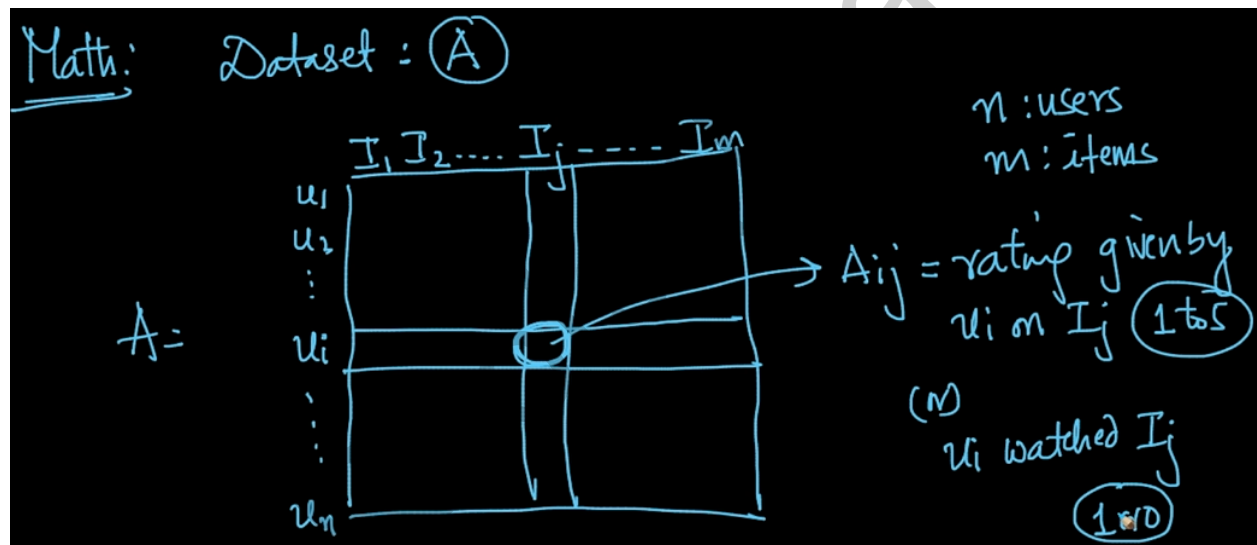


53.1 Problem Formulation: Movie Reviews

Let us look at a new type of problem called a Recommendation problem. One of the best examples of recommendation problems is recommending movies/videos on platforms like Netflix, Amazon Prime and Youtube.

Mathematical Formulation

Let us assume we were given a dataset and it is represented as a matrix 'A'. This dataset contains the data on the ratings given by the different users on different items. Let us assume we have 'n' users and 'm' items, then the matrix looks as below



This matrix contains either binary (or) discrete numerical values. Every element in this matrix is denoted as ' A_{ij} ' which denotes the rating given by the user ' i ' (ie., ' u_i ') on the item ' j ' (ie., ' I_j ').

' A_{ij} ' is binary, if the value just indicates whether the user ' u_i ' has watched the item ' I_j ', whereas ' A_{ij} ' is discrete numerical if the value indicates the rating given by the user ' u_i ' on the item ' I_j '. (This rating could be mostly on the scale of 1 to 5, where 1 being the worst and 5 being the best). These values are given by the users, and in case, if any user hasn't given the value for any item, then we can leave it blank.

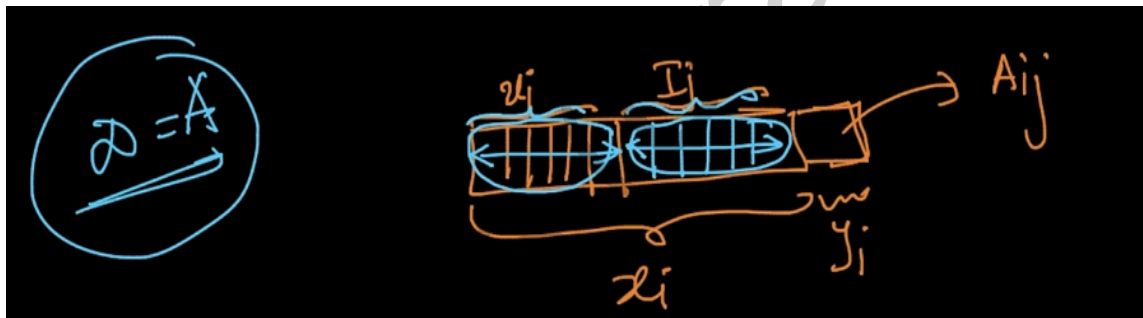
Properties of Matrix of Ratings

- 1) The Matrix of Ratings is very sparse. As there are a huge number of users and a huge number of items, it is not possible for every user to watch and rate every item. Hence we see this matrix sparse in most of the cases.

Sparsity of Matrix = Number of Empty cells/Total number of cells

The main goal of building a recommender system is if a user ' u_i ' has given a rating for a few items, then the recommender system should recommend new items to the user ' u_i ' depending on his/her interests.

- 2) Recommender System problem can also be posed as a Regression or a Classification problem.



The features for ' u_i ' and ' i_j ' are not given directly. We will be given only the matrix ' A ' as the dataset. We have to obtain the features for ' u_i ' and ' i_j ' using Feature Engineering.

Here ' D_{Train} ' = All the data points with non empty cells

' D_{Test} ' = All the data point with empty cells

We need to predict the values for all the non-empty cells. This is the way of posing a recommender system problem as a regression or a classification problem.

- 3) Recommender System problem can also be posed as a Matrix completion problem. Here ' A ' is the given matrix, and only a few values of the matrix ' A ' are given, whereas the remaining cells are empty. The problem is to fill up the empty cells with reasonable values, based on the values in the non-empty cells.

$A := \textcircled{RS}$ \textcircled{RS} as a matrix completion problem

Matrix Completion

$A = \begin{bmatrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{bmatrix}$

→ some values of A_{ij} are given
→ many A_{ij} are empty

✓ Fill up the empty cells with reasonable values based on values in the non-empty cells

