

Matrix Factorization and Matrix Completion

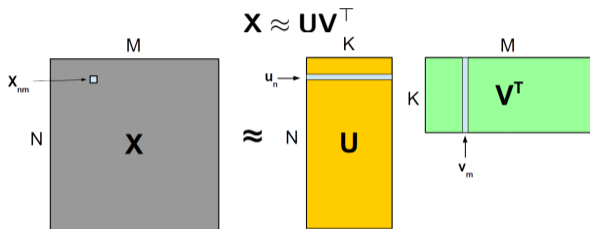
Piyush Rai

Machine Learning (CS771A)

Sept 21, 2016

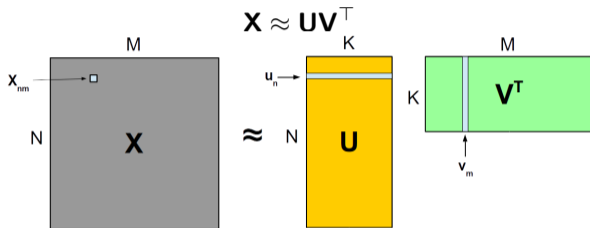
Matrix Factorization

- Given a matrix \mathbf{X} of size $N \times M$, approximate it as a product of two matrices



Matrix Factorization

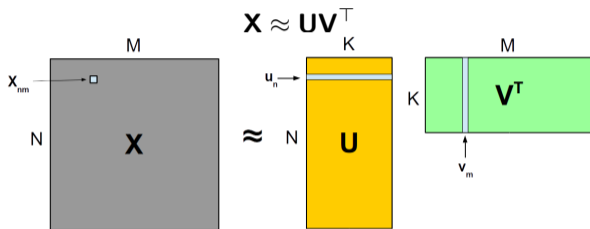
- Given a matrix \mathbf{X} of size $N \times M$, approximate it as a product of two matrices



- \mathbf{U} : $N \times K$ latent factor matrix

Matrix Factorization

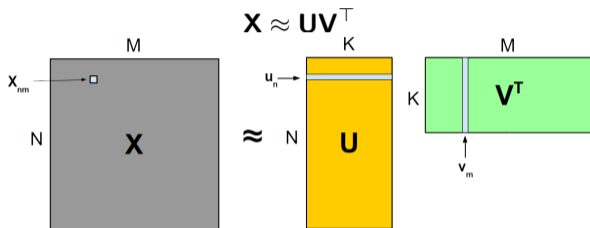
- Given a matrix \mathbf{X} of size $N \times M$, approximate it as a product of two matrices



- \mathbf{U} : $N \times K$ latent factor matrix
 - Each row of \mathbf{U} represents a K -dim latent factor u_n

Matrix Factorization

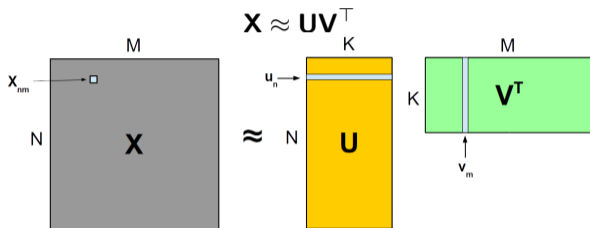
- Given a matrix \mathbf{X} of size $N \times M$, approximate it as a product of two matrices



- \mathbf{U} : $N \times K$ latent factor matrix
 - Each row of \mathbf{U} represents a K -dim latent factor u_n
- \mathbf{V} : $M \times K$ latent factor matrix

Matrix Factorization

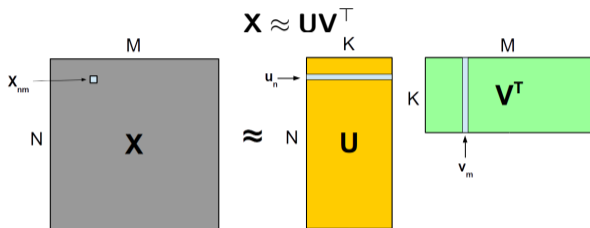
- Given a matrix \mathbf{X} of size $N \times M$, approximate it as a product of two matrices



- \mathbf{U} : $N \times K$ latent factor matrix
 - Each row of \mathbf{U} represents a K -dim latent factor u_n
- \mathbf{V} : $M \times K$ latent factor matrix
 - Each row of \mathbf{V} represents a K -dim latent factor v_n

Matrix Factorization

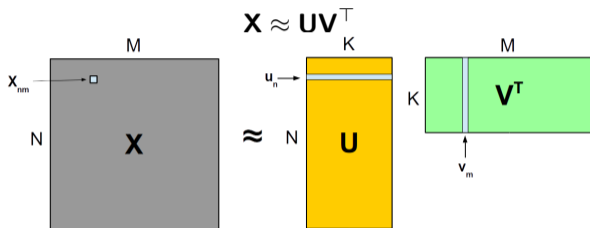
- Given a matrix \mathbf{X} of size $N \times M$, approximate it as a product of two matrices



- \mathbf{U} : $N \times K$ latent factor matrix
 - Each row of \mathbf{U} represents a K -dim latent factor u_n
- \mathbf{V} : $M \times K$ latent factor matrix
 - Each row of \mathbf{V} represents a K -dim latent factor v_n
- Each entry of \mathbf{X} can be written as: $X_{nm} \approx \mathbf{u}_n^T \mathbf{v}_m = \sum_{k=1}^K u_{nk} v_{mk}$

Matrix Factorization

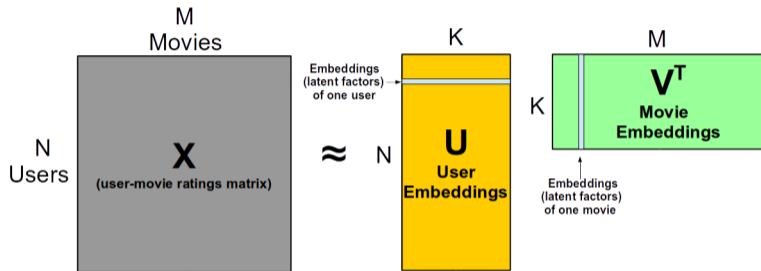
- Given a matrix \mathbf{X} of size $N \times M$, approximate it as a product of two matrices



- \mathbf{U} : $N \times K$ latent factor matrix
 - Each row of \mathbf{U} represents a K -dim latent factor u_n
- \mathbf{V} : $M \times K$ latent factor matrix
 - Each row of \mathbf{V} represents a K -dim latent factor v_n
- Each entry of \mathbf{X} can be written as: $X_{nm} \approx \mathbf{u}_n^T \mathbf{v}_m = \sum_{k=1}^K u_{nk} v_{mk}$
- If X_{nm} is large (small) then u_n and v_m should be similar (dissimilar)

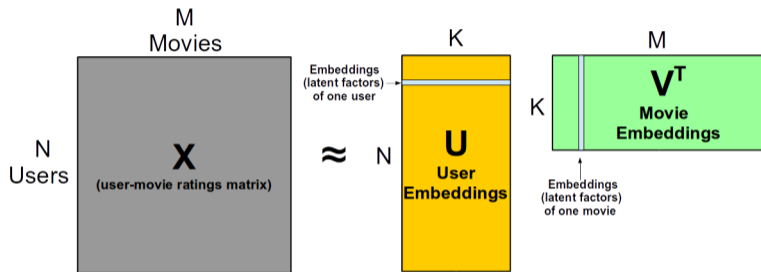
Why Matrix Factorization?

