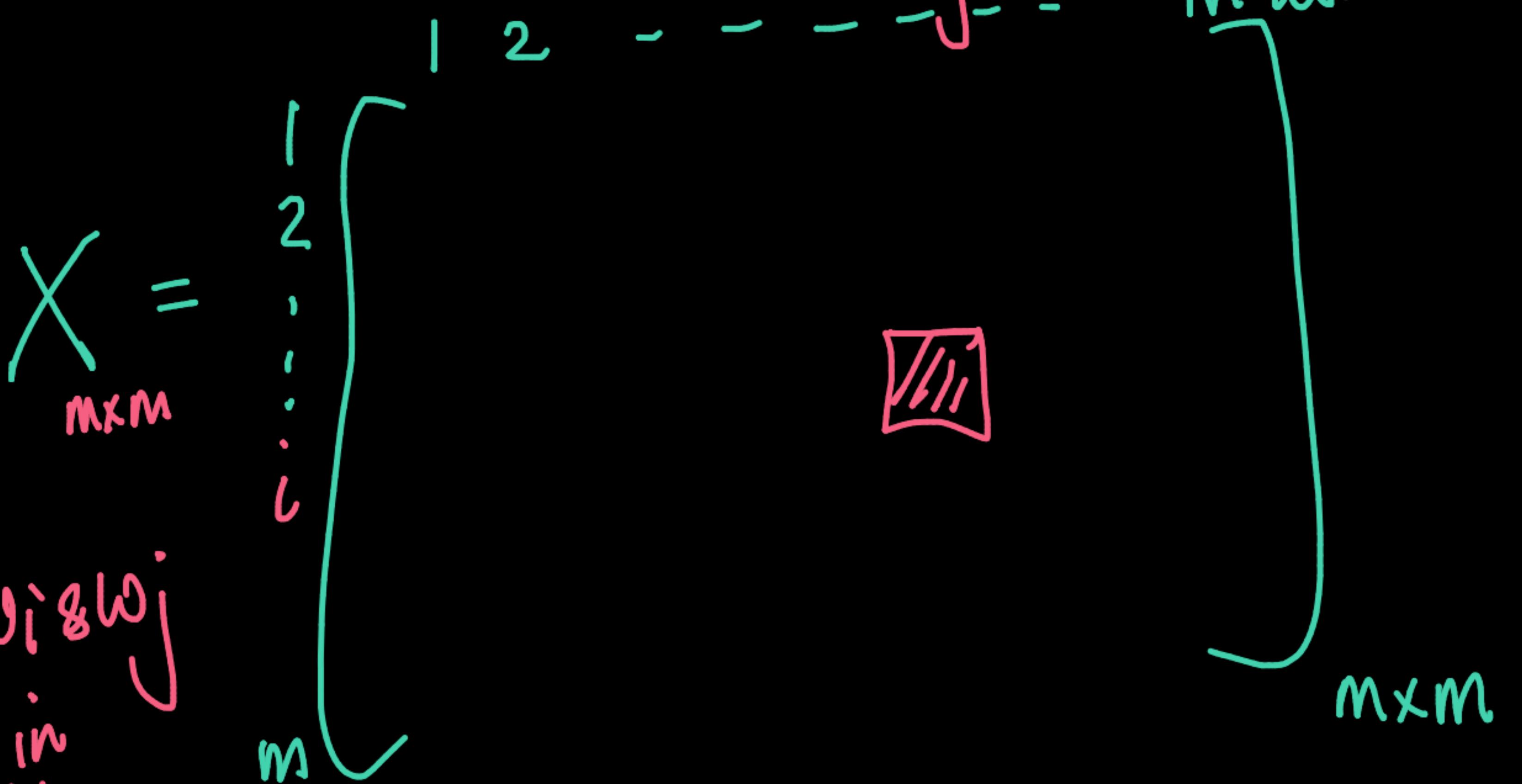






Alt

$x_{ij} = 1$ if $w_i \& w_j$
occurs in
a vicinity
n context



Using prediction algorithms — x | RecSys_MF.ipynb - Colaboratory x | MovieLens 100K Dataset | Grouped x | Recommender-Systems-[Netflix] x | Microsoft Word - BellKorSolutio x | sklearn.decomposition.NMF — x +

scikit-learn.org/stable/modules/generated/sklearn.decomposition.NMF.html

scikit
learn

Prev Up Next

scikit-learn 1.1.2

Other versions

Please cite us if you use the software.

sklearn.decomposition.NMF

Examples using sklearn.decomposition.NMF

See also:

[DictionaryLearning](#)

Find a dictionary that sparsely encodes data.