$A =$

	I_1	I_2	...	I_j	...	I_m
u_i	3	4.5				5

predicted or filled value

fill up these empty cells

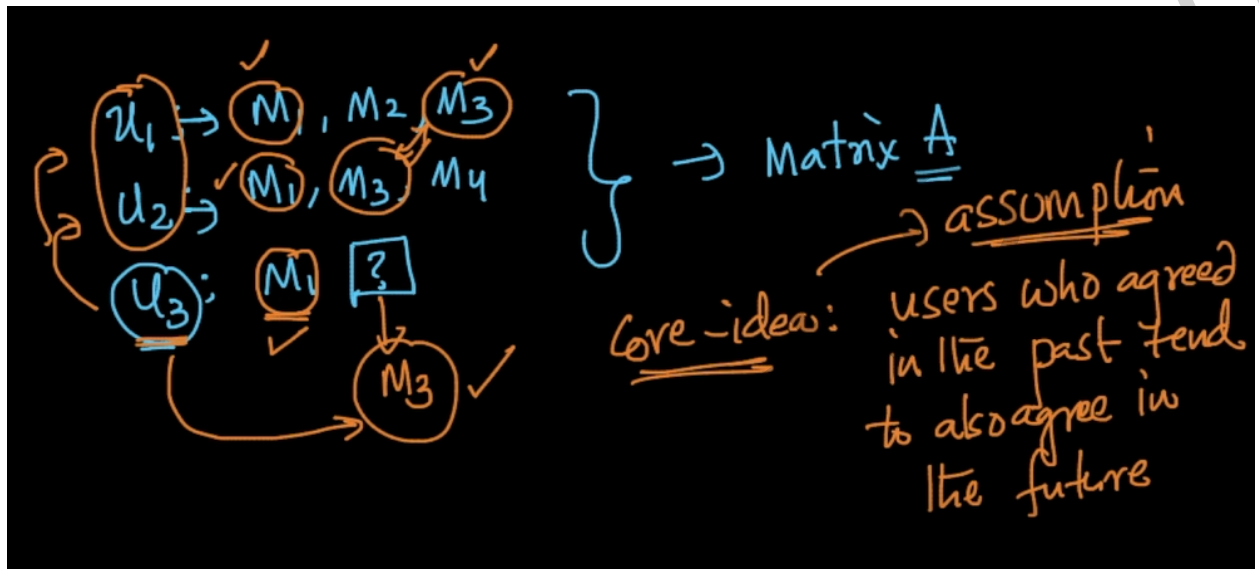
$\hat{A}_{i2} = 4.5$

✓ $\{ u_i \text{ could be recommend } I_2 \}$

53.2 Content Based vs Collaborative Filtering

There are 2 types of recommendation algorithms. They are

a) Collaborative Filtering based recommendation



For example, if we have 3 users ' u_1 ', ' u_2 ' and ' u_3 '. If the user ' u_1 ' has watched the movies ' M_1 ', ' M_2 ' and ' M_3 ', and the user ' u_2 ' has watched the movies ' M_1 ', ' M_3 ' and ' M_4 '. Let us assume the user ' u_3 ' has watched only the movie ' M_1 ', and now we have to recommend more movies to the user ' u_3 ', on the basis of the movies that were watched by the other users, who also have watched the movies watched by ' u_3 '.

So when we look at the above data, we can see that the movie ' M_1 ' which was watched by the user ' u_3 ', was also watched by the users ' u_1 ' and ' u_2 '. So it means both the users ' u_1 ' and ' u_2 ' share the similar taste with the user ' u_3 '. In such a case, it is always more appropriate to recommend the movies to ' u_3 ', that were watched by both ' u_1 ' and ' u_2 ', rather than those watched by only one of them. So here we could see the movie ' M_3 ' was watched by both ' u_1 ' and ' u_2 '. The movie ' M_2 ' was watched only by ' u_1 ' and the movie ' M_4 ' was watched only by ' u_2 '. So it is appropriate to recommend the movie ' M_3 ' first, to the user ' u_3 ', then ' M_2 ' and ' M_4 '.

The core idea/assumption in collaborative filtering is the users who agreed in the past, tend to agree in the future also.

b) Content based recommendation

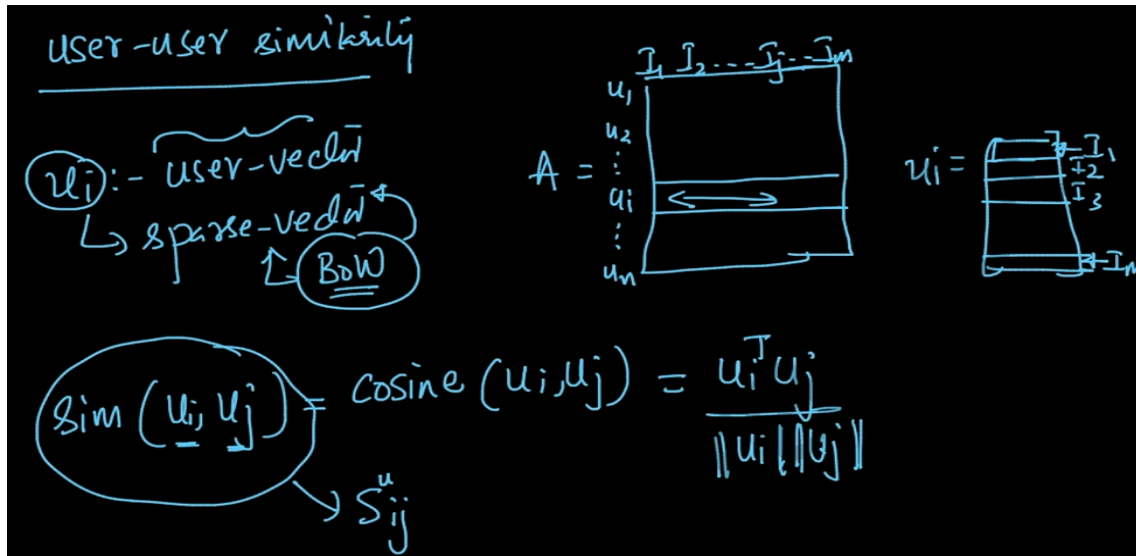
The content based recommendation system doesn't use the ratings ' A_{ij} ' given by the users. This system needs a way to represent the user vector ' u_i ' and the item vector ' I_j '. The point of content based filtering is that we have to know the content of both the users and the items. Usually we construct the user profile and the item profile using the content of the shared attribute space.

In this recommendation system, the engine compares the items that are already positively rated by the user with the items that he didn't rate and looks for similarities. Items similar to the positively rated ones, will be recommended to the user. The content based systems will only recommend the articles/items related to the same categories, and never recommend anything from the other categories.

53.3 Similarity based Algorithms

There are 2 types of Similarity based Recommender. They are

a) User-User Similarity



$u_i \rightarrow$ user vector (It is a sparse vector)

Every ' u_i ' is m-dimensional. It is the rating given by each user for all the 'm' items. So the similarity between two users ' u_i ' and ' u_j ' is given as

$$S_{ij}^u = \text{sim}(u_i, u_j) = \text{cosine-similarity}(u_i, u_j) = \frac{u_i^T u_j}{(\|u_i\| * \|u_j\|)}$$

So we need to build a similarity matrix with all ' S_{ij} '. So ' S^u ' is an $n \times n$ matrix, and ' S_{ij} ' represents how similar two users ' u_i ' and ' u_j ' are. Larger the value of ' S_{ij}^u ', the more similar the users ' u_i ' and ' u_j ' are.

For example, let us assume the users ' u_1 ', ' u_2 ' and ' u_7 ' are the 3 users most similar to the user ' u_{10} ', then we have to pick the items liked/bought by ' u_1 ', ' u_2 ' and ' u_7 ' and not yet bought/liked by the user ' u_{10} ', and finally recommend them to the user ' u_{10} '.

There is one major problem with user-user similarity. That is the user's preferences change over time. In such cases, the chances to build the best recommender systems are low. If the user preferences do not change over time, then we can build a matrix with respect to time. So an alternative to handle this problem is **item-item similarity**.

b) Item-Item Similarity

$I_i \rightarrow$ item vector

Each item is represented by an item vector. Every item vector ' I_i ' is m-dimensional, and it is the rating given for each item. So the similarity between two items is given as

$$S_{ij} = \text{sim}(I_i, I_j) = \text{cosine-similarity}(I_i, I_j) = I_i^T I_j / (||I_i|| * ||I_j||)$$

The major disadvantage with the user-user based similarity is the preferences of the user change over time, but the advantage of item-item similarity is the ratings on the given item do not change significantly.

