Data Mining and Matrices 01 – Introduction

Prof. Dr. Rainer Gemulla

Universität Mannheim

FSS 2018

Outline

- 1. What is data mining?
- 2. What is a matrix?
- 3. Why data mining and matrices?
- 4. Summary

Outline

- 1. What is data mining?
- 2. What is a matrix?
- 3. Why data mining and matrices?
- 4. Summary

"Data mining is the process of discovering knowledge or patterns from massive amounts of data." (Encyclopedia of Database Systems)



Estimated \$100 billion industry around managing and analyzing data.

"Data mining is the process of discovering knowledge or patterns from massive amounts of data." (Encyclopedia of Database Systems)

- Science
 - ► The Sloan Digital Sky Survey gathered 140TB of information
 - ► NASA Center for Climate Simulation stores 32PB of data
 - ► 3B base pairs exist in the human genome
 - ► LHC registers 600M particle collisions per second, 25PB/year
- Social/business data
 - ▶ 1M customer transactions are performed at Walmart per hour
 - 4.6B mobile-phone subscriptions worldwide
 - ▶ 300M active Twitter users write 500M tweets per day
 - 100M Netflix customers view and rate hundreds of thousands of movies
 - 300h of videos uploaded every minute on YouTube
- Government, health care, news, stocks, books, web search, ...
- Often: less data than in these extreme cases, but still "massive"

"Data mining is the process of discovering knowledge or patterns from massive amounts of data." (Encyclopedia of Database Systems)





Outlier detection

"Regnet es am Siebenschläfertag, der Regen sieben Wochen nicht weichen mag." (German folklore)

Pattern mining

"Data mining is the process of discovering knowledge or patterns from massive amounts of data." (Encyclopedia of Database Systems)



Knowledge discovery process

b ai				V.			40		÷	1			ε			V			01	t8						t				0	+
				64			8		t	12	Ц		-tt						23	7.0	7)	X				141	Ŧ			t	Ē
4	0			† 8	3		43				÷		2		ŧ.				τi	2.6	÷	A	П		UT() (e	え		K T	7	Д
1	70	8		T-	53		Ŧ		ŝ	đ	E	t	τ		8				đ		T	U	ł.	8	Uð	T	0		87	X	ŋ
6	÷i	3	5	Ŧŀ	5		V		5	÷	R	Ŷ	6		- F	3			i I	T	R	÷.	Ť.		37	÷	Ŧ		10	S.	Ŧ
÷	00	8	Â	Ľ,			đ		E c		Δ	V							Ň		E	1	T	y.	21	1	а			R	v
a	¥.	εN	4	÷λ	2		n	-	p i	T,	e la	4	¢.		ų Ł	h_{V}		+	2	÷	÷	a			T		Ν	ĉ.		6	Ę
d.		÷	č		M					dia.		110			1							E	ř.	199		P	i.	e,		Ě.	n
	U.	4	n											1			Ŧ					2					2	÷.	ñ	H	Ď
+	÷		п							<u>P</u>							Ľ					Ă					Ť	Ň		H	ñ
t.	2	i u	A	ř	e.	ñ ł	ìñ	54	č r	i e	Ť	2				1		0	Ň	-1	ĕ			×.	ž i	10	-	e i		'n	č
à	1	12		-+ B								K				Y			N		8	7				i i		4	22	Н	Ă
E.			2							•									X			Ň	6		1	12	X				X
		ι.	Υ.	50	14							Ц													77		1	Q.			0
Ą	Ì	5		2	N		E	0	ę.	J,	X	X	A	Q I		L	4			6	٤.	٩L	뷥	4	S.		귀	런	C L		X
オ	I	1+	81	+	5		-1 E		Y	0	E	Ť,	+		7	7			U	ζ.	g	オ	A.	V.	5	ЦĻ	lan.	Ċ,	E.	8	
£0	U	3	3.44	L.	9	3	XC	0	0	1	A	Ļ	1	X,	F	5			8 -	ŧ₹	Ŧ	9		ę,	t.	11	Π	¢,	tτ		
0	γ.		2	7		1	r4	Ψ.	47		А	8	L	0			Ŧ		۰,	17	67	T.		đ	ł	3.6	Π	ŧ.	FQ	Ŧ	
Д			П	U	<i>t</i>	<			د	75	ų,	A	0	П	T	4	t		X,	(đ	8	Ŧ	2	61	N,	13	H	5	s di	2	
Ŧ	X() (t	σC	۲.	1	ŧ٧		ą	20		ð	ŶĴ.	N-	Ŧ.	18	P		1	(8	X	5	8	51	Π			ſ,	11	Ŧ	ę
t	٢,	(t)		sI	3	(70				+	Ŧ	+	7	2	Ŕ	Y			it i	63	0	ć.	N	9 (Sec.	01	2	Π	X
Ċ.	÷	11	Γ.	-	П	П	10		4	10	-	a	a	14	-8	C,	2		1.3	2.5	4	0	Λ^{1}	T	n 4	1.	-4	+-	1.4	-	

Womb



- *mater* (Latin) = *mother*
- matrix (Latin) = pregnant animal
- *matrix* (Late Latin) = *womb* also *source*, *origin*
- Since 1550s: *place or medium where something is developed*
- Since 1640s: *embedding or enclosing mass*

Rectangular arrays of numbers

• "Rectangular arrays" known in ancient China (rod calculus, estimated as early as 300BC)

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

• Term "matrix" coined by <u>J.J. Sylvester</u> in 1850



System of linear equations

Three sheafs of a good crop, two sheafs of a mediocre crop, and one sheaf of a bad crop are sold for 39 dou. Two sheafs of good, three mediocre, and one bad are sold for 34 dou; and one good, two mediocre, and three bad are sold for 26 dou. What is the price received for each sheaf of a good crop, each sheaf of a mediocre crop, and each sheaf of a bad crop? Chiu-chang Suan-shu (Nine Chapters on Arithmetic), $\approx 200BC$

• Systems of linear equations can be written as matrices

$$3x + 2y + z = 39$$

$$2x + 3y + z = 34 \qquad \rightarrow \qquad \begin{pmatrix} 3 & 2 & 1 & | & 39 \\ 2 & 3 & 1 & | & 34 \\ 1 & 2 & 3 & | & 26 \end{pmatrix}$$

$$x + 2y + 3z = 26$$

• and then be solved using linear algebra methods

$$\begin{pmatrix} 3 & 2 & 1 & | & 39 \\ 5 & 1 & | & 24 \\ & & 12 & | & 33 \end{pmatrix} \implies \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 9.25 \\ 4.25 \\ 2.75 \end{pmatrix}$$

Set of data points



Linear maps

• Linear maps from \mathbb{R}^3 to \mathbb{R}

$$f_1(x, y, z) = 3x + 2y + z$$

$$f_2(x, y, z) = 2x + 3y + z$$

$$f_3(x, y, z) = x + 2y + 3z$$

$$f_4(x, y, z) = x$$

• Linear map f₁ written as a matrix

$$\begin{pmatrix} 3 & 2 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = f_1(x, y, z)$$

• Linear map from \mathbb{R}^3 to \mathbb{R}^4

$$\begin{pmatrix} 3 & 2 & 1 \\ 2 & 3 & 1 \\ 1 & 2 & 3 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} f_1(x, y, z) \\ f_2(x, y, z) \\ f_3(x, y, z) \\ f_4(x, y, z) \end{pmatrix}$$

Original data



Rotated and stretched



Graphs



Objects and attributes

Anna, Bob, and Charlie went shopping

- Anna bought butter and bread
- Bob bought butter, bread, and beer
- Charlie bought bread and beer



Many different kinds of data fit this object-attribute viewpoint.

What is a matrix?

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

• • • •

• A means to describe *data*

Rows	Columns	Entries					Attı	
Objects	Attributes	Values		(a ₁₁	a_{12}		a_{1j}	•••\
Equations	Variables	Coefficients		a ₂₁	a 22	• • •	a 2j	• • •
Data points	Axes	Coordinates		:	:	•	1	•
Vertices	Vorticos	Edges		· ·	•			•
vertices	vertices	Luges	Object <i>i</i>	a _{i1}	a i2	• • •	aij	• • •
:	:	:			:			· .
				$\langle : \rangle$				· · /

Linear operators

In data mining, we make use of both viewpoints simultaneously.

bute j

Outline

1. What is data mining?

2. What is a matrix?

3. Why data mining and matrices?

4. Summary

Key tool: Matrix decompositions

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

D = LMR,

where

- **D** is an $m \times n$ data matrix,
- **L** is an $m \times r_1$ matrix,
- **M** is an $r_1 imes r_2$ matrix,
- **R** is an $r_2 \times n$ matrix, and
- r_1 and r_2 are integers ≥ 1 .
- Often $r_1 = r_2 = r \ge 1$

$$d_{ij} = \sum_{k,k'} I_{ik} m_{kk'} r_{k'j}$$



Why matrix decompositions?

- Decompositions as just defined are not really helpful
 - Suppose we set $r = r_1 = r_2 = n$, L = D, $M = R = I_n$ (the $n \times n$ identity matrix)
 - Then $D = LMR = DI_nI_n = D$
 - But this does not provide insight
- To make decompositions useful, we want the decomposition to satisfy certain (carefully chosen) properties or contraints
- For example: we may want r to be small
 - Each object is represented by n numbers in D
 - Each object is represented by r numbers in L (called embedding or distributed representation)
 - If r < n, we performed some form of compression
- Another example: we may want factors to have certain properties
 - Compare: integer factorization
 - ▶ 391 = 17 · 13



Approximate matrix decompositions

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

$\boldsymbol{D} pprox \boldsymbol{L} \boldsymbol{M} \boldsymbol{R} = \hat{\boldsymbol{D}},$

where each matrix has conforming dimensions (as before).

- We often look at approximate decompositions
 - Data is noisy anyway
 - Approximation may remove noise
 - ▶ Allows to focus on global (small *r*) or local (large *r*) patterns
 - Often more insightful
 - More efficient to compute
- \approx defined by some loss function $L(D, \hat{D})$
 - E.g., squared error
 - Low means good approximation, high means bad
 - Finding the best approximation (smallest loss) can be hard
- We often say "matrix decomposition" when we actually mean "approximate matrix decomposition"

Some matrix decompositions

- There are many different decompositions, each enforcing different constraints and serving a different purpose
 - Singular value decomposition (SVD)
 - k-means (crisp or fuzzy)
 - Non-negative matrix factorization (NMF)
 - Semi-discrete decomposition (SDD)
 - Boolean matrix decomposition (BMF)
 - Independent component analysis (ICA)
 - Matrix completion
 - Probabilistic matrix factorization
 - Tensor decompositions
 - ▶ ...
- Picking the right one for the problem (if any) at hand is hard, experience helps
- Decompositions are not always easy (and often hard) to compute

Example: Singular value decomposition



L_{50×2}

 $M_{2\times 2}$

 $R_{2\times 2}$



 $\begin{pmatrix} 11.73 & 0 \\ 0 & 1.71 \end{pmatrix}$



Example: k-Means





Example: Non-negative matrix factorization



Example: Latent Dirichlet allocation

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

R (topic×word)

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services." Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

(doc×topic)

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
- Cluster 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.

^{• ..}

Factor interpretion of matrix decompositions

Assume that M is $r \times r$ and 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 = representation of object via weighted parts, called embedding or distr. representation (r pieces of information)
 d_i^T = ∑_k
- 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")



27 / 36

Example: Weather data (r = 1)1.00 9.05 16.55 26.73 18.75 17.81

		Jan	Apr	Jul	Oct	Year
0.62	Stockholm	-0.70	8.60	21.90	9.90	10.00
0.69	Minsk	-2.10	12.20	23.60	10.20	10.60
0.83	London	7.90	13.30	22.80	15.20	14.80
0.90	Budapest	1.20	16.30	26.50	16.10	15.00
0.88	Paris	6.90	14.70	24.40	15.80	15.50
0.98	Bucharests	1.50	18.00	28.80	18.00	16.50
1.09	Barcelona	12.40	17.60	27.50	21.50	20.00
1.14	Rome	11.90	17.70	30.30	21.40	20.40
1.16	Lisbon	14.80	19.80	27.90	22.50	21.50
1.24	Athens	12.90	20.30	32.60	23.10	22.30
1.21	Valencia	16.10	20.20	29.10	23.60	22.30
1.27	Malta	16.10	20.00	31.50	25.20	23.20

Example: Weather data (r = 1), reconstruction 9.05 16.55 26.73 18.75 17.81 1.00

		Jan	Apr	Jul	Oct	Year
0.62	Stockholm	5.65	10.33	16.68	11.70	11.11
0.69	Minsk	6.21	11.36	18.35	12.87	12.23
0.83	London	7.55	13.80	22.28	15.63	14.85
0.90	Budapest	8.11	14.83	23.94	16.80	15.96
0.88	Paris	7.96	14.56	23.52	16.50	15.67
0.98	Bucharests	8.91	16.30	26.32	18.47	17.54
1.09	Barcelona	9.88	18.06	29.17	20.46	19.44
1.14	Rome	10.28	18.80	30.35	21.30	20.23
1.16	Lisbon	10.47	19.15	30.92	21.70	20.61
1.24	Athens	11.21	20.50	33.11	23.23	22.07
1.21	Valencia	10.92	19.96	32.24	22.62	21.48
1.27	Malta	11.47	20.98	33.88	23.77	22.58