[MiniBatchSparsePCA](#)

Mini-batch Sparse Principal Components Analysis.

[PCA](#)

Principal component analysis.

[SparseCoder](#)

Find a sparse representation of data from a fixed, precomputed dictionary.

[SparsePCA](#)

Sparse Principal Components Analysis.

[TruncatedSVD](#)

Dimensionality reduction using truncated SVD.



References

[1]

["Fast local algorithms for large scale nonnegative matrix and tensor factorizations"](#) Cichocki, Andrzej, and P. H. A. N. Anh-Huy. IEICE transactions on fundamentals of electronics, communications and computer sciences 92.3: 708-721, 2009.

[2]

["Algorithms for nonnegative matrix factorization with the beta-divergence"](#) Fevotte, C., & Idier, J. (2011). Neural Computation, 23(9).

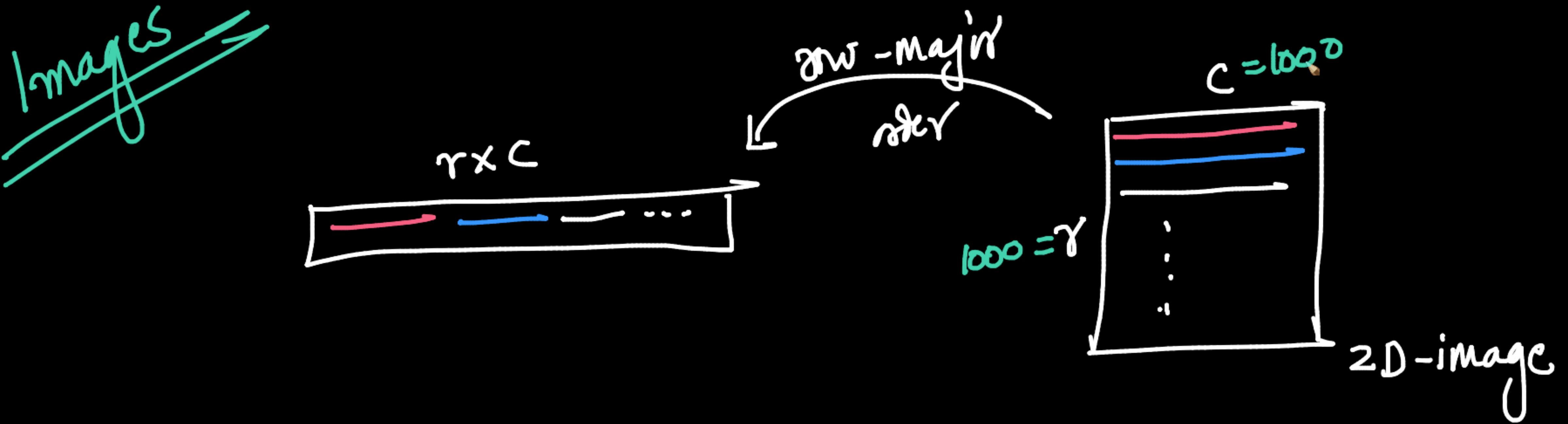
Examples

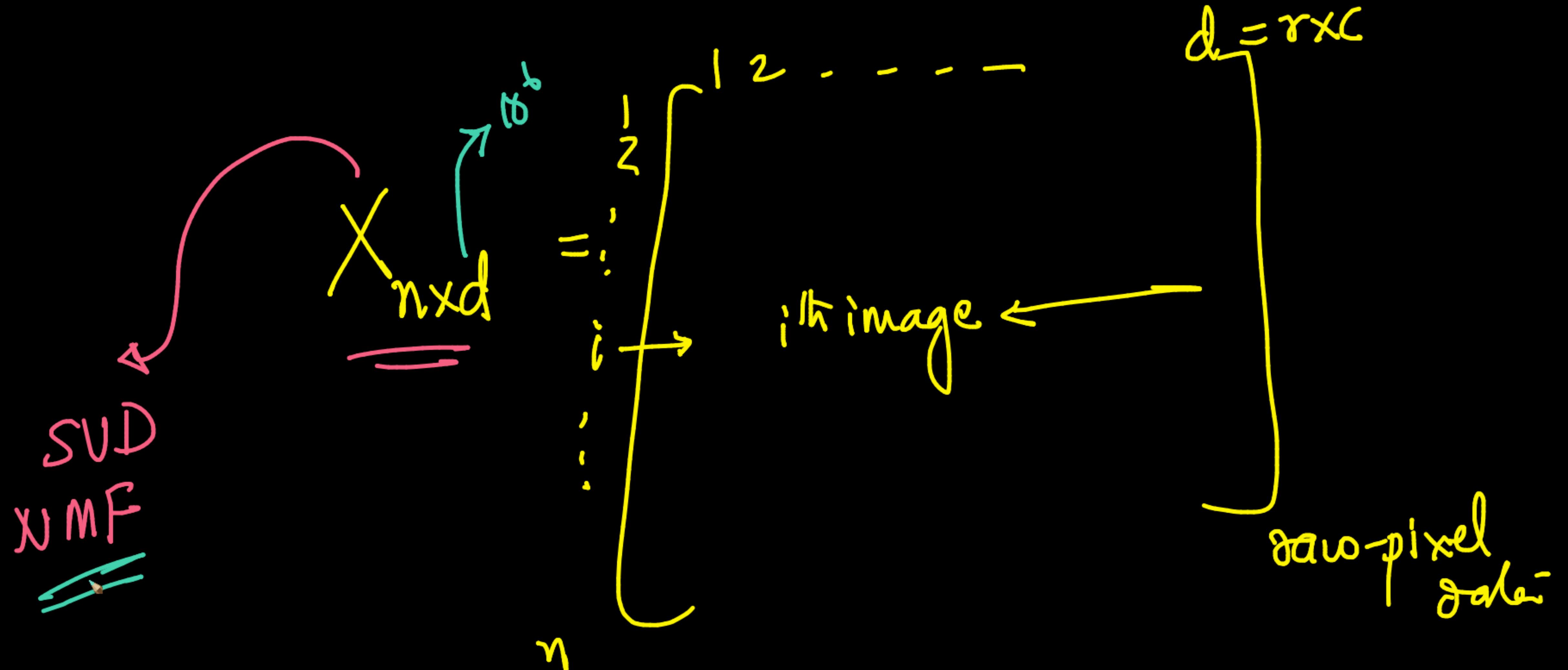


$$X_{n \times n} = B_{n \times d}$$
$$C_{d \times n}$$

B_i^T

~~C^T~~





$$\underset{10^50}{X} \underset{n \times d}{=} \underset{n \times M}{B} \underset{M \times d}{C} \xrightarrow{2D}$$