- The latent factors can be used/interpreted as “embeddings” or “features”



Why Matrix Factorization?

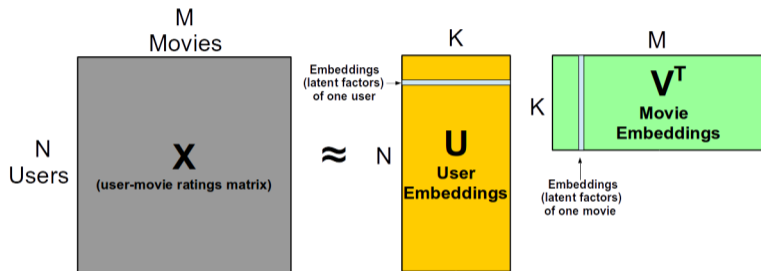
- The latent factors can be used/interpreted as “embeddings” or “features”



- Especially useful for learning good features for “dyadic” or relational data
 - Examples: Users-Movies ratings, Users-Products purchases, etc.

Why Matrix Factorization?

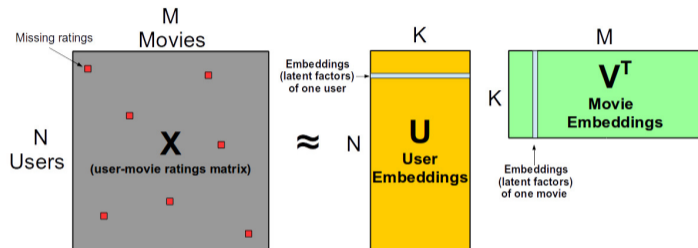
- The latent factors can be used/interpreted as “embeddings” or “features”



- Especially useful for learning good features for “dyadic” or relational data
 - Examples: Users-Movies ratings, Users-Products purchases, etc.
- If $K \ll \min\{M, N\} \Rightarrow$ then can also be seen as dimensionality reduction or a “low-rank factorization” of the matrix \mathbf{X}

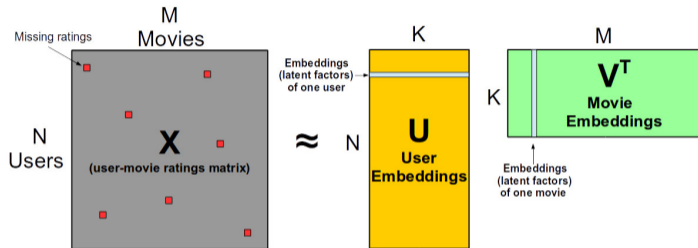
Why Matrix Factorization?

- Can also predict the missing/unknown entries in the original matrix



Why Matrix Factorization?

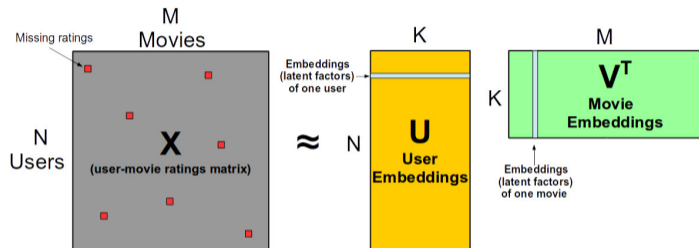
- Can also predict the missing/unknown entries in the original matrix



- Note: The latent factor matrices U and V can be learned even when the matrix X is only **partially observed** (as we will see shortly)

Why Matrix Factorization?

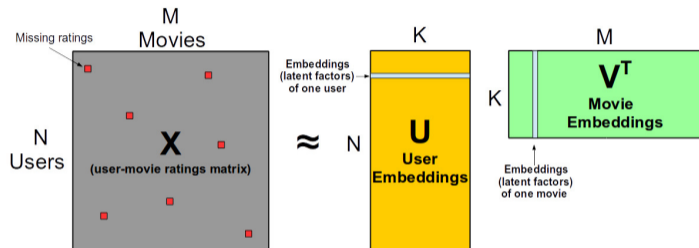
- Can also predict the missing/unknown entries in the original matrix



- Note: The latent factor matrices \mathbf{U} and \mathbf{V} can be learned even when the matrix \mathbf{X} is only **partially observed** (as we will see shortly)
- After learning \mathbf{U} and \mathbf{V} , any missing X_{nm} can be approximated by $\mathbf{u}_n^T \mathbf{v}_m$

Why Matrix Factorization?

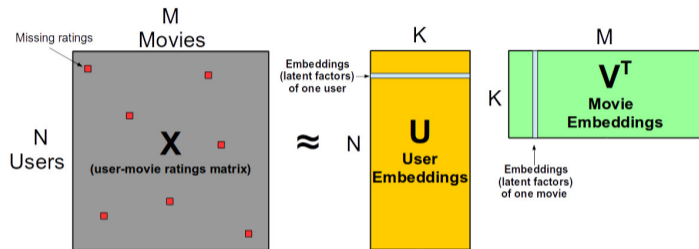
- Can also predict the missing/unknown entries in the original matrix



- Note: The latent factor matrices \mathbf{U} and \mathbf{V} can be learned even when the matrix \mathbf{X} is only **partially observed** (as we will see shortly)
- After learning \mathbf{U} and \mathbf{V} , any missing X_{nm} can be approximated by $\mathbf{u}_n^T \mathbf{v}_m$
- \mathbf{UV}^T is the best low-rank matrix that approximates the full \mathbf{X}

Why Matrix Factorization?