Ô

(RMSE: 2.66) 29/36

Example: Weather data $(r = 2)$, reconstruction											
1.00			9.05	16.55	26.73	18.75	17.81				
	1.00		-4.14	0.27	2.32	-0.89	-0.69				
			Jan	Apr	Jul	Oct	Year				
0.62	1.69	Stockholm	-1.34	10.79	20.59	10.20	9.95				
0.69	2.11	Minsk	-2.52	11.94	23.23	10.99	10.77				
0.83	0.00	London	7.54	13.80	22.28	15.63	14.85				
0.90	1.52	Budapest	1.82	15.24	27.46	15.45	14.91				
0.88	0.30	Paris	6.71	14.65	24.22	16.23	15.46				
0.98	1.59	Bucharests	2.31	16.74	30.02	17.05	16.44				
1.09	-0.66	Barcelona	12.61	17.88	27.64	21.05	19.90				
1.14	-0.31	Rome	11.55	18.71	29.64	21.57	20.44				
1.16	-1.09	Lisbon	15.00	18.85	28.39	22.67	21.36				
1.24	-0.35	Athens	12.65	20.41	32.31	23.54	22.31				
1.21	-1.26	Valencia	16.14	19.62	29.31	23.74	22.36				
1.27	-1.12	Malta	16.10	20.67	31.29	24.76	23.35				
							D				

(RMSE: 0.60) 30/36

Example: Weather data (r = 2), plot



Factor 1

Example: Netflix prize data





Koren et al., 2009

Other interpretions

- Geometric interpretation
 - ► Transformation of *n*-dimensional space in *r*-dimensional space
 - Row of *R* = new axes
 - ▶ Row of *L* = new coordinates (*embeddings*)
- Component interpretation
 - ► **D** is viewed as consisting of r layers (of same shape as **D**)
 - k-th layer described by $\boldsymbol{I}_k m_{kk} \boldsymbol{r}_k^T$
 - $\blacktriangleright \mathbf{D} = \sum_{k} \mathbf{I}_{k} m_{kk} \mathbf{r}_{k}^{T}$
- Graph interpretation
 - ► **D** is thought of as a bipartite graph with object and attribute vertexes
 - ► Edge weights measure association b/w objects and attributes
 - Decomposition thought of as a tripartite graph with row, waypoint, and column vertexes

All interpretations are useful (more later).

Outline

1. What is data mining?

2. What is a matrix?

3. Why data mining and matrices?

4. Summary

Lessons learned

- Data mining = from data to knowledge
 → Prediction, clustering, outlier detection, patterns
- Matrices are common representation of datasets \rightarrow Linear maps, data points, sets, graphs, relational data, \ldots
- Often: rows = objects, columns = attributes
- Matrix decompositions reveal structure in the data \rightarrow *D* \approx *LMR*
- Many different decompositions with different applications \rightarrow SVD, *k*-means, NMF, SDD, BMF, ICA, completion, ...
- Factor interpretation: objects described by "prototypical objects"

Suggested reading

- Skillicorn, Ch. 1: Data Mining
- Skillicorn, Ch. 2: Matrix Decompositions
Data Mining and Matrices 02 – Vectors & Matrices

Prof. Dr. Rainer Gemulla

Universität Mannheim

FSS 2018

Outline

- 1. Vectors
- 2. Matrices
- 3. Summary

Outline

1. Vectors

2. Matrices

3. Summary

Vector

A weater is		Year
A vector is	Stockholm	9.95
 A 1D array of numbers 	Minsk	10.77
• A geometric ontity with magnitude and direction	London	14.85
• A geometric entity with magnitude and direction	Budapest	14.91
 A matrix with exactly one row or column 	Paris	15.46
Called row vector and column vector, resp.	Bucharests	16.44
Transpose \mathbf{v}^T transposes a row vector into a	Barcelona	19.90
column vector and vice versa	Rome	20.44
column vector and vice versa	Lisbon	21.36
 A (latent) object or attribute 	Athens	22.31
	Valencia	22.36
	Malta	\23.35/



	Jan	Apr	Jul	Oct	Year
Stockholm	(-0.70	8.60	21.90	9.90	10.00)

Vector norm

The norm of vector defines its magnitude. Let $\boldsymbol{v} = \begin{pmatrix} v_1 & v_2 & \cdots & v_n \end{pmatrix}^T$.

- Euclidean norm: $\|\mathbf{v}\| = \sqrt{\sum_{i=1}^{n} v_i^2}$
 - Corresponds to intuitive notion of length in Euclidean space

•
$$L_p$$
 norm for $1 \le p \le \infty$: $\|m{v}\|_p = (\sum_{i=1}^n |v_i|^p)^{1/p}$

- L₁ norm = sum of absolute values (Manhattan distance from origin)
- L₂ norm = Euclidean norm (bird-fly distance from origin)
- L_{∞} norm = maximum absolute value
- The L_p norms never increase as p increases, i.e.,

$$\| \boldsymbol{v} \|_{p+a} \leq \| \boldsymbol{v} \|_{p}$$
 for $a \geq 0$

- Properties of vector norms
 - $\blacktriangleright \ \| \boldsymbol{\nu} \| > 0 \text{ when } \boldsymbol{\nu} \neq 0 \text{ and } \| \boldsymbol{\nu} \| = 0 \text{ iff } \boldsymbol{\nu} = \boldsymbol{0}$
 - $||a\mathbf{v}|| = |a| ||\mathbf{v}||$ (absolute scalability)
 - $\blacktriangleright \|\boldsymbol{v}_1 + \boldsymbol{v}_2\| \le \|\boldsymbol{v}_1\| + \|\boldsymbol{v}_2\| \text{ (triangle inequality)}$

 $v = \begin{pmatrix} 4 \\ -3 \end{pmatrix}$

 $\|v\|_{1} = 7$

 $\|v\| = 5$

 $\| v \|_{\infty} = 4$

Norms and distances

The distance between two vectors $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^n$ can be quantified with norm $\|\boldsymbol{u} - \boldsymbol{v}\|$.

		Jan	Apr	Jul	Oct	Year
•	Stockholm, ${m s}=$	(_0.70	8.60	21.90	9.90	10.00)
•	Minsk, m =	(-2.10)	12.20	23.60	10.20	10.60)
	A . I	1000	00.00	22.00	00.10	a a a a

• Athens, $\boldsymbol{a} = ig(12.90 \ \ 20.30 \ \ 32.60 \ \ 23.10 \ \ 22.30ig)$

L_1	S	т	а	L_2	S	m	а
s	0.00	7.60	61.50	s	0.00	4.27	27.60
т	7.60	0.00	56.70	т	4.27	0.00	25.98
а	61.50	56.70	0.00	а	27.60	25.98	0.00

L_{∞}	S	т	а
S	0.00	3.60	13.60
т	3.60	0.00	15.00
а	13.60	15.00	0.00

Dot product (algebraic definition)

The **dot product** of two vectors $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^n$ is given by

$$\boldsymbol{u}\cdot\boldsymbol{v}=\sum_{i=1}^n u_iv_i.$$

- Also known as scalar product or inner product
- We'll often use matrix product notation and write $\boldsymbol{u}^T \boldsymbol{v}$
- Properties (with $a, b \in \mathbb{R}$)

$$\bullet \quad u \cdot v = v \cdot u$$

$$\bullet (au) \cdot v = a(u \cdot v)$$

- $\blacktriangleright (a\mathbf{u} + b\mathbf{v}) \cdot \mathbf{w} = (a\mathbf{u}) \cdot \mathbf{w} + (b\mathbf{v}) \cdot \mathbf{w}$
- Many uses, many interpretations

With dot products, we can ...

• Compute the (squared) Euclidean norm

$$oldsymbol{v}\cdotoldsymbol{v}=\sum_{i=1}^nv_i^2=\|oldsymbol{v}\|^2$$

• Normalize a vector to length 1 (then a unit vector)

$$\hat{oldsymbol{v}} = oldsymbol{v} / \|oldsymbol{v}\|$$

• Determine the value of a coordinate

$$v_i = \boldsymbol{v} \cdot \boldsymbol{e}_i,$$

where e_i denotes the *i*-th standard basis vector (i.e., $[e_i]_j = 1$ if i = j else 0)

• Compute the sum of the elements of a vector

$$\boldsymbol{v}\cdot\boldsymbol{1}_{\boldsymbol{n}}=\sum_{i=1}^{n}v_{i},$$

where $\mathbf{1}_n$ is the all-ones vector of dimensionality n

• ...

Dot product: Weighted sum

The elements of one vector are interpreted as weights for the elements of the other vector.

Example: Anna goes shopping

ltem	Bread	Butter	Pizza
Price/piece	1€	0.50€	3€
Quantity bought	1	2	5

- How much does Anna pay?
- Prices can be interpreted as "weights": $\boldsymbol{p} = \begin{pmatrix} 1 & 0.5 & 3 \end{pmatrix}^T$
- Quantities are $\boldsymbol{n} = \begin{pmatrix} 1 & 2 & 5 \end{pmatrix}^T$
- Total is $p \cdot n = 1 \cdot 1 + 0.5 \cdot 2 + 3 \cdot 5 = 17$
- Similarly: Can interpret quantities as weights for prices

Dot product: Expected value

One vector corresponds to probabilities, the other one to a random variable.

Example: Bob is gambling

Outcome	Jackpot	Win	Loss
Probability	0.1	0.2	0.7
Amount won	5€	1€	-2€

- How much does Bob win in expectation? (Should he play?)
- Probabilities $\boldsymbol{p} = \begin{pmatrix} 0.1 & 0.2 & 0.7 \end{pmatrix}^T$
 - ► A non-negative vector that sums to one (||p||₁ = 1) is called a probability vector
 - Corresponds to a probability distribution over a finite set of outcomes
- Amounts won $\boldsymbol{x} = \begin{pmatrix} 5 & 1 & -2 \end{pmatrix}^T$
 - Corresponds to a random variable; associates a real value with each outcome

• Expected value
$$\boldsymbol{p} \cdot \boldsymbol{x} = 0.1 \cdot 5 + 0.2 \cdot 1 + 0.7 \cdot (-2) = -0.7$$

Dot product: Sample variance

Denote by $\bar{u} = \frac{1}{n} \sum_{i} u_{i}$ the mean of \boldsymbol{u} . If we treat the entries of \boldsymbol{u} as samples from some distribution, then the unbiased sample variance is given by

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (u_{i} - \bar{u})^{2} = \frac{\|\boldsymbol{u} - \bar{\boldsymbol{u}}\|^{2}}{n-1} = \frac{(\boldsymbol{u} - \bar{\boldsymbol{u}}) \cdot (\boldsymbol{u} - \bar{\boldsymbol{u}})}{n-1},$$

where $ar{m{u}}$ denotes the sample mean vector, i.e., $[m{u}]_i = ar{u}$ for

- $1 \leq i \leq n$.
- Example
 - $\bullet \quad \boldsymbol{u} = \begin{pmatrix} 10 & 11 & 12 \end{pmatrix}^T$
 - $\bullet \quad \bar{u} = 11, \ \bar{\boldsymbol{u}} = \begin{pmatrix} 11 & 11 & 11 \end{pmatrix}^T$
 - $\bullet \ \boldsymbol{u} \bar{\boldsymbol{u}} = \begin{pmatrix} -1 & 0 & 1 \end{pmatrix}^T$

•
$$s^2 = 1$$
, $\|\boldsymbol{u}\|^2 = 365$

- Variances are thus closely related to norms; the key difference is centering and averaging
- When we center data before analyzing it, dot products are proportional to variances (*u* · *u*) or covariances (*u* · *v*)

Dot product: Sets and intersections

The indicator vector of a subset T of a set $S = \{s_1, \ldots, s_n\}$ is the vector \mathbf{x} such that $x_i = 1$ if $s_i \in T$ and $x_i = 0$ if $s_i \notin T$. If \mathbf{u} and \mathbf{v} are indicator vectors for subsets $U, V \subseteq S$, resp., then $\mathbf{u} \cdot \mathbf{v} = |U \cap V|$.

- $S = \{ France, Germany, Denmark, Poland \}$
- Anna visited France, Germany, and Poland: $m{u}=egin{pmatrix}1&1&0&1\end{pmatrix}^T$
- Bob visited Germany, Denmark, and Poland: $oldsymbol{v}=egin{pmatrix}0 & 1 & 1\end{pmatrix}^{\mathcal{T}}$
- Number of countries visited by both:

 $u \cdot v = 1 \cdot 0 + 1 \cdot 1 + 0 \cdot 1 + 1 \cdot 1 = 2 = |\{ \text{Germany}, \text{Poland} \}|$

Dot product (geometric definition)

An alternative geometric definition of the dot product of two vectors $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^n$ is

 $\boldsymbol{u}\cdot\boldsymbol{v}=\|\boldsymbol{u}\|\|\boldsymbol{v}\|\cos\theta,$

where $-\pi \leq \theta \leq \pi$ denotes the angle between \boldsymbol{u} and \boldsymbol{v} .