For example, if a user ' u_{10} ' likes the items ' I_1 ', ' I_3 ' and ' I_7 ', then we have to generate 3 different sets with the items similar to ' I_1 ', ' I_3 ' and ' I_7 ' respectively. If any item is present in more than one of these sets, then the chances that it to be liked by the user ' u_{10} ' are more.

Rule of Thumb

When we have more number of users when compared to the items, and if the item ratings do not change much over the time after the initial period, then it is better to use item-item similarity approach over the user-user similarity approach.

53.4 Matrix Factorization: PCA, SVD

To solve the recommendation system problems, we have a technique called **Matrix Factorization**. The core idea of Matrix Factorization is to decompose a given matrix into a product of two or more matrices.

If we can decompose matrix 'A' into a product of two or more matrices, we call it **Multiplicative Matrix Decomposition** (or) **Matrix Factorization**.

Principal Component Analysis (PCA)

We are given a data matrix ' $X_{n \times d}$ ', which is centered. The covariance matrix of 'X' is given as $S_{d \times d} = X^T X / (n-1)$

Let the eigenvalues be denoted as $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_d$

Let the respective eigenvectors be denoted as $w_1, w_2, w_3, \dots, w_d$

Now the covariance matrix $S_{d \times d}$ can be decomposed into a product of two or more matrices as

$$S_{d \times d} = W_{d \times d} \Lambda_{d \times d} W_{d \times d}^T$$

Where ' $W_{d \times d}$ ' is a matrix of the eigenvectors as the columns, and ' $\Lambda_{d \times d}$ ' is a diagonal matrix of the eigenvalues as shown below.

The diagram shows the decomposition of the covariance matrix $S_{d \times d}$ into its eigenvectors W and eigenvalues Λ . The equation $S_{d \times d} = W_{d \times d} \Lambda_{d \times d} W_{d \times d}^T$ is written at the top. Below it, the matrix W is shown as a $d \times d$ matrix with columns $w_1, w_2, w_3, \dots, w_d$. The matrix Λ is shown as a $d \times d$ diagonal matrix with eigenvalues $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_d$ on the diagonal. The matrix S is labeled as the 'Covariance matrix'. The matrix Λ is labeled as the 'diagonal matrix'. The matrix W is labeled as the 'eigen-vec of S'. The matrix W^T is labeled as the 'Transpose of W'.

This whole decomposition ($\mathbf{S}=\mathbf{W}\mathbf{\Lambda}\mathbf{W}^T$) is called **Eigen Decomposition**. PCA is the eigen decomposition of the covariance matrix 'S'.

Singular Value Decomposition (SVD)

SVD is a matrix factorization technique related to PCA. PCA is applied only on the covariance matrix (S) which is square and symmetric. SVD on the other hand can be applied on any rectangular matrix. So we can apply SVD directly on the data matrix which is of the **nxd** shape.

$$\mathbf{X}_{n \times d} = \mathbf{u}_{n \times n} \mathbf{\Sigma}_{n \times d} \mathbf{V}_{d \times d}^T$$

$\mathbf{u}_{n \times n} \rightarrow$ Left singular Vectors

$\mathbf{V}_{d \times d}^T \rightarrow$ Right singular vectors

$\mathbf{\Sigma}_{n \times d} \rightarrow$ Diagonal matrix of singular values of 'X' (ie., $s_1, s_2, s_3, \dots, s_d$)

The relation between the singular values (s_i) and eigenvalues (λ_i) is given by $s_i^2/(n-1) = \lambda_i$

$\lambda_i \rightarrow$ Eigen value

$s_i \rightarrow$ singular value

$\mathbf{u}_i \rightarrow i^{\text{th}}$ column of $\mathbf{u} \rightarrow i^{\text{th}}$ eigenvector of $(\mathbf{X}_{n \times d} \mathbf{X}_{d \times n}^T)_{n \times n}$

$\mathbf{v}_i \rightarrow i^{\text{th}}$ column of $\mathbf{v} \rightarrow i^{\text{th}}$ eigenvector of $(\mathbf{X}_{d \times n}^T \mathbf{X}_{n \times d})_{d \times d} \rightarrow \mathbf{S}$

Note

Eigenvalues and Eigenvectors are defined for square matrices only while singular values and singular vectors are defined for all rectangular matrices.

53.5 Matrix Factorization: Non Negative Matrix Factorization (NMF)

Let us assume we are given a matrix ' $\mathbf{A}_{n \times m}$ ' which has to be decomposed as a product of two matrices ' \mathbf{B} ' and ' \mathbf{C}^T ' such that ' \mathbf{B} ' is a **nxd** matrix and ' \mathbf{C}^T ' is a **dxm** matrix, such that for all $i, j \rightarrow B_{ij} \geq 0, C_{ij} \geq 0$.

$$\mathbf{A}_{n \times m} = \mathbf{B}_{n \times d} (\mathbf{C}^T)_{d \times m}$$

53.6 Matrix Factorization for Collaborative Filtering

Consider a given matrix 'A' which consists of the ratings given by various users for various items.

Let $A_{ij} \rightarrow$ Rating given by the user ' u_i ' on item ' I_j '

Moreover 'A' is a sparse matrix. Now 'A' can be decomposed as a product of two matrices 'B' and ' C^T '.

$$\mathbf{A}_{n \times m} = \mathbf{B}_{n \times d} \mathbf{C}_{d \times m}^T \text{ such that } d > 0 \text{ and } d \leq \min(m, n).$$

Here 'B', 'C' are column vectors.

$n \rightarrow$ number of users

$m \rightarrow$ number of items

So the optimization problem is $\text{argmin}_{\mathbf{B}, \mathbf{C}} \sum_{i,j} (\mathbf{A}_{ij} - \mathbf{B}_i^T \mathbf{C}_j)^2$, such that ' A_{ij} ' is not empty.

The goal is to find 'B' and 'C' matrices. So here we have to perform this summation only where ' A_{ij} ' is not empty.

Using the Stochastic Gradient Descent algorithm, we now have to solve the optimization problem $\text{argmin}_{\mathbf{B}, \mathbf{C}} \sum_{\mathbf{A}_{ij} \text{ is not empty}} (\mathbf{A}_{ij} - \mathbf{B}_i^T \mathbf{C}_j)^2$. It is like minimizing the squared loss. After solving the optimization problem, we'll be having the values of the matrices 'B' and 'C'.

Handwritten diagram illustrating the matrix factorization $\mathbf{A} = \mathbf{B} \mathbf{C}^T$.

Matrix \mathbf{B} is shown as a column vector of size $n \times d$, where $n = \# \text{ users}$ and d is the latent dimension. The elements are B_1, B_2, \dots, B_n . A specific row B_i is highlighted with an arrow pointing to the right.

Matrix \mathbf{C} is shown as a column vector of size $d \times m$, where $m = \# \text{ items}$. The elements are C_1, C_2, \dots, C_m . A specific column C_j is highlighted with an arrow pointing to the left.

The diagram shows the product $\mathbf{B} \mathbf{C}^T$ resulting in the matrix \mathbf{A} .

After obtaining the values of the matrices 'B' and 'C', we have to complete the matrix 'A' by filling the missing values of ' A_{ij} ' using the formula $\mathbf{A}_{ij} = \mathbf{B}_i^T \mathbf{C}_j$. Let us name this product matrix obtained as $\mathbf{A}^{\wedge} = \mathbf{B} \mathbf{C}^T$. The matrix ' \mathbf{A}^{\wedge} ' has no empty cells.