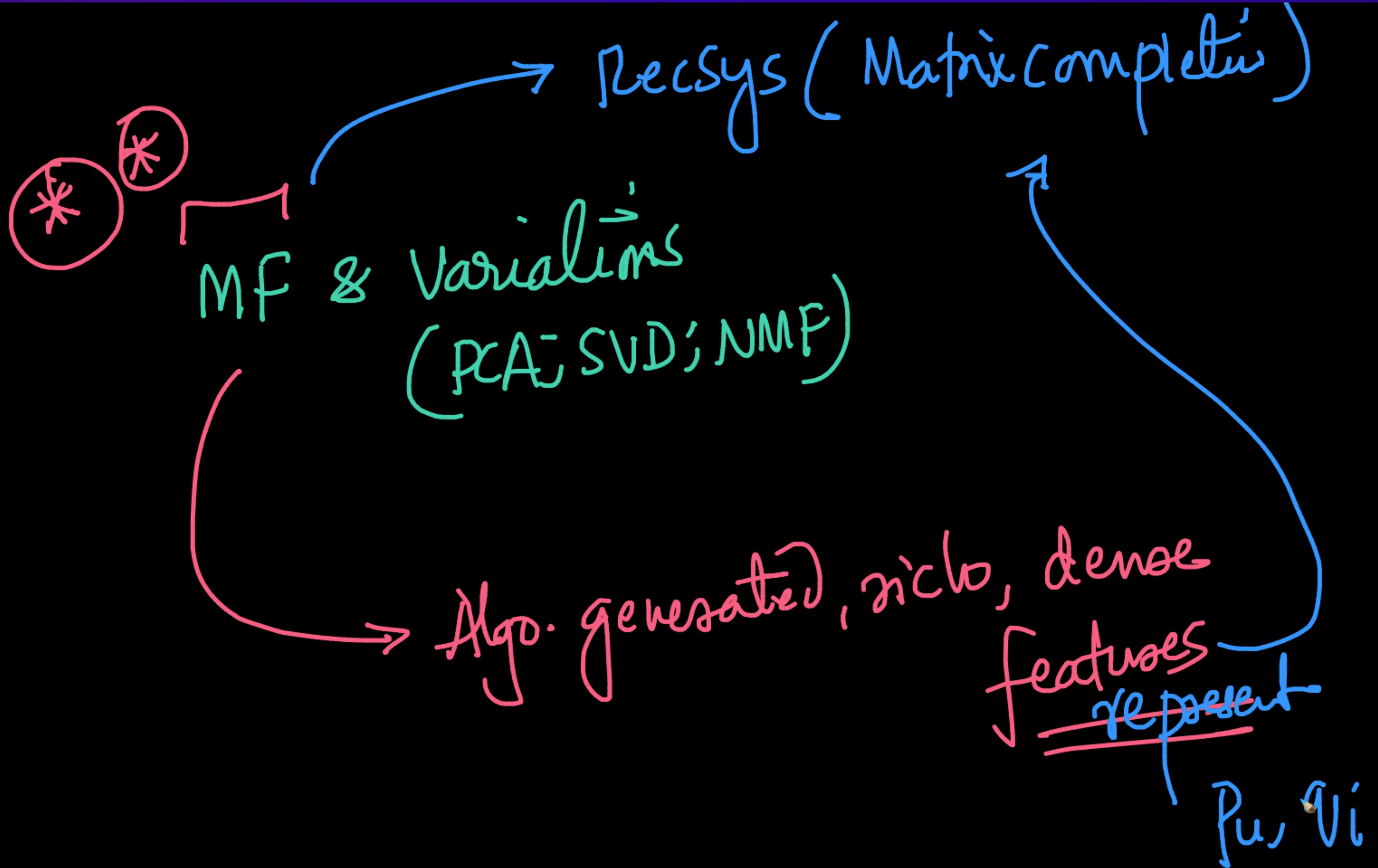
Diagram illustrating the dimensions of matrices B and C in the equation $X = BC$.

