

# Matrix Decomposition Algorithms For Data Mining and Machine Learning

---

2019. 1. 4

**Troy University**  
Ingyu Lee

(Contents are based on Petros Drineas, Saara Hyvonen and Rainer Gemulla's **Materials**)

# Content

---

## **Part 0. Preliminary (30 mins)**

- 0-1. Modern Data
- 0-2. About the Tutorial
- 0-3. References
- 0-4. Suggested Readings

## **Part 1. Linear Algebra Review (1 hour)**

- 1-1. Basics of Linear Algebra
- 1-2. Vector Notations
- 1-3. Matrix Notations

## **Part 2. Matrix Algorithms (1 hour)**

- 2-1. Singular Value Decomposition (SVD)
- 2-2. Non-negative Matrix Factorization (NMF)
- 2-3. Independent Component Analysis (ICA)
- 2-4. Semi-Discrete Decomposition (SDD)

## **Part 3. Applications of Matrix Algorithms (30 mins)**

- 3-1. Examples of Matrix Algorithms Applied on Data Mining Applications

# Part 0: Modern data

---

## Facts:

- Computers make it **easy to collect and store** data.
- Costs of storage are very **low and are dropping** very fast.  
(most laptops have a storage capacity of more than 1TB..)

## When it comes to storing data:

- The current policy typically is “**store everything in case it is needed later**” instead of deciding what could be deleted.

## Data Mining:

- **Extract useful information from the massive amount of available data.**
- “Data mining is **the process of discovering knowledge or patterns** from massive amounts of data.”

# Part 0: About the tutorial

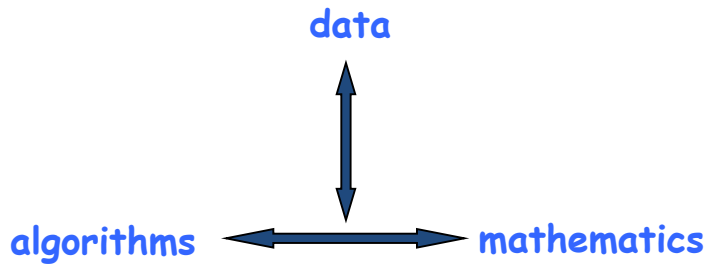
---

## Tools:

- Introduce **matrix algorithms and matrix decompositions** for data mining applications.

## Goals:

- Learn a model for the underlying “physical” system generating the dataset.



**Math** is necessary to design and analyze principled **algorithmic** techniques to **data-mine** the massive datasets that have become ubiquitous in scientific research.

## Part 0: References

---

- **Information Retrieval and Data Mining: A Linear Algebraic Perspective**,  
<http://www.cs.rpi.edu/~drinep/talks.html>  
(or <https://www.ipam.ucla.edu/schedule.aspx?pc=setut> Podcast is available)  
Petros Drineas, Rensselaer Polytechnic Institute
- **Data Mining and Matrices**,  
Rainer Gemulla, and Pauli Miettinen, Max Planck Institute
- **Linear Algebra Methods for Data Mining**,  
<http://www.uta.edu/faculty/rcli/Teaching/math5392/s2010/>  
Saara Hyvonen, University of Helsinki
- **Data Mining with Graphs and Matrices**,  
<http://feiwang03.googlepages.com/sdm-tutorial>  
Fei Wang, Tao Li and Chris Ding, University of Texas at Arlington

## Part 0: Suggested Readings

---

- Lars Elden  
**Matrix Methods in Data Mining and Pattern Recognition**, SIAM, 2007.
- David Skillicorn  
**Understanding Complex Datasets: Data Mining with Matrix Decompositions**, Chapman and Hall, 2007.
- Gene H. Golub and Charles F. Van Loan  
**Matrix Computations**, 3rd ed. Johns Hopkins University Press, 1996
- Carl Meyer  
**Matrix Analysis and Applied Linear Algebra**, SIAM, 2000.  
<http://www.matrixanalysis.com>

# Part 0: Why linear (or multilinear) algebra?

---

**Data are represented by matrices:**

- Numerous modern datasets are in **matrix form**.

**Data are represented by tensors:**

- **Data in the form of tensors** (multi-mode arrays) are becoming very common in the data mining and information retrieval literature in the last few years.
- Linear algebra (and numerical analysis) **provide the fundamental mathematical and algorithmic tools to deal with matrix and tensor computations**.

# Part 0: What is a matrix?

- A means to describe computation
  - Rotation
  - Rescaling
  - Permutation
  - Projection
  - ...



Linear operators

- A means to describe data

Rows	Columns	Entries
<i>Objects</i>	<i>Attributes</i>	<i>Values</i>
Equations	Variables	Coefficients
Data points	Axes	Coordinates
Vertices	Vertices	Edges
⋮	⋮	⋮

$$\begin{array}{c} \text{Object } i \end{array} \begin{pmatrix} A_{11} & A_{12} & \cdots & A_{1j} & \cdots \\ A_{21} & A_{22} & \cdots & A_{2j} & \cdots \\ \vdots & \vdots & \ddots & \vdots & \ddots \\ A_{i1} & A_{i2} & \cdots & A_{ij} & \cdots \\ \vdots & \vdots & \ddots & \vdots & \ddots \end{pmatrix} \begin{array}{c} \text{Attribute } j \end{array}$$

\* In data mining, we make use of both viewpoints simultaneously



# Part 0: Datasets in the form of matrices

We are given  $m$  objects and  $n$  features describing the objects.  
(Each object has  $n$  numeric values describing it.)

## Dataset:

- An  $m$ -by- $n$  matrix  $\mathbf{A}$ ,  $A_{ij}$  shows the “importance” of feature  $j$  for object  $i$ .
- Every row of  $\mathbf{A}$  represents an object.

## Goal:

- We seek to **understand the structure of the data**, e.g. the underlying process generating the data.

	Bread	Butter	Beer
Anna	1	1	0
Bob	1	1	1
Charlie	0	1	1

*Customer transactions*

	Data	Matrix	Mining
Book 1	5	0	3
Book 2	0	0	7
Book 3	4	6	5

*Document-term matrix*

	Avatar	The Matrix	Up
Alice		4	2
Bob	3	2	
Charlie	5		3

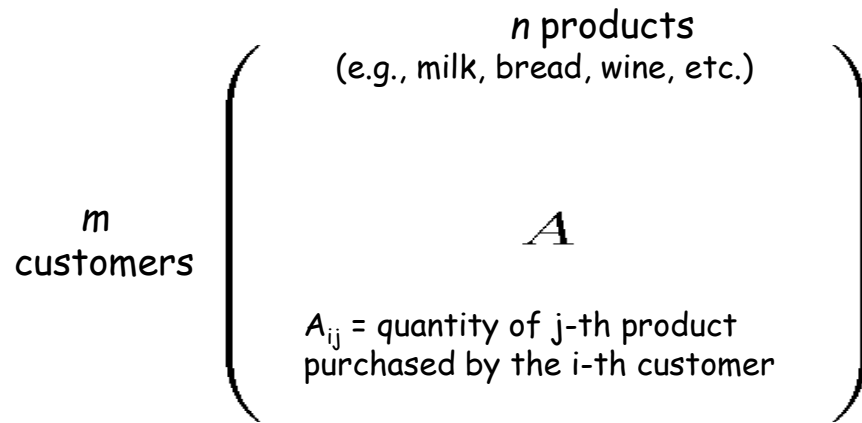
*Incomplete rating matrix*

	Jan	Jun	Sep
Saarbrücken	1	11	10
Helsinki	6.5	10.9	8.7
Cape Town	15.7	7.8	8.7

*Cities and monthly temperatures*

# Part 0: Market basket matrices

- Common representation for association rule mining.



## Data mining tasks

- Find association rules

E.g., customers who buy product  $x$  buy product  $y$  with probability 89%.

- Such rules are used to make item display decisions, advertising decisions, etc.

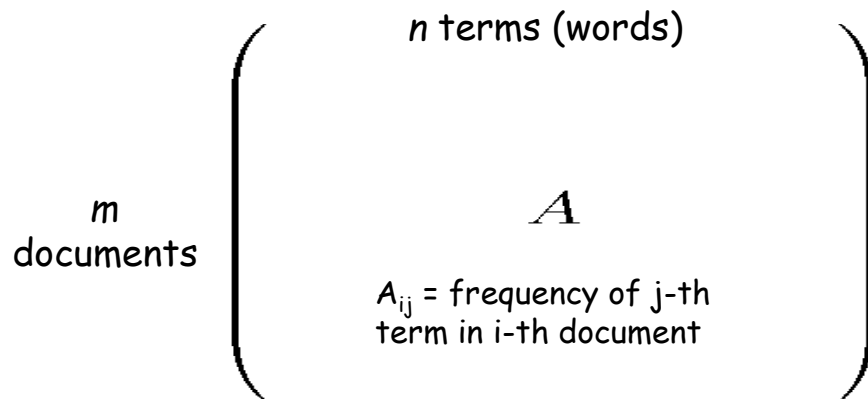
## Market Basket Example



Image source: deepclimate.org

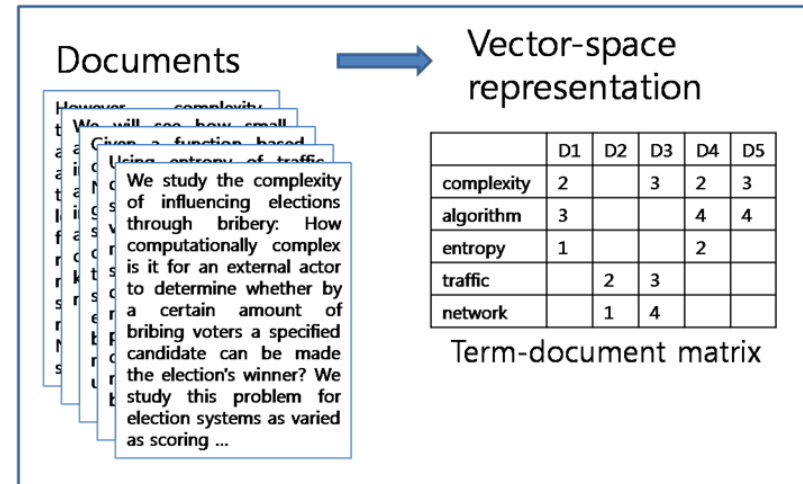
# Part 0: Document-term matrices

- A collection of documents is represented by an m-by-n matrix (bag-of-words model).



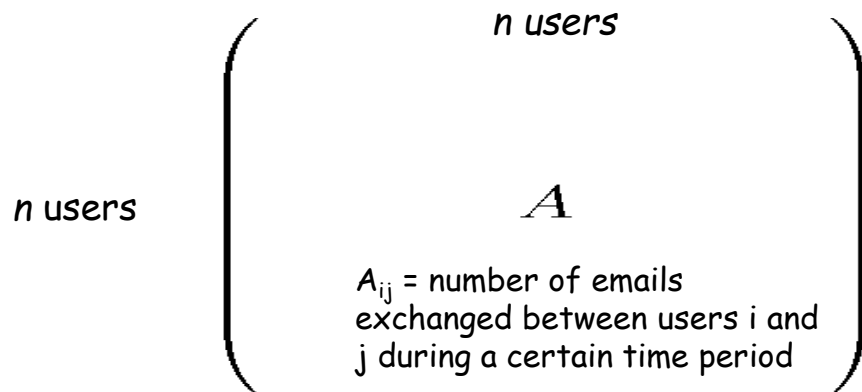
## Data mining tasks

- Cluster or classify documents
- Find "nearest neighbors"
- Feature selection: find a subset of terms that (accurately) clusters or classifies documents.



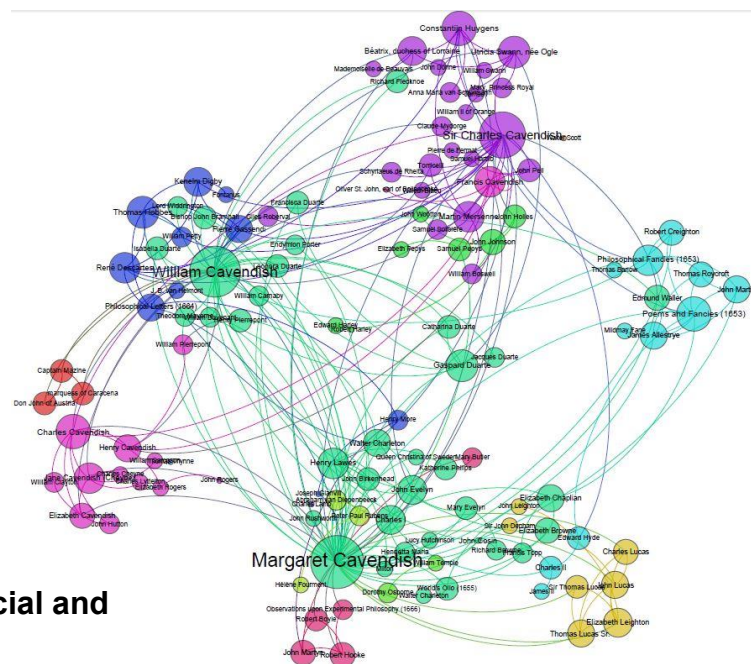
# Part 0: Social networks (e-mail graph)

- Represent the email communications between groups of users.



## Data mining tasks

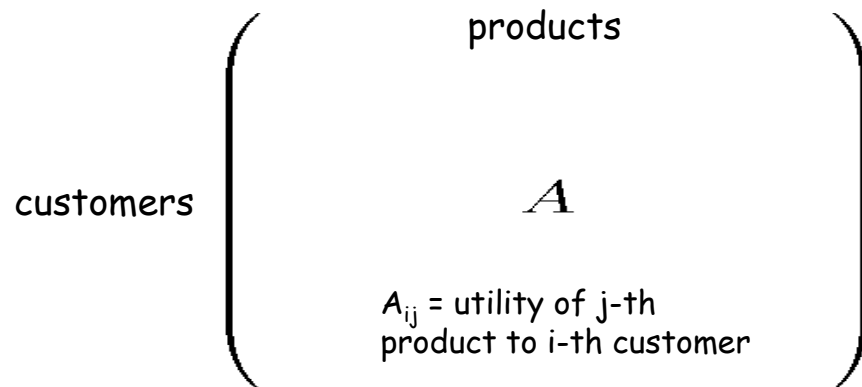
- cluster the users
- identify “dense” networks of users (dense subgraphs)



Margaret Cavendish's social and some textual networks.

# Part 0: Recommendation systems

- The  $m$ -by- $n$  matrix  $A$  represents  $m$  customers and  $n$  products.



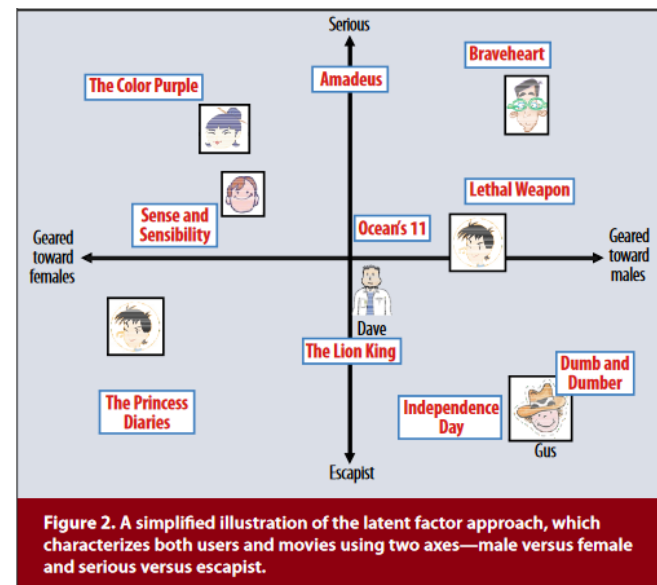
## MATRIX FACTORIZATION TECHNIQUES FOR RECOMMENDER SYSTEMS

Yehuda Koren, Yahoo Research

Robert Bell and Chris Volinsky, AT&T Labs—Research

### Data mining task

- Given a few samples from  $A$ , recommend high utility products to customers.