The key thing now is $\mathbf{A} - \mathbf{\hat{A}}$ has to be very small. We should ensure that for all i, j , the value of ' $\mathbf{\hat{A}}_{ij}$ ' should be very close to ' \mathbf{A}_{ij} '.

Now ' $\mathbf{\hat{A}}_{ij}$ ' is a complete form of ' \mathbf{A} ' using matrix factorization. The values in ' $\mathbf{\hat{A}}$ ' are the predicted ratings, whereas the values in ' \mathbf{A} ' are the actual ratings.

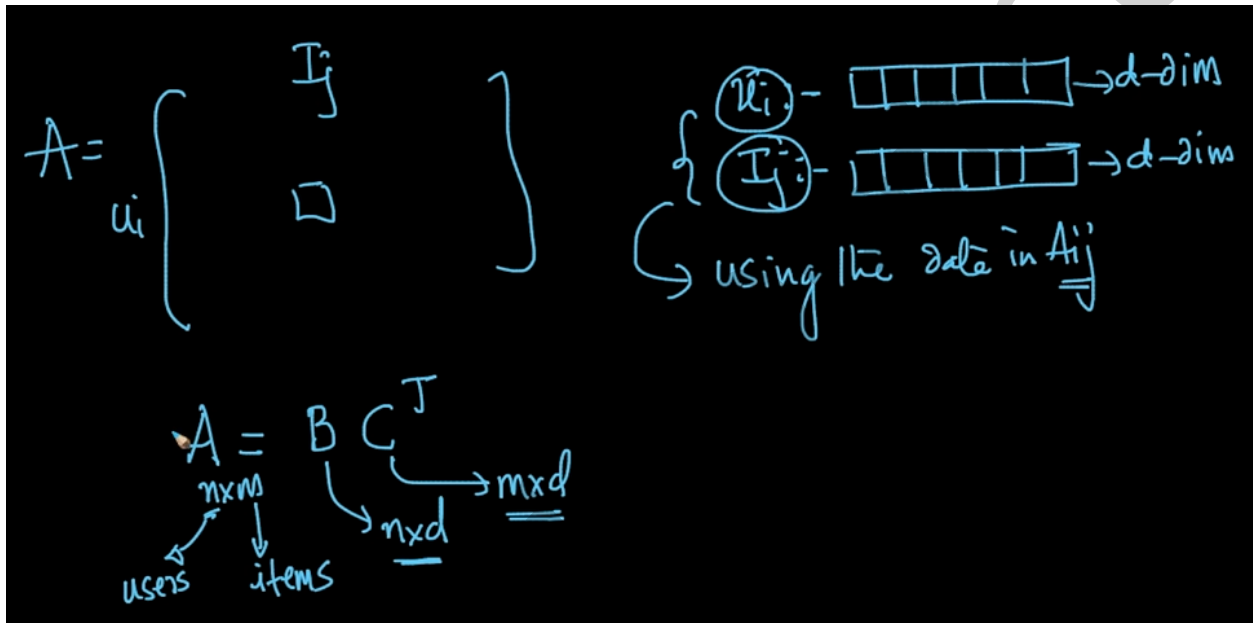
53.7 Matrix Factorization for Feature Engineering

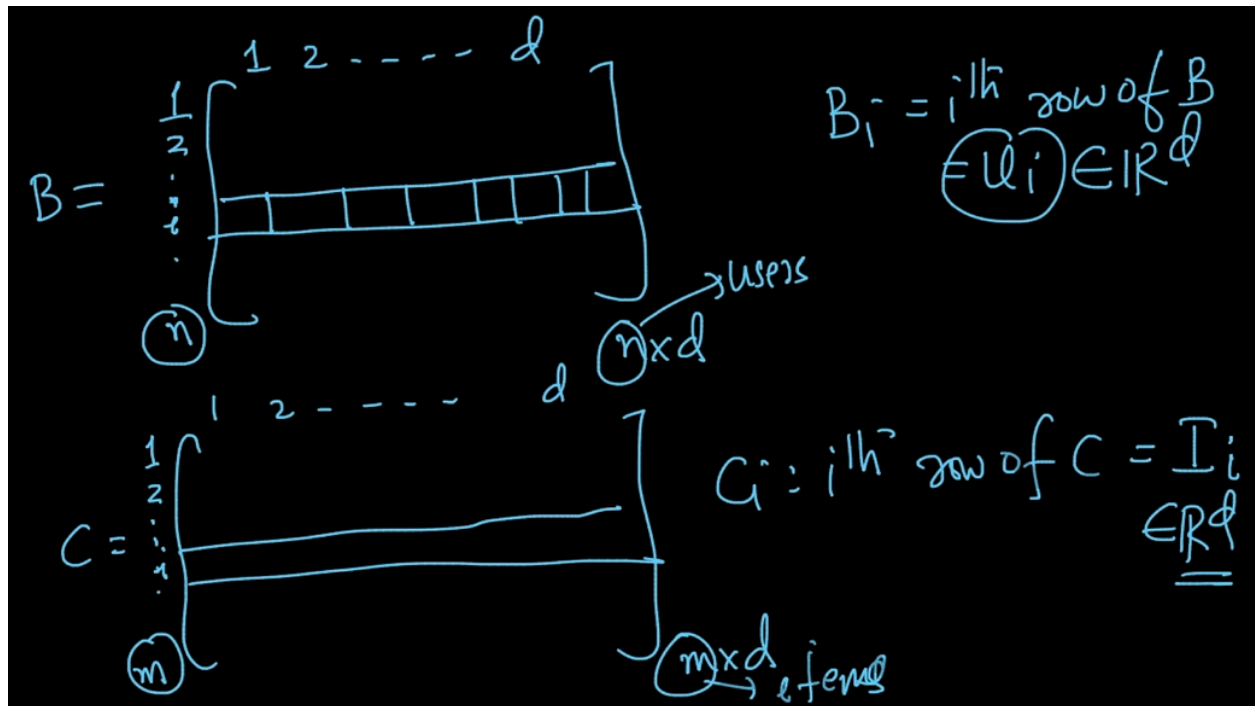
Let us assume we are given a matrix 'A'. Let us now decompose the matrix 'A' into a product of two matrices.

$$\mathbf{A}_{n \times m} = \mathbf{B}_{n \times d} \mathbf{C}_{m \times d}^T$$

Where $\mathbf{B}_{n \times d} \rightarrow$ Matrix of user vectors

$\mathbf{C}_{m \times d}^T \rightarrow$ Matrix of item vectors





As long as $d > 0$ and $d \leq \min(m, n)$, we can represent both 'B' and 'C' in d -dimensional form. So by leveraging the data in ' A_{ij} ', we have to represent the user and the item vectors in d -dimensional form. We arrive at d -dimensional representation using the matrix factorization.

For example, if ' u_i ', ' u_j ' are similar in their ratings of products, then the distance between ' u_i ' and ' u_j ' is very small. Similarly, if ' I_i ' and ' I_j ' are similar on the basis of the ratings, then the distance between ' I_i ' and ' I_j ' is very small.

$B_i \rightarrow i^{\text{th}} \text{ row of 'B'} = u_i \in \mathbb{R}^d$

$C_i \rightarrow i^{\text{th}} \text{ row of 'C'} = I_i \in \mathbb{R}^d$

In the later chapters, we'll learn more in detail on how the matrix factorization can be used in creating the word vectors.

