### **Chapter IX: Matrix factorizations**

Information Retrieval & Data Mining Universität des Saarlandes, Saarbrücken Winter Semester 2011/12

### **Chapter IX: Matrix factorizations\***

- 1. The general idea
- 2. Matrix factorization methods
  - 2.1. Eigendecompositions
  - 2.2. SVD
  - **2.3. PCA**
  - 2.4. Nonnegative matrix factorization
  - **2.5. Some other matrix factorizations**
- 3. Latent topic models
- 4. Dimensionality reduction

\*Zaki & Meira, Ch. 8; Tan, Steinbach & Kumar, App. B; Manning, Raghavan & Schütze, Ch. 18 Extra reading: Golub & Van Loan: *Matrix computations*. 3rd ed., JHU press, 1996

## IX.1: The general idea

#### 1. The general definition

- 1.1. Matrix factorizations we've seen so far
- **1.2. Matrices as data and functions**
- **1.3. Matrix distances and types of matrices**
- 2. Very quick recap of linear algebra
- **3. Why matrix factorizations**

#### The general definition

- Given *n*-by-*m* matrix *X*, represent it as a product of two (or more) factor matrices *A* and *B* 
  - -X = AB
  - We are more interested in *approximate* matrix factorizations  $X \approx AB$
  - -Matrix *A* is *n*-by-*k*; matrix *B* is *k*-by-*m* ( $k \le \min(n, m)$ )
    - For more factor matrices, their *inner dimension* must match
- The distance between *X* and *AB* is the *representation error* of (approximate) factorization  $-E.g. \|X - AB\|_{F}^{2} = \sum_{i=1}^{n} \sum_{j=1}^{m} (x_{ij} - (AB)_{ij})^{2}$

#### Variations

- We can change the distance measure
  - Squared element-wise error
  - Absolute element-wise error
- We can restrict the matrices involved
  - Types of values
    - Non-negative
    - Binary
  - Types of factor matrices
    - Upper triangular
    - Diagonal
    - Orthogonal
- We can have more factor matrices
- We can change the matrix multiplication

#### Matrix factorizations we've seen so far

• Clustering:  $\|\mathbf{X} - \mathbf{C}\mathbf{M}\|_2^2$ 

-C has to be cluster assignment matrix

• Co-clustering:  $\|\mathbf{X} - \mathbf{R}\mathbf{M}\mathbf{C}^{\mathsf{T}}\|_{2}^{2}$ 

-R and C are cluster assignment matrices

- Linear regression:  $\|\mathbf{y} \mathbf{X}\mathbf{\beta}\|_2$ 
  - -y is vector, as is  $\beta$
  - -"decomposes" y but also is X is known
- Singular value decomposition (SVD) and eigendecomposition
  - -Have been mentioned earlier

#### Two views of a matrix: data or function

- In IR & DM (and most CS) a matrix is a way to write down data
  - -A two-dimensional flat database
  - Items and transactions, documents and terms, ...
- In linear algebra, a matrix is a linear function between vector spaces
  - *n*-by-*m* matrix maps *m*-dimensional vectors to
     *n*-dimensional ones
  - If y = Mx, then  $y_i = \sum_j m_{ij} x_j$
- Different views motivate different techniques

#### Matrix distances and norms

• Frobenius norm  $||X||_F = (\sum_{i,j} x_{ij}^2)^{1/2}$ 

-Corresponds to Euclidean norm of vectors

- Sum of absolute values  $|X| = \sum_{i,j} x_{ij}$ - Corresponds to  $L_1$ -norm of vectors
- The above elementwise norms are sometimes (imprecisely) called  $L_2$  and  $L_1$  norms