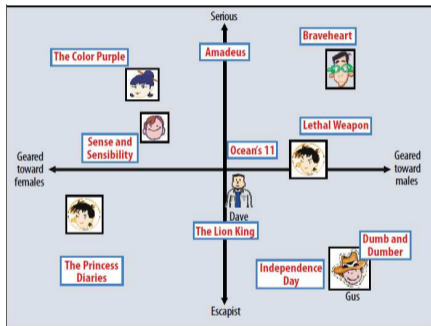
- Can also predict the missing/unknown entries in the original matrix



- Note: The latent factor matrices \mathbf{U} and \mathbf{V} can be learned even when the matrix \mathbf{X} is only **partially observed** (as we will see shortly)
- After learning \mathbf{U} and \mathbf{V} , any missing X_{nm} can be approximated by $\mathbf{u}_n^T \mathbf{v}_m$
- \mathbf{UV}^T is the best low-rank matrix that approximates the full \mathbf{X}
- Note: The “[Netflix Challenge](#)” was won by a matrix factorization method

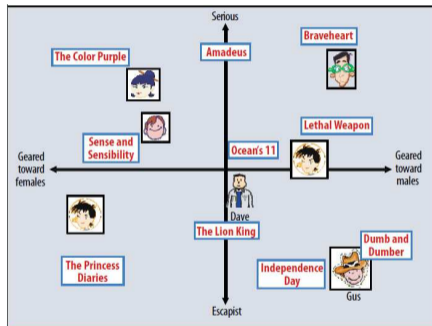
Interpreting the Embeddings/Latent Factors

- Embeddings/latent factors can often be interpreted. E.g., as “genres” if \mathbf{X} represents a user-movie rating matrix. A cartoon with $K = 2$ shown below



Interpreting the Embeddings/Latent Factors

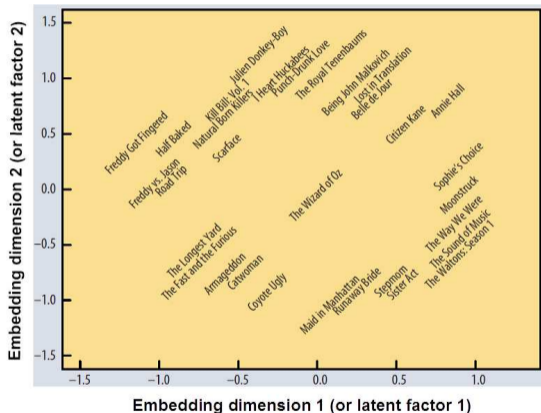
- Embeddings/latent factors can often be interpreted. E.g., as “genres” if \mathbf{X} represents a user-movie rating matrix. A cartoon with $K = 2$ shown below



- Similar things (users/movies) get embedded nearby in the embedding space (two things will be deemed similar if their embeddings are similar). Thus useful for [computing similarities](#) and/or [making recommendations](#)

Interpreting the Embeddings/Latent Factors

- Another illustration of 2-D embeddings of the movies only

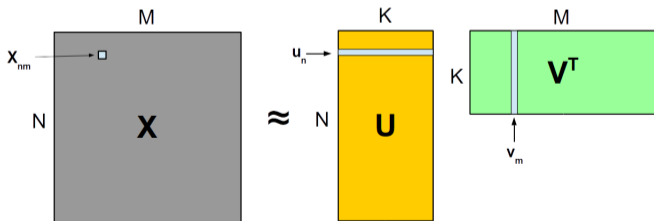


- Similar movies will be embedded at nearby locations in the embedding space

Solving Matrix Factorization

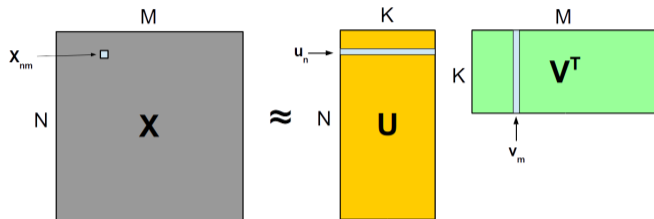
Matrix Factorization

- Recall our matrix factorization model: $\mathbf{X} \approx \mathbf{UV}^T$
- Goal: learn \mathbf{U} and \mathbf{V} , given a subset Ω of entries in \mathbf{X} (let's call it \mathbf{X}_Ω)
- Some notations:
 - $\Omega = \{(n, m)\}$: X_{nm} is observed
 - Ω_{r_n} : column indices of observed entries in row n of \mathbf{X}
 - Ω_{c_m} : row indices of observed entries in column m of \mathbf{X}



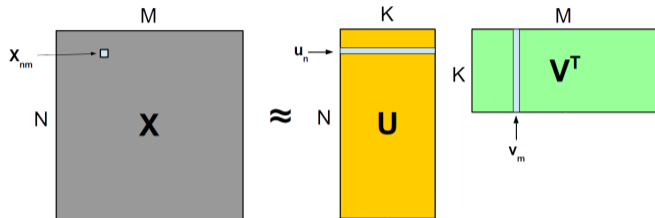
Matrix Factorization

- We want \mathbf{X} to be as close to \mathbf{UV}^T as possible



Matrix Factorization

- We want \mathbf{X} to be as close to \mathbf{UV}^T as possible

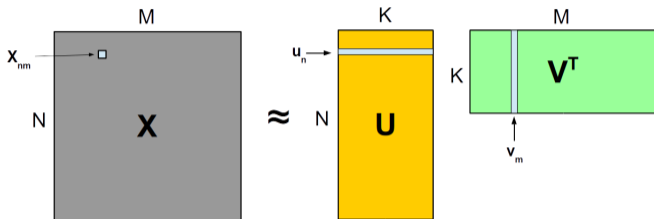


- Let's define a squared "loss function" over the observed entries in \mathbf{X}

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^T \mathbf{v}_m)^2$$

Matrix Factorization

- We want \mathbf{X} to be as close to \mathbf{UV}^T as possible



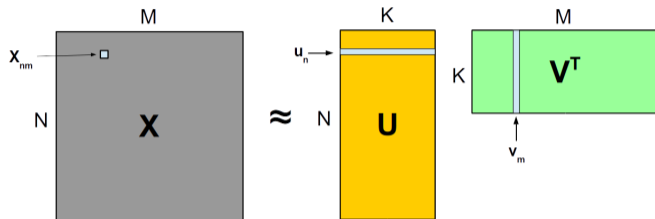
- Let's define a squared “loss function” over the observed entries in \mathbf{X}

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^T \mathbf{v}_m)^2$$

- Here the latent factors $\{\mathbf{u}_n\}_{n=1}^N$ and $\{\mathbf{v}_m\}_{m=1}^M$ are the **unknown parameters**

Matrix Factorization

- We want \mathbf{X} to be as close to \mathbf{UV}^T as possible



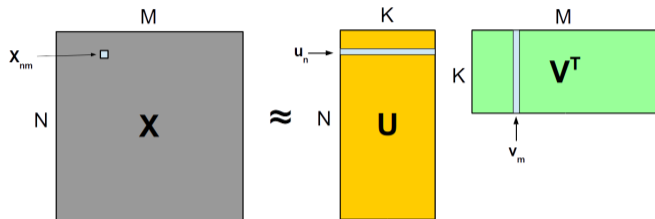
- Let's define a squared "loss function" over the observed entries in \mathbf{X}