Why is this?

Let's focus on the 2D case. Recall the law of cosines:

$$c^2 = a^2 + b^2 - 2ab\cos\theta.$$

Now set $\boldsymbol{u} = B - C$ and $\boldsymbol{v} = A - C$ and observe that $\boldsymbol{v} - \boldsymbol{u} = A - B$.

$$\cos \theta = \frac{a^2 + b^2 - c^2}{2ab} = \frac{\|\boldsymbol{u}\|^2 + \|\boldsymbol{v}\|^2 - \|\boldsymbol{v} - \boldsymbol{u}\|^2}{2\|\boldsymbol{u}\| \|\boldsymbol{v}\|}$$
$$= \frac{\boldsymbol{u} \cdot \boldsymbol{u} + \boldsymbol{v} \cdot \boldsymbol{v} - (\boldsymbol{v} - \boldsymbol{u}) \cdot (\boldsymbol{v} - \boldsymbol{u})}{2\|\boldsymbol{u}\| \|\boldsymbol{v}\|}$$
$$= \frac{\boldsymbol{u} \cdot \boldsymbol{u} + \boldsymbol{v} \cdot \boldsymbol{v} - \boldsymbol{v} \cdot \boldsymbol{v} + 2\boldsymbol{u} \cdot \boldsymbol{v} - \boldsymbol{u} \cdot \boldsymbol{u}}{2\|\boldsymbol{u}\| \|\boldsymbol{v}\|}$$
$$= \frac{\boldsymbol{u} \cdot \boldsymbol{v}}{\|\boldsymbol{u}\| \|\boldsymbol{v}\|}$$



Dot product: Test for orthogonality

Two nonzero vectors $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^n$ are orthogonal iff $\boldsymbol{u} \cdot \boldsymbol{v} = 0$.

- Since $0 = \boldsymbol{u} \cdot \boldsymbol{v} = \|\boldsymbol{u}\| \|\boldsymbol{v}\| \cos \theta$ and $\|\boldsymbol{u}\|, \|\boldsymbol{v}\| > 0$, we have $\cos \theta = 0$
- And this means that the angle is 90 degrees



Dot product: Cosine similarity (1)

The angle between u and v is another way to measure the similarity between two vectors. The cosine similarity of u and v is given by

$$\cos(\boldsymbol{u},\boldsymbol{v})=\frac{\boldsymbol{u}\cdot\boldsymbol{v}}{\|\boldsymbol{u}\|\|\boldsymbol{v}\|}.$$

•
$$-1 \leq \cos(u, v) \leq 1$$

- Vectors that point in roughly the same direction \rightarrow small angle \rightarrow cosine similarity ≈ 1
- Vectors that point in roughly opposite directions ightarrow large angle ightarrow cosine similarity pprox -1
- Vectors that are roughly orthogonal \rightarrow roughly right angle \rightarrow cosine similarity ≈ 0
- Popular in IR to determine the similarity between a document and a query

Dot product: Cosine similarity (2)



Dot product: Pearson correlation

The sample Pearson correlation coefficient is a measure of linear correlation. It is given by

$$r_{\mathbf{x},\mathbf{y}} = rac{(\mathbf{x} - ar{\mathbf{x}}) \cdot (\mathbf{y} - ar{\mathbf{y}})}{\|\mathbf{x} - ar{\mathbf{x}}\| \|\mathbf{y} - ar{\mathbf{y}}\|}.$$

- Numerator proportional to the sample covariance
- Denominator proportional to sample standard deviations
- Closely related to cosine similarity but performs centering
 - This is sometimes desired
 - And sometimes a bad idea (e.g., example last slide)



Dot product: Similarity

The dot product itself can also be seen as a measure of similarity or compatibility. Recall

 $\boldsymbol{u} \cdot \boldsymbol{v} = \|\boldsymbol{u}\| \|\boldsymbol{v}\| \cos \theta.$

Example: Shopping transactions

- Like in previous example, vectors \boldsymbol{u} and \boldsymbol{v} correspond to persons
- Elements correponds to frequencies of buying each product
- We can think of the direction of a vector as "preference"
 - Which products are being bought?
 - $\cos \theta$ large when **u** and **v** have similar interest
- We can think of the magnitude of a vector as "strength"
 - How much is being bought?
 - $\|\boldsymbol{u}\| \|\boldsymbol{v}\|$ large when both persons buy a lot
- If $\boldsymbol{u} \cdot \boldsymbol{v}$ is large, \boldsymbol{u} and \boldsymbol{v} have similar shopping behaviour and buy a lot

Dot product: Projection

The vector projection of v onto u is given by



Outline

1. Vectors

2. Matrices

3. Summary

Notation

Let $\boldsymbol{A} \in \mathbb{R}^{m imes n}$ be a real m imes n matrix. We write

- a_{ij} or A_{ij} (both scalars) for the value of entry (i, j)
- \boldsymbol{a}_j or \boldsymbol{A}_{*j} (both column vectors) for the *j*-th column of \boldsymbol{A}
- *a_i* (column vector) or *A_{i*}* (row vector) for the *i*-th row of *A* Thus

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_{11} & \mathbf{a}_{12} & \cdots & \mathbf{a}_{1n} \\ \mathbf{a}_{21} & \mathbf{a}_{22} & \cdots & \mathbf{a}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_{m1} & \mathbf{a}_{m2} & \cdots & \mathbf{a}_{mn} \end{pmatrix}$$
$$= (\mathbf{A}_{*1} \ \mathbf{A}_{*2} \ \cdots \ \mathbf{A}_{*n}) = \begin{pmatrix} \mathbf{A}_{1*} \\ \mathbf{A}_{2*} \\ \vdots \\ \mathbf{A}_{m*} \end{pmatrix}$$

Full matrix ring (addition)

The set of all matrices in $\mathbb{R}^{n \times n}$ form a ring, the **full matrix ring**.

• Addition and substraction are element-wise

$$[\mathbf{A} + \mathbf{B}]_{ij} = a_{ij} + b_{ij}$$

 $[\mathbf{A} - \mathbf{B}]_{ij} = a_{ij} - b_{ij}$

- The additive identity is the $n \times n$ zero matrix $\mathbf{0}_{n \times n}$
- The additive inverse is $-\mathbf{A}$ with $[-\mathbf{A}]_{ij} = -a_{ij}$
- In general $[c A]_{ij} = c a_{ij}$ for $c \in \mathbb{R}$ (scalar multiplication)
- Addition is associative and commutative

Full matrix ring (multiplication)

• For multiplication, we take dot products

$$[\boldsymbol{A}\boldsymbol{B}]_{ij}=\boldsymbol{a}_i\cdot\boldsymbol{b}_j=\sum_{k=1}^na_{ik}b_{kj}$$

• The multiplicative identity is the *n* × *n* identity matrix *I*_n

$$\boldsymbol{I}_n = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$



- Multiplication is associative, but not commutative $(AB \neq BA \text{ in general})$
- Multiplication distributes over addition (A(B + C) = AB + AC and (B + C)A = BA + CA)
- Multiplication does not always have an inverse (division)

Rectangular matrices

- We generally have rectangular matrices $oldsymbol{A} \in \mathbb{R}^{m imes n}$
- We can only add and substract matrices of the same dimensions $(A_{m \times n} + B_{m \times n})$
- We can only multiply matrices with a matching inner dimension
 - We can multiply $\mathbf{A} \in \mathbb{R}^{m \times r}$ with $\mathbf{B} \in \mathbb{R}^{r \times n}$ (inner dimension is r)
 - Gives an $m \times n$ matrix (outer dimensions)

•
$$[\boldsymbol{A}\boldsymbol{B}]_{ij} = \boldsymbol{a}_i \cdot \boldsymbol{b}_j = \sum_{k=1}^r a_{ik} b_{kj}$$



Interpretation for matrix multiplication (1)

When we multiply A and B, we compute all dot products between rows of A and columns of B.

- We can apply any of the interpretations of the dot product
- E.g., weighted sum
 - m supermarkets, r products, n persons
 - a_{ik} = price of product k at supermarket i
 - b_{kj} = quantity of product k bought by person i
 - [AB]_{ij} = how much the j-th person would pay when buying at the i-th supermarket
- E.g., covariance
 - ► If all columns of $\mathbf{A}_{m \times n}$ are centered $(\sum_{k} a_{ik} = 0)$, then $\frac{1}{m-1} \mathbf{A}^T \mathbf{A} \in \mathbb{R}^{n \times n}$ is the sample covariance matrix
 - $\left[\frac{1}{m-1} \mathbf{A}^T \mathbf{A}\right]_{ii}$ holds the sample variance of column *i*
 - $\left[\frac{1}{m-1} \mathbf{A}^T \mathbf{A}\right]_{ij}$ holds the sample covariance between columns *i* and *j*



Interpretation for matrix multiplication (2)

We can also interpret rows i of AB as a linear combination of the rows of B with the coefficients coming from A

$$[\boldsymbol{A}\boldsymbol{B}]_{i*} = a_{i1}\boldsymbol{B}_{1*} + a_{i2}\boldsymbol{B}_{2*} + \cdots + a_{ir}\boldsymbol{B}_{r*},$$

and, similarly, the columns of $\boldsymbol{A}\boldsymbol{B}$ as linear combinations of the columns of \boldsymbol{A}

$$[\boldsymbol{A}\boldsymbol{B}]_{*j} = b_{1j}\boldsymbol{A}_{*1} + b_{2j}\boldsymbol{A}_{*2} + \cdots + b_{rj}\boldsymbol{A}_{*r}.$$



Interpretation for matrix multiplication (3)

We can view matrix AB as the sum of the *r* component matrices obtained by multiplying the *k*-th column of A and the *k*-th row of B:

$\boldsymbol{A}\boldsymbol{B} = \boldsymbol{A}_{*1}\boldsymbol{B}_{1*} + \boldsymbol{A}_{*2}\boldsymbol{B}_{2*} + \dots + \boldsymbol{A}_{*r}\boldsymbol{B}_{r*}$

- Components $A_{*k}B_{k*}$ are outer products ($m \times n$ matrices)
- Note: when $\boldsymbol{u} \in \mathbb{R}^m$ and $\boldsymbol{v} \in \mathbb{R}^n$, the matrix product
 - $\boldsymbol{u}^T \boldsymbol{v}$ corresponds to a dot product (a scalar), m = n required
 - uv^T corresponds to an outer product (an $m \times n$ matrix)
- In our supermarket example
 - Components correspond to products
 - Entry (i, j) of k-th component indicates how much the j-th person would pay for product k when buying at the i-th supermarket



Transposes

The matrix transpose \mathbf{A}^{T} switches rows and columns, i.e.,

$$[\mathbf{A}^{\mathsf{T}}]_{ij} = a_{ji}.$$

The following properties hold

- $(\boldsymbol{A}^T)^T = \boldsymbol{A}$
- $(\boldsymbol{A} + \boldsymbol{B})^T = \boldsymbol{A}^T + \boldsymbol{B}^T$
- $(c\mathbf{A})^T = c\mathbf{A}^T$
- $(\boldsymbol{A}\boldsymbol{B})^T = \boldsymbol{B}^T \boldsymbol{A}^T$

Summing and scaling

Let $A \in \mathbb{R}^{m \times n}$. Denote by $\mathbf{1}_n$ the all-ones vector of dimensionality n. For $s \in \mathbb{R}^n$, denote by diag (s) the $n \times n$ matrix with the entries of s on the main diagonal:

$$\operatorname{diag}(\boldsymbol{s}) = \begin{pmatrix} s_1 & 0 & \cdots & 0 \\ 0 & s_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_n \end{pmatrix}$$

- A1_n computes the row sums of A
- $\mathbf{1}_m^T \mathbf{A}$ computes the column sums of \mathbf{A}
- $A \operatorname{diag}(\boldsymbol{c})$ scales each column j of A by $c_j, \boldsymbol{c} \in \mathbb{R}^n$
- diag (\mathbf{r}) \mathbf{A} scales each row i of \mathbf{A} by r_i , $\mathbf{r} \in \mathbb{R}^m$

Matrices as linear maps

- A matrix $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ is a linear map from \mathbb{R}^n to \mathbb{R}^m
 - If $x \in \mathbb{R}^n$, then $y = Ax \in \mathbb{R}^m$ is the image of x
 - $\mathbf{y} = \sum_{j=1}^{n} \mathbf{a}_{j} x_{j}$, i.e., a linear combination of the columns of \mathbf{A}
- If $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$, then AB maps from \mathbb{R}^n to \mathbb{R}^m
 - Product AB corresponds to composition of linear maps A and B
- Square matrix $A \in \mathbb{R}^{n \times n}$ is invertible (= nonsingular) iff there is matrix $B \in \mathbb{R}^{n \times n}$ such that AB = I
 - ▶ Matrix **B** is the inverse of **A**, denoted **A**⁻¹
 - If **A** is invertible, then $AA^{-1} = A^{-1}A = I$