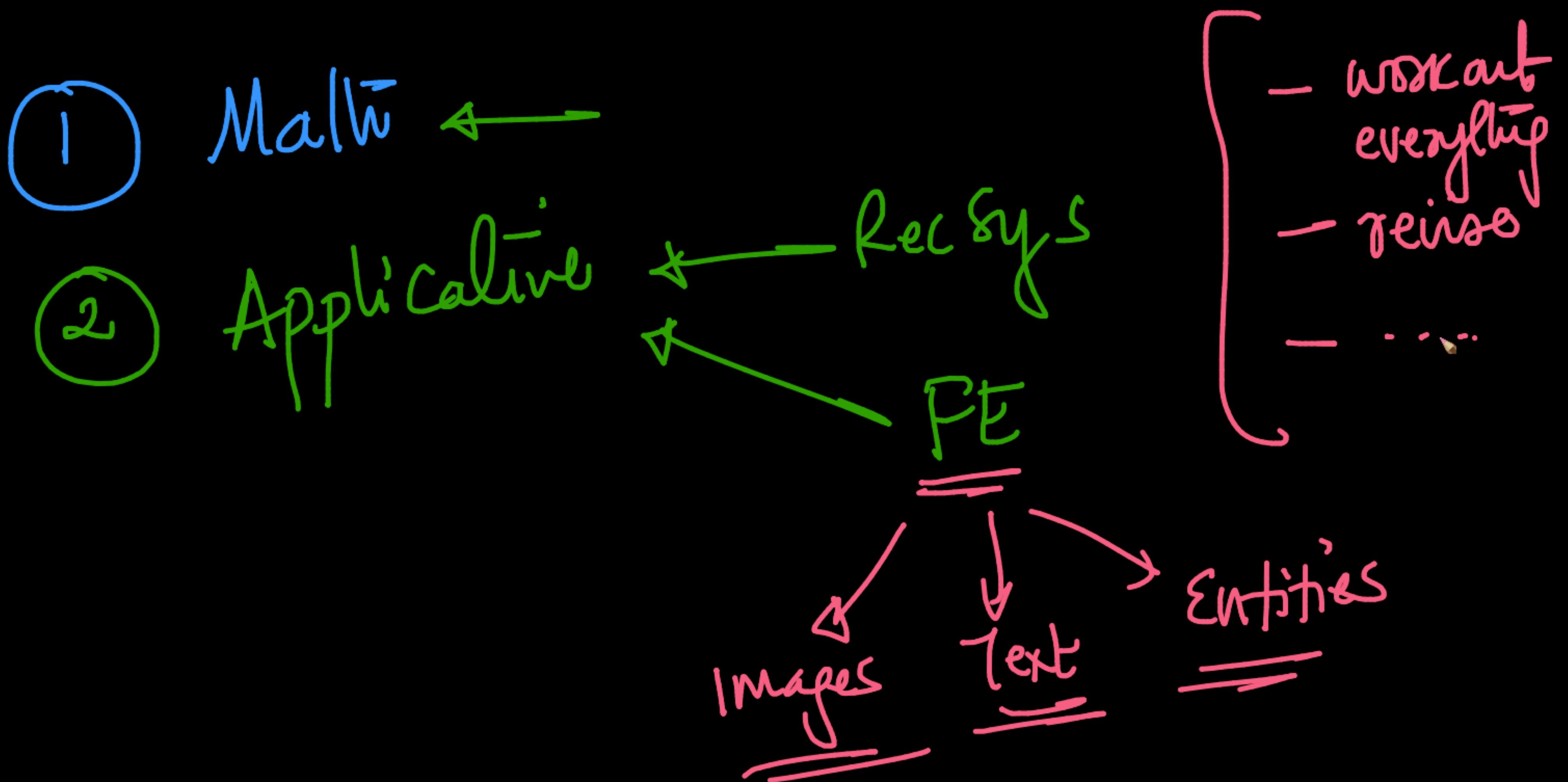
The matrix B has dimensions $n \times M$. The matrix C has dimensions $M \times d$. The resulting matrix X has dimensions $n \times d$, indicated by the pink arrow.