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^T \mathbf{v}_m)^2$$

- Here the latent factors $\{\mathbf{u}_n\}_{n=1}^N$ and $\{\mathbf{v}_m\}_{m=1}^M$ are the **unknown parameters**
- Squared loss chosen only for simplicity; other loss functions can be used

Matrix Factorization

- We want \mathbf{X} to be as close to \mathbf{UV}^T as possible



- Let's define a squared "loss function" over the observed entries in \mathbf{X}

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^T \mathbf{v}_m)^2$$

- Here the latent factors $\{\mathbf{u}_n\}_{n=1}^N$ and $\{\mathbf{v}_m\}_{m=1}^M$ are the **unknown parameters**
- Squared loss chosen only for simplicity; other loss functions can be used
- How do we learn $\{\mathbf{u}_n\}_{n=1}^N$ and $\{\mathbf{v}_m\}_{m=1}^M$?

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

- A **non-convex** problem. Difficult to optimize w.r.t. \mathbf{u}_n and \mathbf{v}_m jointly.

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

- A **non-convex** problem. Difficult to optimize w.r.t. \mathbf{u}_n and \mathbf{v}_m jointly.
- One way is to solve for \mathbf{u}_n and \mathbf{v}_m in an **alternating fashion**, e.g.,

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

- A **non-convex** problem. Difficult to optimize w.r.t. \mathbf{u}_n and \mathbf{v}_m jointly.
- One way is to solve for \mathbf{u}_n and \mathbf{v}_m in an **alternating fashion**, e.g.,
 - $\forall n$, fix all variables except \mathbf{u}_n and solve the optim. problem w.r.t. \mathbf{u}_n

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

- A **non-convex** problem. Difficult to optimize w.r.t. \mathbf{u}_n and \mathbf{v}_m jointly.
- One way is to solve for \mathbf{u}_n and \mathbf{v}_m in an **alternating fashion**, e.g.,

- $\forall n$, fix all variables except \mathbf{u}_n and solve the optim. problem w.r.t. \mathbf{u}_n

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

- $\forall m$, fix all variables except \mathbf{v}_m and solve the optim. problem w.r.t. \mathbf{v}_m

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

- A **non-convex** problem. Difficult to optimize w.r.t. \mathbf{u}_n and \mathbf{v}_m jointly.
- One way is to solve for \mathbf{u}_n and \mathbf{v}_m in an **alternating fashion**, e.g.,

- $\forall n$, fix all variables except \mathbf{u}_n and solve the optim. problem w.r.t. \mathbf{u}_n

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

- $\forall m$, fix all variables except \mathbf{v}_m and solve the optim. problem w.r.t. \mathbf{v}_m

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

- Iterate until not converged

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

- A **non-convex** problem. Difficult to optimize w.r.t. \mathbf{u}_n and \mathbf{v}_m jointly.
- One way is to solve for \mathbf{u}_n and \mathbf{v}_m in an **alternating fashion**, e.g.,

- $\forall n$, fix all variables except \mathbf{u}_n and solve the optim. problem w.r.t. \mathbf{u}_n

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

- $\forall m$, fix all variables except \mathbf{v}_m and solve the optim. problem w.r.t. \mathbf{v}_m

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

- Iterate until not converged
- Each of these subproblems has a simple, convex objective function

Alternating Optimization

- We will use an ℓ_2 regularized version of the squared loss function

$$\mathcal{L} = \sum_{(n,m) \in \Omega} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \sum_{n=1}^N \lambda_U \|\mathbf{u}_n\|^2 + \sum_{m=1}^M \lambda_V \|\mathbf{v}_m\|^2$$

- A **non-convex** problem. Difficult to optimize w.r.t. \mathbf{u}_n and \mathbf{v}_m jointly.
- One way is to solve for \mathbf{u}_n and \mathbf{v}_m in an **alternating fashion**, e.g.,

- $\forall n$, fix all variables except \mathbf{u}_n and solve the optim. problem w.r.t. \mathbf{u}_n

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

- $\forall m$, fix all variables except \mathbf{v}_m and solve the optim. problem w.r.t. \mathbf{v}_m

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

- Iterate until not converged
- Each of these subproblems has a simple, convex objective function
- Convergence properties of such methods have been studied extensively

The Solutions

- Easy to show that the problem

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^T \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

The Solutions

- Easy to show that the problem

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

.. has the solution

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

The Solutions

- Easy to show that the problem

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

.. has the solution

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

- Likewise, the problem

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

The Solutions

- Easy to show that the problem

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

.. has the solution

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

- Likewise, the problem

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

.. has the solution

$$\mathbf{v}_m = \left(\sum_{n \in \Omega_{c_m}} \mathbf{u}_n \mathbf{u}_n^\top + \lambda_V \mathbf{I}_K \right)^{-1} \left(\sum_{n \in \Omega_{c_m}} X_{nm} \mathbf{u}_n \right)$$

The Solutions

- Easy to show that the problem

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

.. has the solution

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

- Likewise, the problem

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

.. has the solution

$$\mathbf{v}_m = \left(\sum_{n \in \Omega_{c_m}} \mathbf{u}_n \mathbf{u}_n^\top + \lambda_V \mathbf{I}_K \right)^{-1} \left(\sum_{n \in \Omega_{c_m}} X_{nm} \mathbf{u}_n \right)$$

- Note that this is very similar to (regularized) least squares regression

The Solutions

- Easy to show that the problem

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$$

.. has the solution

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

- Likewise, the problem

$$\arg \min_{\mathbf{v}_m} \sum_{n \in \Omega_{c_m}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_V \|\mathbf{v}_m\|^2$$

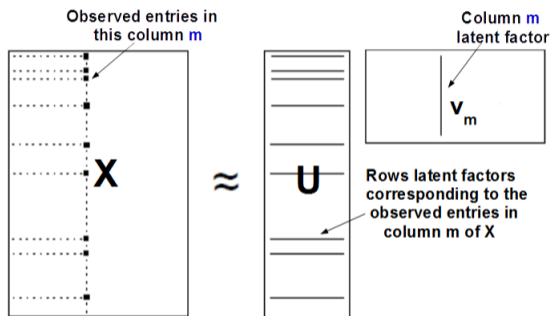
.. has the solution

$$\mathbf{v}_m = \left(\sum_{n \in \Omega_{c_m}} \mathbf{u}_n \mathbf{u}_n^\top + \lambda_V \mathbf{I}_K \right)^{-1} \left(\sum_{n \in \Omega_{c_m}} X_{nm} \mathbf{u}_n \right)$$