53.8 Clustering as MF

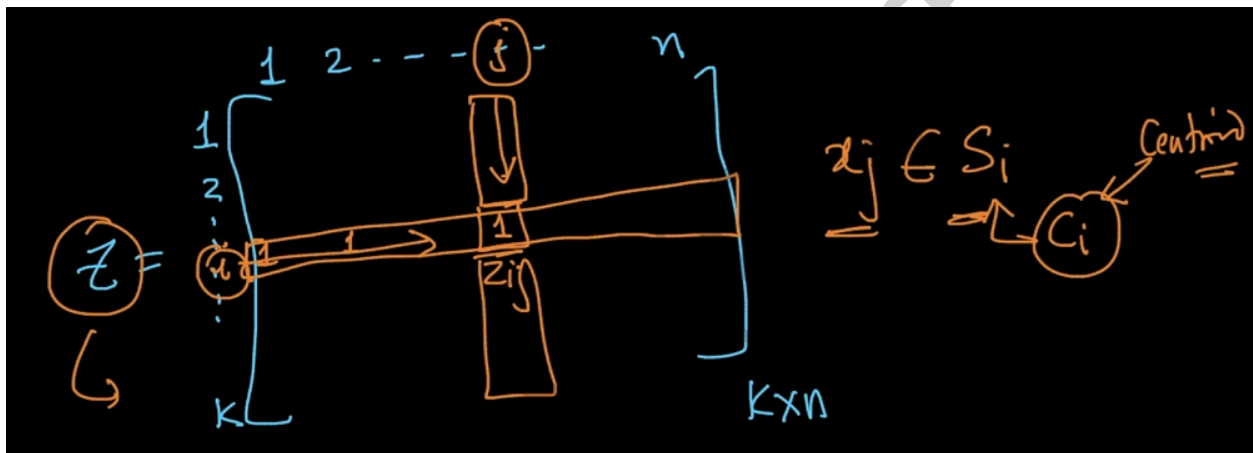
Let us take the example of the KMeans algorithm. The given dataset is denoted as $D = \{x_{ij}\}$. The task is to compute centroids $\{C_1, C_2, C_3, \dots, C_k\}$ and the cluster sets $\{S_1, S_2, S_3, \dots, S_k\}$ such that $S_i \cap S_j = \emptyset \forall i, j$.

The optimization problem we have to solve is

$$\min_{C_i} \sum_{i=1}^k \sum_{x \in S_i} \|x - C_i\|^2 \text{ where } x \in \mathbb{R}^d$$

Let us define a matrix 'Z' (of shape $k \times n$) such that

$$Z_{ij} = \begin{cases} 1, & \text{if } x_j \in S_i \\ 0, & \text{otherwise} \end{cases}$$



Here 'Z' is a binary matrix. So if $Z_{ij}=1$, then it means the point ' x_j ' belongs to the cluster ' S_i ' with the centre ' C_i '. The actual optimization problem was $\min_{C_i} \sum_{i=1}^k \sum_{x \in S_i} \|x - C_i\|^2$ where $x \in \mathbb{R}^d$. Now this problem becomes

$$\min_{C_i, Z_{ij}} \sum_{i=1}^k \sum_{j=1}^n Z_{ij} \|x_j - C_i\|^2$$

We have to convert the above problem into a matrix factorization problem. Let us take the matrices 'X' (of shape $d \times n$) and 'C' (of shape $d \times k$) where 'X' consists of all the data points, whereas 'C' consists of all the centroids. The optimization problem now becomes $\min_{C, Z} \|X - CZ\|^2$

Diagram illustrating matrix factorization $X = CZ$.

Matrix X is of size $d \times n$, with columns $x_1, x_2, \dots, x_j, \dots, x_n$.

Matrix C is of size $d \times k$, with columns $c_1, c_2, \dots, c_i, \dots, c_k$.

Matrix Z is defined as $z_{ij} = \begin{cases} 1 & \text{if } x_j \in S_i \\ 0 & \text{otherwise} \end{cases}$.

So if we want to write the KMeans clustering problem as a recommender system, then we have to decompose it as $X=CZ$, such that $Z_{ij}=1$ (or) 0 (this is called **0-1 constraint**) and in every column ' Z_j ' of the matrix ' Z ', only one value should be 1, and the remaining all values should be 0. (ie., sum of all the elements of the column vector $Z_j = 1$). There is no chance for the existence of negative values in the matrix ' Z '.

But finding the accurate solution for KMeans is NP-hard. Hence we go with approximation. So

KMeans = Matrix Factorization + (0-1 constraint) + column constraint

Without these constraints, KMeans is the same as Matrix Factorization.

Note

In the case of general matrix factorization, we do not have any constraints on the values of the matrices ' C ' and ' Z '.

In the case of Non negative matrix factorization, we have the constraints $C_{ij} \geq 0$ and $Z_{ij} \geq 0$.

In the case of KMeans, the constraint is $Z_{ij} = 0$ or 1 , and $\sum_{i=1}^k Z_{ij} = 1$.

53.9 Hyperparameter Tuning

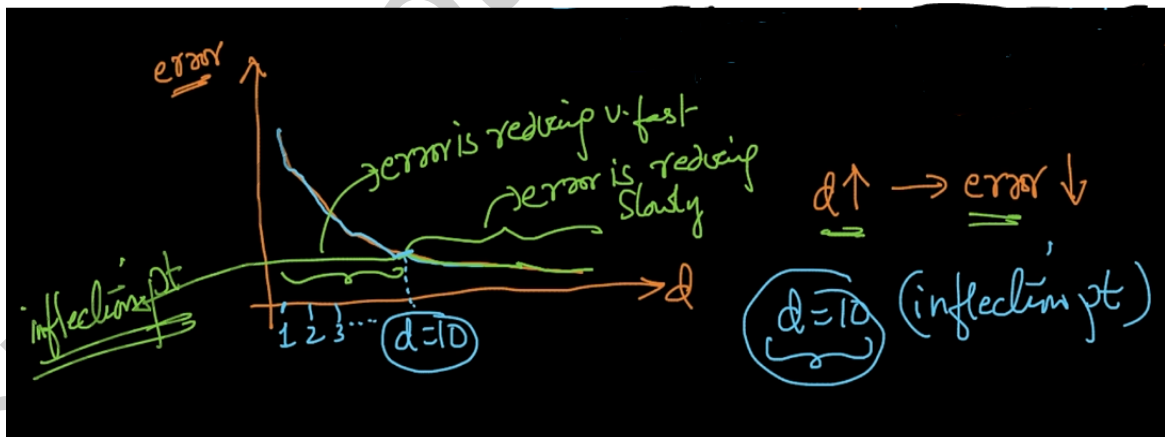
Let 'A' be a 'n x m' matrix and is decomposed into a product of two matrices.

$$\mathbf{A}_{n \times m} = \mathbf{B}_{n \times d} \mathbf{C}_{m \times d}^T \quad (\text{where } d > 0, \text{ and } d \leq \min(m, n))$$

Here 'd' is the hyperparameter.

How to select the value of 'd'?