 $AA^{-1}x = A^{-1}Ax = x$

- Non-square matrices do not have (general) inverses but can have left or right inverses: AR = I or LA = I
- The transpose of $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a linear map $\mathbf{A}^T : \mathbb{R}^m \to \mathbb{R}^n$
 - $\blacktriangleright (\boldsymbol{A}^{T})_{ij} = \boldsymbol{A}_{ji}$
 - Generally, transpose is **not** the inverse $(\mathbf{A}\mathbf{A}^T \neq \mathbf{I})$

Matrix norms

- Matrix norms measure the magnitude of the matrix
 - Magnitude of the values in the matrix
 - Magnitude of the image
- Operator norms measure how large the image of an unit vector can get
 - Induced by a vector norm
 - ► For $p \ge 1$, $\|\boldsymbol{A}\|_p = \max\{\|\boldsymbol{A}\boldsymbol{x}\|_p \mid \|\boldsymbol{x}\|_p = 1\}$

 - $\left\| \boldsymbol{A} \right\|_{\infty} =$ maximum sum of absolute values of a row
 - Spectral norm: $\|\mathbf{A}\|_2 = \text{largest singular value of } \mathbf{A} \text{ (more later)}$
- The Frobenius norm is the vector- L_2 norm applied to matrices (treating them as a vector)

$$\bullet \|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$$

Note: ||A||_F ≠ ||A||₂ (but sometimes Frobenius norm is referred to as L₂ norm)

Matrix rank and linear independence

- A vector *u* ∈ ℝⁿ is linearly dependent on set of vectors
 V = {*v_i*} ⊂ ℝⁿ if *u* can be expressed as a linear combination of vectors in *V*
 - $\boldsymbol{u} = \sum_i a_i \boldsymbol{v}_i$ for some $a_1, \ldots, a_n \in \mathbb{R}$
 - Set V is linearly dependent if some v_i ∈ V is linearly dependent on V \ {v_i}
 - ► If V is not linearly dependent, it is linearly independent
- The column rank of matrix **A** is the maximum number of linearly independent columns of **A**
- The row rank of **A** is the maximum number of linearly independent rows of **A**
- The Schein rank of A is the least integer r such that A = LR for some $L \in \mathbb{R}^{m \times r}$ and $R \in \mathbb{R}^{r \times n}$
 - Equivalently, the least r such that A is a sum of r vector outer products
- All these ranks are equivalent
 - E.g., matrix has rank 1 iff it is an outer product of two (non-zero) vectors

Matrices as systems of linear equations

• A matrix can hold the coefficients of a system of linear equations (c.f. Chinese *Nine Chapters on Arithmetic*)

$$\begin{array}{c} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m \end{array} \Leftrightarrow \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}$$

- If the coefficient matrix **A** is invertible, the system has exact solution $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$
- If *m* < *n* the system is **underdetermined** and can have an infinite number of solutions
- If m > n the system is **overdetermined** and (usually) does not have an exact solution
- The least-squares solution is the vector \boldsymbol{x} that minimizes $\|\boldsymbol{A}\boldsymbol{x} \boldsymbol{b}\|_2^2$ (cf. linear regression)

Special types of matrices

- The diagonals of matrix A go from top-left to bottom-right
 - The main diagonal contains the elements a_{ii}
 - The k-th upper diagonal contains the elements $a_{i,(i+k)}$
 - The k-th lower diagonal contains the elements $a_{(i+k),i}$
 - The anti-diagonals go from top-right to bottom-left
- Matrix is **diagonal** if all its non-zero values are in a diagonal (typically main diagonal)
 - Bi-diagonal matrices have values in two diagonals, etc.
- Matrix **A** is **upper (right) triangular** if all of its non-zeros are in or above the main diagonal
 - Lower (left) triangular matrices have all non-zeros in or below main diagonal
 - Upper left and lower right triangular matrices: replace diagonal with anti-diagonal
- A square matrix *P* is **permutation matrix** if each row and each column of *P* has exactly one 1 and rest are 0s
 - ► If *P* is a permutation matrix, *PA* permutes the order of the rows and *AP* the order of the columns of *A*

Orthogonal matrices

- A set V = {v_i} ⊂ ℝⁿ is orthogonal if all vectors in V are mutually orthogonal
 - $\boldsymbol{v} \cdot \boldsymbol{u} = 0$ for all $\boldsymbol{v} \neq \boldsymbol{u} \in V$
 - ► If all vectors in V also have unit norm (||v||₂ = 1), V is orthonormal
- A square matrix **A** is orthogonal if its columns are a set of orthonormal (!) vectors or equivalently
 - Its rows are orthonormal
 - $\bullet \mathbf{A}^{\mathsf{T}}\mathbf{A} = \mathbf{I}_{\underline{n}}$
 - $\bullet \mathbf{A}^{-1} = \mathbf{A}^{\mathsf{T}}$
- An *m* × *n* matrix **A** is
 - ► column-orthogonal if columns are a set of orthonormal vectors (only possible if m ≥ n); then A^T is left inverse (A^TA = I_n)
 - ▶ row-orthogonal if rows are a set of orthonormal vectors (only possible if $m \le n$); then \mathbf{A}^T is right inverse $(\mathbf{A}\mathbf{A}^T = \mathbf{I}_m)$
- If **A** and **B** are orthogonal, so is **AB**
 - Similarly: column-orthogonality and row-orthogonality is preserved
Outline

- 1. Vectors
- 2. Matrices
- 3. Summary

Lessons learned

- Many uses, many interpretations
 - Vectors
 - Matrices
 - Dot products
 - Matrix products
- Magnitudes and distances are measured by norms
- Basic concepts of linear algebra
- Special types of matrices: diagonal, triangular, orthogonal

Suggested reading

• Any (elementary) linear algebra text book

 Carl Meyer Matrix Analysis and Applied Linear Algebra Society for Industrial and Applied Mathematics, 2000 http://www.matrixanalysis.com

- <u>Wolfram MathWorld</u> articles
- Wikipedia articles

Data Mining and Matrices 03 – Singular Value Decomposition

Prof. Dr. Rainer Gemulla

Universität Mannheim

FSS 2018

The SVD is the Swiss Army knife of matrix decompositions.

—Diane O'Leary, 2006

Outline

1. Definition

- 2. Properties of the SVD
- 3. Interpreting SVD
- 4. Using the SVD
- 5. Computing the SVD
- 6. Wrap-Up

Definition

Theorem

For each $\mathbf{A} \in \mathbb{R}^{m \times n}$, there are orthogonal matrices $\mathbf{U}_{m \times m}$, $\mathbf{V}_{n \times n}$, and a diagonal matrix $\mathbf{\Sigma}_{m \times n}$ with values $\sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_{\min(m,n)} \ge 0$ on the main diagonal such that $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$.

- $U\Sigma V^T$ is called the singular value decomposition (SVD) of **A**
- Values σ_i are the singular values of **A**
- Columns of **U** are the left singular vectors of **A**
- Columns of **V** are the right singular vectors of **A**



Outline

1. Definition

- 2. Properties of the SVD
- 3. Interpreting SVD
- 4. Using the SVD
- 5. Computing the SVD
- 6. Wrap-Up

Characterization of the four fundamental subspaces

The fundamental theorem of linear algebra states that every matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ induces four fundamental subspaces:

- The column space (range, image) of dimension rank $(\mathbf{A}) = r$
 - ► The set of all possible linear combinations of columns of **A**
- The nullspace (kernel) of dimension n r
 - The set of all vectors $\pmb{x} \in \mathbb{R}^n$ for which $\pmb{A}\pmb{x} = \pmb{0}$
- The row space (coimage) of dimension r
- The left nullspace (cokernel) of dimension m r

Explicit bases for these subspaces can be obtained from the SVD:

- Column space: the first r columns of U
- Null space: the last (n r) columns of **V**
- Row space: the first r columns of V
- Left null space: the last (m r) columns of **U**

The four fundamental subspaces



The action of **A**. Row space to column space, nullspace to zero.

Pseudo-inverse

Problem.

Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, find $\mathbf{x} \in \mathbb{R}^n$ minimizing $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$.

- If **A** is invertible, the solution is $\mathbf{A}^{-1}\mathbf{A}\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} \Leftrightarrow \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$
- A pseudo-inverse A^+ captures some properties of inverse A^{-1}
- The Moose-Penrose pseudo-inverse of A is a matrix $\mathbf{A}^+ \in \mathbb{R}^{n \times m}$ satisfying the following criteria
 - $AA^+A = A$ (but it is possible that $AA^+ \neq I$)

 - ► $\mathbf{A}^+ \mathbf{A} \mathbf{A}^+ = \mathbf{A}^+$ (cf. above) ► $(\mathbf{A} \mathbf{A}^+)^T = \mathbf{A} \mathbf{A}^+$ ($\mathbf{A} \mathbf{A}^+$ is symmetric)
 - $(\mathbf{A}^{+}\mathbf{A})^{T} = \mathbf{A}^{+}\mathbf{A} \quad (as is \mathbf{A}^{+}\mathbf{A})$
- If $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{T}$ is the SVD of \mathbf{A} , then $|\mathbf{A}^{+} = \mathbf{V} \mathbf{\Sigma}^{+} \mathbf{U}^{T}|$
 - Σ^+ replaces each $\sigma_i > 0$ by $1/\sigma_i$ and transposes

Theorem.

An optimum solution for the above problem is $\mathbf{x} = \mathbf{A}^+ \mathbf{b}$.

Pseudo inverse (illustration)



The inverse of **A** is (where possible) the pseudo-inverse A^+ .

Truncated SVD

- The rank of the matrix is the number of its non-zero singular values
 - Easy to see by writing $\mathbf{A} = \sum_{j=1}^{\min\{n,m\}} \sigma_j \mathbf{u}_j \mathbf{v}_j^T$
- The truncated SVD only takes the first k columns of U and V and the main k × k submatrix of Σ
 - $\mathbf{A}_{k} = \sum_{j=1}^{k} \sigma_{j} \mathbf{u}_{j} \mathbf{v}_{j}^{\mathsf{T}} = \mathbf{U}_{k} \mathbf{\Sigma}_{k} \mathbf{V}_{k}^{\mathsf{T}}$

• rank
$$(\boldsymbol{A}_k) = k$$
 (if $\sigma_k > 0$)

- ► U_k and V_k are not orthogonal anymore, but they are column-orthogonal
- If $k = \min\{m, n\}$, then $A_k = A$; called thin SVD
- If $k = \operatorname{rank}(A)$, then $A_k = A$; called compact SVD
- If $k < \operatorname{rank}(A)$, then A_k is low-rank approximation of A



SVD and matrix norms

Let
$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{T}$$
 be the SVD of \mathbf{A} . Then

•
$$\|\mathbf{A}\|_{F}^{2} = \sum_{i=1}^{\min\{n,m\}} \sigma_{i}^{2}$$

•
$$\|\bm{A}\|_2 = \sigma_1$$

• Remember: $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{\min\{n,m\}} \geq 0$

- Therefore $\| \boldsymbol{A} \|_2 \leq \| \boldsymbol{A} \|_F \leq \sqrt{n} \, \| \boldsymbol{A} \|_2$
- Sq. Frobenius norm of truncated SVD is $\|\mathbf{A}_k\|_F^2 = \sum_{i=1}^k \sigma_i^2$
 - And of the "approximation error" $\|\boldsymbol{A} \boldsymbol{A}_k\|_F^2 = \sum_{i=k+1}^{\min\{n,m\}} \sigma_i^2$

The Eckart–Young theorem

Let A_k be the rank-k truncated SVD of A. Then A_k is the rank-k matrix closest to A in the Frobenius sense. That is

$$\|\boldsymbol{A}-\boldsymbol{A}_k\|_F \leq \|\boldsymbol{A}-\boldsymbol{B}\|_F$$
 for all rank-k matrices \boldsymbol{B} .