# Part 0: Why matrix decompositions?

## Matrix decompositions:

(e.g. SVD, QR, SDD, CX and CUR, NMF, ICA, etc.)

- They use the relationships between the available data in order to **identify components** of the underlying physical system generating the data.
- Some assumptions on the relationships between the underlying components are necessary.
- **Very active area of research**: some matrix decompositions are more than one century old, whereas others are very recent.

- A matrix decomposition of a data matrix **D** is given by three matrices **L**, **M**, **R** such that

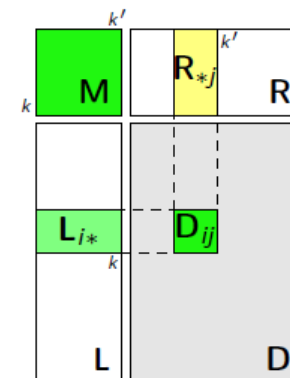
$$\mathbf{D} = \mathbf{L}\mathbf{M}\mathbf{R},$$

where

- **D** is an  $m \times n$  data matrix,
- **L** is an  $m \times r$  matrix,
- **M** is an  $r \times r$  matrix,
- **R** is an  $r \times n$  matrix, and
- $r$  is an integer  $\geq 1$

There are many different kinds of matrix decompositions, each putting certain constraints on matrices **L**, **M**, **R** (which may not be easy to find).

$$D_{ij} = \sum_{k,k'} L_{ik} M_{kk'} R_{k'j}$$



# Part 0: What can we do with matrix decompositions?

---

- Separate data from multiple processes
- Remove noise from the data
- Remove redundancy from the data
- Reveal latent structure and similarities in the data
- Fill in missing entries
- Find local patterns
- Reduce space consumption
- Reduce computational cost
- Aid visualization

Matrix decompositions can [make data mining algorithms more effective](#). They may also [provide insight into the data](#) by themselves.

# Part 0: Factor interpretation of matrix decompositions?

- Assume that  $M$  is diagonal. Consider object  $i$

■ Row of  $R$  = part (or piece), called *latent factor* ("latent object")

■ Entry of  $M$  = weight of corresponding part

Row of  $MR$  = weighted part

■ Row of  $L$  = "view" of corresponding row of  $D$   
in terms of the weighted parts  
( $r$  pieces of information)

$r$  forces "compactness" (often  $r < n$ )

Each object can be viewed as a combination of  $r$  (weighted) "latent objects" (or "prototypical objects"). Similarly, each attribute can be viewed as a combination of  $r$  (weighted) "latent attributes."

(e.g., latent attribute = "body size"; latent object relates body size to real attributes such as "height", "weight", "shoe size")

$$D_{i*} = \sum_k L_{ik} M_{kk} R_{k*}$$

The diagram shows three matrices arranged in a block matrix structure. Matrix  $M$  is a square matrix with a diagonal line and a single pink square at the top-left corner. Matrix  $R$  is a matrix with a single green row at the top. Matrix  $L$  is a matrix with a single yellow square at the top-left corner. The resulting matrix  $D$  is a matrix with a single grey row at the top, labeled  $D_{i*}$ . The equation  $D_{i*} = \sum_k L_{ik} M_{kk} R_{k*}$  is shown above the matrices.



## Part 0: Lessons learned

---

- Data mining = from data to knowledge  
→ Prediction, clustering, outlier detection, local patterns
- Many different data types can be represented with a matrix  
→ Linear equations, data points, maps, graphs, relational data, ...
- Common interpretation: rows = objects, columns = attributes
- Matrix decompositions reveal structure in the data  
→  **$D = LMR$**
- Many different decompositions with different applications exist  
→ QR, SVD, NMF, SDD, ICA, etc.
- Factor interpretation: objects described by “latent attributes”

# Content

---

## Part 0. Preliminary (30 mins)

- 0-5. Why linear algebra?
- 0-6. Why matrix decompositions?
- 0-7. Datasets in the form of matrices
- 0-8. Examples

## Part 1. Linear Algebra Review (1 hour)

- 1-1. Basics of Linear Algebra
- 1-2. Vector Notations
- 1-3. Matrix Notations

## Part 2. Matrix Algorithms (1 hour)

- 2-1. Singular Value Decomposition (SVD)
- 2-2. Non-negative Matrix Factorization (NMF)
- 2-3. Semi-Discrete Decomposition (SDD)
- 2-4. Independent Component Analysis (ICA)

## Part 3. Applications of Matrix Algorithms (30 mins)

- 3-1. Examples of Matrix Algorithms Applied on Data Mining Applications

# Part 1: Vector

---

- A **vector** is
  - a 1D array of numbers
  - a geometric entity with magnitude and direction
  - a matrix with exactly one row or column → row and column vectors
- A **transpose**  $\mathbf{x}^T$  transposes a row vector into a column vector and vice versa
- The **norm** of vector defines its magnitude
  - 1-norm:  $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$
  - Euclidean or  $L_2$  norm:  $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
  - General  $L_p$ :  $\|\mathbf{x}\|_p = \left(\sum_{i=1}^n x_i^p\right)^{1/p}$
- Generally, a vector norm is a mapping  $R^n \rightarrow R$  with the properties
  - $\|\mathbf{x}\| \geq 0$  for all  $x$ ,
  - $\|\mathbf{x}\| = 0$  if and only if  $x = 0$ ,
  - $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$ , for all  $\alpha \in R$ ,
  - $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ , the triangular equality.

## Part 1: How to measure distance between vectors?

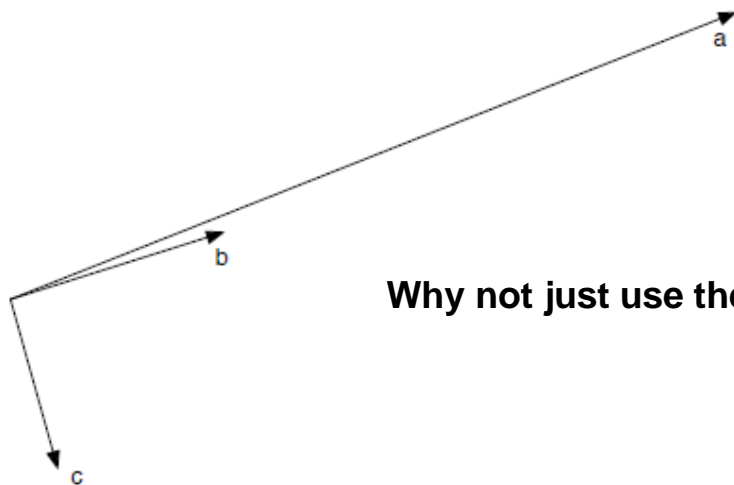
---

- Obvious answer: the distance between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is  $||\mathbf{x} - \mathbf{y}||$ , where  $|| \cdot ||$  is some vector norm.
- Frequently one measures the distance by Euclidean norm  $||\mathbf{x} - \mathbf{y}||_2$ .
- So usually, if the index is dropped, this is what is meant.
- Alternative: use the angle between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  to measure the distance between them.
- How to calculate the angle between two vectors?

## Part 1: Angle between vectors

---

- The inner product between two vectors is defined by  $(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$ .
- This is associated with the Euclidean norm:  $\|\mathbf{x}\|_2 = (\mathbf{x}^T \mathbf{x})^{1/2}$ .
- The angle  $\theta$  between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is  $\cos\theta = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$ .
- The cosine of the angle between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  can be used to measure the **similarity** between the two vectors:
  - If  $\mathbf{x}$  and  $\mathbf{y}$  are close, the angle between them is small, and  $\cos\theta \approx 1$ .
  - $\mathbf{x}$  and  $\mathbf{y}$  are **orthogonal**, if  $\theta = \frac{\pi}{2}$ , i. e.  $\mathbf{x}^T \mathbf{y} = 0$



Why not just use the Euclidean distance?

## Part 1: Example: term-document matrix

---

- Each entry tells how many times a term appears in the document:

	Doc1	Doc2	Doc3
Term1	10	1	0
Term2	10	1	0
Term3	0	0	1

- Using the Euclidean distance Documents 1 and 2 look dissimilar, and Documents 2 and 3 look similar. This is just due to the length of the documents!!!
- Using the cosine of the angle between document vectors Documents 1 and 2 are similar to each other and dissimilar to Document 3.

# Part 1: Matrix Norms

---

- Let  $\|\cdot\|$  be a vector norm and  $A \in R^{m \times n}$ .
- The corresponding matrix norm is  $\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}$ .
- $\|A\|_2 = (\max_{1 \leq i \leq n} \lambda_i(A^T A))^{1/2}$  = square root of the largest eigenvalue of  $A^T A$ . Heavy to compute!!!
- $\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$  (maximum over rows)
- $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$  (maximum over columns)
- $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$  **Frobenius** norm: does not correspond to any vector norm. Still, related to Euclidean vector norm.

# Part 1: Linear Independence

---

- Given a set of vectors  $(\mathbf{v}_j)^n$  in  $R^m, m \geq n$ , consider the set of linear combinations  $y = \sum_{j=1}^n \alpha_j \mathbf{v}_j$  for arbitrary coefficients  $\alpha_j$ .
- The vectors  $(\mathbf{v}_j)^n$  are linearly independent, if  $\sum_{j=1}^n \alpha_j \mathbf{v}_j = 0$  if and only if  $\alpha_j = 0$  for all  $j = 1, \dots, n$ .
- A set of  $m$  linearly independent vectors of  $R^m$  is called a **basis** in  $R^m$ : any vector in  $R^m$  can be expressed as a linear combination of the basis vectors.
- Example: The column vectors of the matrix are not linearly independent

$$[\mathbf{v}_1 \ \mathbf{v}_2 \ \mathbf{v}_3 \ \mathbf{v}_4] = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}$$

$$\alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \alpha_3 \mathbf{v}_3 + \alpha_4 \mathbf{v}_4 = 0$$

holds for  $\alpha_1 = \alpha_3 = 1, \alpha_2 = \alpha_4 = -1$ .



## Part 1: Rank of a matrix

---

- The **rank** of a matrix is the maximum number of linearly independent column vectors.
- A square matrix  $A \in R^{n \times n}$  with rank  $n$  is called **nonsingular**, and it has an inverse  $A^{-1}$  satisfying  $AA^{-1} = A^{-1}A = I$ .
- The (outer product) matrix  $\mathbf{xy}^T$  has rank 1:
- All columns are linearly dependent (and so are all the rows).
- Example

The  $4 \times 4$  matrix

$$\begin{pmatrix} 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}$$

has rank 3.

## Part 1: Example

---

- Consider a  $m \times n$  term-document matrix  $\mathbf{A} = [a_1, a_2 \cdots a_n]$ , where  $a_i \in R^m$  are the documents.
- If  $\mathbf{A}$  has rank 3, then all the documents can be expressed as a linear combination of only three vectors  $\mathbf{v}_1, \mathbf{v}_2$  and  $\mathbf{v}_3 \in R^m$ :

$$a_j = w_{1j}\mathbf{v}_1 + w_{2j}\mathbf{v}_2 + w_{3j}\mathbf{v}_3, \quad j = 1, \dots, n.$$

- The term-document matrix can be written as

$$\mathbf{A} = \mathbf{V}\mathbf{W}$$

$$\text{where } \mathbf{V} = (\mathbf{v}_1 \mathbf{v}_2 \mathbf{v}_3) \in R^{m \times 3} \text{ and } \mathbf{W} = (w_{ij}) \in R^{3 \times n}.$$

# Part 1: Orthogonality

---

- Two vectors  $\mathbf{x}$  and  $\mathbf{y}$  are orthogonal, if  $\mathbf{x}^T \mathbf{y} = 0$ .
- Let  $q_j, j = 1, \dots, n$  be orthogonal, i.e.  $q_i^T q_j = 0, i \neq j$ . Then, they are linearly independent.
- Let the set of orthogonal vectors  $q_j, j = 1, \dots, n$  in  $R^n$  be normalized,  $\|q_j\| = 1$ . Then they are **orthonormal**, and constitute an **orthonormal basis** in  $R^n$ .
- A matrix  $R^{m \times m} \ni Q = [q_1 q_2 \dots q_m]$  with orthonormal column is called an **orthogonal matrix**.
- Example

$$\begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix} \text{ to be } \mathbf{v}_1 = \begin{pmatrix} 0.52 \\ 0.85 \end{pmatrix}, \quad \mathbf{v}_2 = \begin{pmatrix} 0.85 \\ -0.52 \end{pmatrix}.$$

The vectors  $\mathbf{v}_1$  and  $\mathbf{v}_2$  are orthogonal:

$$\mathbf{v}_1^T \mathbf{v}_2 = 0.52 \cdot 0.85 + 0.85 \cdot (-0.52) = 0.$$

This is no coincidence!

# Part 1: Why we like orthogonal matrices

---

- An orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{m \times m}$  has rank  $m$  (since its columns are linearly independent).
- $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ .  $\mathbf{Q} \mathbf{Q}^T = \mathbf{I}$ .
- The inverse of an orthogonal matrix  $\mathbf{Q}$  is  $\mathbf{Q}^{-1} = \mathbf{Q}^T$ .
- The Euclidean length of a vector is invariant under an orthogonal transformation  $\mathbf{Q}$ :  $\|\mathbf{Q}\mathbf{x}\|^2 = (\mathbf{Q}\mathbf{x})^T \mathbf{Q}\mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|^2$ .
- The product of two orthogonal matrices  $\mathbf{Q}$  and  $\mathbf{P}$  is orthogonal:  
 $\mathbf{X}^T \mathbf{X} = (\mathbf{P}\mathbf{Q})^T (\mathbf{P}\mathbf{Q}) = \mathbf{Q}^T \mathbf{P}^T \mathbf{P} \mathbf{Q} = \mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ .

# Part 1: Eigenvalues and eigenvectors of a symmetric matrix

- Let  $\mathbf{A}$  be a  $n \times n$  matrix. The vector  $\mathbf{v} \neq 0$  that satisfies  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$  for some scalar  $\lambda$  is called the eigenvector of  $\mathbf{A}$  and  $\lambda$  is the eigenvalue corresponding to the eigenvector  $\mathbf{v}$ .
- The eigenvectors of a symmetric matrix are mutually orthogonal and its eigenvalues are real.
- A symmetric matrix  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$  where the columns of  $\mathbf{U}$  are the eigenvectors of  $\mathbf{A}$  and  $\mathbf{\Lambda}$  is a diagonal matrix, the diagonal elements of which are the corresponding eigenvalues of  $\mathbf{A}$ . Note, that  $\mathbf{U}$  is orthogonal. This is called the eigendecomposition or symmetric Schur decomposition of  $\mathbf{A}$ .

$$\mathbf{A}\mathbf{v} = \begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix} \mathbf{v} = \lambda\mathbf{v}.$$

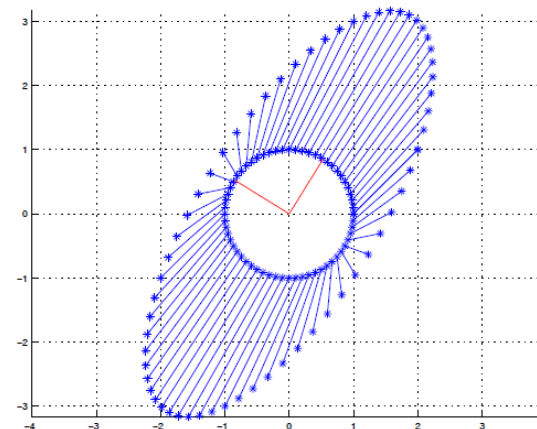
$$\det(\mathbf{A} - \lambda\mathbf{I}) = \begin{vmatrix} 2 - \lambda & 1 \\ 1 & 3 - \lambda \end{vmatrix} = (2 - \lambda)(3 - \lambda) - 1 = 0$$

$$\lambda_1 = 3.62$$

$$\lambda_2 = 1.38$$

$$\mathbf{v}_1 = \begin{pmatrix} 0.52 \\ 0.85 \end{pmatrix}$$

$$\mathbf{v}_2 = \begin{pmatrix} 0.85 \\ -0.52 \end{pmatrix}$$



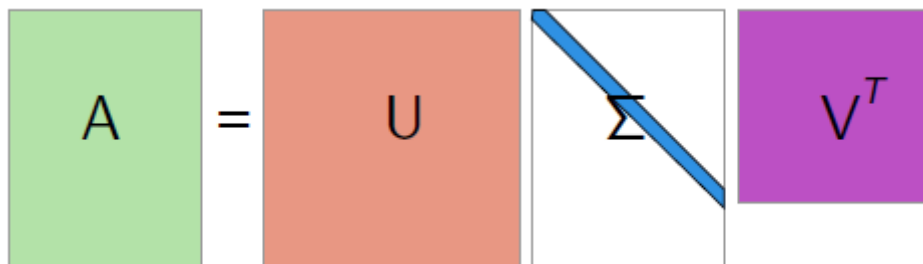
## Part 1: What if the matrix is not symmetric?

---

- If  $\mathbf{A}$  is symmetric, we can write  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ , where  $\mathbf{U}$  and  $\mathbf{\Lambda}$  contain the eigenvectors and corresponding eigenvalues of  $\mathbf{A}$ .
- What if  $\mathbf{A}$  Definitely not symmetric!!!
- Then we can use the singular value decomposition (SVD)

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T,$$

where  $\mathbf{U}$  and  $\mathbf{V}$  are orthogonal, and  $\mathbf{\Sigma}$  is diagonal.



## Part 1: Example

customer \ day	Wed	Thu	Fri	Sat	Sun
ABC Inc.	1	1	1	0	0
CDE Co.	2	2	2	0	0
FGH Ltd.	1	1	1	0	0
NOP Inc.	5	5	5	0	0
Smith	0	0	0	2	2
Brown	0	0	0	3	3
Johnson	0	0	0	1	1

$$\begin{pmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{pmatrix} \times \begin{pmatrix} 9.64 & 0 \\ 0 & 5.29 \end{pmatrix} \times \begin{pmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{pmatrix}$$

# Part 1: Lessons learned

---

- Vector is
  - A 1D array of numbers, entity with magnitude and direction.
- Matrix is
  - A rectangular form representing observations.
- Orthogonal matrix
  - Has full rank
  - Has an inverse  $\mathbf{Q}^T$
  - Is invariant
- Symmetric matrix has eigendecomposition as in  $\mathbf{A}\mathbf{v}=\lambda\mathbf{v}$ 
  - Eigenvector  $\mathbf{v}$ , eigenvalue  $\lambda$
- Unsymmetric matrix has singular value and singular vectors as in  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ 
  - Left singular vector  $\mathbf{U}$
  - Right singular vector  $\mathbf{V}$
  - Singular value  $\Sigma$



# Content

---

## Part 0. Preliminary (30 mins)

- 0-1. Modern Data
- 0-2. About the Tutorial
- 0-3. References
- 0-4. Suggested Readings

## Part 1. Linear Algebra Review (1 hour)

- 1-1. Basics of Linear Algebra
- 1-2. Vector Notations
- 1-3. Matrix Notations

## **Part 2. Matrix Algorithms (1 hour)**

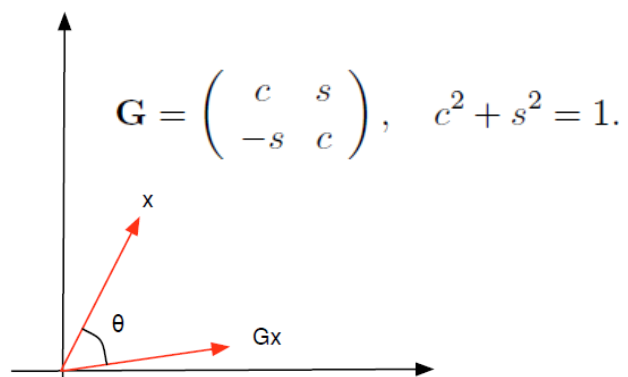
- 2-1. Singular Value Decomposition (SVD)
- 2-2. Non-negative Matrix Factorization (NMF)
- 2-3. Semi-Discrete Decomposition (SDD)
- 2-4. Independent Component Analysis (ICA)

## Part 3. Applications of Matrix Algorithms (30 mins)

- 3-1. Examples of Matrix Algorithms Applied on Data Mining Applications

## Part 2: How to compute matrix decompositions?

- We have had a brief glimpse at eigenvalue decomposition and singular value decomposition.
- Before taking a closer look: how can we compute these?
- Answer 1: use LAPACK, MATLAB, Mathematica, ...they are implemented everywhere!
- Answer 2: we need more linear algebra tools!!, Givens rotation, Householder transforms



$$P = I - \frac{2}{v^T v} v v^T;$$

$$Px = x - \frac{2v^T x}{v^T v} v = x - v = y,$$

$$P_2 A^{(1)} = P_2 \begin{pmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \end{pmatrix} = \begin{pmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \end{pmatrix} =: A^{(2)}.$$

## Part 2: Matrix decompositions

---

- We wish to decompose the matrix  $\mathbf{A}$  by writing it as a product of two or more matrices

$$\mathbf{A}_{m \times n} = \mathbf{B}_{m \times k} \mathbf{C}_{k \times n}, \quad \mathbf{A}_{m \times n} = \mathbf{B}_{m \times k} \mathbf{C}_{k \times r} \mathbf{D}_{r \times n}$$

- This is done in such a way that the **right side of the equation yields some useful information or insight** to the nature of the data matrix  $\mathbf{A}$ .
- Or is in other ways useful for solving the problem at hand.
- There are numerous examples of useful matrix decompositions
  - **Eigendecomposition**:  $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ ,  $\mathbf{A}$  symmetric,  $\mathbf{U}$  and  $\mathbf{\Lambda}$  eigenvectors and eigenvalues.
  - **Singular value decomposition**:  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ ,  $\mathbf{U}, \mathbf{V}$  orthogonal,  $\mathbf{\Sigma}$  diagonal.
  - Matrix factorization is the same thing as matrix decomposition (e.g. **NMF** = nonnegative matrix factorization,  $\mathbf{V} = \mathbf{W} \mathbf{H}$ , all elements nonnegative.)

## Part 2: How to compute these?

---

- Roughly

(1) Manipulate **A** by multiplying it by intelligently chosen, fairly simple (orthogonal) matrices from both sides:

$$\mathbf{V}_k \dots \mathbf{V}_2 \mathbf{V}_1 \mathbf{A} \mathbf{W}_1 \dots \mathbf{W}_s = \mathbf{B}, \quad \text{until } \mathbf{B} \text{ is "nice".}$$

(2) Denote  $\mathbf{V} = \mathbf{V}_k \dots \mathbf{V}_2 \mathbf{V}_1$ ,  $\mathbf{W} = \mathbf{W}_1 \dots \mathbf{W}_s$ . Now  $\mathbf{A} = \mathbf{V} \mathbf{B} \mathbf{W}^T$ .

- But how to choose  $\mathbf{V}_1, \dots, \mathbf{V}_k$ ,  $\mathbf{W}_1, \dots, \mathbf{W}_s$ ?
- Needed: linear algebra tools for transforming matrices in an orderly fashion: Givens!

## Part 2: QR decomposition

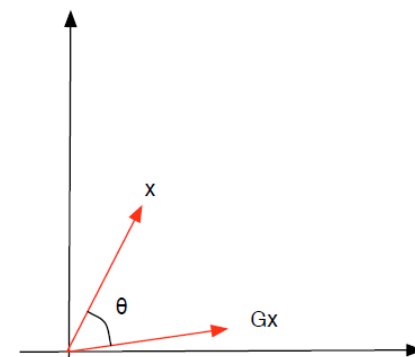
- Any matrix  $\mathbf{A} \in R^{m \times n}, m \geq n$ , can be transformed to upper triangular form by an orthogonal matrix

$$\mathbf{A} = \mathbf{Q} \begin{pmatrix} \mathbf{R} \\ 0 \end{pmatrix}$$

- If the columns of  $\mathbf{A}$  are linearly independent, then  $\mathbf{R}$  is non-singular.
- Needed: linear algebra tools for transforming matrices in an orderly fashion: Givens!

Let  $\mathbf{x}$  be a vector. The parameters  $c$  and  $s$ ,  $c^2 + s^2 = 1$ , can be chosen so that multiplication of  $\mathbf{x}$  by

$$\mathbf{G} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & c & 0 & s \\ 0 & 0 & 1 & 0 \\ 0 & -s & 0 & c \end{pmatrix}$$



will zero the element 4 in vector  $\mathbf{x}$  by a rotation in plane (2,4). How?

## Part 2: What is QR good for?

---

- It will come up when we discuss the computation of some matrix decompositions
- It can also be used for solving the **least squares problem**.
- Other applications exist.

### "Skinny" QR decomposition

Partition  $Q = (Q_1 \ Q_2)$ , where  $Q_1 \in \mathbb{R}^{m \times n}$ :

$$A = (Q_1 \ Q_2) \begin{pmatrix} R_1 \\ 0 \end{pmatrix} = Q_1 R_1.$$

## Part 2: Least squares for linear regression

---

- Consider problems of the following form: given a number of measurements  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]$  and an outcome  $\mathbf{y}$ , build a linear model

$$\hat{\mathbf{y}} = b_0 + \sum_{j=1}^n \mathbf{x}_j b_j,$$

- which uses the measurements  $\mathbf{X}$  to predict the outcome  $\mathbf{y}$ .
- $\mathbf{X}$  = clinical measures of patients,  $\mathbf{y}$  = level of cancer specific antigen
- $\mathbf{X}$  = atmospheric measurements of each day,  $\mathbf{y}$  = occurrence of spontaneous particle formation.
- The model  $\hat{\mathbf{y}} = b_0 + \sum_{j=1}^n \mathbf{x}_j b_j$  can be written in the form
  - $\hat{\mathbf{y}} = \mathbf{X}^T \mathbf{b}$ ,  $\mathbf{b} = (b_1 \dots b_n b_0)^T$ .
    - To do this, one must append  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n \mathbf{x}_{n+1}]$ , where  $\mathbf{x}_{n+1}$  is a vector of ones.
- Fitting a linear model to data is usually done using the method of least squares.
- Note.  $\mathbf{X}$  need not be a square matrix !!!

## Part 2: Example

- We have measurement data

$x$	1	2	3	4	5
$y$	7.97	10.2	14.2	16.0	21.2

We wish to find  $\alpha$  and  $\beta$  such that  $\alpha x + \beta = y$ . Thus

$$\alpha + \beta = 7.97$$

$$2\alpha + \beta = 10.2$$

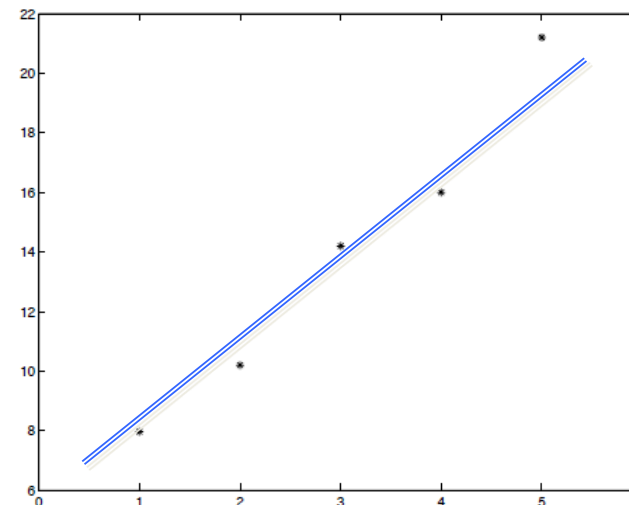
$$3\alpha + \beta = 14.2$$

$$4\alpha + \beta = 16.0$$

$$5\alpha + \beta = 21.2$$

In matrix form:

$$\begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \\ 4 & 1 \\ 5 & 1 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} 7.97 \\ 10.2 \\ 14.2 \\ 16.0 \\ 21.2 \end{pmatrix}$$





## Part 2: Example

---

- Make the residual vector  
 $\mathbf{r} = \mathbf{b} - \mathbf{Ax}$  as small as possible. But how?
- Make  $\mathbf{r}$  orthogonal to the columns of  $\mathbf{A}$ :  
 $\mathbf{r}^T(\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n) = \mathbf{r}^T\mathbf{A} = 0.$
- Now write  $\mathbf{r} = \mathbf{b} - \mathbf{Ax}$  to get the normal equations
  - $\mathbf{A}^T\mathbf{Ax} = \mathbf{A}^T\mathbf{b}$ ; solve for  $\mathbf{x}$
- If the column vectors of  $\mathbf{A}$  are linearly independent, then the **normal equations**  $\mathbf{A}^T\mathbf{Ax} = \mathbf{A}^T\mathbf{b}$  are non-singular and have a unique solution.
- But the normal equations have two significant drawbacks:
  - Forming  $\mathbf{A}^T\mathbf{A}$  leads to loss of information
  - The condition number of  $\mathbf{A}^T\mathbf{A}$  is the square of that of  $\mathbf{A}$

## Part 2: Solving least squares problem using QR

---

- $||\mathbf{r}||^2 = ||\mathbf{b} - \mathbf{Ax}||^2 = ||\mathbf{b} - \mathbf{Q}\begin{pmatrix} \mathbf{R} \\ 0 \end{pmatrix}\mathbf{x}||^2 = ||\mathbf{Q}(\mathbf{Q}^T\mathbf{b} - \begin{pmatrix} \mathbf{R} \\ 0 \end{pmatrix}\mathbf{x})||^2 = ||\mathbf{Q}^T\mathbf{b} - \begin{pmatrix} \mathbf{R} \\ 0 \end{pmatrix}\mathbf{x}||^2$
- Partition  $\mathbf{Q}=(\mathbf{Q}_1 \ \mathbf{Q}_2)$ , where  $\mathbf{Q}_1 \in \mathbf{R}^{m \times n}$ , and denote
  - $\mathbf{Q}^T\mathbf{b} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix} := \begin{pmatrix} \mathbf{Q}_1^T\mathbf{b} \\ \mathbf{Q}_2^T\mathbf{b} \end{pmatrix}$
- Then
  - $||\mathbf{r}||^2 = \left\| \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix} - \begin{pmatrix} \mathbf{R} \\ 0 \end{pmatrix}\mathbf{x} \right\|^2 = ||\mathbf{b}_1 - \mathbf{Rx}||^2 + ||\mathbf{b}_2||^2$
- Minimize  $||\mathbf{r}||$  by making the first term equal to zero: i.e. solve  $\mathbf{Rx}=\mathbf{b}_1$
- Theorem. Let  $\mathbf{A} \in \mathbf{R}^{m \times n}$  have full column rank and a thin QR-decomposition  $\mathbf{A} = \mathbf{Q}_1\mathbf{R}$ . Then the least squares problem

$$\min_x ||\mathbf{Ax} - \mathbf{b}||_2$$

has the unique solution  $\mathbf{x} = \mathbf{R}^{-1}\mathbf{Q}_1^T\mathbf{b}$ .