– Matrix  $L_1$  and  $L_2$  norms are something different altogether

- Operator norm  $||X||_p = \max_{y \neq 0} ||Xy||_p / ||y||_p$ 
  - Largest norm of an image of a unit norm vector  $||\mathbf{Y}||_{\mathbf{V}} \leq \sqrt{(\operatorname{rank}(\mathbf{Y}))} ||\mathbf{Y}||_{\mathbf{V}}$
  - $-||X||_{2} \le ||X||_{F} \le \sqrt{(\operatorname{rank}(X))} ||X||_{2}$

# Types of matrices

- Diagonal *n*-by-*n* matrix
  Identity matrix *I<sub>n</sub>* is a diagonal *n*-by-*n* matrix with 1s in diagonal
- Upper triangular matrix
  - -Lower triangular is the transpose
  - -If diagonal is full of 0s, matrix is *strictly triangular*
- Permutation matrix
  - -Each row and column has exactly one 1, rest are 0



## Very quick recap of linear algebra

- An *n*-by-*m* matrix *X* can be represented exactly as a product of *n*-by-*k* and *k*-by-*m* matrices *A* and *B* if and only if rank of *X* is at most *k* 
  - $-\operatorname{rank}(AB) \le \min(\operatorname{rank}(A), \operatorname{rank}(B))$
  - If rank(X) =  $n \le m$ , we can set  $A = I_n$  and B = X
  - In general, if  $n \le m$ , columns of A are linearly independent *basis vectors* for the subspace spanned by X and columns of B tell the linear combinations of these vectors needed to get the original columns of X
- If *X* is rank-*k*, it can be written as a sum of *k* rank-1 matrices, but no fewer
  - Another way to define rank
  - In general,  $rank(A + B) \le rank(A) + rank(B)$

#### Spaces

- Let X be an n-by-m (real-valued) matrix
  - -Set  $\{u \in \mathbb{R}^n : Xv = u, v \in \mathbb{R}^m\}$  is the *column space* of X •Image of X
  - -Set { $v \in \mathbb{R}^m : X^T u = v, u \in \mathbb{R}^n$ } is the *row space* of X•Image of  $X^T$
  - -Set { $v \in \mathbb{R}^m : Xv = 0$ } is the *null space* of *X*
  - -Set { $u \in \mathbb{R}^n : X^T u = 0$ } is the *left null space* of X

## Orthogonality and orthonormality

- Two vectors x and y are **orthogonal** if their inner product  $\langle x, y \rangle$  is 0
  - -Vectors are orthonormal if they have unit norm, ||x|| = ||y|| = 1
- A square matrix *X* is *orthogonal* if its rows and columns are orthonormal
  - -Equivalently,  $X^T = X^{-1}$
  - -Yet equivalently,  $XX^T = X^TX = I$

#### Why matrix factorizations?

- A general way of writing many problems
  - Makes easier to see similarities & differences
  - -May help finding new approaches and tools
- A method to remove noise
  - -"True" matrix A is low-rank
  - -Observed matrix  $\tilde{A}$  has some noise  $A + \varepsilon$  and has full rank
  - -Finding a low-rank approximation of  $\tilde{A}$  helps remove the noise and leave only the original matrix A
  - -Here we're interested in the representation of A
- Alternatively we can be interested on the factors...

#### Factors and dimensionality reduction

- Let X be *n*-by-*m*, A be *n*-by-*k*, B be *k*-by-*m*, and  $X \approx AB$ 
  - -Rows of A are k-dimensional representations of rows of X
  - -Columns of B are k-dimensional representations of columns of X
  - We can project rows of X to k-dimensional subspace  $XB^T$ 
    - Columns of X are projected with  $A^{T}X$
- Low-dimensional views allow
  - -Direct study of factors
    - By hand, plotting, etc.
  - Avoidance of *curse of dimensionality* (more on this later)
  - -Better scalability / avoidance of noise

- 10-dimensional data
- Clustered using *k*-means in 3 clusters
- Want to visualize the clusters
  - Are they "natural"?
- Project the data to first two principal components:



17 January 2012

#### **IX.2 Matrix factorization methods**

- 1. Eigendecomposition
- 2. Singular value decomposition (SVD)
- 3. Principal component analysis (PCA)
- 4. Non-negative matrix factorization
- **5.** Other matrix factorization methods
  - **5.1. CX matrix factorization**
  - 5.2. Boolean matrix factorization
  - 5.3. Regularizers
  - **5.4. Matrix completion**

### Eigendecomposition

- If X is an *n*-by-*n* matrix and *v* is a vector such that  $Xv = \lambda v$  for some scalar  $\lambda$ , then
  - $-\lambda$  is an **eigenvalue** of *X*
  - -v is an **eigenvector** of *X* associated to  $\lambda$
- Matrix X has to *diagonalizable*-PXP<sup>-1</sup> is a diagonal matrix for some invertible matrix P
- Matrix X has to have n linearly independent eigenvectors
- The eigendecomposition of X is  $X = QAQ^{-1}$

-Columns of Q are the eigenvectors of X

 $-\Lambda$  is a diagonal matrix with eigenvalues in the diagonal

#### Some useful facts

- Not all matrices have eigendecomposition
  - -Not all invertible matrices have eigendecomposition
  - -Not all matrices that have eigendecomposition are invertible
  - If X is invertible and has eigendecomposition, then  $X^{-1} = QA^{-1}Q^{-1}$
- If X is symmetric and invertible (and real), then X has eigendecomposition  $X = QAQ^T$

#### How to find eigendecomposition, part 1

- Recall the *power method* for computing the stationary distribution of a Markov chain
  - $-\boldsymbol{v}_{t+1} = \boldsymbol{v}_t \boldsymbol{P}$
  - -Computes the *dominant* eigenvalue and eigenvector
    - Can't be used to find the full eigendecomposition
- Similar iterative idea is usually used:
  - -Let  $X_0 = X$  and find orthogonal  $Q_t$  such that  $X_t = Q_t^T X_{t-1} Q_t$  is "more diagonal" than  $X_{t-1}$
  - -When  $X_t$  is diagonal enough, set  $\Lambda = X_t$  and  $Q = Q_t Q_{t-1} Q_{t-2} \cdots Q_1$

#### The Jacobi method for symmetric matrix

- We assume that X is symmetric *n*-by-*n*
- The idea is to reduce the quantity off(X) =  $\sqrt{\sum_{i,j:i \neq j} x_{ij}^2}$
- The Jacobi rotations are matrices of form

$$J(p,q,\theta) = \begin{pmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \cdots & c & \cdots & s & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & & \vdots \\ 0 & \cdots & -s & \cdots & c & \cdots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix} q \qquad s = \sin(\theta)$$

D

q

#### Basic Jacobi step

- 1. Choose index pair (p,q) s.t.  $1 \le p \le q \le n$
- 2. Compute  $c = cos(\theta)$  and  $s = sin(\theta) s.t.$

 $\begin{pmatrix} y_{pp} & y_{pq} \\ y_{qp} & y_{qq} \end{pmatrix} = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}^{\mathsf{T}} \begin{pmatrix} x_{pp} & x_{pq} \\ x_{qp} & x_{qq} \end{pmatrix} \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$ is diagonal  $(y_{pq} = y_{qp} = 0)$ 

3. Overwrite X with  $Y = J^T X J$  where  $J = J(p,q,\theta)$ .

Each Jacobi step reduces off-diagonal values of the 2-by-2 matrix by  $off(\mathbf{Y})^2 = off(\mathbf{X})^2 - 2x_{pq}^2$ 

# Basic Jacobi step

- 1. Choose index pair (p,q) s.t. 1  $\leq$
- 2. Compute  $c = cos(\theta)$  and  $s = sin(\theta) s.t.$

$$\begin{pmatrix} y_{pp} & y_{pq} \\ y_{qp} & y_{qq} \end{pmatrix} = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}^{\mathsf{T}} \begin{pmatrix} x_{pp} & x_{pq} \\ x_{qp} & x_{qq} \end{pmatrix} \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$
  
is diagonal  $(y_{pq} = y_{qp} = 0)$ 

3. Overwrite X with  $Y = J^T X J$  where  $J = J(p,q,\theta)$ .

Each Jacobi step reduces off-diagonal values of the 2-by-2 matrix by  $off(\mathbf{Y})^2 = off(\mathbf{X})^2 - 2x_{pq}^2$ 