- 1) It is problem specific. Here it depends on the task we are solving. For example, if we are working on the problem of predicting the ratings given by a user for an item, then if we do not need larger dimensional (or) smaller dimensional representation for both the user and the item vectors, then the domain expert could select some moderate value like 100. It could be specified by the domain expert.
- 2) The optimization problem we are solving for matrix factorization is $\min_{\mathbf{B}, \mathbf{C}} \sum (\mathbf{A}_{ij} - \mathbf{B}_i^T \mathbf{C}_j)^2$. We can find the 'd' value by minimizing the loss. Whichever value of 'd' gives the least amount of loss, that would be the optimal value.



Plot the values of 'd' on the 'X' axis and the error on the 'Y' axis. As the 'd' value increases, the error falls. Upto the value of $d < 10$, the error is reducing fast, and after crossing the point $d = 10$, the error is reducing slowly. So here $d = 10$ is the inflection point, because we could see a huge change in the behaviour of the loss with the increasing number of 'd' values.

- 3) The other way to find out the 'd' value is by plotting the 'd' values on the 'X' axis and the rate of change of the loss/error (ie., slope) on the 'Y' axis, and picking the inflection point.

AppliedRoots (Draft Copy)

53.10 Matrix Factorization for Recommender Systems: Netflix Prize Solution

Let $r_{ui} \rightarrow$ Rating given by the user 'u' on the item 'i'.

$q_i \rightarrow$ Item Vector

$p_u \rightarrow$ User Vector

So the optimization problem is

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

In the above optimization problem,

$(r_{ui} - q_i^T p_u)^2 \rightarrow$ Squared Loss

$(\|q_i\|^2 + \|p_u\|^2) \rightarrow L_2$ regularization to avoid overfitting

Solution 1

Using SGD, we have to compute $\partial L / \partial p_u$, $\partial L / \partial q_i$. But it takes more time, as there will be a huge number of users and items.

Solution 2 (Alternative Least Squares (ALS))

Step (a)

Assuming ' p_u ' as constant and running the gradient descent to find the best ' q_i '.

Step(b)

Assuming ' q_i ' as constant and running the gradient descent to find the best ' p_u '.

This way, we get the optimal values for both ' p_u ' and ' q_i ' by alternating. Both steps (a) and (b) are executed alternatively, till it converges to minima.

ALS is faster than SGD. Also ALS is a well suited algorithm when working with the least squares.

Let us assume 3 scalars. They are

$b_u \rightarrow$ Bias term for each user

$b_i \rightarrow$ Bias term for each item

$\mu \rightarrow$ Mean term

The optimization problem now becomes

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

Where $\mu \rightarrow$ Average ratings of all the users across all the items.
(ie., average of all the non empty cell ratings)

Here ' μ ' is a constant, whereas ' b_u ' and ' b_i ' are variables.

The main advantage of using Matrix Factorization in recommender systems is incorporation of problem specific changes is easy (like here we have incorporated the scalars ' μ ', ' b_u ' and ' b_i '). This advantage is not directly available in item-item and user-user similarities.

Note

Refer to the research paper in the link given below

[https://datajobs.com/data-science-repo/Recommender-Systems-\[Netflix\].pdf](https://datajobs.com/data-science-repo/Recommender-Systems-[Netflix].pdf)

53.11 Cold Start Problem

In case if we have a new user 'u' coming into the ecosystem, then we do not have any ratings from this user, as he is new to the ecosystem. So this new user vector will have empty cells. Similarly if a new item 'i' is added to the inventory, then we do not have any ratings for that item. Here this new item vector will have empty cells. Such a problem is called a **cold-start** problem.

Case 1 (When we have a new user added in the ecosystem)

In this case, we recommend top items based on the metadata. For example, if the new user is from San Francisco, and is using an ipad with safari browser, then those items that are used/purchased by the people in the same area, with the same devices are recommended to this new user. This will be a content based recommendation system.

Case 2 (When we have a new item added to the inventory)

In this case, it checks to which category this item belongs to (like electronics, clothing, sports, provisions, etc), cost of the product, description of the item and then recommends this new item to those users who bought similar items from the same categories. This also will be a content based recommendation system.

These 2 cases come under a cold-start problem.

If a new user creates an account (or) if a new item is added to the inventory, for such users/items, the system doesn't have enough information to make a decision.

Cold-Start problem is that problem, where the system is not able to recommend the items to the users. For every recommender system, it is required to build a user profile by considering the user's activities and behaviours with the system. On the basis of the user's previous history, the system makes the decisions and recommends the items consequently.

If a new user has not rated some items and not yet viewed/visited some items, then it would be difficult for the system to build a model on that basis. The Cold-Start problem arises in 3 different situations.

- a) For new users
- b) For new items
- c) For new community

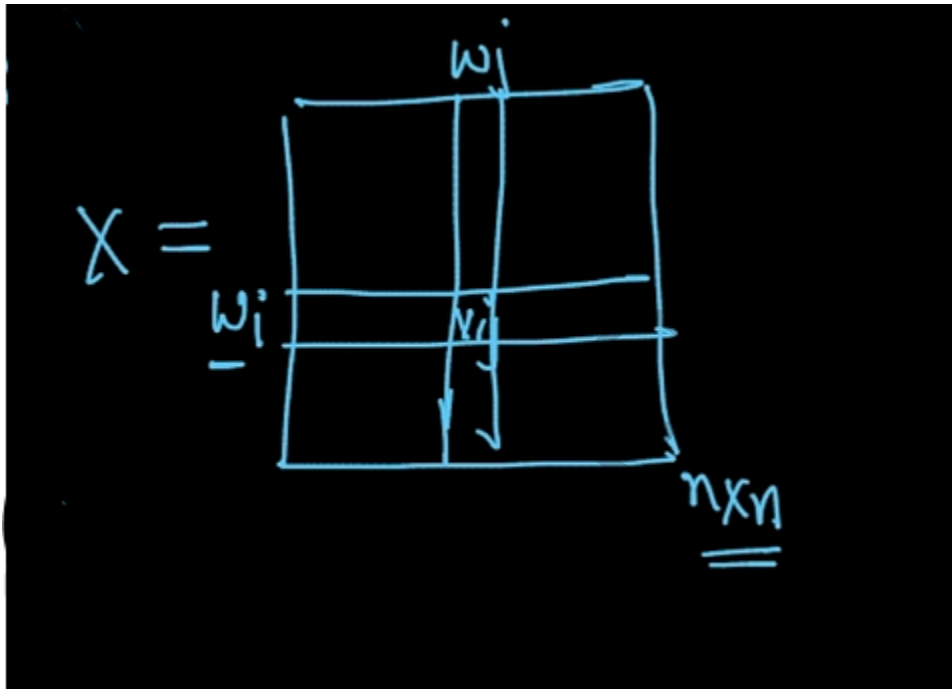
So till we get sufficient data to recommend the items, we use content based recommendations. Once after getting enough data, we can go for collaborative recommendations.

53.12 Word Vectors as MF

So far we have obtained Word Vectors from models like Word2Vec, GLoVe, PLSA, LDA, etc. Here we will use matrix factorization techniques like SVD to compute the word vectors.

1) Co-occurrence Matrix (In Text Data)

Let us denote the co-occurrence matrix by 'X'.