## Part 2: Updating the solution of the LS problem

---

Assume we have reduced the matrix and the right hand side

$$(\mathbf{A} \quad \mathbf{b}) \rightarrow \mathbf{Q}^T (\mathbf{A} \quad \mathbf{b}) = \begin{pmatrix} \mathbf{R} & \mathbf{Q}_1^T \mathbf{b} \\ \mathbf{0} & \mathbf{Q}_2^T \mathbf{b} \end{pmatrix}.$$

From this the solution of the LS problem is readily available.

Assume we have not saved  $\mathbf{Q}$ .

We then get a new observation  $(\mathbf{a} \quad b)$ ,  $\mathbf{a} \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ .

Do we have to recompute the whole solution?

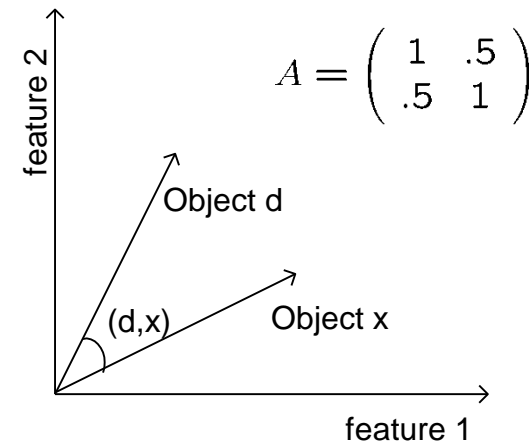
No. Instead, write the reduction on the previous slide as

$$\begin{pmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{a}^T & b \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{Q}^T & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \begin{pmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{a}^T & b \end{pmatrix} = \begin{pmatrix} \mathbf{R} & \mathbf{Q}_1^T \mathbf{b} \\ \mathbf{0} & \mathbf{Q}_2^T \mathbf{b} \\ \mathbf{a}^T & b \end{pmatrix}.$$

And reduce this to triangular form using plane rotations.

## Part 2: The Singular Value Decomposition (SVD)

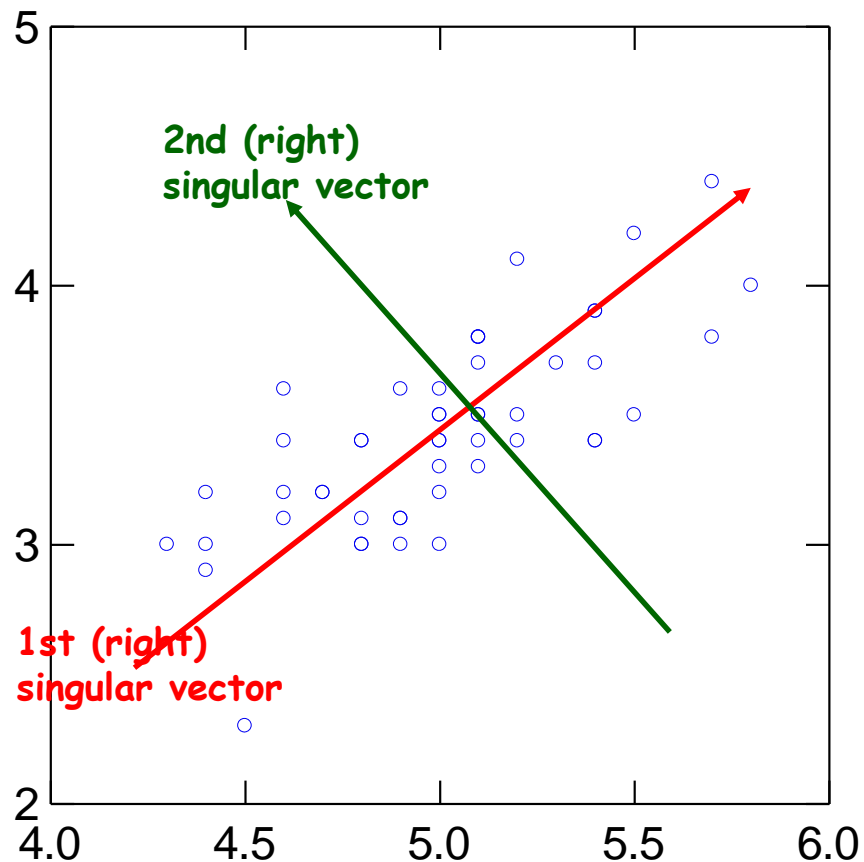
- Recall: data matrices have  $m$  rows (one for each object) and  $n$  columns (one for each feature).
- Matrix rows: points (vectors) in a Euclidean space,
  - e.g., given 2 objects ( $x$  &  $d$ ), each described with respect to two features, we get a 2-by-2 matrix.
- Two objects are “close” if the angle between their corresponding vectors is small.



$$A = U \Sigma V^T$$

- I.e. every  $A$  has decomposition  $A = U \Sigma V^T$ 
  - The **singular value decomposition** (SVD)
- The values  $\sigma_i$  are the **singular values** of  $A$
- Columns of  $U$  are the **left singular vectors** and columns of  $V$  the **right singular vectors** of  $A$

## Part 2: Singular Value Decomposition (SVD) intuition



Let the **blue circles** represent  $m$  data points in a 2-D Euclidean space.

Then, the SVD of the  $m$ -by-2 matrix of the data will return ...

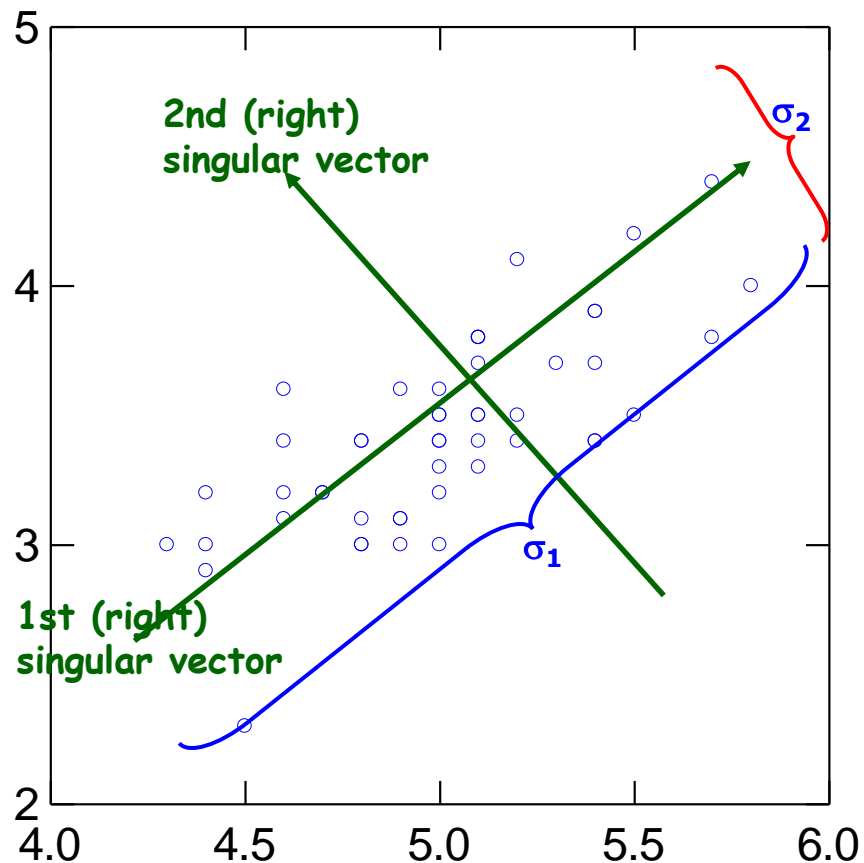
- 1st (right) singular vector:

direction of maximal variance,

- 2nd (right) singular vector:

direction of maximal variance, after **removing the projection of the data** along the first singular vector.

## Part 2: Singular Values



$\sigma_1$ : measures how much of the data variance is explained by the first singular vector.

$\sigma_2$ : measures how much of the data variance is explained by the second singular vector.

## Part 2: SVD: formal definition

---

$$\begin{pmatrix} A \\ m \times n \end{pmatrix} = \begin{pmatrix} U \\ m \times \rho \end{pmatrix} \cdot \begin{pmatrix} \text{ } & \text{ } & 0 \\ \text{ } & \Sigma & \text{ } \\ 0 & \text{ } & \text{ } \end{pmatrix} \cdot \begin{pmatrix} V \\ \rho \times n \end{pmatrix}^T$$

$\rho$ : rank of **A**

**U** (**V**): orthogonal matrix containing the left (right) singular vectors of **A**.

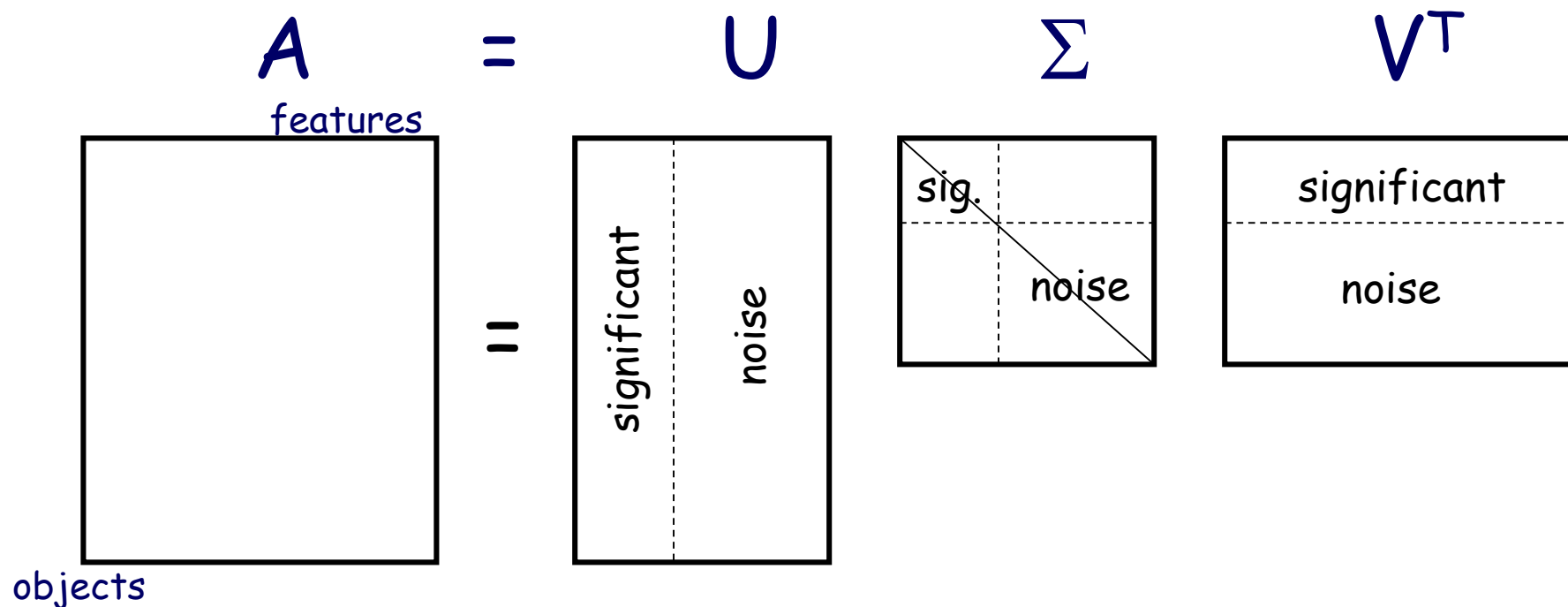
$\Sigma$ : diagonal matrix containing the singular values of **A**.

Let  $\sigma_1, \sigma_2, \dots, \sigma_\rho$  be the entries of  $\Sigma$ .

Exact computation of the SVD takes  $O(\min\{mn^2, m^2n\})$  time.

The top  $k$  left/right singular vectors/values can be computed faster using Lanczos/Arnoldi methods.

## Part 2: Rank-k approximations via the SVD



$$\begin{pmatrix} A_k \\ m \times n \end{pmatrix} = \begin{pmatrix} U_k \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} \Sigma_k \\ k \times k \end{pmatrix} \cdot \begin{pmatrix} V_k^T \\ k \times n \end{pmatrix}$$

$U_k$  ( $V_k$ ): orthogonal matrix containing the top  $k$  left (right) singular vectors of  $A$ .  
 $\Sigma_k$ : diagonal matrix containing the top  $k$  singular values of  $A$ .