Symmetric 2-by-2

eigendecomposition

#### Basic Jacobi step

- 1. Choose index pair (p,q) s.t.  $1 \le p \le q \le n$
- 2. Compute  $c = cos(\theta)$  and  $s = sin(\theta) s.t.$

 $\begin{pmatrix} y_{pp} & y_{pq} \\ y_{qp} & y_{qq} \end{pmatrix} = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}^{\mathsf{T}} \begin{pmatrix} x_{pp} & x_{pq} \\ x_{qp} & x_{qq} \end{pmatrix} \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$ is diagonal  $(y_{pq} = y_{qp} = 0)$ 

3. Overwrite X with  $Y = J^T X J$  where  $J = J(p,q,\theta)$ .

Each Jacobi step reduces off-diagonal values of the 2-by-2 matrix by  $off(\mathbf{Y})^2 = off(\mathbf{X})^2 - 2x_{pq}^2$ 

#### How to select c and s

- We want to have  $c = \cos(\theta)$  and  $s = \sin(\theta)$  s.t.  $0 = y_{pq} = x_{pq}(c^2 - s^2) + (x_{pp} - x_{qq})cs$
- If  $x_{pq} = 0$ , set c = 1 and s = 0
- Else set  $\tau = (x_{qq} x_{pp})/(2x_{pq})$
- If  $\tau \ge 0$ , set  $t = 1/(\tau + \sqrt{(1 + \tau^2)})$ - Else set  $t = -1/(-\tau + \sqrt{(1 + \tau^2)})$
- Set  $c = 1/\sqrt{(1 + t^2)}$  and s = tc

#### How to select *p* and *q*

- In *Classical Jacobi* select (*p*,*q*) such that |*x<sub>pq</sub>*| = max<sub>*i≠j*</sub> |*x<sub>ij</sub>*|
  -Finding this value takes O(*n*<sup>2</sup>) time
- In *Cyclic Jacobi* go thru the off-diagonal elements in a fixed order
  - $-E.g.(p,q) = (1,2), (1,3), (1,4), (2,3), (2,4), (3,4), (1,2), \dots$

#### Jacobi in nutshell

- 1. Set  $V = I_n$ ;  $eps = tol \times ||X||_F$ ; Y = X
- 2. while off(*Y*) > *eps* 
  - 2.1.Choose (p,q) so  $|x_{pq}| = \max_{i \neq j} |x_{ij}|$  (or use cyclic order) 2.2.Compute cosine—sin pair (c,s)2.3. $Y = J(p,q,\theta)^T Y J(p,q,\theta)$ 2.4. $V = V J(p,q,\theta)$
- 3. end while
- 4. return  $\Lambda = Y$  and  $Q = V(X \approx Q\Lambda Q^T)$

#### Some notes

- The quality (and running time) depends on parameter tol > 0
- Jacobi method is easy to parallellize
   Split the update in non-conflicting steps
- Other methods exist
  - -Symmetric QR algorithm
  - -Tri-diagonal methods
    - Bisecting algorithm
    - Divide-and-conquer
- Numerical stability is an issue with all these methods

### Singular value decomposition (SVD)

- Not every matrix has eigendecomposition, but: **Theorem.** If *X* is *n*-by-*m* real matrix, there exists *n*-by-*n* orthogonal matrix *U* and *m*-by-*m* orthogonal matrix *V* such that  $U^TXV$  is *n*-by-*m* matrix  $\Sigma$  with values  $\sigma_1, \sigma_2, ..., \sigma_{\min(n,m)}, \sigma_1 \ge \sigma_2 \ge ... \ge \sigma_{\min(n,m)} \ge 0$ , in its diagonal.
  - In other words,  $X = U\Sigma V^T$
  - -Values  $\sigma_i$  are the **singular values** of *X*
  - -Columns of *U* are the left singular vectors and columns of *V* the right singular vectors of *X*