X_{ij} = number of times the word ' w_j ' occurs in the context of ' w_i '. Let us assume we have a review ' r_{10} ', and it consists of the words

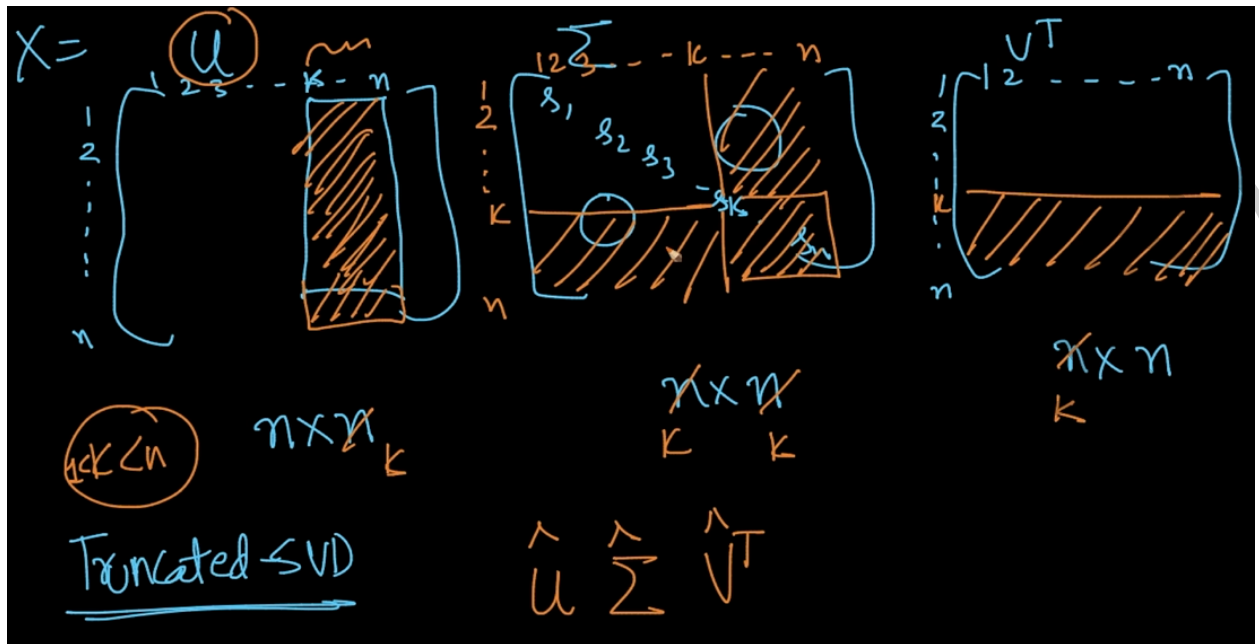
$$r_{10} = w_6 w_j w_1 w_2 w_i w_1 w_3 w_4$$

We can say that the word ' w_j ' occurs in the context of ' w_i ', if ' w_i ' and ' w_j ' are within a neighborhood.

Let us assume the neighborhood distance of 10, then if ' w_i ' and ' w_j ' are separated by less than or equal to 10 words, then we can say that both these words are in the neighborhood of each other.

2) Decompose 'X' as $U\Sigma V^T$

$$X_{n \times n} = U_{n \times n} \Sigma_{n \times n} V_{n \times n}^T$$



Here now we are considering only the first 'K' columns from the matrix 'U', and are discarding the remaining 'n-K' columns. Let us denote it as ' U^\wedge ' and its shape is $n \times k$. Here the matrix ' U^\wedge ' denotes the top left 'K' singular vectors.

We are considering only the first 'K' rows and first 'K' columns from the matrix ' Σ ' and discard the remaining 'n-K' rows and columns. Let us denote it as ' Σ^\wedge ' and its shape is $k \times k$. Here the matrix ' Σ^\wedge ' denotes the top 'K' singular values in diagonal order.

Similarly we are considering only the first 'K' rows from the matrix ' V^T ' and are discarding the remaining 'n-K' rows. Let us denote it as ' V^{T^\wedge} ' and its shape is $k \times n$. Here the matrix ' V^{T^\wedge} ' denotes the top right 'K' singular vectors.

In PCA, we take the top 'K' eigenvalues and the top 'K' eigenvectors. In TruncatedSVD, we take the top 'K' singular values, top 'K' left singular vectors, and the top 'K' right singular vectors.

So in TruncatedSVD(K) \rightarrow 'K' is a hyperparameter

$$X_{n \times n} \rightarrow U_{n \times k}^\wedge \Sigma_{k \times k}^\wedge V_{k \times n}^{T^\wedge}$$

- 3) Take the rows of matrix U^{\wedge} as $u_1, u_2, u_3, \dots, u_n$. Each row ' u_i ' is the vector of the word ' w_i '.

The problem with the combined usage of Cooccurrence matrix + TruncatedSVD is that, if the number of words is more in co-occurrence matrix is large, then applying TruncatedSVD becomes highly expensive.

One hack to handle this kind of problem is that instead of using all the words, it is recommended to use the top ' m ' words. The top ' m ' words are picked on the basis of the IDF scores in TF-IDF.

53.13 Eigen-Faces

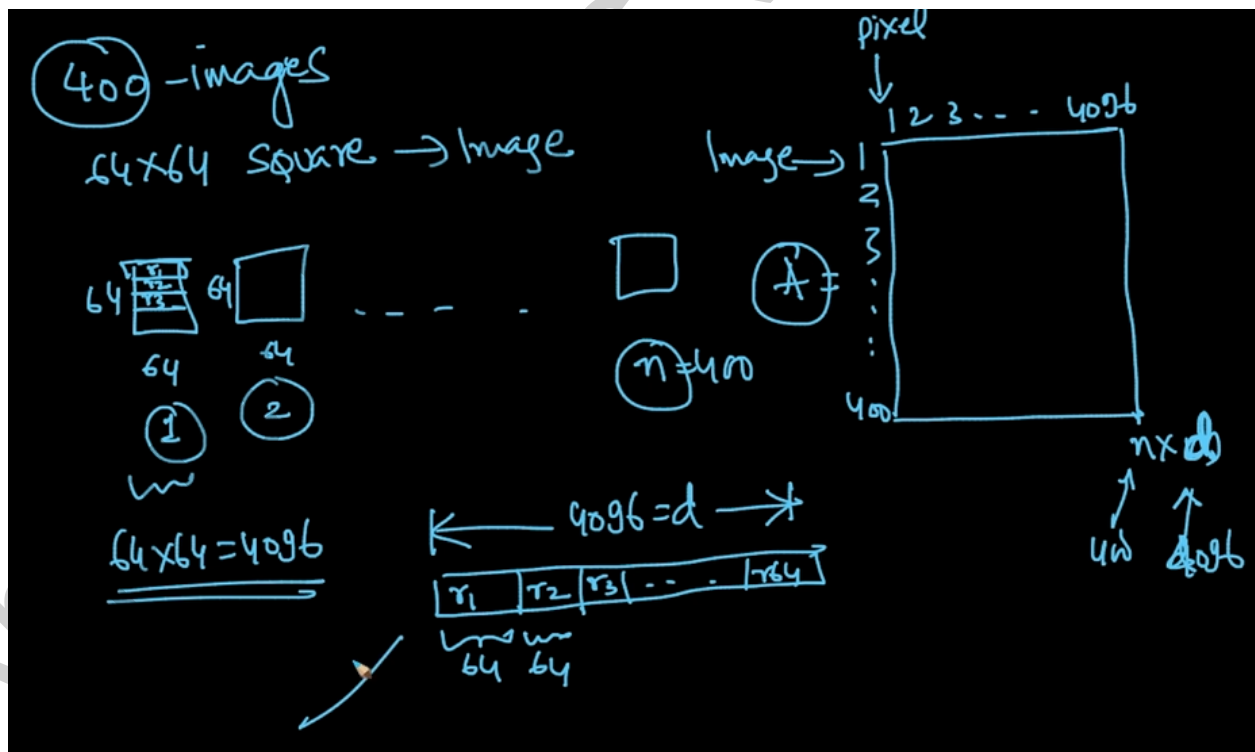
So far, for Word Vector (using co-occurrence matrix), we had used matrix factorization (using SVD). Similarly for the image data, we have a technique called **EigenFaces** (using PCA). This is a very old technique, and today we use the **Convolutional Neural Networks (CNN)**.