## Part 2: SVD is useful for

---

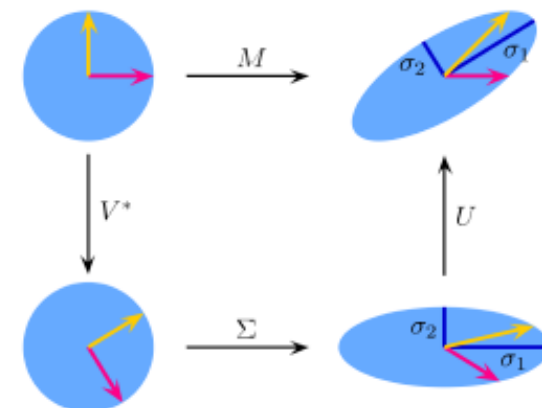
- Compression
- Noise reduction
- Finding “concepts” or “topics” (text mining/LSI)
- Data exploration and visualizing data (e.g. spatial data/PCA)
- Classification (of e.g. handwritten digits)
- SVD appears under different names:
  - Principal component analysis (PCA)
  - Latent semantic indexing (LSI) / Latent Semantic Analysis (LSA)
  - Karhunen-Loeve expansion / Hotelling transform (in image processing)

## Part 2: Interpreting SVD

- The most common way to interpret SVD is to consider the columns of  $\mathbf{U}$  (or  $\mathbf{V}$ )
  - Let  $\mathbf{A}$  be objects-by-attributes and  $\mathbf{U}\Sigma\mathbf{V}^T$  its SVD
  - If two columns have similar values in a row of  $\mathbf{V}^T$ , these attributes are somehow similar (have strong correlation)
  - If two rows have similar values in a column of  $\mathbf{U}$ , these uses are somehow similar

- Geometric interpretation

- SVD shows that every linear mapping  $\mathbf{y}=\mathbf{M}\mathbf{x}$  can be considered as a **series of rotation, stretching, and rotation operations**
- Matrix  $\mathbf{V}^T$  performs the first rotation  $\mathbf{y}_1=\mathbf{V}^T\mathbf{x}$
- Matrix  $\Sigma$  performs the stretching  $\mathbf{y}_2=\Sigma \mathbf{y}_1$
- Matrix  $\mathbf{U}$  performs the second rotation  $\mathbf{y}=\mathbf{U}\mathbf{y}_2$



$$M = U \cdot \Sigma \cdot V^*$$

- Dimension of largest variance
  - The singular vectors give the dimensions of the variance in the data
    - The first singular vector is the dimension of the largest variance
    - The second singular vector is the orthogonal dimension of the second...
- Component interpretations: sums of rank-1 layers
  - The first layer explains the most
  - The second corrects that by adding and removing smaller values

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \sum_{i=1}^r \mathbf{A}_i$$

## Part 2: Eigenvalue decomposition vs. SVD

---

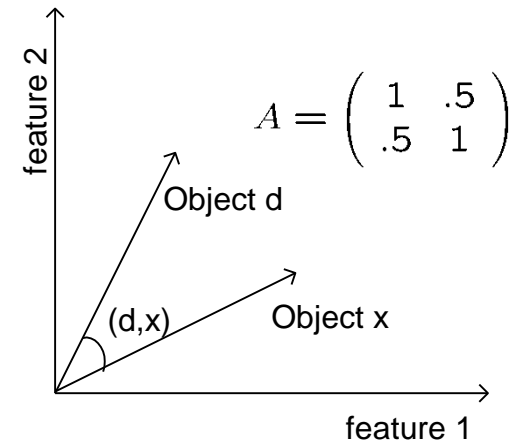
- For symmetric **A** the singular value decomposition is closely related to the eigendecomposition  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ , **U** and **Λ** eigenvectors and eigenvalues.
- Computation of both the eigendecomposition and the **SVD** follow the same pattern
- Computation of the eigenvalue decomposition
  - Use givens to transform **A** into tridiagonal form
  - Use QR iteration with Wilkinson shift  $\mu$  to transform tridiagonal form to diagonal form: Repeat until converged
    - $\mathbf{QR} = \mathbf{T}_k - \mu_k \mathbf{I}$
    - $\mathbf{T}_{k+1} = \mathbf{RQ} + \mu_k \mathbf{I}$
  - Once an eigenvalue is found, forget it, reduce (deflate) the problem and repeat the same step.
- Computing the **SVD**
  - Use Householder transformations to transform **A** into bidiagonal form **B**. Now  $\mathbf{B}^T\mathbf{B}$  is tridiagonal
  - Use **QR** iteration with Wilkinson shift to transform tridiagonal **B** to diagonal form (without forming  $\mathbf{B}^T\mathbf{B}$  implicitly!)

## Part 2: PCA and SVD

Principal Components Analysis (PCA) essentially amounts to the computation of the Singular Value Decomposition (SVD) of a covariance matrix. SVD is the algorithmic tool behind MultiDimensional Scaling (MDS) and Factor Analysis.

SVD is “the Rolls-Royce and the Swiss Army Knife of Numerical Linear Algebra.”\*

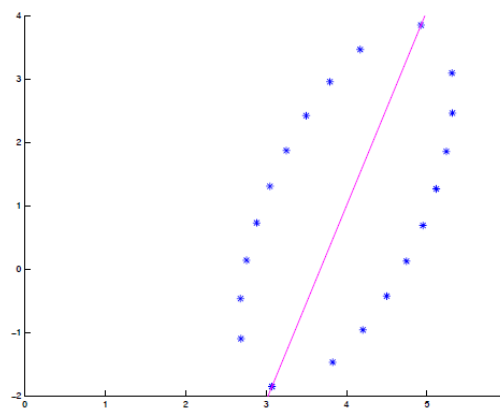
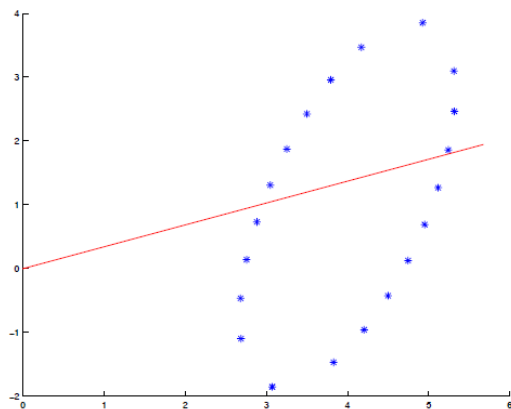
\*Dianne O’Leary, MMDS ’06



- **PCA is SVD done on centered data.**
- PCA looks for such a direction that the data projected onto it has maximal variance.
- When found, PCA continues by seeking the next direction, which is orthogonal to all the previously found directions, and which explains as much of the remaining variance in the data as possible.
- Principal components are uncorrelated.
- PCA is useful for
  - Data exploration, Visualizing data, Compressing data, Outlier detection, Ratio rules.

## Part 2: How to compute the PCA

- Data matrix  $\mathbf{A}$ , rows=data points, columns=variables (attributes, parameters).
1. Center the data by subtracting the mean of each column.
  2. Compute the SVD of the centered matrix  $\hat{\mathbf{A}}$  (or the  $k$  first singular values and vectors)  
$$\hat{\mathbf{A}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$
  3. The principal components are the columns of  $\mathbf{V}$ , the coordinates of the data in the basis defined by the principal components are  $\mathbf{U}\mathbf{\Sigma}$ .



## Part 2: LSI: $A_k$ for document-term matrices

### Latent Semantic Indexing (LSI)

Replace  $A$  by  $A_k$ ; apply clustering/classification algorithms on  $A_k$ .

$$\begin{array}{c} m \\ \text{documents} \end{array} \left( \begin{array}{c} n \text{ terms (words)} \\ \\ \\ A \\ \\ A_{ij} = \text{frequency of } j\text{-th} \\ \text{term in } i\text{-th document} \end{array} \right)$$

- The data: a term-document matrix  $A$ 
  - The values are (weighted) term frequencies
  - Typically TF/IDF values
- The truncated SVD  $A_k = U_k \Sigma_k V_k^T$  of  $A$  is computed
  - Matrix  $U_k$  associates documents to topics
  - Matrix  $V_k$  associates topics to terms
  - If two rows of  $U_k$  are similar, the corresponding documents “talk about the same things”
- A query  $q$  can be answered by considering its term vector  $q$ 
  - $q$  is projected to  $q_k = qV\Sigma^{-1}$
  - $q_k$  is compared to rows of  $U$  and most similar rows are returned

### Pros

- Less storage for small  $k$ .  
 $O(km+kn)$  vs.  $O(mn)$
- Improved performance.  
Documents are represented in a “concept” space.

### Cons

- $A_k$  destroys sparsity.
- Interpretation is difficult.
- Choosing a good  $k$  is tough.

## Part 2: $A_k$ and k-means clustering

---

### k-means clustering

A standard objective function that measures cluster quality.

(Often denotes an iterative algorithm that attempts to optimize the  $k$ -means objective function.)

### k-means objective

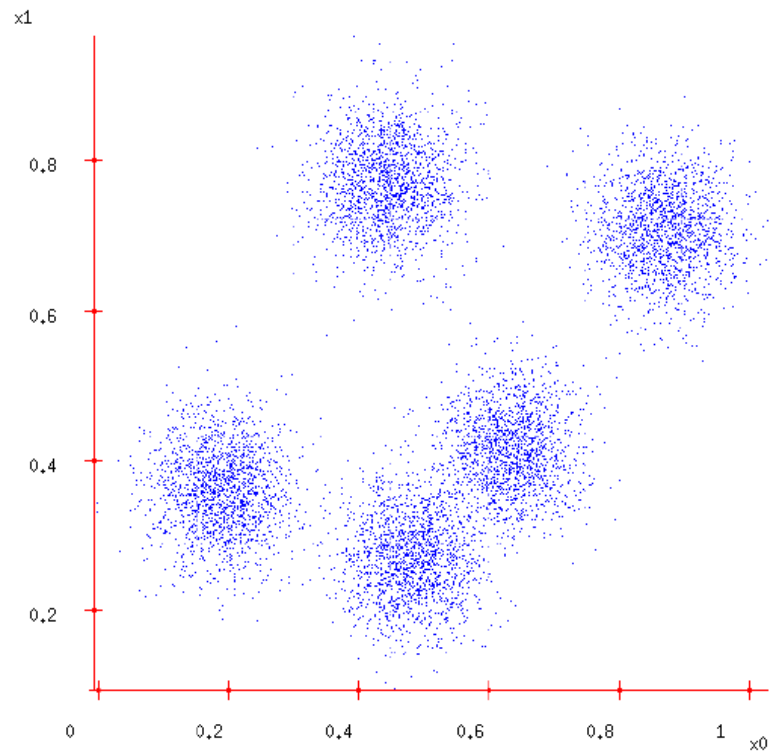
**Input:** set of  $m$  points in  $R^n$ , positive integer  $k$

**Output:** a partition of the  $m$  points to  $k$  clusters

Partition the  $m$  points to  $k$  clusters in order to minimize the sum of the squared Euclidean distances from each point to its cluster centroid.

## Part 2: k-means, cont' d

---

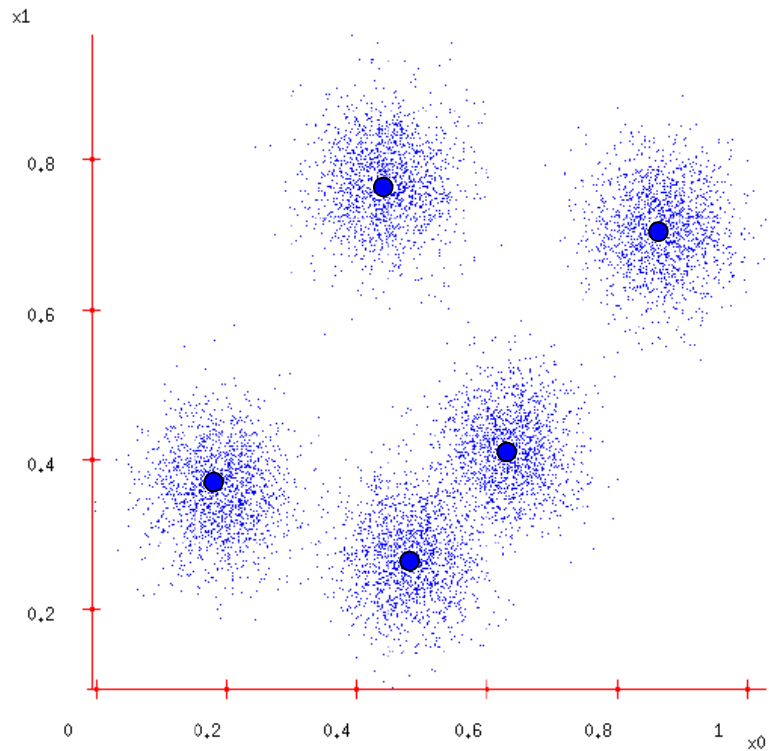


We seek to split the input points in 5 clusters.



## Part 2: k-means, cont' d

---



We seek to split the input points in 5 clusters.

The cluster centroid is the “average” of all the points in the cluster.

## Part 2: k-means: a matrix formulation

---

Let  $A$  be the  $m$ -by- $n$  matrix representing  $m$  points in  $\mathbb{R}^n$ . Then, we seek to

$$\min_{X \in \mathbb{R}^{m \times k}} \|A\|_F^2 - \|X^T A\|_F^2 \quad \text{or} \quad \max_{X \in \mathbb{R}^{m \times k}} \|X^T A\|_F^2$$

$X$  is a **special "cluster membership"** matrix:  $X_{ij}$  denotes if the  $i$ -th point belongs to the  $j$ -th cluster.

clusters

points

$\left( \begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right)$

$X$

$m \times k$

- Columns of  $X$  are **normalized** to have unit length.  
(We divide each column by the square root of the number of points in the cluster.)
- **Every row of  $X$  has at most one non-zero element.**  
(Each element belongs to at most one cluster.)
- $X$  is an orthogonal matrix, i.e.,  $X^T X = I$ .

## Part 2: SVD and $k$ -means

---

If we only require that  $X$  is an orthogonal matrix and remove the condition on the number of non-zero entries per row of  $X$ , then

$$\min_{X \in \mathbb{R}^{m \times k}} \|A\|_F^2 - \|X^T A\|_F^2 \quad \text{or} \quad \max_{X \in \mathbb{R}^{m \times k}} \|X^T A\|_F^2$$

is easy to minimize! The solution is  $X = U_k$ .

### Using SVD to solve $k$ -means

- We can get a 2-approximation algorithm for  $k$ -means.  
(Drineas, Frieze, Kannan, Vempala, and Vinay '99, '04)
- We can get heuristic schemes to assign points to clusters.  
(Zha, He, Ding, Simon, and Gu '01)
- There exist PTAS (based on random projections) for the  $k$ -means problem.  
(Ostrovsky and Rabani '00, '02)
- Deeper connections between SVD and clustering in Kannan, Vempala, and Vetta '00, '04.

### Hypertext Induced Topic Selection (HITS)

A link analysis algorithm that rates Web pages for their authority and hub scores.

**Authority score:** an estimate of the value of the content of the page.

**Hub score:** an estimate of the value of the links from this page to other pages.

These values can be used to rank Web search results.

**Authority:** a page that is pointed to by many pages with high hub scores.

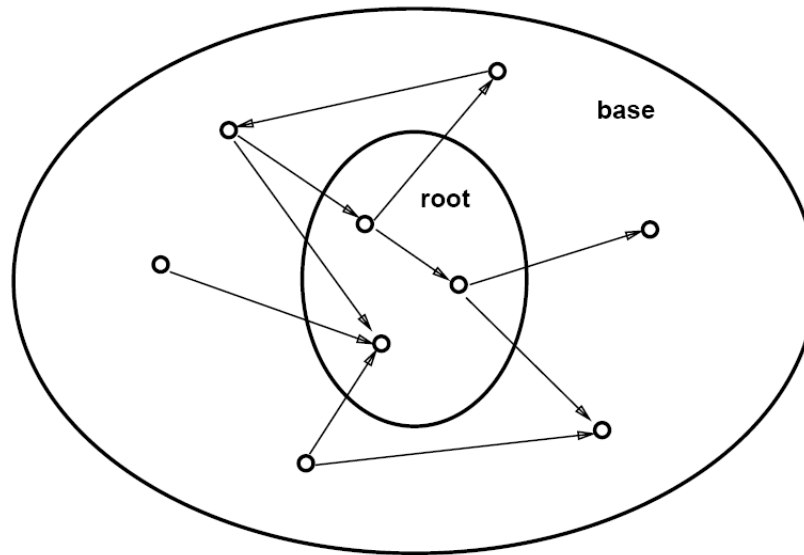
**Hub:** a page pointing to many pages that are good authorities.

Recursive definition; notice that each node has two scores.

## Part 2: $A_k$ and Kleinberg's HITS algorithm (Kleinberg '98, '99)

### Phase 1:

Given a query term (e.g., "jaguar"), find all pages containing the query term (**root set**).  
Expand the resulting graph by one move forward and backward (**base set**).