- Note that this is very similar to (regularized) least squares regression
- Thus matrix factorization can be also seen as [a sequence of regression problems](#) (one for each latent factor)

Matrix Factorization as Regression

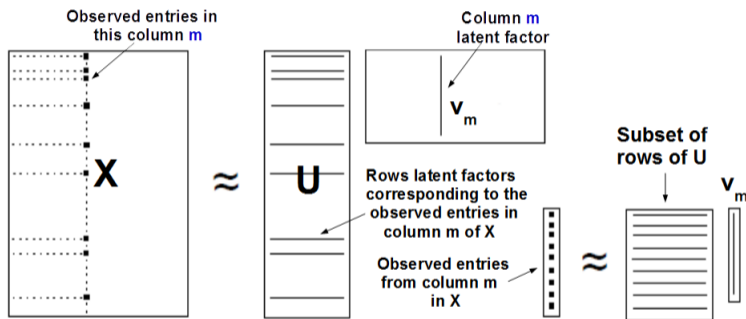
Suppose we are solving for \mathbf{v}_m (with \mathbf{U} and all other \mathbf{v}_m 's fixed)



..1

Matrix Factorization as Regression

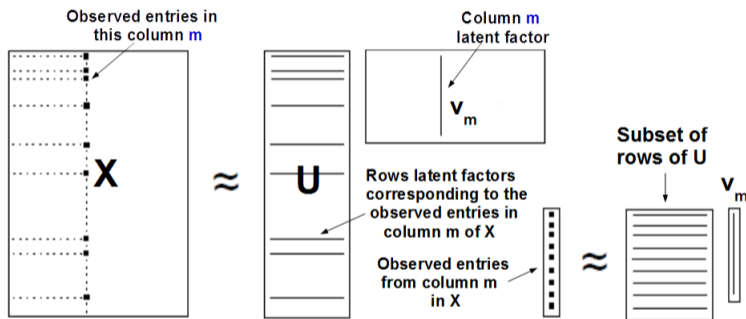
Suppose we are solving for \mathbf{v}_m (with \mathbf{U} and all other \mathbf{v}_m 's fixed)



Now becomes a least-squares type problem for solving for \mathbf{v}_m

Matrix Factorization as Regression

Suppose we are solving for \mathbf{v}_m (with \mathbf{U} and all other \mathbf{v}_m 's fixed)



Now becomes a least-squares type problem for solving for \mathbf{v}_m

Can think of solving for \mathbf{u}_n (with \mathbf{V} and all other \mathbf{u}_n 's fixed) in the same way

Matrix Factorization as Regression

- A very useful way to understand matrix factorization

Matrix Factorization as Regression

- A very useful way to understand matrix factorization
- Can modify the regularized least-squares like objective

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \mathbf{u}_n^\top \mathbf{u}_n$$

.. using other **loss functions** and **regularizers**

Matrix Factorization as Regression

- A very useful way to understand matrix factorization
- Can modify the regularized least-squares like objective

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \mathbf{u}_n^\top \mathbf{u}_n$$

.. using other **loss functions** and **regularizers**

- Some possible modifications:

Matrix Factorization as Regression

- A very useful way to understand matrix factorization
- Can modify the regularized least-squares like objective

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^T \mathbf{v}_m)^2 + \lambda_U \mathbf{u}_n^T \mathbf{u}_n$$

.. using other **loss functions** and **regularizers**

- Some possible modifications:
 - If entries in the matrix \mathbf{X} are binary, counts, etc. then we can replace the squared loss function by some other loss function (e.g., logistic or Poisson)

Matrix Factorization as Regression

- A very useful way to understand matrix factorization
- Can modify the regularized least-squares like objective

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \mathbf{u}_n^\top \mathbf{u}_n$$

.. using other **loss functions** and **regularizers**

- Some possible modifications:
 - If entries in the matrix \mathbf{X} are binary, counts, etc. then we can replace the squared loss function by some other loss function (e.g., logistic or Poisson)
 - Can impose other constraints on the latent factors, e.g., non-negativity, sparsity, etc. (by changing the regularizer)

Matrix Factorization as Regression

- A very useful way to understand matrix factorization
- Can modify the regularized least-squares like objective

$$\arg \min_{\mathbf{u}_n} \sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \mathbf{u}_n^\top \mathbf{u}_n$$

.. using other **loss functions** and **regularizers**

- Some possible modifications:
 - If entries in the matrix \mathbf{X} are binary, counts, etc. then we can replace the squared loss function by some other loss function (e.g., logistic or Poisson)
 - Can impose other constraints on the latent factors, e.g., non-negativity, sparsity, etc. (by changing the regularizer)
 - Can think of this also as a probabilistic model (a likelihood function on X_{nm} and priors on the latent factors $\mathbf{u}_n, \mathbf{v}_m$) and do MLE/MAP

Matrix Factorization: The Complete Algorithm

- Input: Partially complete matrix \mathbf{X}_Ω

Matrix Factorization: The Complete Algorithm

- Input: Partially complete matrix \mathbf{X}_Ω
- Initialize the latent factors $\mathbf{v}_1, \dots, \mathbf{v}_M$ randomly

Matrix Factorization: The Complete Algorithm

- Input: Partially complete matrix \mathbf{X}_Ω
- Initialize the latent factors $\mathbf{v}_1, \dots, \mathbf{v}_M$ randomly
- Iterate until not converged

Matrix Factorization: The Complete Algorithm

- Input: Partially complete matrix \mathbf{X}_Ω
- Initialize the latent factors $\mathbf{v}_1, \dots, \mathbf{v}_M$ randomly
- Iterate until not converged
 - Update each row latent factor \mathbf{u}_n , $n = 1, \dots, N$ (can be in parallel)

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

Matrix Factorization: The Complete Algorithm

- Input: Partially complete matrix \mathbf{X}_Ω
- Initialize the latent factors $\mathbf{v}_1, \dots, \mathbf{v}_M$ randomly
- Iterate until not converged
 - Update each row latent factor \mathbf{u}_n , $n = 1, \dots, N$ (can be in parallel)

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

- Update each column latent factor \mathbf{v}_m , $m = 1, \dots, M$ (can be in parallel)

$$\mathbf{v}_m = \left(\sum_{n \in \Omega_{c_m}} \mathbf{u}_n \mathbf{u}_n^\top + \lambda_V \mathbf{I}_K \right)^{-1} \left(\sum_{n \in \Omega_{c_m}} X_{nm} \mathbf{u}_n \right)$$