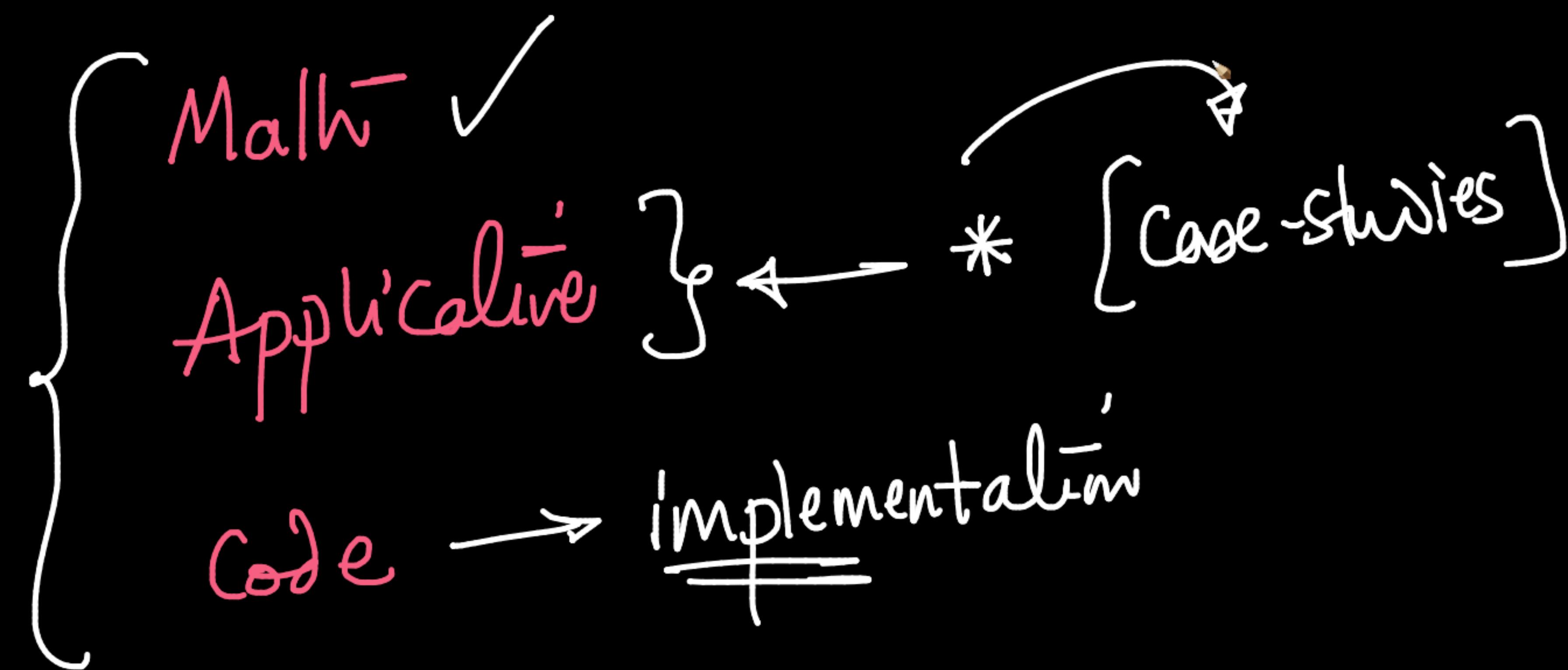
A yellow bracket labeled i, k indicates the columns of B . A yellow bracket labeled j, h indicates the rows of C . A yellow bracket labeled $1, 2, \dots, m$ indicates the columns of C and the rows of X .

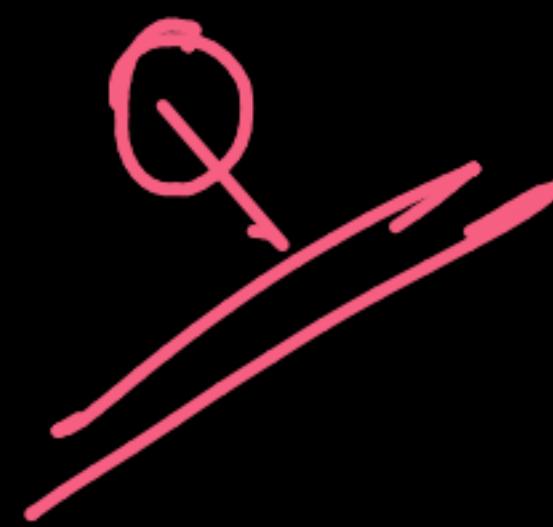
A green arrow points from the matrix B to a diagram showing a sequence of M boxes, each containing a vector. This diagram represents the M columns of B .

A green arrow points from the matrix C to a green bracket labeled "dense - Vec w", indicating the result of multiplying B and C .