## Part 2: $A_k$ and Kleinberg's HITS algorithm (Kleinberg '98, '99)

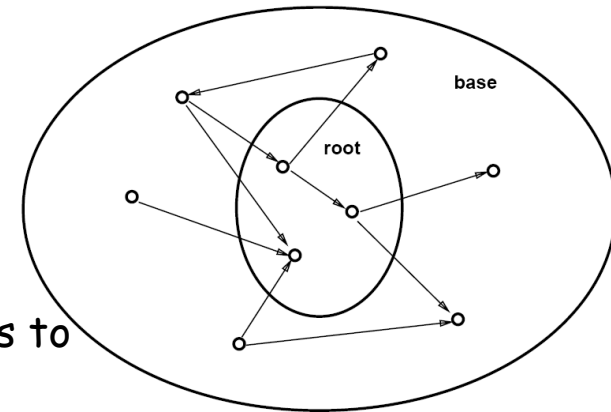
### Phase 2:

Let  $A$  be the adjacency matrix of the (directed) graph of the **base set**.

Let  $h, a \in \mathbb{R}^n$  be the vectors of hub (authority) scores.

Then,  $h = Aa$  and  $a = A^T h$

$$h = AA^T h \text{ and } a = A^T Aa.$$



Thus, the top left (right) singular vector of  $A$  corresponds to hub (authority) scores.

### What about the rest?

They provide a natural way to extract **additional densely linked collections of hubs and authorities** from the base set.

See the "jaguar" example in Kleinberg '99.

## Part 2: SVD Lessons learned

---

- SVD is the Swiss Army knife of (numerical) linear algebra  
→ ranks, kernels, norms, ...
- SVD is also very useful in data analysis  
→ noise removal, visualization, dimensionality reduction, ...
- Selecting the correct rank for the truncated SVD is still a problem

## Part 2: Non-Negative Datasets

---

- Some datasets are intrinsically non-negative
  - Counters (e.g., no. occurrences of each word in a text document)
  - Quantities (e.g., amount of each ingredient in a chemical experiment)
  - Intensities (e.g., intensity of each color in an image)
- The corresponding data matrix **D** has only non-negative values
  - Decompositions such as **SVD** and **SDD** may involve negative values in factors and components
  - Negative values describe the absence of something
  - Often no natural interpretation

Can we find a decomposition that is more natural to non-negative data?



## Part 2: Example (SVD)

Consider the following "bridge" matrix and its truncated SVD:

$$\begin{array}{ccccc}
 \begin{array}{|c|c|c|c|c|}
 \hline
 1 & 1 & 1 & 1 & 1 \\
 \hline
 0 & 1 & 0 & 1 & 0 \\
 \hline
 0 & 1 & 0 & 1 & 0 \\
 \hline
 \end{array} & = & \begin{array}{|c|c|}
 \hline
 0.8 & 0.6 \\
 \hline
 0.5 & -0.5 \\
 \hline
 0.5 & -0.5 \\
 \hline
 \end{array} & \begin{array}{|c|c|}
 \hline
 1.5 & 0 \\
 \hline
 0 & 1.3 \\
 \hline
 \end{array} & \begin{array}{|c|c|c|c|c|}
 \hline
 0.3 & 0.6 & 0.3 & 0.6 & 0.3 \\
 \hline
 0.5 & -0.3 & 0.5 & -0.3 & 0.5 \\
 \hline
 \end{array} \\
 \mathbf{D} & & \mathbf{U} & \Sigma & \mathbf{V}^T
 \end{array}$$

Here are the corresponding components:

$$\begin{array}{ccccc}
 \begin{array}{|c|c|c|c|c|}
 \hline
 1 & 1 & 1 & 1 & 1 \\
 \hline
 0 & 1 & 0 & 1 & 0 \\
 \hline
 0 & 1 & 0 & 1 & 0 \\
 \hline
 \end{array} & = & \begin{array}{|c|c|c|c|c|}
 \hline
 0.8 & 1.3 & 0.6 & 1.3 & 0.8 \\
 \hline
 0.3 & 0.6 & 0.3 & 0.6 & 0.3 \\
 \hline
 0.3 & 0.6 & 0.3 & 0.6 & 0.3 \\
 \hline
 \end{array} & + & \begin{array}{|c|c|c|c|c|}
 \hline
 0.4 & -0.3 & 0.4 & -0.3 & 0.4 \\
 \hline
 -0.3 & 0.2 & -0.3 & 0.2 & -0.3 \\
 \hline
 -0.3 & 0.2 & -0.3 & 0.2 & -0.3 \\
 \hline
 \end{array} \\
 \mathbf{D} & & \mathbf{U}_{*1} \mathbf{D}_{11} \mathbf{V}_{*1}^T & & \mathbf{U}_{*2} \mathbf{D}_{22} \mathbf{V}_{*2}^T
 \end{array}$$

Negative values make interpretation unnatural or difficult.

## Part 2: Example (NMF)

Consider the following “bridge” matrix and its rank-2 NMF:

$$\mathbf{D} = \mathbf{L} \mathbf{R}$$

1	1	1	1	1
0	1	0	1	0
0	1	0	1	0

1	0
0	1
0	1

1	1	1	1	1
0	1	0	1	0

Here are the corresponding components:

$$\mathbf{D} = \mathbf{L}_{*1} \mathbf{R}_{1*} + \mathbf{L}_{*2} \mathbf{R}_{2*}$$

1	1	1	1	1
0	1	0	1	0
0	1	0	1	0

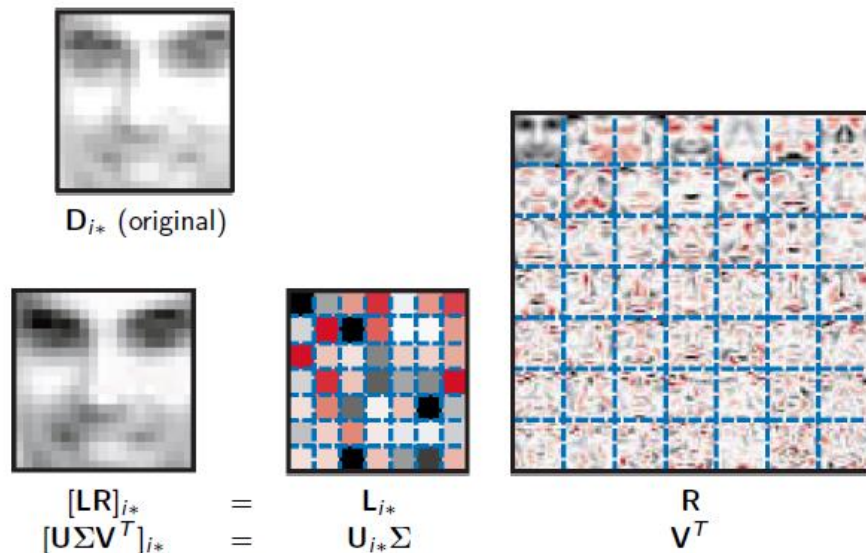
1	1	1	1	1
0	0	0	0	0
0	0	0	0	0

0	0	0	0	0
0	1	0	1	0
0	1	0	1	0

Non-negative matrix decomposition encourage a more natural, part-based representation and (sometimes) sparsity.

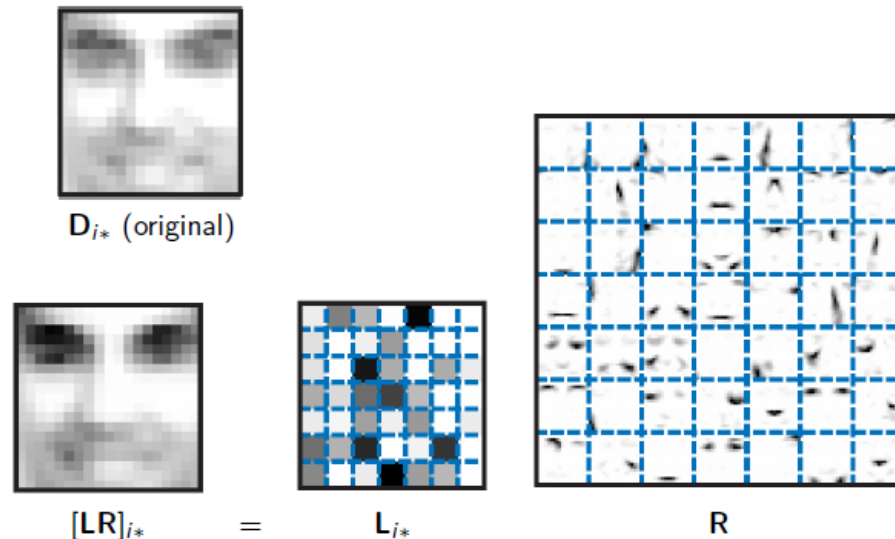
# Part 2: Examples

## Decomposing faces (PCA)



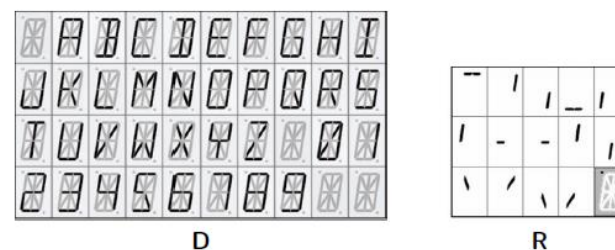
PCA factors are hard to interpret.

## Decomposing faces (NMF)



NMF factors correspond to parts of faces.

## Decomposing digits (NMF)



NMF factors correspond to parts of digits and "background".

## Part 2: Decompositions that respect the data

- Given a nonnegative matrix  $\mathbf{A} \in \mathbf{R}^{m \times n}$ , we wish to express the matrix as a product of two nonnegative matrices  $\mathbf{W} \in \mathbf{R}^{m \times k}$  and  $\mathbf{H} \in \mathbf{R}^{k \times n}$  such that

$$\mathbf{A} \approx \mathbf{WH}$$

- Why require nonnegativity?
  - Nonnegativity is natural in many applications
  - Term-document
  - Market basket
  - Etc.

### Non-negative matrix factorization

(Lee and Seung '00)

Assume that the  $A_{ij}$  are non-negative for all  $i, j$ .

$$\min_{\Phi \in \mathcal{R}^{m \times k}, X \in \mathcal{R}^{k \times n}} \left\| \begin{pmatrix} A \\ m \times n \end{pmatrix} - \begin{pmatrix} \Phi \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} X \\ k \times n \end{pmatrix} \right\|_F^2$$

Constrain  $\Phi$  and  $X$  to have only non-negative entries as well.

This should respect the structure of the data better than  $A_k = U_k \Sigma_k V_k^T$  which introduces a lot of (difficult to interpret) negative entries.

## Part 2: Example

---

3	2	4	0	0	0	0
1	0	1	0	0	0	0
0	1	2	0	0	0	0
0	0	0	1	1	1	1
0	0	0	1	0	1	1
0	0	0	3	2	2	3

Rows: customers, columns: products they buy.

PC=			W =		H' =	
			0	2.3515	0	1.1206
-0.3751	0.4383	-0.8047	0	0.5073	0	0.7992
-0.2728	0.2821	0.4166	0	0.7955	0	1.7347
-0.5718	0.4372	0.4073	0.8760	0	1.1918	0
0.3940	0.4311	0.0413	0.6937	0	0.7532	0
0.2537	0.3341	0.0907	2.4179	0	0.8510	0
0.2883	0.2322	-0.0378			1.1918	0
0.3940	0.4311	0.0413				

Rows of  $W$  are customers, rows of  $H^T$  are products.

Principal components. Each column corresponds to a product.

## Part 2: The Non-negative Matrix Factorization

---

$$\min_{\substack{\Phi \in \mathcal{R}^{m \times k}, X \in \mathcal{R}^{k \times n} \\ \forall i, j, \Phi_{ij}, X_{ij} \geq 0}} \left\| \begin{pmatrix} A \\ m \times n \end{pmatrix} - \begin{pmatrix} \Phi \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} X \\ k \times n \end{pmatrix} \right\|_F^2$$

It has been extensively applied to:

1. Image mining (Lee and Seung '00)
2. Enron email collection (Berry and Brown '05)
3. Other text mining tasks (Berry and Plemmons '04)

Algorithms for NMF:

1. Multiplicative update rules (Lee and Seung '00, Hoyer '02)
2. Gradient descent (Hoyer '04, Berry and Plemmons '04)
3. Alternating least squares (dating back to Paatero '94)

## Part 2: The Non-negative Matrix Factorization

---

### How to compute NMF: Multiplicative algorithm

```
W=rand(m,k);  
H=rand(k,n);  
for i=1:maxiter  
    H=H.*(W'*A)./(W'*W*H+epsilon);  
    W=W.*(A*H')./(W*H*H'+epsilon);  
end
```

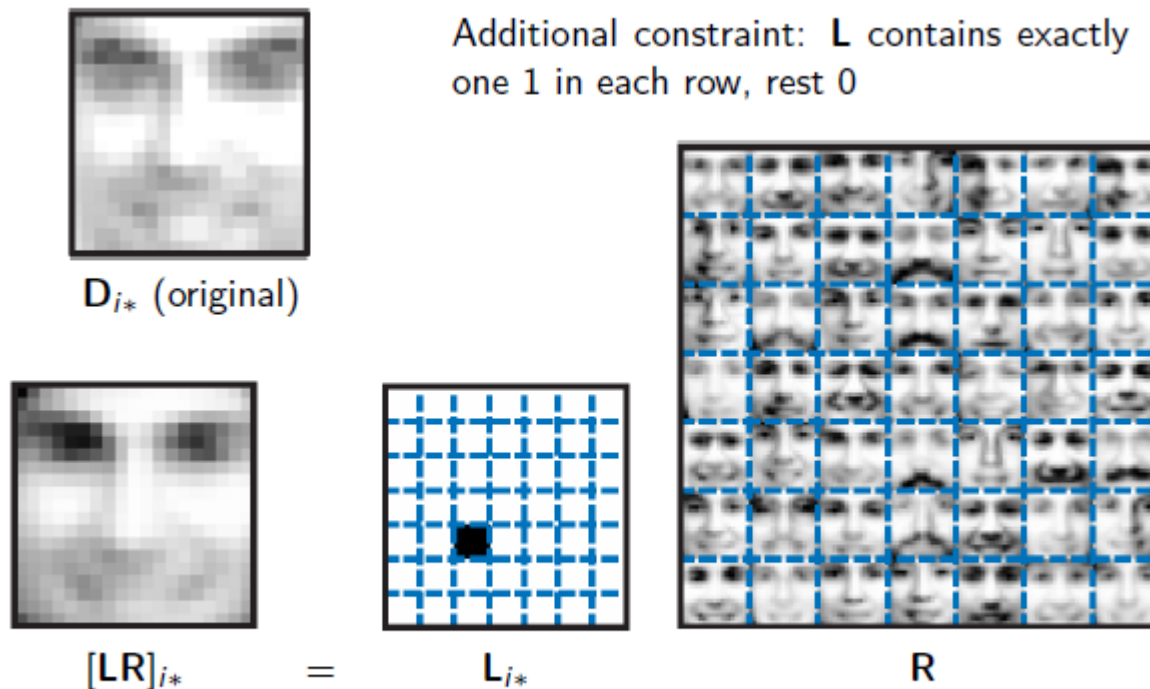
### How to compute NMF: ALS

```
W = rand(m, k);  
for i=1:maxiter
```

Solve for  $\mathbf{H}$  in equation  $\mathbf{W}^T \mathbf{W} \mathbf{H} = \mathbf{W}^T \mathbf{A}$   
Set all negative elements in  $\mathbf{H}$  to 0.  
Solve for  $\mathbf{W}$  in equation  $\mathbf{H} \mathbf{H}^T \mathbf{W}^T = \mathbf{H} \mathbf{A}^T$ .  
Set all negative elements of  $\mathbf{W}$  to zero.

```
end
```

## Part 2: K-means can be seen as a variant of NMF



k-Means factors correspond to prototypical faces.