Eigendecompositions

- An eigenvector of a square matrix **A** is a vector **v** such that **A** only changes the magnitude of **v**
 - I.e. $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ for some $\lambda \in \mathbb{R}$
 - Such λ is an **eigenvalue** of **A**
 - ► <u>Try it!</u>
- The eigendecomposition of **A** is $\mathbf{A} = \mathbf{Q} \Delta \mathbf{Q}^{-1}$
 - The columns of Q are the eigenvectors of A
 - Matrix $oldsymbol{\Delta}$ is a diagonal matrix with the eigenvalues
- Not every (square) matrix has eigendecomposition
 - ► If **A** is of form **BB**^T, it always has eigendecomposition
- The SVD of A is closely related to the eigendecompositions of AA^T and A^TA
 - ► The left singular vectors are the eigenvectors of **AA**^T_{_}
 - The right singular vectors are the eigenvectors of $\mathbf{A}^T \mathbf{A}$
 - ► The singular values are the square roots of the eigenvalues of both AA^T and A^TA

Outline

- 1. Definition
- 2. Properties of the SVD
- $3. \ Interpreting \ SVD$
- 4. Using the SVD
- 5. Computing the SVD
- 6. Wrap-Up

Factor interpretation

- The most common way to interpret SVD is to consider the columns of \boldsymbol{U} (or \boldsymbol{V})
 - Let **A** be objects-by-attributes and $U\Sigma V^{T}$ its SVD
 - ► If two rows have similar values in a column of *U*, these objects are somehow similar
 - ► If two columns have similar values in a row of V^T, these attributes are somehow similar
 - \blacktriangleright In both cases, first entries often matter most \rightarrow truncated SVD



Figure 3.2. The first two factors for a dataset ranking wines. Skillicorn, p. 55

- Example: people's ratings of different wines
- Scatterplot of first and second column of U
 - left: likes wine
 - right: doesn't like
 - up: prefers red wine
 - bottom: prefers white vine
- Conclusion: winelovers like red and white, others care 14/52

Example: Weather data ($k = 2$)										
1.00			9.05	16.55	26.73	18.75	17.81			
	1.00		-4.14	0.27	2.32	-0.89	-0.69			
				•		<u> </u>				
			Jan	Apr	Jul	Oct	Year			
0.62	1.69	Stockholm	-0.70	8.60	21.90	9.90	10.00			
0.69	2.11	Minsk	-2.10	12.20	23.60	10.20	10.60			
0.83	0.00	London	7.90	13.30	22.80	15.20	14.80			
0.90	1.52	Budapest	1.20	16.30	26.50	16.10	15.00			
0.88	0.30	Paris	6.90	14.70	24.40	15.80	15.50			
0.98	1.59	Bucharests	1.50	18.00	28.80	18.00	16.50			
1.09	-0.66	Barcelona	12.40	17.60	27.50	21.50	20.00			
1.14	-0.31	Rome	11.90	17.70	30.30	21.40	20.40			
1.16	-1.09	Lisbon	14.80	19.80	27.90	22.50	21.50			
1.24	-0.35	Athens	12.90	20.30	32.60	23.10	22.30			
1.21	-1.26	Valencia	16.10	20.20	29.10	23.60	22.30			
1.27	-1.12	Malta	16.10	20.00	31.50	25.20	23.20			

D

(RMSE: 0.60) 15/52

Thin SVD of Weather data (U)

		, 1	2	3	4	5 、
	Stockholm	0.18	0.41	0.61	0.28	-0.32
	Minsk	0.19	0.51	0.08	-0.54	0.40
	London	0.24	0.00	0.20	-0.15	0.04
	Budapest	0.25	0.37	-0.39	0.18	-0.10
U =	Paris	0.25	0.07	0.05	-0.25	-0.22
	Bucharests	0.28	0.39	-0.49	0.30	0.08
	Barcelona	0.31	-0.16	-0.01	0.33	-0.26
	Rome	0.32	-0.07	0.30	0.10	0.07
	Lisbon	0.33	-0.27	-0.23	-0.27	-0.46
	Athens	0.35	-0.08	0.10	-0.23	-0.09
	Valencia	0.34	-0.31	-0.12	-0.21	0.17
	Malta	\0.36	-0.27	0.12	0.37	0.59 /

Thin SVD of Weather data (Σ)

		_ 1	2	3	4	5 ्
Σ =	1/	[′] 147.55	0.00	0.00	0.00	0.00
	2	0.00	20.09	0.00	0.00	0.00
	3	0.00	0.00	4.25	0.00	0.00
	4	0.00	0.00	0.00	1.77	0.00
	5 \	0.00	0.00	0.00	0.00	0.32/

Thin SVD of Weather data (V)

		1	2	3	4	5
V =	Jan	0.22	-0.85	0.31	-0.30	0.21
	Apr	0.40	0.06	-0.74	-0.52	0.17
	Jul	0.64	0.47	0.54	-0.16	0.21
	Oct	0.45	-0.18	-0.25	0.78	0.30
	Year	\0.43	-0.14	-0.03	0.05	-0.8/

Example: Weather data ($k = 2$), SVD										
147.54			0.22	0.40	0.64	0.45	0.43			
	20.09		-0.85	0.06	0.47	-0.18	-0.14			
			Jan	Apr	Jul	Oct	Year			
0.18	0.41	Stockholm	-0.70	8.60	21.90	9.90	10.00			
0.19	0.51	Minsk	-2.10	12.20	23.60	10.20	10.60			
0.24	0.00	London	7.90	13.30	22.80	15.20	14.80			
0.25	0.37	Budapest	1.20	16.30	26.50	16.10	15.00			
0.25	0.07	Paris	6.90	14.70	24.40	15.80	15.50			
0.28	0.39	Bucharests	1.50	18.00	28.80	18.00	16.50			
0.31	-0.16	Barcelona	12.40	17.60	27.50	21.50	20.00			
0.32	-0.07	Rome	11.90	17.70	30.30	21.40	20.40			
0.33	-0.27	Lisbon	14.80	19.80	27.90	22.50	21.50			
0.35	-0.08	Athens	12.90	20.30	32.60	23.10	22.30			
0.34	-0.31	Valencia	16.10	20.20	29.10	23.60	22.30			
0.36	-0.27	Malta	16.10	20.00	31.50	25.20	23.20			

D

(RMSE: 0.60) 19/52

Example: Weather data (r = 2), SVD plot



Orthogonal matrices and rotations

- Orthogonal matrices are rotation matrices
- Consider orthogonal matrix $oldsymbol{Q}$
- Inner products are retained: $(\mathbf{Q}\mathbf{x})^T (\mathbf{Q}\mathbf{y}) = \mathbf{x}^T \mathbf{Q}^T \mathbf{Q}\mathbf{y} = \mathbf{x}^T \mathbf{y}$
- Thus Euclidean norms also retained: $\|\boldsymbol{Q}\boldsymbol{x}\| = \|\boldsymbol{x}\|$
- Implies that all angles are retained
- In 2D: $\boldsymbol{Q}_{\theta} = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$ • Consider vector $\begin{pmatrix} x \\ y \end{pmatrix} = x \boldsymbol{e}_1 + y \boldsymbol{e}_2$ • $\boldsymbol{Q}_{\theta} \begin{pmatrix} x \\ y \end{pmatrix} = x [\boldsymbol{Q}_{\theta}]_{*1} + y [\boldsymbol{Q}_{\theta}]_{*2}$
- Thus: the columns of Q form "new axes" for rotation Qv
- Similarly: rows of *Q* form "new axes" for rotation *v^TQ* (rotates backwards)



Geometric interpretation (1)



Geometric interpretation (2)



- Let **UΣV**^T be the SVD of M
- SVD shows that every linear mapping y = Mxcan be considered as a series of rotation, stretching, and rotation operations
 - Matrix V^T performs the first rotation y₁ = V^Tx
 - Matrix Σ performs the stretching y₂ = Σy₁
 - Matrix U performs the second rotation y = Uy₂

Dimension of largest variance (1)



- The singular vectors give the dimensions of the variance in the data
 - The first singular vector is the dimension of the largest variance

Dimension of largest variance (2)



- 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 largest variance
 - First two dimensions span a hyperplane
- From Eckart-Young we know that if we project the data to the spanned hyperplanes, the distance of the projection is minimized

SVD and linear regression (1)

Consider the 2-dimensional case.

- When fitting a linear regression model, we want to predict the value of one distinguished variable (the response variable) given the others (the explanatory variable)
- Say y is the response variable, x the explanatory variable
- We obtain a coefficient vector $\boldsymbol{\beta} = \begin{pmatrix} \beta_0 & \beta_1 \end{pmatrix}^T \in \mathbb{R}^2$, consisting of offset β_0 and slope β_1
- Prediction is $\hat{y} = \beta \cdot \begin{pmatrix} 1 & x \end{pmatrix}^T$, the error is $(y - \hat{y})^2$ (least squares)
- Goal is to minimize this error, i.e., $\| \boldsymbol{y} \hat{\boldsymbol{y}} \|_2$
- Set $oldsymbol{X} = oldsymbol{\left(1 \ x
 ight)}$, then $oldsymbol{eta} = oldsymbol{X}^+ oldsymbol{y}$



SVD and linear regression (2)

- Contrast this to the rank-1 truncated SVD of $m{A} = egin{pmatrix} m{x} & m{y} \end{pmatrix}$
- We obtain a vector $oldsymbol{u}_1$, a scaling coefficent σ_1 , and a vector $oldsymbol{v}_1$
- Let's look at the line described by \mathbf{v}_1 (roughly corresponds to eta)
- Reconstructed data is $\mathbf{A}_1 = \mathbf{u}_1 \sigma_1 \mathbf{v}_1^T$; all points lie on the line
- There is no distinguished response variable; we minimize distance to line instead; i.e., ||A - A₁||_F.
- This is different from regression



SVD and linear regression (3)



Component interpretation

- Recall that we can write $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \sum_{i=1}^r \mathbf{A}_i$ • $\mathbf{A}_i = \sigma_i \mathbf{v}_i \mathbf{u}_i^T$
- This explains the data as a sums of (rank-1) layers
 - The first layer explains the most
 - The second corrects that by adding and removing smaller values
 - The third corrects that by adding and removing even smaller values
 ...
- The layers don't have to be very intuitive

Outline

- 1. Definition
- 2. Properties of the SVD
- 3. Interpreting SVD
- 4. Using the SVD
- 4.1 How many factors?
- 4.2 Using SVD: Data preprocessing and visualization
- 5. Computing the SVD
- 6. Wrap-Up

Outline

- 1. Definition
- 2. Properties of the SVD
- 3. Interpreting SVD
- 4. Using the SVD
- 4.1 How many factors?
- 4.2 Using SVD: Data preprocessing and visualization
- 5. Computing the SVD
- 6. Wrap-Up

Problem

- Most data mining applications do not use full SVD, but truncated SVD
 - To concentrate on "the most important parts"
- But how to select the rank k of the truncated SVD?
 - What is important, what is unimportant?
 - What is structure, what is noise?
 - Too small rank: all subtlety is lost
 - Too big rank: all smoothing is lost
- Typical methods rely on singular values in a way or another

Guttman-Kaiser criterion and captured energy

- Perhaps the oldest method is the Guttman-Kaiser criterion:
 - Select k so that for all i > k, $\sigma_i < 1$
 - Motivation: all components with singular value less than unit are uninteresting
- Another common method is to select enough singular values such that the sum of their squares is 90% of the total sum of the squared singular values
 - ▶ The exact percentage can be different (80%, 95%)
 - Motivation: The resulting matrix "explains" 90% of the (sq.) Frobenius norm of the matrix
- **Problem:** Both of these methods are based on arbitrary thresholds and do not consider the "shape" of the data
Cattell's Scree test

- The scree plot plots the squared singular values in decreasing order
 - The plot looks like a side of the hill, hence the name
- The scree test is a subjective decision on the rank based on the shape of the scree plot
- The rank should be set to a point where
 - there is a clear drop in the magnitudes of the values values; or
 - the values start to even out
- **Problem:** Scree test is subjective, and many data don't have any clear shapes to use (or have many)
 - Automated methods have been developed to detect the shapes from the scree plot





Entropy-based method

• Consider the relative contribution of each singular value to the overall (sq.) Frobenius norm

• Relative contribution of σ_k is $f_k = \sigma_k^2 / \sum_i \sigma_i^2$

• We can treat these as probabilities and define the (normalized) entropy of the singular values as

$$E = -\frac{1}{\log\left(\min\{n, m\}\right)} \sum_{i=1}^{\min\{n, m\}} f_i \log f_i$$

- The basis of the logarithm doesn't matter
- We assume that $0 \cdot \infty = 0$
- Low entropy (close to 0): the first singular value has almost all mass
- ▶ High entropy (close to 1): the singular values are almost equal
- The rank is selected to be the smallest k such that $\sum_{i=1}^{k} f_i \ge E$
- Problem: Why entropy?

Random flip of signs

- Multiply every element of the data ${\pmb A}$ randomly with either 1 or -1 to get $\tilde{\pmb A}$
 - The Frobenius norm doesn't change $(\|\mathbf{A}\|_F = \|\mathbf{\tilde{A}}\|_F)$
 - The spectral norm does change $(\|\mathbf{A}\|_2 \neq \|\mathbf{\tilde{A}}\|_2)$
 - ► How much this changes depends on how much "structure" **A** has
- We try to select k such that the residual matrix contains only noise
 - ► The residual matrix contains the last m k columns of U, min{n,m} - k singular values, and last n - k rows of V^T
 - If \mathbf{A}_{-k} is the residual matrix of \mathbf{A} after rank-k truncated SVD, and $\tilde{\mathbf{A}}_{-k}$ is obtained from \mathbf{A}_{-k} by randomly flipping signs, we select rank k to be such that $(\|\mathbf{A}_{-k}\|_2 \|\tilde{\mathbf{A}}_{-k}\|_2)/\|\mathbf{A}_{-k}\|_F$ is small
- Problem: How small is small?