### Properties of SVD, part 1

- rank(X) = r iff X has exactly r non-zero singular values ( $\sigma_1 \ge \sigma_2 \ge ... \ge \sigma_r > \sigma_{r+1} = ... = \sigma_{\min(n,m)} = 0$ )
- Vectors  $u_1, u_2, ..., u_r$  are a basis for the column space of X
- Vectors  $u_{r+1}$ ,  $u_{r+2}$ , ...,  $u_n$  are a basis for the left null space of X
- Vectors  $v_1, v_2, ..., v_r$  are a basis for the row space of X
- Vectors  $v_{r+1}$ ,  $v_{r+2}$ , ...,  $v_m$  are a basis for the null space of X

## Properties of SVD, part 2

• If X is rank-r, then  $X = \sum_{i=1}^{r} \sigma_i u_i v_i^T$ 

-X is a sum of r rank-1 matrices scaled with singular values

- $\|\mathbf{X}\|_{\mathrm{F}}^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_{\min(n,m)}^2$
- $\bullet \|\mathbf{X}\|_2 = \sigma_1$
- Eckart–Young theorem. Let *X* be of rank-*r* and let  $U\Sigma V^T$  be its SVD. Denote by  $U_k$  the first *k* columns of *U*, by  $V_k$  the first *k* columns of *V* and by  $\Sigma_k$  the upperleft *k*-by-*k* corner of  $\Sigma$ . Then  $X_k = U_k \Sigma_k V_k^T$  is the best rank-*k* approximation of *X* in the sense that  $\|X X_k\|_F \leq \|X Y\|_F$  and  $\|X X_k\|_2 \leq \|X Y\|_2$  for any rank-*k* matrix *Y*.

#### SVD and pseudo-inverse

- Recall that if X is *n*-by-*m* with rank(X) =  $m \le n$ , the *pseudo-inverse* of X is  $X^{\dagger} = (X^T X)^{-1} X^T$
- If rank(X) = r and  $X = U\Sigma V^T$ , then we can define  $X^{\dagger} = V\Sigma^{\dagger}U^T$ 
  - $-\Sigma^{\dagger}$  is a diagonal matrix with  $1/\sigma_i$  in its *i*th position -More general than the above definition
- This gives the least-squares solution to the following problem: given A and X, find Y s.t.  $||A XY||_F^2$  is minimized

-Setting  $Y = X^{\dagger}A$  minimizes the squared Frobenius

#### SVD and eigendecomposition

- Let X be *n*-by-*m* and  $X = U\Sigma V^T$  its SVD
- Recall that the *Gram matrix* of the columns of X is  $X^T X$ -For the rows it is  $XX^T$
- Now  $X^T X = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T U^T U \Sigma V^T$ =  $V \Sigma^T \Sigma V^T = V \Sigma_m^2 V^T$ 
  - $-\Sigma_m^2$  is an *m*-by-*m* diagonal matrix with  $\sigma_i^2$  in its *i*th position
- Similarly  $XX^T = U\Sigma_n^2 U^T$
- Therefore
  - -Columns of U are the eigenvectors of  $XX^T$
  - -Columns of V are the eigenvectors of  $X^T X$
  - Singular values are square roots of the associated eigenvalues

# Computing the SVD

- Simple idea: Compute the eigendecompositions of  $XX^T$  and  $X^TX$ 
  - -Bad for numerical stability
- We can adapt the Jacobi method:
  - At each step find a Jacobi rotation  $J(p,q,\theta)$  such that columns p and q of  $XJ(p,q,\theta)$  are orthogonal
    - Corresponds to zeroing (p,q) and (q,p) in  $X^T X$
    - The product of this sequence of Jacobi rotations gives orthogonal V
    - Rest follows by  $AV = U\Sigma$
  - -This is called one-sided Jacobi

### Principal component analysis (PCA)

- Let rows of matrix denote observations and columns denote variables
- In **principal component analysis** (PCA) we want to find new variables (dimensions) that capture the variance of the data
  - -First variable has as much variance as possible
  - Second variable is orthogonal to the first and captures as much as possible of the remaining variance
  - -Third variable ...