Let us assume we have 400 images, and each image is of the shape 64×64 squares. Then for each image, the number of pixels = $64 \times 64 = 4096$. As we get the information for every image in a format of 4096 pixels, we have to flatten this pixel matrix into an array of 4096 dimensions. Now this pixel data is represented as a 4096 dimensional vector for every image.

As each row has 64 pixels, and there are 64 rows, then these rows should be concatenated, and we get a 4096 dimensional vector.

Step 1

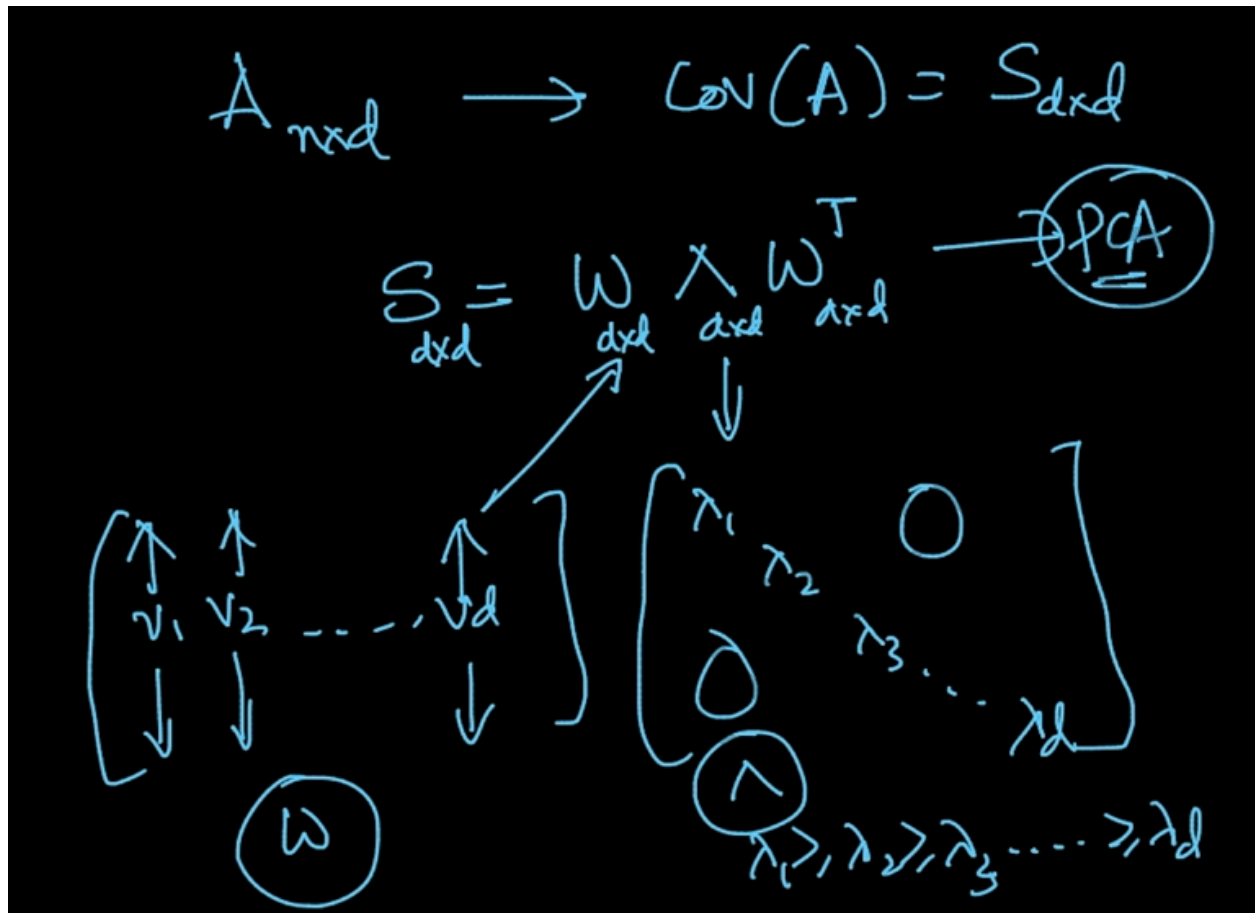
So now we have to construct a data matrix for all these 400 images. Let us denote it as $\mathbf{A}_{n \times d}$. (here $n = 400$, $d = 4096$)



For the data matrix ' $\mathbf{A}_{n \times d}$ ', we have to compute the covariance matrix ' $\mathbf{S}_{d \times d}$ '.

Step 2

$$S_{d \times d} = W_{d \times d} \Lambda_{d \times d} W_{d \times d}^T$$



Step 3

Let us now reduce the dimensionality from 'd' to 'k' to pick the top 'k' eigenvalues and eigenvectors. (here $k < d$).

$$\begin{array}{c}
 d\text{-dim} \rightarrow k\text{-dim} \quad \underline{k < d} \\
 \uparrow \text{Dimensionality Reduction} \\
 \left\{ \begin{array}{c} \text{top } k\text{-eigen val} \\ \downarrow \\ \text{eig-vec} \end{array} \right\} \\
 \underline{W_k} = \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ v_1 & v_2 & \dots & v_k \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix}_{d \times k} \quad v_i \rightarrow v_i \in \underline{\mathbb{R}^d}
 \end{array}$$

Step 4

Now in order to get our images in k-dimensional space, we have to multiply ' $A_{n \times d}$ ' with ' $W_{d \times k}^k$ '. Let us denote the product as ' $A_{n \times k}^{\sim}$ '

$$A_{n \times d} W_{d \times k}^k = A_{n \times k}^{\sim}$$

This ' $A_{n \times k}^{\sim}$ ' denotes the k-dimensional representation of the images.

Note

As the video lecture 53.14 is based on the code sample discussion, we are not providing any notes for it. We suggest you to go through the video lecture and the reference link for it is given below. For any queries, you can post them on the comments wall (or) you can mail us at mentors.diploma@appliedroots.com

Reference Links for Video Lectures 53.13 and 53.14

http://scikit-learn.org/stable/auto_examples/decomposition/plot_faces_decomposition.html#sphx-glr-auto-examples-decomposition-plot-faces-decomposition-py

<https://github.com/sumitjha4321/Notebooks/blob/master/Notebooks/EigenFaces.ipynb>

<http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.NMF.html>

<http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html>

http://scikit-learn.org/stable/auto_examples/decomposition/plot_faces_decomposition.html#sphx-glr-auto-examples-decomposition-plot-faces-decomposition-py

AppliedRoots (Draft Copy)

AppliedRoots (Draft Copy)