Matrix Factorization: The Complete Algorithm

- Input: Partially complete matrix \mathbf{X}_Ω
- Initialize the latent factors $\mathbf{v}_1, \dots, \mathbf{v}_M$ randomly
- Iterate until not converged
 - Update each row latent factor \mathbf{u}_n , $n = 1, \dots, N$ (can be in parallel)

$$\mathbf{u}_n = \left(\sum_{m \in \Omega_{r_n}} \mathbf{v}_m \mathbf{v}_m^\top + \lambda_U \mathbf{I}_K \right)^{-1} \left(\sum_{m \in \Omega_{r_n}} X_{nm} \mathbf{v}_m \right)$$

- Update each column latent factor \mathbf{v}_m , $m = 1, \dots, M$ (can be in parallel)

$$\mathbf{v}_m = \left(\sum_{n \in \Omega_{c_m}} \mathbf{u}_n \mathbf{u}_n^\top + \lambda_V \mathbf{I}_K \right)^{-1} \left(\sum_{n \in \Omega_{c_m}} X_{nm} \mathbf{u}_n \right)$$

- Final prediction for any (missing) entry: $X_{nm} = \mathbf{u}_n^\top \mathbf{v}_m$

A Faster Algorithm via SGD

- Alternating optimization is nice but can be slow (note that it requires matrix inversion with cost $O(K^3)$ for updating each latent factor $\mathbf{u}_n, \mathbf{v}_m$)

A Faster Algorithm via SGD

- Alternating optimization is nice but can be slow (note that it requires matrix inversion with cost $O(K^3)$ for updating each latent factor $\mathbf{u}_n, \mathbf{v}_m$)
- An alternative is to use stochastic gradient descent (SGD). In each round, select a randomly chosen entry X_{nm} with $(n, m) \in \Omega$

A Faster Algorithm via SGD

- Alternating optimization is nice but can be slow (note that it requires matrix inversion with cost $O(K^3)$ for updating each latent factor $\mathbf{u}_n, \mathbf{v}_m$)
- An alternative is to use stochastic gradient descent (SGD). In each round, select a randomly chosen entry X_{nm} with $(n, m) \in \Omega$
- Consider updating \mathbf{u}_n . For loss function $\sum_{m \in \Omega_n} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$, the **stochastic gradient** w.r.t. \mathbf{u}_n using this randomly chosen entry X_{nm} is

$$-(X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{v}_m + \lambda_U \mathbf{u}_n$$

A Faster Algorithm via SGD

- Alternating optimization is nice but can be slow (note that it requires matrix inversion with cost $O(K^3)$ for updating each latent factor $\mathbf{u}_n, \mathbf{v}_m$)
- An alternative is to use stochastic gradient descent (SGD). In each round, select a randomly chosen entry X_{nm} with $(n, m) \in \Omega$
- Consider updating \mathbf{u}_n . For loss function $\sum_{m \in \Omega_n} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$, the **stochastic gradient** w.r.t. \mathbf{u}_n using this randomly chosen entry X_{nm} is

$$-(X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{v}_m + \lambda_U \mathbf{u}_n$$

- Thus the SGD update for \mathbf{u}_n will be

$$\mathbf{u}_n = \mathbf{u}_n - \eta(\lambda_U \mathbf{u}_n - (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{v}_m)$$

A Faster Algorithm via SGD

- Alternating optimization is nice but can be slow (note that it requires matrix inversion with cost $O(K^3)$ for updating each latent factor $\mathbf{u}_n, \mathbf{v}_m$)
- An alternative is to use stochastic gradient descent (SGD). In each round, select a randomly chosen entry X_{nm} with $(n, m) \in \Omega$
- Consider updating \mathbf{u}_n . For loss function $\sum_{m \in \Omega_n} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$, the **stochastic gradient** w.r.t. \mathbf{u}_n using this randomly chosen entry X_{nm} is

$$-(X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{v}_m + \lambda_U \mathbf{u}_n$$

- Thus the SGD update for \mathbf{u}_n will be

$$\mathbf{u}_n = \mathbf{u}_n - \eta(\lambda_U \mathbf{u}_n - (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{v}_m)$$

- Likewise, the SGD update for \mathbf{v}_m will be

$$\mathbf{v}_m = \mathbf{v}_m - \eta(\lambda_V \mathbf{v}_m - (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{u}_n)$$

A Faster Algorithm via SGD

- Alternating optimization is nice but can be slow (note that it requires matrix inversion with cost $O(K^3)$ for updating each latent factor $\mathbf{u}_n, \mathbf{v}_m$)
- An alternative is to use stochastic gradient descent (SGD). In each round, select a randomly chosen entry X_{nm} with $(n, m) \in \Omega$
- Consider updating \mathbf{u}_n . For loss function $\sum_{m \in \Omega_{r_n}} (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m)^2 + \lambda_U \|\mathbf{u}_n\|^2$, the **stochastic gradient** w.r.t. \mathbf{u}_n using this randomly chosen entry X_{nm} is

$$-(X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{v}_m + \lambda_U \mathbf{u}_n$$

- Thus the SGD update for \mathbf{u}_n will be

$$\mathbf{u}_n = \mathbf{u}_n - \eta(\lambda_U \mathbf{u}_n - (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{v}_m)$$

- Likewise, the SGD update for \mathbf{v}_m will be

$$\mathbf{v}_m = \mathbf{v}_m - \eta(\lambda_V \mathbf{v}_m - (X_{nm} - \mathbf{u}_n^\top \mathbf{v}_m) \mathbf{u}_n)$$

- The SGD algorithm chooses a random entry X_{nm} in each iteration, updates $\mathbf{u}_n, \mathbf{v}_m$, and repeats until convergence (\mathbf{u}_n 's, \mathbf{v}_m 's randomly initialized).

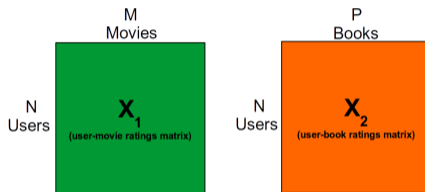
Some Other Extensions of Matrix Factorization

Joint Matrix Factorization

- Can do joint matrix factorization of more than one matrices

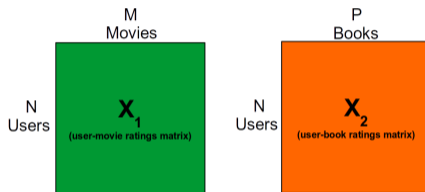
Joint Matrix Factorization

- Can do joint matrix factorization of more than one matrices
- Consider two “ratings” matrices with the N users shared in both



Joint Matrix Factorization

- Can do joint matrix factorization of more than one matrices
- Consider two “ratings” matrices with the N users shared in both

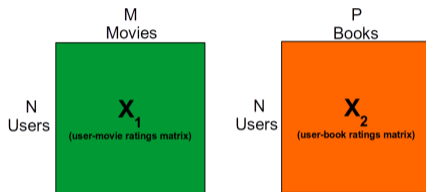


- Can assume the following matrix factorization

$$\mathbf{X}_1 \approx \mathbf{UV}_1^T \quad \text{and} \quad \mathbf{X}_2 \approx \mathbf{UV}_2^T$$