## Part 2: Algorithmic challenges for NMF

---

$$\min_{\substack{\Phi \in \mathcal{R}^{m \times k}, X \in \mathcal{R}^{k \times n} \\ \forall i, j, \Phi_{ij}, X_{ij} \geq 0}} \left\| \begin{pmatrix} A \\ m \times n \end{pmatrix} - \begin{pmatrix} \Phi \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} X \\ k \times n \end{pmatrix} \right\|_F^2$$

Algorithmic challenges for the NMF:

1. NMF (as stated above) is **convex given  $\Phi$  or  $X$** , but not if both are unknown.
2. **No unique solution**: many matrices  $\Phi$  and  $X$  that minimize the error.
3. **Other optimization objectives** could be chosen (e.g., spectral norm, etc.)
4. NMF becomes harder if **sparsity constraints** are included (e.g.,  $X$  has a small number of non-zeros).
5. For the multiplicative update rules there exists some theory proving that they **converge to a fixed point**; this might be a local optimum or a saddle point.
6. Little theory is known for the other algorithms.

# Part 2: NMF Challenges

## NMF is not unique

- Factors are not "ordered"

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

- One way of ordering: decreasing Frobenius norm of components (i.e., order by  $\|\mathbf{L}_{*k} \mathbf{R}_{k*}\|_F$ )

- Factors/components are not unique

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix} \\ = \begin{bmatrix} 1 & 0.5 & 1 & 0.5 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0.5 & 0 & 0.5 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix} \\ = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

Additional constraints or regularization can encourage uniqueness.

## General framework

- Gradient descent generally slow
- Stochastic gradient descent inappropriate
- Key approach: alternating minimization
  - Pick starting point  $\mathbf{L}_0$  and  $\mathbf{R}_0$
  - while** not converged **do**
  - Keep  $\mathbf{R}$  fixed, optimize  $\mathbf{L}$
  - Keep  $\mathbf{L}$  fixed, optimize  $\mathbf{R}$
  - end while**
- Update steps 3 and 4 easier than full problem

## NMF is not hierarchical

- Rank-1 NMF

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix} \approx \begin{bmatrix} 0.6 & 1.3 & 0.6 & 1.3 & 0.6 \\ 0.3 & 0.8 & 0.3 & 0.8 & 0.3 \\ 0.3 & 0.8 & 0.3 & 0.8 & 0.3 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 0.5 \\ 0.5 \end{bmatrix} \begin{bmatrix} 0.7 & 1.5 & 0.7 & 1.5 & 0.7 \end{bmatrix}$$

- Rank-2 NMF

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix} \\ = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

- Best rank- $k$  approximation may differ significantly from best rank- $(k - 1)$  approximation
- Rank influences sparsity, interpretability, and statistical fidelity
- Optimum choice of rank is not well-studied (often requires experimentation)

## Part 2: NMF Summary

---

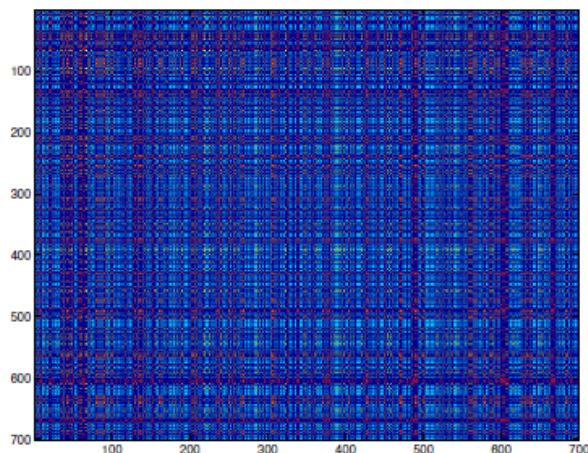
- Given a nonnegative matrix  $\mathbf{A} \in \mathbf{R}^{m \times n}$  find nonnegative matrices  $\mathbf{W} \in \mathbf{R}^{m \times k}$  and  $\mathbf{H} \in \mathbf{R}^{k \times n}$  so that

$$\min_{\mathbf{W}, \mathbf{H}} ||\mathbf{A} - \mathbf{WH}||$$

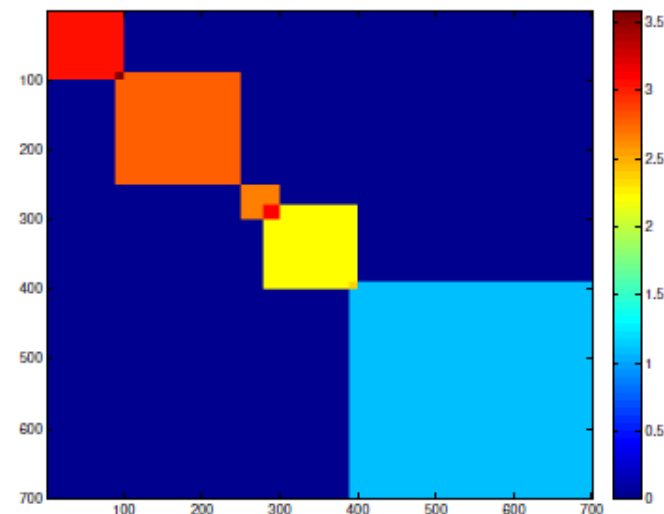
- Algorithms exist, both basic (easy to implement) and more advanced (implementing e.g. sparsity constraints)
- Interpretability
- Applications
  - Text mining
  - Email surveillance
  - Music transcription
  - Bioinformatics
  - Source separation
  - Spatial data analysis
  - Etc.

# Part 2: SemiDiscrete Decomposition (SDD)

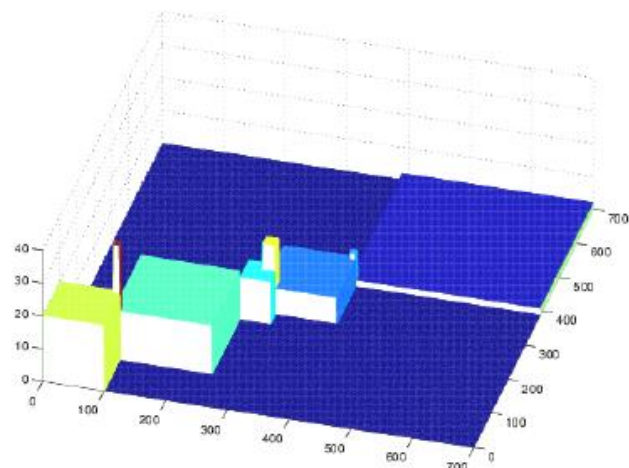
An example data



The data



The data after permuting rows and columns



The data in a 3D view

Can we find the bump in the picture automatically (from unpermuted data)?

What is a bump?

A submatrix of a matrix  $A \in \mathbb{R}^{m \times n}$  contains some rows of  $A$  and some columns of those row

$$\mathbf{A} = \begin{pmatrix} 3 & 1 & 3 \\ 2 & 3 & 1 \\ 3 & 2 & 3 \end{pmatrix} \quad \mathbf{x} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad \mathbf{A} \circ \mathbf{xy}^T = \begin{pmatrix} 3 & 0 & 3 \\ 0 & 0 & 0 \\ 3 & 0 & 3 \end{pmatrix}$$

## Part 2: SemiDiscrete Decomposition (SDD)

---

$$\begin{pmatrix} A_k \\ m \times n \end{pmatrix} = \begin{pmatrix} U_k \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} \Sigma_k \\ k \times k \end{pmatrix} \cdot \begin{pmatrix} V_k^T \\ k \times n \end{pmatrix}$$

$$\begin{pmatrix} A_{\text{SDD}} \\ m \times n \end{pmatrix} = \begin{pmatrix} X_k \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} D_k \\ k \times k \end{pmatrix} \cdot \begin{pmatrix} Y_k^T \\ k \times n \end{pmatrix}$$

$D_k$ : diagonal matrix

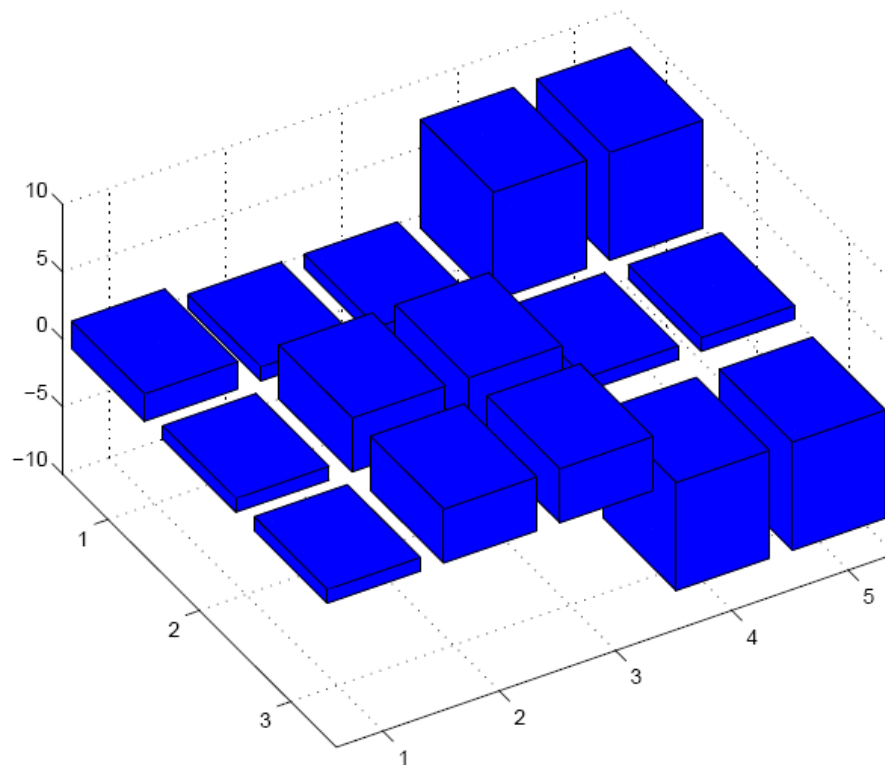
$X_k, Y_k$ : all entries are in  $\{-1, 0, +1\}$

SDD identifies regions of the matrix that have **homogeneous density**.

- Assume we know how to find the largest bump of a matrix
- To find another bump, we can subtract the found bump from the matrix and find the largest bump of the residual matrix

## Part 2: SemiDiscrete Decomposition (SDD)

$$\begin{bmatrix} 2 & 1 & 1 & 8 & 8 \\ 1 & 4 & 4 & 1 & 1 \\ 1 & 4 & 4 & -8 & -8 \end{bmatrix}$$



SDD looks for blocks of similar height towers and similar depth holes: "bump hunting". Applications include image compression and text mining.

O'Leary and Peleg '83, Kolda and O'Leary '98, '00, O'Leary and Roth '06  
The figures are from D. Skillkorn's book on Data Mining with Matrix Decompositions.

## Part 2: Properties of SDD

---

- The columns of  $\mathbf{X}_k$  and  $\mathbf{Y}_k$  do not need to be linearly independent
  - The same column can be even repeated multiple times
- The dimension  $k$  might need to be large for accurate approximation (compared to SVD)
- SDD factors are local
  - Only affect a certain submatrix, typically not every element
  - SVD factors typically change every value
- Storing an  $k$ -dimensional SDD takes less space than storing rank- $k$  truncated SVD
  - $\mathbf{X}_k$  and  $\mathbf{Y}_k$  are ternary and often sparse
- For every rank-1 layer of an SDD, all non-zero values in the layer have the same magnitude ( $d_i$  for layer  $i$ )

## Part 2: Interpretation of SDD

---

- The factor interpretation is not very useful as the factors are not independent
  - A later factor can change just a subset of values already changed by an earlier factor
- The SDD can be interpreted as a form of bi-clustering
  - Every layer (bump) defines a group of rows and columns with homogeneous values in the residual matrix
- The component interpretation is natural to SDD
  - The SDD is a sum of local bumps
  - SDD does not model global phenomena (e.g. noise) well



## Part 2: The outline of the SDD algorithm

---

- ① **Input:** Matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , non-negative integer  $k$
- ② **Output:**  $k$ -dimensional SDD of  $\mathbf{A}$ , i.e. matrices  $\mathbf{X}_k \in \{-1, 0, 1\}^{m \times k}$ ,  $\mathbf{Y}_k \in \{-1, 0, 1\}^{n \times k}$ , and diagonal  $\mathbf{D}_k \in \mathbb{R}_+^{k \times k}$
- ③  $\mathbf{R}_1 \leftarrow \mathbf{A}$
- ④ **for**  $i = 1, \dots, k$ 
  - ① Select  $\mathbf{y}_i \in \{-1, 0, 1\}^n$
  - ② **while** not converged
    - ① Compute  $\mathbf{x}_i \in \{-1, 0, 1\}^m$  given  $\mathbf{y}_i$  and  $\mathbf{R}_i$
    - ② Compute  $\mathbf{y}_i$  given  $\mathbf{x}$  and  $\mathbf{R}_i$
  - ③ **end while**
  - ④ Set  $d_i$  to the average of  $\mathbf{R}_i \circ \mathbf{x}_i \mathbf{y}_i^T$  over the non-zero locations of  $\mathbf{x} \mathbf{y}^T$
  - ⑤ Set  $\mathbf{x}_i$  as the  $i$ th column of  $\mathbf{X}_i$ ,  $\mathbf{y}_i$  the  $i$ th column of  $\mathbf{Y}_i$ , and  $d_i$  the  $i$ th value of  $\mathbf{D}_i$
  - ⑥  $\mathbf{R}_{i+1} \leftarrow \mathbf{R}_i - d_i \mathbf{x}_i \mathbf{y}_i^T$
- ⑤ **end for**
- ⑥ **return**  $\mathbf{X}_k$ ,  $\mathbf{Y}_k$ , and  $\mathbf{D}_k$

## Part 2: SDD Clustering

---

- SDD performs a type of bi-clustering of the matrix
  - Every bump  $\mathbf{dxy}^T$  gives a cluster of rows  $x_i \neq 0$  and a cluster of columns  $y_j \neq 0$ , together with 'centroid'  $d$
  - This is **not** a partition clustering: the clusters can overlap and not every row or column has to belong to some cluster
- We can impose an ordering of the bumps based on the values of  $d_i$ 
  - The algorithm usually returns the bumps in that order
- This ordering can be used to obtain a hierarchical clustering
  - First column of  $\mathbf{X}$  clusters the rows of  $\mathbf{A}$  into three sets (-1,0,1) and same for the first column of  $\mathbf{Y}$  and columns of  $\mathbf{A}$
  - The second column of  $\mathbf{X}$  splits the previous clusters again into three sets
  - And so on and so forth
- Distance between two objects in the hierarchical clustering is not the usual dendrogram depth, but depends on whether we use the 0 or +/-1 branches