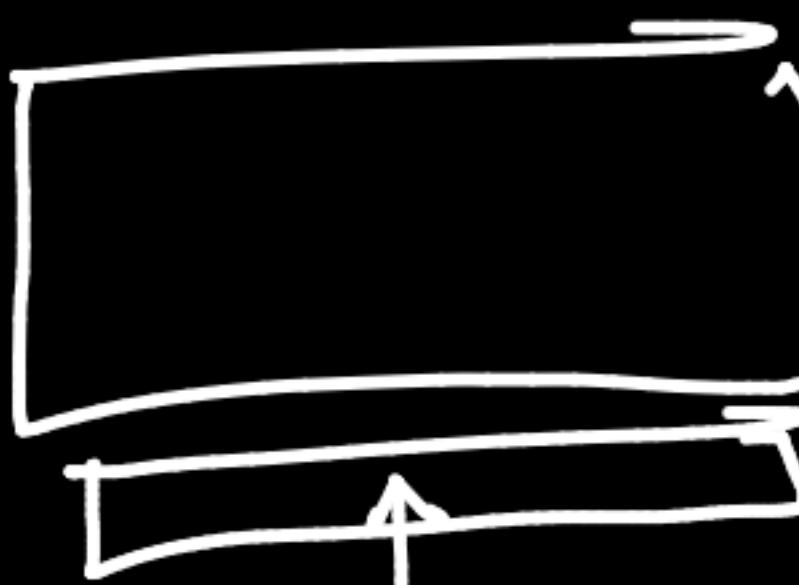






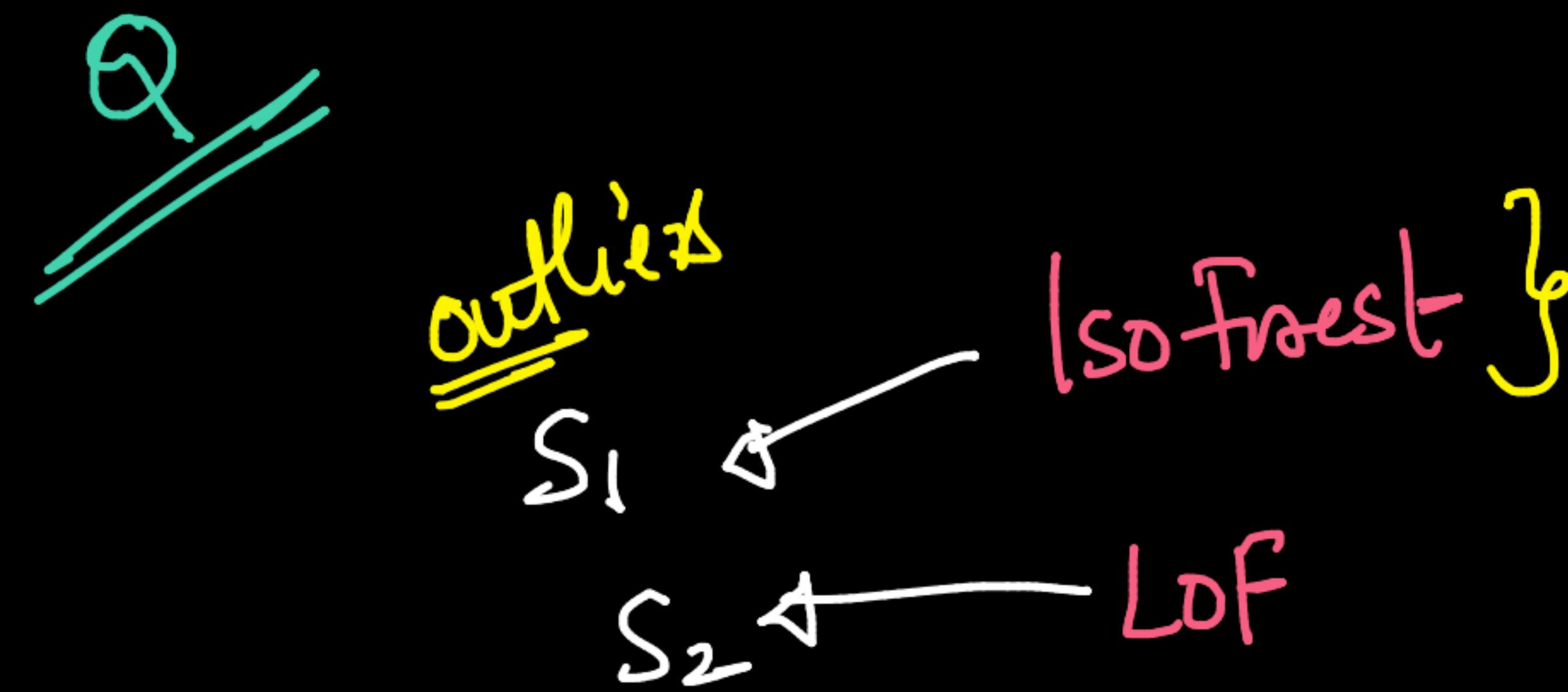
Search - KW

↓ { Webpage → KW
Video, ts → KW

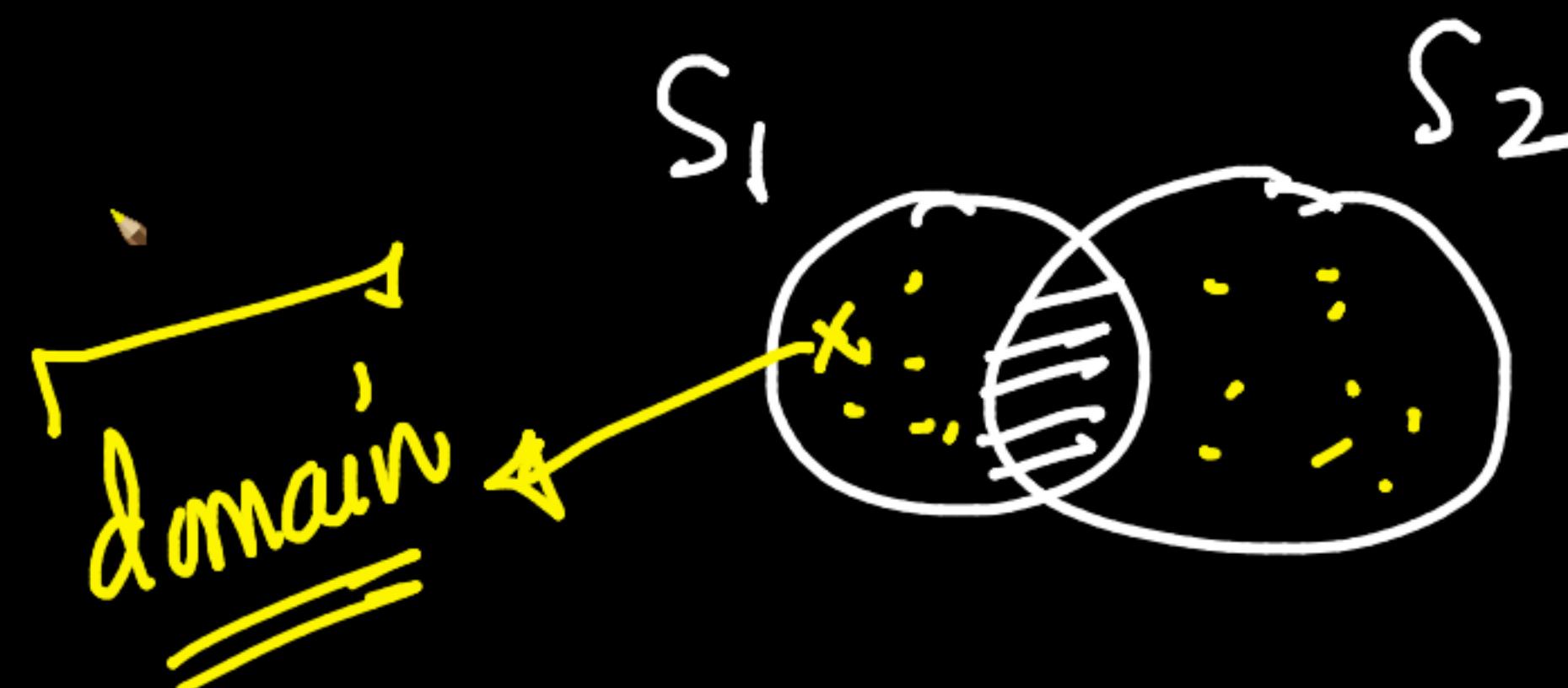


timestamp

DL ↗
Audio → text
↖



~~ground truth - given~~
NOT





InterviewBit Software Services | Live | ML, RS-3 (Basket Analysis) | +

scaler.com/meetings/i/ml-rs-3-basket-analysis-association-rules/live

You have left the meeting

We get frequent requests for your notes!

Notes written by you helps in understanding the topic better. You can upload the notes in two simple steps mentioned below

- 1 Scan the QR code with your iPad**
Scanner should be present in the top menu on your iPad
- 2 Upload Notes on the generated link**
All notes uploaded will be visible in the saved version of this session

OR

Drag and drop files or [click here to upload](#)

Files Uploaded from your computer appear here