Joint Matrix Factorization

- Can do joint matrix factorization of more than one matrices
- Consider two “ratings” matrices with the N users shared in both



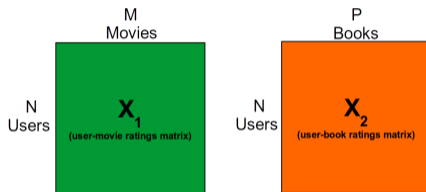
- Can assume the following matrix factorization

$$\mathbf{X}_1 \approx \mathbf{U}\mathbf{V}_1^T \quad \text{and} \quad \mathbf{X}_2 \approx \mathbf{U}\mathbf{V}_2^T$$

- Note that the user latent factor matrix \mathbf{U} is shared in both factorizations

Joint Matrix Factorization

- Can do joint matrix factorization of more than one matrices
- Consider two “ratings” matrices with the N users shared in both



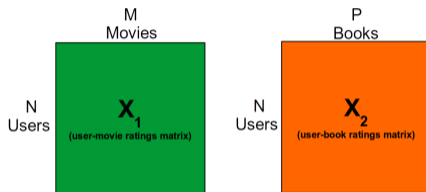
- Can assume the following matrix factorization

$$\mathbf{X}_1 \approx \mathbf{U}\mathbf{V}_1^T \quad \text{and} \quad \mathbf{X}_2 \approx \mathbf{U}\mathbf{V}_2^T$$

- Note that the user latent factor matrix \mathbf{U} is shared in both factorizations
- Gives a way to [learn features by combining multiple data sets](#) (2 in this case)

Joint Matrix Factorization

- Can do joint matrix factorization of more than one matrices
- Consider two “ratings” matrices with the N users shared in both



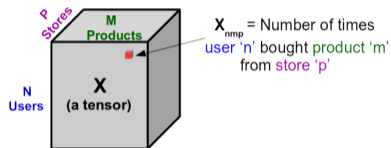
- Can assume the following matrix factorization

$$X_1 \approx UV_1^T \quad \text{and} \quad X_2 \approx UV_2^T$$

- Note that the user latent factor matrix U is shared in both factorizations
- Gives a way to [learn features by combining multiple data sets](#) (2 in this case)
- Can use the alternating optimization to solve for U , V_1 and V_2

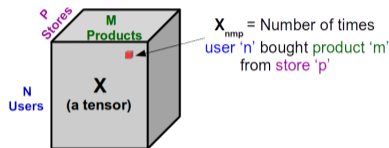
Tensor Factorization

- A “tensor” is a generalization of a matrix to more than two dimensions
- Consider a 3-dim (or 3-mode or 3-way) tensor \mathbf{X} of size $N \times M \times P$



Tensor Factorization

- A “tensor” is a generalization of a matrix to more than two dimensions
- Consider a 3-dim (or 3-mode or 3-way) tensor \mathbf{X} of size $N \times M \times P$

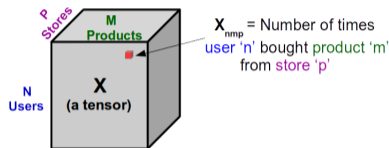


- We can model each entry of tensor \mathbf{X} as

$$X_{nmp} \approx \mathbf{u}_n \odot \mathbf{v}_m \odot \mathbf{w}_p = \sum_{k=1}^K u_{nk} v_{mk} w_{pk}$$

Tensor Factorization

- A “tensor” is a generalization of a matrix to more than two dimensions
- Consider a 3-dim (or 3-mode or 3-way) tensor \mathbf{X} of size $N \times M \times P$



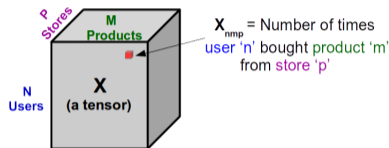
- We can model each entry of tensor \mathbf{X} as

$$X_{nmp} \approx \mathbf{u}_n \odot \mathbf{v}_m \odot \mathbf{w}_p = \sum_{k=1}^K u_{nk} v_{mk} w_{pk}$$

- Can learn $\{\mathbf{u}_n\}_{n=1}^N, \{\mathbf{v}_m\}_{m=1}^M, \{\mathbf{w}_p\}_{p=1}^P$ using alternating optimization

Tensor Factorization

- A “tensor” is a generalization of a matrix to more than two dimensions
- Consider a 3-dim (or 3-mode or 3-way) tensor \mathbf{X} of size $N \times M \times P$



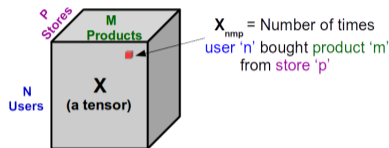
- We can model each entry of tensor \mathbf{X} as

$$X_{nmp} \approx \mathbf{u}_n \odot \mathbf{v}_m \odot \mathbf{w}_p = \sum_{k=1}^K u_{nk} v_{mk} w_{pk}$$

- Can learn $\{\mathbf{u}_n\}_{n=1}^N, \{\mathbf{v}_m\}_{m=1}^M, \{\mathbf{w}_p\}_{p=1}^P$ using alternating optimization
- These K -dim. “embeddings” can be used as features for other tasks (e.g., tensor completion, computing similarities, etc.)

Tensor Factorization

- A “tensor” is a generalization of a matrix to more than two dimensions
- Consider a 3-dim (or 3-mode or 3-way) tensor \mathbf{X} of size $N \times M \times P$



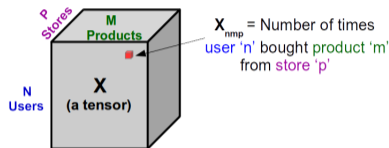
- We can model each entry of tensor \mathbf{X} as

$$X_{nmp} \approx \mathbf{u}_n \odot \mathbf{v}_m \odot \mathbf{w}_p = \sum_{k=1}^K u_{nk} v_{mk} w_{pk}$$

- Can learn $\{\mathbf{u}_n\}_{n=1}^N, \{\mathbf{v}_m\}_{m=1}^M, \{\mathbf{w}_p\}_{p=1}^P$ using alternating optimization
- These K -dim. “embeddings” can be used as features for other tasks (e.g., tensor completion, computing similarities, etc.)
- The model also be easily extended to tensors having than 3 dimensions

Tensor Factorization

- A “tensor” is a generalization of a matrix to more than two dimensions
- Consider a 3-dim (or 3-mode or 3-way) tensor \mathbf{X} of size $N \times M \times P$



- We can model each entry of tensor \mathbf{X} as

$$X_{nmp} \approx \mathbf{u}_n \odot \mathbf{v}_m \odot \mathbf{w}_p = \sum_{k=1}^K u_{nk} v_{mk} w_{pk}$$

- Can learn $\{\mathbf{u}_n\}_{n=1}^N, \{\mathbf{v}_m\}_{m=1}^M, \{\mathbf{w}_p\}_{p=1}^P$ using alternating optimization
- These K -dim. “embeddings” can be used as features for other tasks (e.g., tensor completion, computing similarities, etc.)
- The model also be easily extended to tensors having than 3 dimensions
- Several specialized algorithms for tensor factorization (CP/Tucker decomposition, etc.)