## Part 2: SDD Other applications

---

- Image compression (O'Leary and Peleg, 1983)
  - A grayscale image is compressed using its SDD
  - The original application, modern image compression techniques are better
- Latent topic models (Kolda and O'Leary, 1998)
  - Used similarly as SVD is used to compute LSA
  - Compute the SDD of the term-document matrix
- SDD to correlation matrices
  - A bump in the correlation matrix  $\mathbf{A}\mathbf{A}^T$  corresponds to rows of  $\mathbf{A}$  with similar values
  - A bump in  $\mathbf{A}^T\mathbf{A}$  corresponds to columns with similar values

## Part 2: SVD+SDD General approaches

---

- Most common way to combine SVD and SDD is to first use SVD to denoise the data and then to compute the SDD on the clean data
  - Clean data = truncated  $\text{SVD}(\mathbf{A}_k)$
  - SVD is good at finding global structure, SDD at finding local structure
- Another option is to first compute the  $\mathbf{A}_k$  with SVD and then apply SDD to  $\mathbf{A}_k\mathbf{A}_k^T$  (or  $\mathbf{A}_k^T\mathbf{A}_k$ )
  - SDD finds the objects with similar values
- The results can be visualized using the first 2-3 columns of  $\mathbf{U}$  from SVD and the first layers of the hierarchical clustering from SDD

## Part 2: SDD Lessons learned

---

- SDD finds the local areas of values with uniform magnitude  
→ Easier interpretation, 'orthogonal' view to SVD
- Finding SDD is hard and requires a heuristic
- Together SVD and SDD provide a strong analysis toolset

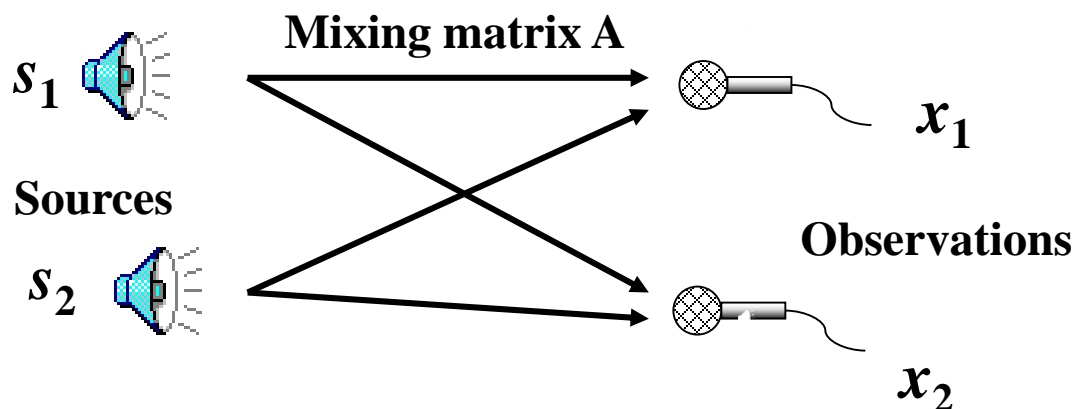
## Part 2: Independent Component Analysis

- Consider the cocktail-party problem: in a room, you have two people speaking simultaneously, and two microphones recording the mixture of these speech signals.
- Each recording is a weighted sum of speech signals  $s_1(t)$  and  $s_2(t)$

$$x_1(t) = a_{11}s_1 + a_{12}s_2$$

$$x_2(t) = a_{21}s_1 + a_{22}s_2$$

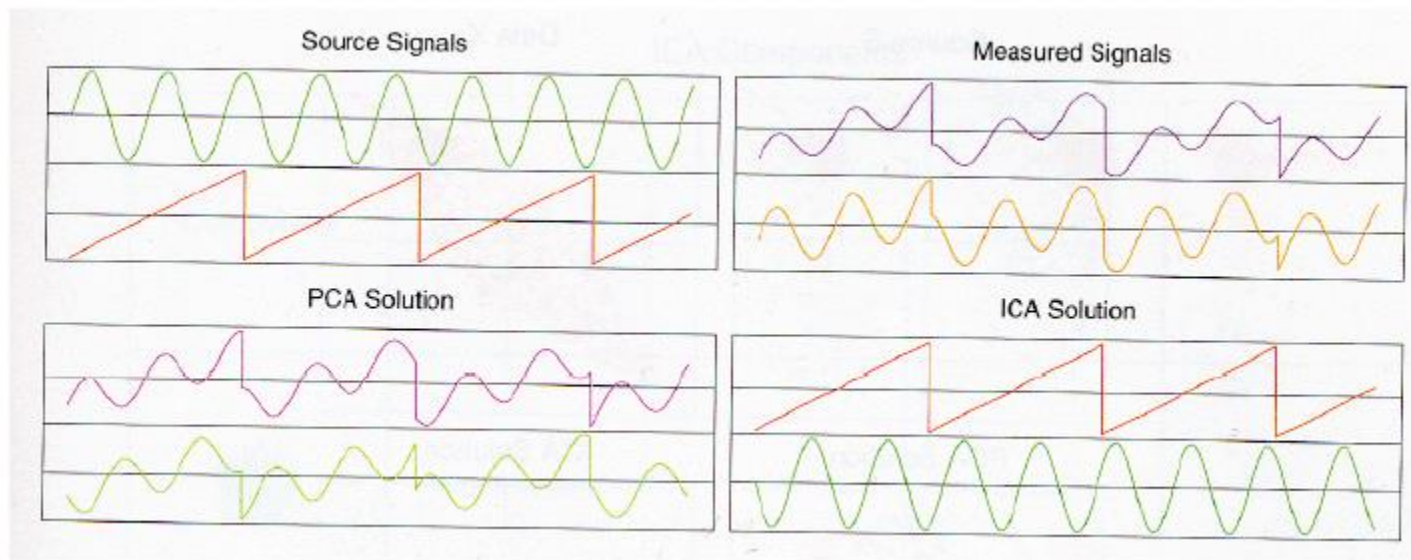
- where  $a_{ij}$  are some parameters depending on the distances of the microphones from the speakers.
- How to recover the original signals  $s_1$  and  $s_2$  from the recorded signals  $x_1$  and  $x_2$ ?



$$\mathbf{x} = \mathbf{A}\mathbf{s}$$

**Task: estimate  $\mathbf{A}$  and  $\mathbf{s}$  using only the observable random vector  $\mathbf{x}$**

## Part 2: Independent Component Analysis



- PCA gives uncorrelated components. In the cocktail-party problem this is not the right answer.
- ICA gives statistically independent components.
- Variables  $y_1$  and  $y_2$  are independent, if information on the value of  $y_1$  does not give any information on the value of  $y_2$ , and vice versa.
- Note: data must be nongaussian.

## Part 2: Example: Image separation. Mixtures of images

---





## Part 2: Example: Image separation. Mixtures of images

---

ICA produced the following images:



## Part 2: Lessons learned

---

- SVD is the Swiss Army knife of (numerical) linear algebra  
→ ranks, kernels, norms, ...
- SVD is also very useful in data analysis  
→ noise removal, visualization, dimensionality reduction, ...
- Selecting the correct rank for truncated SVD is still a problem
- SDD finds the local areas of values with uniform magnitude  
→ easier interpretation, 'orthogonal' view to SVD
- Finding SDD is hard and requires a heuristic
- Non-negative matrix factorization (NMF) appears natural for non-negative data
- NMF encourages parts-based decomposition, interpretability, and sparseness
- Usually solved via alternating minimization algorithms.

## Part 0. Preliminary

- 0-1. Modern Data
- 0-2. About the Tutorial
- 0-3. References
- 0-4. Suggested Readings

## Part 1. Linear Algebra Review

- 1-1. Basics of Linear Algebra
- 1-2. Vector Notations
- 1-3. Matrix Notations

## Part 2. Matrix Algorithms

- 2-1. Singular Value Decomposition (SVD)
- 2-2. Non-negative Matrix Factorization (NMF)
- 2-3. Semi-Discrete Decomposition (SDD)
- 2-4. Independent Component Analysis (ICA)

## **Part 3. Applications of Matrix Algorithms**

- 3-1. Examples of Matrix Algorithms Applied on Data Mining Applications

## Part 3: Problem definition and motivation

---

- In many applications (e.g., [statistical data analysis](#) and [scientific computation](#)), one has [n observations](#) of the form:

$$y_i = y(t_i), i = 1, \dots, n$$

Model  $y(t)$  (unknown) as a linear combination of [d basis functions](#):

$$y(t) \approx x_1 \phi_1(t) + \dots + x_d \phi_d(t)$$

$A$  is an  $n \times d$  “design matrix” ( $n \gg d$ ):

$$A_{ij} = \phi_j(t_i)$$

In matrix-vector notation,

$$y \approx Ax$$

## Part 3: Least-norm approximation problems

---

- Recall a **linear measurement model**:

$$y = Ax + \varepsilon \quad \begin{cases} y \text{ are the measurements} \\ x \text{ is the unknown} \\ \varepsilon \text{ is an error process} \end{cases}$$

- In order to **estimate  $x$** , solve:

$$\hat{x} = \arg \min \|y - Ax\|$$

- First application: Astronomy**

Predicting the orbit of the asteroid Ceres (in 1801!).

Gauss (1809) -- see also Legendre (1805) and Adrain (1808).

First application of “least squares optimization” and runs in  $O(nd^2)$  time!

- Data analysis**: Fit parameters of a **biological, chemical, economical, physical (astronomical), social, internet**, etc. model to experimental data.

## Part 3: Norms of common interest

---

Let  $y = b$  and define the residual:  $r = Ax - b \in R^n$

Least-squares approximation:

$$\text{minimize: } \|Ax - b\|_2^2 = r_1^2 + r_2^2 + \cdots + r_n^2$$

Chebyshev or mini-max approximation:

$$\text{minimize: } \|Ax - b\|_\infty = \max\{|r_1|, \dots, |r_n|\}$$

Sum of absolute residuals approximation:

$$\text{minimize: } \|Ax - b\|_1 = |r_1| + |r_2| + \cdots + |r_n|$$

## Part 3: $L_p$ norms and their unit balls

Recall the  $L_p$  norm for  $z \in R^n$

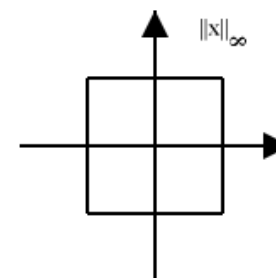
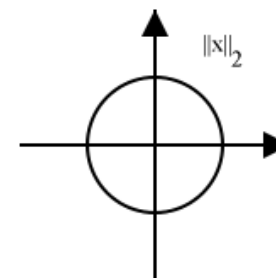
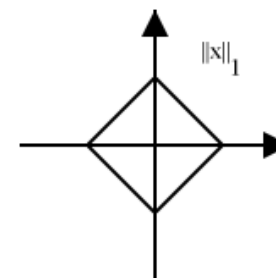
$$\|z\|_p = \left( \sum_{i=1}^n |z_i|^p \right)^{1/p}, \quad p \in [1, \infty)$$

$$\|z\|_\infty = \max_i |z_i|$$

$$\|z\|_2^2 = \sum_i z_i^2 = z^T z$$

Some inequality relationships include:

$$\frac{1}{\sqrt{n}} \|z\|_2 \leq \|z\|_\infty \leq \|z\|_2 \leq \|z\|_1 \leq \sqrt{n} \|z\|_2$$



## Part 3: $L_p$ regression problems

---

$$\begin{aligned} \mathcal{Z}_p &= \min_{x \in \mathbb{R}^d} \|b - Ax\|_p \\ &= \|b - A\hat{x}\|_p \end{aligned} \quad \rightarrow \quad \begin{matrix} \begin{pmatrix} & & \\ & A & \\ & & \end{pmatrix}_{n \times d, \, n \gg d} \begin{pmatrix} \hat{x} \end{pmatrix} \approx \begin{pmatrix} b \end{pmatrix} \end{matrix}$$

We are interested in **over-constrained  $L_p$  regression problems**,  $n \gg d$ .

Typically, there is no  $x$  such that  $Ax = b$ .

Want to find the “best”  $x$  such that  $Ax \approx b$ .

$L_p$  regression problems are convex programs (or better).

There exist **poly-time algorithms**.

We want to **solve them faster**.



# Part 3: Exact solution to $L_2$ regression

## Cholesky Decomposition:

If  $A$  is full rank and well-conditioned,  
decompose  $A^T A = R^T R$ , where  $R$  is upper triangular, and  
solve the normal equations:  $R^T R x = A^T b$ .

## QR Decomposition:

Slower but numerically stable, esp. if  $A$  is rank-deficient.  
Write  $A = QR$ , and solve  $Rx = Q^T b$ .

## Singular Value Decomposition:

Most expensive, but best if  $A$  is very ill-conditioned.  
Write  $A = U \Sigma V^T$ , in which case:  $\mathbf{x}_{\text{OPT}} = \mathbf{A}^+ \mathbf{b} = \mathbf{V} \Sigma^{-1} \mathbf{U}^T \mathbf{b}$ .

Complexity is  $O(nd^2)$ , but constant factors differ.

$$\begin{aligned} \mathcal{Z}_2 &= \min_{x \in \mathbb{R}^d} \|b - Ax\|_2 \\ &= \|b - A\hat{x}\|_2 \end{aligned}$$

Projection of  $b$  on the  
subspace spanned by the  
columns of  $A$

$$\begin{aligned} \mathcal{Z}_2^2 &= \|b\|_2^2 - \|AA^+b\|_2^2 \\ \hat{x} &= A^+b \end{aligned}$$

Pseudoinverse of  $A$

## Part 3: Questions ...

---

$$\mathcal{Z}_p = \min_{x \in \mathbb{R}^d} \|b - Ax\|_p = \|b - A\hat{x}\|_p$$

### Approximation algorithms:

Can we approximately solve  $L_p$  regression faster than “exact” methods?

### Core-sets (or induced sub-problems):

Can we find a small set of constraints such that solving the  $L_p$  regression on those constraints gives an approximation to the original problem?

## Part 3: Lessons learned

---

- Least square problems  $\mathbf{Ax}=\mathbf{b}$  can be solved by
  - Normal equation,  $\mathbf{A}^T\mathbf{Ax}=\mathbf{A}^T\mathbf{b} \rightarrow \mathbf{R}^T\mathbf{Rx} = \mathbf{A}^T\mathbf{b}$
  - QR decomposition,  $\mathbf{Rx} = \mathbf{Q}^T\mathbf{b}$
  - Singular value decomposition,  $\mathbf{x}_{\text{OPT}} = \mathbf{A}^+\mathbf{b} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{b}$
- If the matrix size is huge
  - Approximate with random algorithm

# Conclusions

---

- Linear Algebraic techniques (e.g. matrix decompositions and regression) are fundamental in data mining.
- Matrix decompositions reveal structure in the data  
→  **$D = LMR$**
- Many different decompositions with different applications exist  
→ QR, SVD, k-means, NMF, SDD, ICA, etc.
- SVD is the Swiss Army knife of (numerical) linear algebra  
→ ranks, kernels, norms, ...
- SVD is also very useful in data analysis  
→ noise removal, visualization, dimensionality reduction, ...
- Non-negative matrix factorization (NMF) appears natural for non-negative data
- NMF encourages parts-based decomposition, interpretability, and sparseness
- Randomized algorithms for linear algebra computations contribute novel results and ideas, both from a theoretical as well as an applied perspective.