Outline

- 1. Definition
- 2. Properties of the SVD
- 3. Interpreting SVD
- 4. Using the SVD
- 4.1 How many factors?
- 4.2 Using SVD: Data preprocessing and visualization
- 5. Computing the SVD
- 6. Wrap-Up

Normalization

- SVD is sensitive to data scaling
- Data should usually be normalized before SVD is applied
 - If one attribute is height in meters and another weights in grams, weight seems to carry much more importance in data about humans
 - If data is all positive, the first singular vector just explains where in the positive quadrant the data is
- The *z*-scores are attributes whose values are transformed by
 - Centering them to 0
 - Remove the mean of the attribute's values from each value
 - Normalizing the magnitudes
 - Divide every value with the standard deviation of the attribute
- Notice that the *z*-scores assume that
 - all attributes are equally important
 - attribute values are approximately normally distributed
- Values that have larger magnitude than importance can also be normalized by first taking logarithms (from positive values) or cubic roots
- The effects of normalization should always be considered

Relationship to PCA (1)

- Truncated SVD can also be used to battle the **curse of dimensionality**
 - All points are far from each other in very high dimensional spaces
 - High dimensionality slows down data mining algorithms
 - If we use the truncated SVD, every object is represented by its row in U_k (k values instead of n)
 - If $k \ll n$, we performed dimensionality reduction
- SVD is closely related to principal components analysis (PCA)
 - Key difference is "centering" of the data
 - 1. Center each column of \boldsymbol{A} to obtain \boldsymbol{M}
 - 2. Compute the sample covariance matrix $\boldsymbol{S} = \boldsymbol{M}^T \boldsymbol{M} / (m-1)$
 - 3. Compute the eigendecomposition $\boldsymbol{S} = \boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^{T}$ s.t. \boldsymbol{Q} orthogonal
 - ► The columns of **Q** are called principal components
 - \blacktriangleright The corresponding eigenvalues in Λ are the component weights
 - ► <u>Try it!</u>

Relationship to PCA (2)

- Relationship between SVD and PCA
 - $\boldsymbol{Q} = \text{eigenvectors of } \boldsymbol{M}^T \boldsymbol{M} / (m-1), \boldsymbol{\Lambda} = \text{corresponding eigenvalues}$
 - Let $\boldsymbol{M} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T$ be the SVD of \boldsymbol{M}
 - The SVD of $\boldsymbol{M}/\sqrt{m-1}$ is then $\boldsymbol{U}(\boldsymbol{\Sigma}/\sqrt{m-1})\boldsymbol{V}^{T}$
 - From slide 12, we know that $\boldsymbol{Q} = \boldsymbol{V}$
 - From slide 12, we know that $\Sigma^2/(m-1) = \Lambda$
- PCA associates each data object with a set of scores
 - One per principal component
 - $m \times n$ "score matrix" given by MQ
 - We have: $\boldsymbol{M}\boldsymbol{Q} = (\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T)\boldsymbol{V} = \boldsymbol{U}\boldsymbol{\Sigma}$
 - Usually we only take first k-components: $MQ_k = U_k \Sigma_k$
 - Also known as the Karhunen–Loève transform (KLT) of the rows of M

Relationship to PCA (3)



 $\mathbf{MQ}_{2} = \mathbf{U}_{2}\mathbf{\Sigma}_{2}$

What is the difference between the PCA and the SVD of $A \ (\neq M)$?

Relationship to PCA (4)

- Assumptions of PCA
 - Linearity: The data set is a linear combination of the (retained) components.
 - Large variances indicate important structure: PCA scores may or may not be suitable for subsequent tasks (e.g., classification).
 - Components are orthogonal
- Recall that [*M^TM/(m-1)*]_{ij} holds the sample covariance between attributes *i* and *j* (i.e., columns *A_{*i}* and *A_{*j}*)
 - ▶ When each row of *M* is an i.i.d. observation of some underlying distribution, then *M^TM/(m 1)* is an *unbiased estimator* of the covariance matrix of that distribution

Relationship to PCA (5)

- For our example data, \boldsymbol{M} is a 50 imes 2 matrix
- Here is its sample covariance matrix

$$\boldsymbol{M}^{T}\boldsymbol{M}/(m-1) = \begin{bmatrix} 2.1 & 1.16 \\ 1.16 & 0.72 \end{bmatrix}$$

• And here the sample covariance matrix of the PCA scores

$$(\boldsymbol{M}\boldsymbol{Q})^{T}(\boldsymbol{M}\boldsymbol{Q})/(m-1) = egin{array}{ccc} 2.76 & 0 \ 0 & 0.06 \ \end{array}$$

• This is "by construction"

Relationship to PCA (6)

- Only information used by PCA is sample means and covariances
 - > Data assumed i.i.d. from a multivariate normal distribution
 - Mean and covariance parameters are estimated from data (and ordered)
 - ► We obtain distribution N(0, Λ) after rotation (with correct estimates and assuming normality assumption holds)
- We already established that $f_k = \sigma_k^2 / \sum_i \sigma_i^2$ can be seen as the relative contribution of a singular value to the sq. Frobenius norm
- We also know that $\sigma_k^2/(m-1) = \lambda_k$
- Implies that f_k can also be seen as the fraction of the total variance explained by the *k*-th principal component
 - Interpretation only valid when data is centered (!)
 - ▶ In a Scree plot, one may plot f_k (or their cumulative sum) instead of the σ_k^2
- To make sure PCA is adequate, ensure that data "looks" normally distributed (or transform it to look normal)

Removing noise

- Very common application of SVD is to remove the noise from the data
- This works simply by taking the truncated SVD from the (normalized) data
 - The big problem is to select the rank of the truncated SVD
- Example:



- Original data
 - Looks like 1-dimensional with some noise
- The right singular vectors show the directions
 - The first looks like the data direction
 - The second looks like the noise direction
- The singular values confirm this

Visualization

. . .

- Truncated SVD with k = 2,3 allows us to visualize the data
 - We can plot the projected data points after 2D or 3D PCA
 - Or we can plot the scatter plot of two or three singular vectors
 - Or we color data points based on their entries in a singular vector





Figure 3.2. The first two factors for a dataset ranking wines.

Latent semantic analysis

- The latent semantic analysis (LSA) is an information retrieval method that uses SVD
- The data: a term-document matrix **A**
 - Values are (weighted) term frequencies
 - Typically tf-idf values (the frequency of the term in the document divided by the global frequency of the term)
- The truncated SVD $\boldsymbol{A}_k = \boldsymbol{U}_k \boldsymbol{\Sigma}_k \boldsymbol{V}_k^T$ of \boldsymbol{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 same things"
- A query q can be answered by considering its term vector \boldsymbol{q}
 - **q** is projected to $\mathbf{q}_k = \mathbf{q} \mathbf{V} \mathbf{\Sigma}^{-1}$ (called: fold in)
 - \boldsymbol{q}_k is compared to rows of \boldsymbol{U} and most similar rows are returned
- Later more (when we consider NMF and LDA)

Outline

- 1. Definition
- 2. Properties of the SVD
- 3. Interpreting SVD
- 4. Using the SVD
- 5. Computing the SVD
- 6. Wrap-Up

Algorithms for SVD

- In principle, the SVD of *A* can be computed by computing the eigendecomposition of *AA*^T
 - > This gives us left singular vectors and squares of singular values
 - Right singular vectors can be solved: $\boldsymbol{V}^{T} = \boldsymbol{\Sigma}^{-1} \boldsymbol{U}^{T} \boldsymbol{A}$
 - Bad for <u>numerical stability</u>!
- Full SVD can be computed in time $O(nm\min\{n, m\})$
 - Matrix A is first reduced to a bidiagonal matrix
 - The SVD of the bidiagonal matrix is computed using iterative methods (similar to eigendecompositions)
- Methods that are fast in practice exist
 - Especially for truncated SVD
- Efficient implementation of an SVD algorithm requires considerable work and knowledge
 - Luckily (almost) all numerical computation packages and programs implement SVD

Outline

- 1. Definition
- 2. Properties of the SVD
- 3. Interpreting SVD
- 4. Using the SVD
- 5. Computing the SVD
- 6. Wrap-Up

Lessons learned

- SVD is the Swiss Army knife of (numerical) linear algebra
 → Ranks, kernels, norms, inverses,
- SVD is also very useful in data analysis \rightarrow Dimensionality reduction, noise removal, visualization, \ldots
- Truncated SVD is best low-rank factorization of the data in terms of Frobenius norm
- Selecting the correct rank for truncated SVD is still a problem
- Interpretation of results can be challenging

Suggested reading

- Skillicorn, Ch. 3
- Meyer, Ch. 5.12
- Gene H. Golub & Charles F. Van Loan: *Matrix Computations*, 3rd ed. Johns Hopkins University Press, 1996
 - ► Excellent source for the algorithms and theory, but very dense

Data Mining and Matrices 04 – Matrix Completion

Prof. Dr. Rainer Gemulla

Universität Mannheim

FSS 2018

Recap: Singular Value Decomposition

- SVD is useful in data analysis
 - \rightarrow Noise removal, visualization, dimensionality reduction, \ldots
- Provides a means to understand the hidden structure in the data

We may think of A_k and its factor matrices as a low-rank model of the data:

- Used to capture the important aspects of the data (cf. principal components)
- Ignores the rest
- Truncated SVD is best low-rank factorization of the data in terms of Frobenius norm
- Truncated SVD $\mathbf{A}_k = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T$ of \mathbf{A} thus satisfies

$$\|\boldsymbol{A}-\boldsymbol{A}_k\|_F = \min_{\mathrm{rank}(\boldsymbol{B})=k} \|\boldsymbol{A}-\boldsymbol{B}\|_F$$

Incomplete data

- If all entries of the data matrix are available, computing a low-rank model is easy
- We now look at the case of partially observed data
 - Many data mining methods cannot handle missing data well
 - Need for imputation
- Our goal is to build a low-rank model...
 - …and ultimately approximately recover the full matrix
 - Cannot be done in general; assumptions needed
 - Ongoing research topic, large body of results
- Some reasons for missing entries
 - Failure in data acquisition processes
 - Expensive to obtain all entries
 - Some entries cannot be measured

	Name	Sex	Age	Height	Weight
1	Alfred	м	14	69	112.5
2	Alice		13	56.5	84
3	Barbara	F	13	65.3	98
4	Carol	F	14		102.5
5	Henry	м	14	63.5	102.5
6	James	м	12	57.3	83
7	Jane	F	12		84.5
8	Janet	F	15	62.5	112.5
9	Jeffrey	м	13	62.5	
10	John	м	12	59	99.5
11	Joyce	F	11	51.3	50.5
12	Judy	F	14	64.3	90
13	Louise	F	12	56.3	77
14		F	15	66.5	112
15	Philip	м	16	72	150
	Robert	M	12	64.8	128

Outline

- 1. Collaborative Filtering
- 2. Matrix Completion
- 3. Algorithms
- 4. Variants
- 5. Summary

Outline

1. Collaborative Filtering

- 2. Matrix Completion
- 3. Algorithms
- 4. Variants
- 5. Summary

Recommender systems

- Problem
 - Set of users
 - Set of items (movies, books, jokes, products, stories, ...)
 - ► Feedback (ratings, purchase, click-through, tags, ...)
 - ► Sometimes: metadata (user profiles, item properties, ...)
- Goal: Predict preferences of users for items
- Ultimate goal: Create item recommendations for each user
- Example

	Avatar	The Matrix	Up
Alice (?	4	2
Bob	3	2	?
<i>Charlie</i> \	5	?	3/

Collaborative filtering

- Key idea: Make use of past user behavior
- No domain knowledge required
- No expensive data collection needed
- Allows discovery of complex and unexpected patterns
- Widely adopted: Amazon, TiVo, Netflix, Microsoft
- Key techniques: neighborhood models, latent factor models

$$\begin{array}{c} \text{Avatar} \quad \text{The Matrix} \quad Up\\ \text{Alice} \left(\begin{array}{cc} ? & 4 & 2\\ Bob \\ 3 & 2 & ?\\ \text{Charlie} & 5 & ? & 3 \end{array}\right)$$

Leverage past behavior of other users and/or on other items.

A simple baseline

- *m* users, *n* items, $m \times n$ rating matrix **D**
- Revealed entries $\Omega = \{ (i, j) | \text{ rating } d_{ij} \text{ is revealed } \}, N = |\Omega|$
- Baseline predictor: $b_{ij} = \mu + b_i + b_j$
 - $\mu = \frac{1}{N} \sum_{(i,j) \in \Omega} d_{ij}$ is the overall average rating
 - b_i is a user bias (user's tendency to rate low/high)
 - b_j is an item bias (item's tendency to be rated low/high)

• Least squares estimates: $\operatorname{argmin}_{b_*}\sum_{(i,j)\in\Omega}(d_{ij}-\mu-b_i-b_j)^2$

When does a user like an item?