Tensor Factorization: An Application

- Many data sets can be naturally represented as tensors

Tensor Factorization: An Application

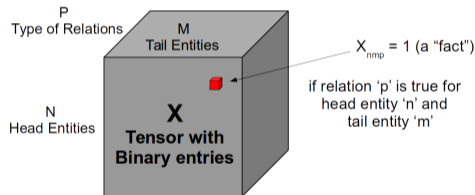
- Many data sets can be naturally represented as tensors
- Knowledge-Graphs (KG) of Knowledge-Bases (KB) is one such example

Tensor Factorization: An Application

- Many data sets can be naturally represented as tensors
- Knowledge-Graphs (KG) of Knowledge-Bases (KB) is one such example
- A KG/KB consists of “facts” in form of triplets (e.g. Modi-PM-India)

Tensor Factorization: An Application

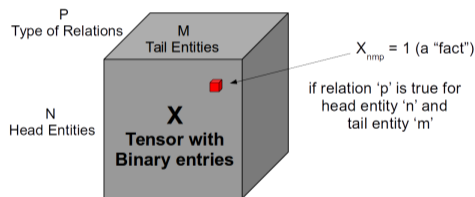
- Many data sets can be naturally represented as tensors
- Knowledge-Graphs (KG) or Knowledge-Bases (KB) is one such example
- A KG/KB consists of “facts” in form of triplets (e.g. Modi-PM-India)



- KGs are highly incomplete. One goal is to “complete” the KG, i.e., generate new valid facts from existing facts

Tensor Factorization: An Application

- Many data sets can be naturally represented as tensors
- Knowledge-Graphs (KG) of Knowledge-Bases (KB) is one such example
- A KG/KB consists of “facts” in form of triplets (e.g. Modi-PM-India)



- KGs are highly incomplete. One goal is to “complete” the KG, i.e., generate new valid facts from existing facts
- We can applying tensor factorization to learn features/embeddings of the entities and relations. Can use the embeddings to predict the unknown facts

Matrix/Tensor Factorization: Another View

- We modeled each matrix/tensor entry as an inner product of latent factors

Matrix/Tensor Factorization: Another View

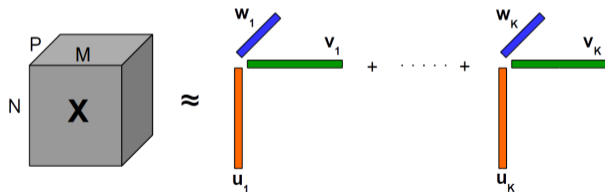
- We modeled each matrix/tensor entry as an inner product of latent factors
- Can also model matrix (or tensor) as a whole, as **sum of rank-1 components**

Matrix/Tensor Factorization: Another View

- We modeled each matrix/tensor entry as an inner product of latent factors
- Can also model matrix (or tensor) as a whole, as **sum of rank-1 components**
- E.g., an $N \times M \times P$ tensor \mathbf{X} as a sum of outer products of column vectors

$$\mathbf{X} \approx \sum_{k=1}^K \mathbf{u}_k \otimes \mathbf{v}_k \otimes \mathbf{w}_k \quad (\text{tensor SVD view; also generalizes to more than 3 dims})$$

$\mathbf{u}_k \in \mathbb{R}^N$, $\mathbf{v}_k \in \mathbb{R}^M$ and $\mathbf{w}_k \in \mathbb{R}^P$ denote the k -th columns of \mathbf{U} , \mathbf{V} , \mathbf{W} resp.

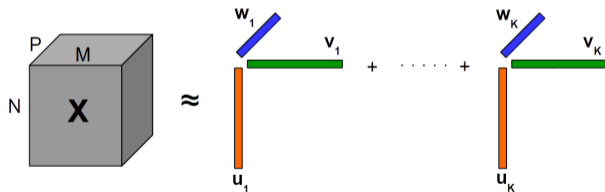


Matrix/Tensor Factorization: Another View

- We modeled each matrix/tensor entry as an inner product of latent factors
- Can also model matrix (or tensor) as a whole, as **sum of rank-1 components**
- E.g., an $N \times M \times P$ tensor \mathbf{X} as a sum of outer products of column vectors

$$\mathbf{X} \approx \sum_{k=1}^K \mathbf{u}_k \otimes \mathbf{v}_k \otimes \mathbf{w}_k \quad (\text{tensor SVD view; also generalizes to more than 3 dims})$$

$\mathbf{u}_k \in \mathbb{R}^N$, $\mathbf{v}_k \in \mathbb{R}^M$ and $\mathbf{w}_k \in \mathbb{R}^P$ denote the k -th columns of \mathbf{U} , \mathbf{V} , \mathbf{W} resp.



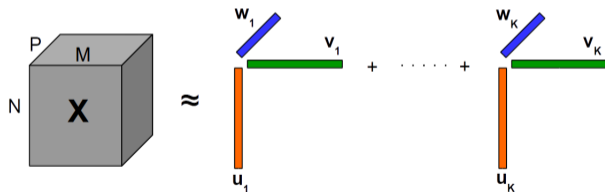
- The matrix case is similar (only 2 dimensions)

Matrix/Tensor Factorization: Another View

- We modeled each matrix/tensor entry as an inner product of latent factors
- Can also model matrix (or tensor) as a whole, as **sum of rank-1 components**
- E.g., an $N \times M \times P$ tensor \mathbf{X} as a sum of outer products of column vectors

$$\mathbf{X} \approx \sum_{k=1}^K \mathbf{u}_k \otimes \mathbf{v}_k \otimes \mathbf{w}_k \quad (\text{tensor SVD view; also generalizes to more than 3 dims})$$

$\mathbf{u}_k \in \mathbb{R}^N$, $\mathbf{v}_k \in \mathbb{R}^M$ and $\mathbf{w}_k \in \mathbb{R}^P$ denote the k -th columns of \mathbf{U} , \mathbf{V} , \mathbf{W} resp.



- The matrix case is similar (only 2 dimensions)
- With this view, in alternating optimization, can update one column at a time