# Computing the PCA

• First, data is centered

– The mean of each column is subtracted from the column

- Then, the *m*-by-*m* covariance matrix **S** is computed  $-s_{ij}$  is the covariance between *i*th and *j*th column (variable) - For centered data  $X, S = 1/n X^T X$
- The first principal vector is given by the eigenvector of S associated with the highest eigenvalue  $\lambda_1$

 $-\lambda 1$  gives the amount of variance explained

- The second principal vector is given by the second eigenvector, etc.
- The total variance of the data is  $\lambda_1 + \lambda_2 + \ldots + \lambda_m$

### PCA and SVD

- Alternatively, we can just compute the SVD of centered data X'
  - -Now the principal vectors are columns of *V* -Therefore, PCA is *SVD done with centered data*
- We can project the data X' into its principal space by X'V



#### How many principal vectors?

- Rule of thumb: keep 90% of variance - Select *k* s.t.  $(\lambda_1 + \lambda_2 + ... + \lambda_k)/(\lambda_1 + \lambda_2 + ... + \lambda_m) \ge 0.9$ - Same as  $(\sigma_1^2 + \sigma_2^2 + ... + \sigma_k^2)/(\sigma_1^2 + \sigma_2^2 + ... + \sigma_m^2) \ge 0.9$
- But if you want to do plotting, you need less...



17 January 2012

#### Nonnegative matrix factorization (NMF)

- Eigenvectors and singular vectors can have negative entries even if the data is non-negative
  - This can make the factor matrices hard to interpret in the context of the data
- In **nonnegative matrix factorization** we assume the data is nonnegative and we require the factor matrices to be nonnegative
  - -Factors have parts-of-whole interpretation
    - Data is represented as a sum of non-negative elements
  - -Models many real-world processes

#### Definition

- Given a nonnegative *n*-by-*m* matrix X (i.e. x<sub>ij</sub> ≥ 0 for all *i* and *j*) and a positive integer k, find an *n*-by-k nonnegative matrix W and a k-by-m nonnegative matrix H s.t. ||X WH||<sub>F</sub><sup>2</sup> is minimized.
  - If  $k = \min(n,m)$ , we can do W = X and  $H = I_m$  (or vice versa)
  - -Otherwise the complexity of the problem is unknown
- If either *W* or *H* is fixed, we can find the other factor matrix in polynomial time
  - Which gives us our first algorithm...

## The alternating least squares (ALS)

- Let's forget the nonnegativity constraint for a while
- The alternating least squares algorithm is the following:
  - -Intialize W to a random matrix
  - -repeat
    - Fix W and find H s.t.  $||X WH||_{F^2}$  is minimized
    - Fix **H** and find **W** s.t.  $||X WH||_{F^2}$  is minimized
  - -until convergence
- For *unconstrained least squares* we can use  $H = W^{\dagger}X$ and  $W = XH^{\dagger}$
- ALS will typically converge to *local optimum*

#### NMF and ALS

- With the nonnegativity constraint pseudo-inverse doesn't work
  - The problem is still *convex* with either of the factor matrices fixed (but not if both are free)
  - -We can use *constrained convex optimization* 
    - In theory, polynomial time
    - In practice, often too slow
- Poor man's nonnegative ALS:
  - -Solve *H* using pseudo-inverse
  - -Set all  $h_{ij} < 0$  to 0
  - -Repeat for W