• Neighborhood models (kNN):

When the user likes similar items

- Find the top-k most similar items the user has rated
- Combine the ratings of these items (e.g., average)
- ▶ Requires a similarity measure (e.g., Pearson correlation coefficient)



is similar to



Bob rated 4

- Unrated by Bob \rightarrow predict 4
- Latent factor models (LFM): When *similar* users like similar items
 - More holistic approach
 - Users and items are placed in the same "latent factor space"
 - Position of a user and an item related to preference (via dot products)



Intuition behind latent factor models (1)



Intuition behind latent factor models (2)

- Does user **u** like item **v**?
- Quality: measured via direction from origin $(\cos \angle(u, v))$
 - Same direction ightarrow attraction: $\cos \angle (oldsymbol{u},oldsymbol{v}) pprox 1$
 - ▶ Opposite direction ightarrow repulsion: cos $\angle(\textit{u},\textit{v}) pprox -1$
 - ▶ Orthogonal direction \rightarrow oblivious: cos∠($m{u},m{v}$) pprox 0
- Strength: measured via distance from origin $(\|\boldsymbol{u}\| \|\boldsymbol{v}\|)$
 - ▶ Far from origin \rightarrow strong relationship: $\|\boldsymbol{u}\| \|\boldsymbol{v}\|$ large
 - \blacktriangleright Close to origin \rightarrow weak relationship: $\| \textbf{\textit{u}} \| \, \| \textbf{\textit{v}} \|$ small
- Overall preference: measured via dot product $(\boldsymbol{u} \cdot \boldsymbol{v})$

$$\boldsymbol{u} \cdot \boldsymbol{v} = \|\boldsymbol{u}\| \|\boldsymbol{v}\| \frac{\boldsymbol{u} \cdot \boldsymbol{v}}{\|\boldsymbol{u}\| \|\boldsymbol{v}\|} = \|\boldsymbol{u}\| \|\boldsymbol{v}\| \cos \angle (\boldsymbol{u}, \boldsymbol{v})$$

- ▶ Same direction, far out → strong attraction: $\boldsymbol{u} \cdot \boldsymbol{v}$ large positive
- $\blacktriangleright \ \ \mathsf{Opposite \ direction, \ far \ out } \rightarrow \mathsf{strong \ repulsion:} \ \ \boldsymbol{u}\cdot\boldsymbol{v} \ \ \mathsf{large \ negative}$
- Orthogonal direction, any distance ightarrow oblivious: : ${m u}\cdot {m v}pprox 0$

But how to select dimensions and where to place items and users? Key idea: Pick dimensions that explain the known data well.

Latent factor models (simple form)

(1.98)

Bob

(1.21)

Charlie

(2.30)

L

(4.4)

3

(2.7)

(5.2)

5

• Given rank r, find $m \times r$ matrix L and $r \times n$ matrix R such that

$$d_{ij} \approx [LR]_{ij} \quad \text{for } (i,j) \in \Omega$$
• Least squares formulation
$$\min_{L,R} \sum_{(i,j)\in\Omega} (d_{ij} - [LR]_{ij})^2$$
• Example $(r = 1)$

$$R$$

$$\frac{Avatar \quad The \; Matrix \quad Up}{(2.24) \quad (1.92) \quad (1.18)}$$

$$Alice \quad ? \qquad \mathbf{4} \qquad \mathbf{2}$$

(3.8)

(2.3)

2

?

(4.4)

(2.3)

(1.4)

3

(2.7)



R

SVD and missing values

Input data



10% of input data

Rank-10 truncated SVD



Rank-10 truncated SVD



Latent factor models and missing values



Input data

Rank-10 LFM



Rank-10 LFM



Example: Netflix prize data





Koren et al., 2009

Latent factor models (summation form)

- Least squares formulation prone to overfitting
- More general summation form:

$$L = \sum_{(i,j)\in\Omega} L_{ij}(\boldsymbol{I}_i, \boldsymbol{r}_j) + R(\boldsymbol{L}, \boldsymbol{R}),$$

- L is global loss
- ► *I_i* and *r_j* are user and item parameters, resp.
- L_{ij} is local loss, e.g., $L_{ij} = (d_{ij} [LR]_{ij})^2$
- *R* is regularization term, e.g., $R = \lambda (\|\boldsymbol{L}\|_F^2 + \|\boldsymbol{R}\|_F^2)$
- As before, goal is to solve $\min_{L,R} L(L, R)$
- Loss function can be more sophisticated
 - Improved predictors (e.g., include user and item bias)
 - Additional feedback data (e.g., time, implicit feedback)
 - Regularization terms (e.g., weighted depending on amount of feedback)
 - Available metadata (e.g., demographics, genre of a movie)



Example: Netflix prize data



Root mean square error of predictions
Outline

- 1. Collaborative Filtering
- 2. Matrix Completion
- 3. Algorithms
- 4. Variants
- 5. Summary

The matrix completion problem

Complete these matrices!



Matrix completion is impossible without additional assumptions!

Let's assume that underlying full matrix is "simple" (here: rank 1).

/1	1	1	1	1	/1	1	1	1	1\
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1/	1	1	1	1	1/

When/how can we recover a low-rank matrix from a sample of its entries?

Rank minimization

Definition (rank minimization problem)

Given an $n \times n$ data matrix **D** and an index set Ω of revealed entries. The *rank minimization problem* is

$$\begin{array}{ll} \text{minimize} & \operatorname{rank}(\boldsymbol{X}) \\ \text{subject to} & d_{ij} = x_{ij} \\ \boldsymbol{X} \in \mathbb{R}^{n \times n}. \end{array} (i, j) \in \Omega$$

- Seeks for "simplest explanation" fitting the data
- If unique and sufficient samples, recovers \boldsymbol{D} (i.e., $\boldsymbol{X} = \boldsymbol{D}$)
- NP-hard

Time complexity of existing rank minimization algorithms double exponential in n (and also slow in practice).

Nuclear norm minimization

- Denote by $\boldsymbol{\sigma} = \begin{pmatrix} \sigma_1 & \sigma_2 & \dots & \sigma_n \end{pmatrix}^T$ the singular values of \boldsymbol{X}
- $\operatorname{rank}(X) = |\{\sigma_k > 0\}| = \sum_{k=1}^n I_{\sigma_k > 0} = \|\sigma\|_0$
- Nuclear norm: $\|\boldsymbol{X}\|_* = \sum_{k=1}^n \sigma_k = \|\boldsymbol{\sigma}\|_1$
- Nuclear norm can be seen as approximation to rank

Definition (nuclear norm minimization)

Given an $n \times n$ data matrix **D** and an index set Ω of revealed entries. The *nuclear minimization problem* is

$$\begin{array}{ll} \text{minimize} & \| \boldsymbol{X} \|_{*} \\ \text{subject to} & d_{ij} = x_{ij} \\ & \boldsymbol{X} \in \mathbb{R}^{n \times n}. \end{array} (i,j) \in \Omega$$

- A convex relaxation of rank minimization
- Nuclear norm is convex function (thus local optimum is global optimum)

Can be optimized (more) efficiently via semidefinite programming.

Why nuclear norm minimization? (1)

Let's look at two-dimensional vectors first.

- Border of blue area $= L_1$ unit ball
- Red lines = L_0 "unit ball" \rightarrow axes except 0
- L_0 unit ball intersects L_1 unit ball at extreme points



Why nuclear norm minimization? (2)

Let's find a solution to the problem $\Phi x = y$.

- Underdetermined system with infinitely many solution
- We usually pick one that has certain structure
- E.g. sparsity: $\min \|\boldsymbol{x}\|_0$ s.t. $\Phi \boldsymbol{x} = y$
 - $\blacktriangleright \ \, {\sf Example \ \, solution: \ \, x_1=y/\phi_1, \ \, x_2=0 \rightarrow L_0=1}$
- Approximation: minimize L₁
 - Recall: $\|\bm{x}\|_1 = \sum |x_i|$
 - Increasing L₁ norm can be seen as "inflating" the L₁ unit ball
 - Minimum L₁ norm
 minimum inflation
 - Achieved at intersection with x₁ or x₂ axis (whatever is smaller)



 X_1

 X_2

 $\Phi_{X=y}$

Why nuclear norm minimization (3)

Figure 1. Unit ball of the nuclear norm for symmetric 2 × 2 matrices. The red line depicts a random one-dimensional affine space. Such a subspace will generically intersect a sufficiently large nuclear norm ball at a rank one matrix.



- Consider SVD of $\boldsymbol{D} = \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^T$
- Unit nuclear norm ball = convex combination (σ_k) of rank-1 matrices of unit Frobenius norm (U_{*k}V^T_{*k})
- Extreme points have low rank (in figure: rank-1 matrices of unit Frobenius norm)
- Nuclear norm minimization: inflate unit ball as little as possible to reach d_{ij} = x_{ij}
- Solution lies at extreme point of inflated ball \rightarrow (hopefully) low rank

Relationship to LFMs

• Recall regularized LFM (\boldsymbol{L} is $m \times r$, \boldsymbol{R} is $r \times n$):

$$\min_{\boldsymbol{L},\boldsymbol{R}}\sum_{(i,j)\in\Omega}(d_{ij}-[\boldsymbol{L}\boldsymbol{R}]_{ij})^2+\lambda(\|\boldsymbol{L}\|_F^2+\|\boldsymbol{R}\|_F^2)$$

• View as matrix completion problem by enforcing $d_{ij} = [LR]_{ij}$:

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \left(\|\boldsymbol{L}\|_{F}^{2} + \|\boldsymbol{R}\|_{F}^{2} \right) \\ \text{subject to} & d_{ij} = x_{ij} \quad (i,j) \in \Omega \\ & \boldsymbol{L}\boldsymbol{R} = \boldsymbol{X}. \end{array}$$

- One can show: for *r* chosen larger than rank of nuclear norm optimum, equivalent to nuclear norm minimization
- For some intuition, suppose $\boldsymbol{X} = \boldsymbol{L}\boldsymbol{R} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{T}$ at optimum \boldsymbol{L} and \boldsymbol{R} : $\frac{1}{2}\left(\|\boldsymbol{L}\|_{F}^{2} + \|\boldsymbol{R}\|_{F}^{2}\right) \leq \frac{1}{2}\left(\|\boldsymbol{U}\boldsymbol{\Sigma}^{1/2}\|_{F}^{2} + \|\boldsymbol{\Sigma}^{1/2}\boldsymbol{V}^{T}\|_{F}^{2}\right)$ $= \frac{1}{2}\sum_{i=1}^{n}\sum_{k=1}^{r}(u_{ik}^{2}\sigma_{k} + v_{ik}^{2}\sigma_{k})$ $= \sum_{k=1}^{r}\sigma_{k} = \|\boldsymbol{X}\|_{*}$

When can we hope to recover D? (1)

Assume **D** is the 5×5 all-ones matrix (rank 1).



Sampling strategy and sample size matter.

When can we hope to recover D? (2)

Consider the following rank-1 matrices and assume **50% revealed** entries.



When can we hope to recover D? (3)

Exact conditions under which matrix completion "works" is active research area:

- Which sampling schemes? (e.g., random, WR/WOR, active)
- Which sample size?
- Which matrices? (e.g., "incoherent" matrices)
- Noise (e.g., independent, normally distributed noise)

Theorem (Candès and Recht, 2009)

Let $\boldsymbol{D} = \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^{\mathsf{T}}$. If \boldsymbol{D} is incoherent in that

$$\max_{ij} u_{ij}^2 \leq \frac{\mu_B}{n} \qquad and \qquad \max_{ij} v_{ij}^2 \leq \frac{\mu_B}{n}$$

for some $\mu_B = O(1)$, and if rank $(\mathbf{D}) \le \mu_B^{-1} n^{1/5}$, then $O(n^{6/5} r \log n)$ random samples without replacement suffice to recover \mathbf{D} exactly with high probability.

Outline

- 1. Collaborative Filtering
- 2. Matrix Completion
- 3. Algorithms
- 4. Variants
- 5. Summary

Overview

Latent factor models in practice

- Millions of users and items
- Billions of ratings
- Sometimes quite complex models

Many algorithms have been applied to large-scale problems

- Gradient descent and quasi-Newton methods
- Coordinate-wise gradient descent
- Stochastic gradient descent
- Alternating least squares

Continuous gradient descent

- Find minimum θ^* of function L
- Pick a starting point $heta_0$
- Compute gradient $L'(\theta_0)$
- Walk downhill
- Differential equation

$$\frac{\partial \boldsymbol{\theta}(t)}{\partial t} = -L'(\boldsymbol{\theta}(t))$$

with boundary cond. $\theta(0) = \theta_0$

• Under certain conditions

 $oldsymbol{ heta}(t) o oldsymbol{ heta}^*$



Discrete gradient descent

- Find minimum θ^* of function L
- Pick a starting point $heta_0$
- Compute gradient $L'(\theta_0)$
- Jump downhill
- Difference equation

 $\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \epsilon_n L'(\boldsymbol{\theta}_n)$

• Under certain conditions, approximates CGD in that

 $\theta^n(t) = \theta_n +$ "steps of size t"

satisfies the ODE as $n
ightarrow \infty$



Recap: Gradient computation

• You know: gradient computations for functions on one input; e.g.,

$$f(x) = x^2$$
$$\nabla_x f(x) = 2x$$

• For functions with multiple inputs, there are multiple partial derivatives

$$f(l, r) = (d - lr)^2$$

$$\nabla_l f(l, r) = -2(d - lr)r$$

$$\nabla_r f(l, r) = -2(d - lr)l$$

Matrix calculus (1)

- We focus on functions from $\mathbb{R}^n \to \mathbb{R}$
- Let $\boldsymbol{\theta} = \begin{pmatrix} I & r \end{pmatrix}^T$
- Then

$$f(\theta) = (d - lr)^2$$

$$\nabla_l f(\theta) = -2(d - lr)r$$

$$\nabla_r f(\theta) = -2(d - lr)l$$

• We can write this in matrix form

$$\nabla_{\theta^{T}} f = \begin{pmatrix} -2(d-lr)r & -2(d-lr)l \end{pmatrix} = -2(d-lr) \begin{pmatrix} r & l \end{pmatrix}$$

• The resulting matrix is called the Jacobian matrix of f

$$\boldsymbol{J}_f = -2(d-lr)\begin{pmatrix}r & l\end{pmatrix}$$

• Note: The Jacobian is a matrix of functions of heta

Matrix calculus (2)

• To recap

$$\boldsymbol{J}_{\boldsymbol{f}} \stackrel{\text{def}}{=} \nabla_{\boldsymbol{x}^{\mathsf{T}}} \boldsymbol{f} \stackrel{\text{def}}{=} \begin{pmatrix} \nabla_{x_1} \boldsymbol{f} & \nabla_{x_2} \boldsymbol{f} & \cdots & \nabla_{x_n} \boldsymbol{f} \end{pmatrix}$$

• The following rules hold

$$\begin{aligned} \nabla_{\mathbf{x}^T} & \mathbf{c} = \begin{pmatrix} 0 & 0 & \cdots & 0 \end{pmatrix} \\ \nabla_{\mathbf{x}^T} & \mathbf{c}^T \mathbf{x} = \mathbf{c}^T \\ \nabla_{\mathbf{x}^T} & \mathbf{x}^T \mathbf{c} = \mathbf{c}^T \\ \nabla_{\mathbf{x}^T} & \mathbf{x}^T \mathbf{x} = 2\mathbf{x}^T \\ \nabla_{\mathbf{x}^T} & \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T), \end{aligned}$$

where constants c, c, and A do not depend on x

• Also: multiplicative rule, product rule, chain rule, ...

Gradient computation for LFMs Set $\theta = (\boldsymbol{L}, \boldsymbol{R}), \ L(\theta) = L(\boldsymbol{L}, \boldsymbol{R}), \ \text{and write}$ $L(\boldsymbol{\theta}) = \sum L_{ij}(\boldsymbol{I}_i, \boldsymbol{r}_j)$ $(i,j)\in\Omega$ $L'(\boldsymbol{\theta}) = \sum L'_{ij}(\boldsymbol{I}_i, \boldsymbol{r}_j)$ $(i,i) \in \Omega$ $\nabla_{l_{ik}} L(\boldsymbol{\theta}) = \sum \nabla_{l_{ik}} L_{i'j}(\boldsymbol{I}_{i'}, \boldsymbol{r}_j)$ $(i',j)\in\Omega$ $= \sum \nabla_{l_{ik}} L_{ij}(\boldsymbol{I}_i, \boldsymbol{r}_j)$ $i \in \{ i' | (i,j') \in \Omega \}$



since $\nabla_{I_{ik}} L_{i'j}(\boldsymbol{I}_{i'}, \boldsymbol{r}_j) = 0$ for $i \neq i'$

Local gradient of entry $(i,j) \in \Omega$ nonzero only for row I_i and column r_j .

Example gradient computation

Simplest form of LFM (unregularized)

$$L_{ij}(\boldsymbol{I}_i,\boldsymbol{r}_j)=(d_{ij}-\boldsymbol{I}_i^T\boldsymbol{r}_j)^2$$

Gradient computation

$$\nabla_{I_{ik}} L_{ij}(\boldsymbol{I}_{i}, \boldsymbol{r}_{j}) = -2(d_{ij} - \boldsymbol{I}_{i}^{T} \boldsymbol{r}_{j})r_{kj}$$
$$\nabla_{r_{kj}} L_{ij}(\boldsymbol{I}_{i}, \boldsymbol{r}_{j}) = -2I_{ik}(d_{ij} - \boldsymbol{I}_{i}^{T} \boldsymbol{r}_{j})$$
$$\nabla_{\boldsymbol{I}_{i}^{T}} L_{ij}(\boldsymbol{I}_{i}, \boldsymbol{r}_{j}) = -2(d_{ij} - \boldsymbol{I}_{i}^{T} \boldsymbol{r}_{j})\boldsymbol{r}_{j}^{T}$$
$$\nabla_{\boldsymbol{r}_{j}^{T}} L_{ij}(\boldsymbol{I}_{i}, \boldsymbol{r}_{j}) = -2\boldsymbol{I}_{i}^{T}(d_{ij} - \boldsymbol{I}_{i}^{T} \boldsymbol{r}_{j})$$

Gradient descent for LFMs

GD epoch

- 1. Compute gradient
 - Initialize zero matrices $\textit{L}^{
 abla}$ and $\textit{R}^{
 abla}$
 - ► For each entry $(i,j) \in \Omega$, update gradients

$$\begin{split} \boldsymbol{I}_{i}^{\nabla} \leftarrow \boldsymbol{I}_{i}^{\nabla} + \nabla_{\boldsymbol{I}_{i}} \boldsymbol{L}_{ij}(\boldsymbol{I}_{i}, \boldsymbol{r}_{j}) \\ \boldsymbol{r}_{j}^{\nabla} \leftarrow \boldsymbol{r}_{j}^{\nabla} + \nabla_{\boldsymbol{r}_{j}} \boldsymbol{L}_{ij}(\boldsymbol{I}_{i}, \boldsymbol{r}_{j}) \end{split}$$



After this step

$$\boldsymbol{I}_{i}^{\nabla} = \sum_{j \in \{j' \mid (i,j') \in \Omega\}} \nabla_{\boldsymbol{I}_{i}} L_{ij}(\boldsymbol{I}_{i}, \boldsymbol{r}_{j})$$

as desired ($\mathbf{r}_{j}^{
abla}$ analog)

2. Update parameters

$$\boldsymbol{L} \leftarrow \boldsymbol{L} - \epsilon_n \boldsymbol{L}^{\nabla}$$
$$\boldsymbol{R} \leftarrow \boldsymbol{R} - \epsilon_n \boldsymbol{R}^{\nabla}$$

Stochastic gradient descent

- Find minimum θ^* of function L
- Pick a starting point $heta_0$
- Approximate gradient $\hat{L}'(\boldsymbol{\theta}_0)$
- Jump "approximately" downhill
- Stochastic difference equation

$$\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \epsilon_n \hat{L}'(\boldsymbol{\theta}_n)$$

 Under certain conditions, asymptotically approximates (continuous) gradient descent



Stochastic gradient descent for LFMs

• Set
$$\boldsymbol{\theta} = (\boldsymbol{L}, \boldsymbol{R})$$
 and use

$$L(\boldsymbol{\theta}) = \sum_{(i,j)\in\Omega} L_{ij}(\boldsymbol{I}_i, \boldsymbol{r}_j)$$
$$L'(\boldsymbol{\theta}) = \sum_{(i,j)\in\Omega} L'_{ij}(\boldsymbol{I}_i, \boldsymbol{r}_j)$$
$$\hat{L}'(\boldsymbol{\theta}, \boldsymbol{z}) = NL'_{i,i}(\boldsymbol{I}_{i,i}, \boldsymbol{r}_{i,i}),$$



where $\textit{N} = |\Omega|$ and $\textit{z} = (\textit{i}_z, \textit{j}_z) \in \Omega$

- SGD epoch (with or without replacement)
 - 1. Pick a random entry $z \in \Omega$ (with or without replacement)
 - 2. Compute approximate gradient $\hat{L}'(\theta, z)$
 - 3. Update parameters

$$\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \epsilon_n \hat{L}'(\boldsymbol{\theta}_n, z)$$

4. Repeat N times

SGD step affects only current row and column.

Comparison

• Per epoch, assuming O(r) gradient computation per element

	GD	SGD
Algorithm	Deterministic	Randomized
Gradient computations	1	Ν
Gradient types	Exact	Approximate
Parameter updates	1	Ν
Time	O(rN)	O(rN)
Space	O((m+n)r)	O((m+n)r)

- Why stochastic?
 - Fast convergence to vicinity of optimum
 - Randomization may help escape local minima
 - Exploitation of "repeated structure"

Example: Netflix data, unregularized



epoch

GD/SGD in practice (1)

Step size (or learning rate) sequence $\{\epsilon_n\}$ needs to be chosen carefully; widely studied, many options:

- Large \rightarrow good initially (move quickly), bad later on (juggle around optimum)
- Keep step size throughout or reduce it gradually
 - E.g., constant (useful for online learning)
 - E.g., $\epsilon_n = a/(b+n)$ for some constants a, b
 - ▶ E.g., pick $\epsilon_n \leq 1/L(\nabla f)$ if f has bounded gradient
- Bold driver heuristic: After every epoch
 - Increase step size slightly when loss decreased (by, say, 5%)
 - Decrease step size sharply when loss increased (by, say, 50%)
 - Not provably correct, but works well in practice
- Pick initial step size based on sample
- Line search: optimize the step size directly

$$\epsilon_n = \operatorname*{argmin}_{\epsilon} L(\boldsymbol{\theta}_n - \epsilon L'(\boldsymbol{\theta}_n))$$

GD/SGD in practice (2)

- SGD is a common learning algorithm
 - E.g., training neural networks
 - Related to incremental gradient descent and online learning
- Many variants exist; e.g.,
 - Use more than one example per step (mini-batch)
 - Polyak averaging
 - ► <u>Momentum</u>
 - Adaptive, per-parameter step sizes (<u>AdaGrad</u>, <u>RMSprop</u>, AdaDelta, Adam)
 - ▶ ...
- And it can (often) be parallelized; e.g.,
 - Hogwild
 - Vowpal Wabbit (for regression problems)
 - $\label{eq:scalar} \blacktriangleright \frac{\text{DSGD}{++}}{\text{E.g.: 10M}{\times}1\text{M}, \ \text{10B} \ \text{observed entries}, \approx 1\text{h on 8 machines}}$

▶ ...

Outline

- 1. Collaborative Filtering
- 2. Matrix Completion
- 3. Algorithms

4. Variants

5. Summary

Variants of latent factor models

- In practice, the basic latent factor model is often modified in an application-dependent way
- We discuss two variants here
 - How to handle implicit feedback? (All observed entries are 1.)
 - How to make use of additional contextual information? (E.g., attributes for users and items.)
- There are many other variants; e.g.
 - Bayesian variants (partly discussed later)
 - Combination of explicit and implicit feedback
 - Additional constraints (e.g., non-negativity of factors)
 - ▶ ...
- Picking the right model requires thought, experience, and experimentation

Implicit feedback

- Implicit feedback is an indicator for preference
 - Which product pages do users look at?
 - Which movies do users watch?
 - How much time do users spend on a news article?
- Absence of implicit feedback is **not** necessarily an indicator for non-preference
- In the basic latent factor model, we used explicit feedback, which contains
 - Positive evidence (a user gave an item a high rating)
 - Negative evidence (a user gave an item a low rating)
- If we only have implicit feedback
 - We may think of the implicit feedback as positive evidence
 - But we do not have negative evidence
 - This implies that the basic latent factor model won't work

Implicit feedback and the basic LFM

• A matrix of implicit feedback (1=present, ?=absent)

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

• Its completion by an (unregularized) LFM

• That's not helpful!

How to handle implicit feedback

- One option is to build a different model
 - E.g., a nearest-neighbor model
- We can also learn variant of an LFM
 - Good approach if low-rank assumption applies
- How? Move from "predicting values" to "ranking values"
 - ► Ranking is per row (e.g, the items for each user) or per column
 - ► Goal is not to predict the preference of a user for an item, but to rank items such that high-preference items appear early on
- Key idea is to find a ranking that
 - Tends to rank items with implicit feedback before items without
 - Is simple (e.g., can be described with a low-rank matrix)

Bayesian personalized ranking (intuition)

- Goal is to learn a personalized ordering $>_i$ for each user i
 - Consider a user *i*, an item j_1 and an item j_2
 - Idea: $j_1 >_i j_2$ if user *i* prefers item j_1 over item j_2
 - Assumption: If user *i* provided implicit feedback for j_1 but not for j_2 , then he prefers j_1 over j_2
 - We want to learn $>_i$ for all pairs of items
 - Ordering should be total, antisymmetric, and transitive

The BPR model

• Associate each triple (i, j_1, j_2) with a score $\hat{x}_{ij_1j_2} \in \mathbb{R}$

- Score is $\begin{cases} positive & \text{if } j_1 \text{ is preferred by user } i \\ negative & \text{if } j_2 \text{ is preferred by user } i \\ 0 & \text{if oblivious} \end{cases}$
- Model probability that *i* prefers j_1 over j_2 as

$$p(j_1 >_i j_2) = \sigma(\hat{x}_{ij_1j_2}),$$

where $\sigma(x) = 1/(1 + \exp(-x)) \in [0, 1]$ is the logistic function

BPR for latent factor models

• Use low-rank matrix factorization to model scores

$$\hat{x}_{ij_1j_2} = \boldsymbol{I}_i^T \boldsymbol{r}_{j_1} - \boldsymbol{I}_i^T \boldsymbol{r}_{j_2},$$

where as before \boldsymbol{L} and \boldsymbol{R} are factor matrices

- Construct database $D \subseteq [m] \times [n] \times [n]$ of "observed orderings"
 - ▶ $(i, j^+, j^-) \in D$ if user *i* has implicit feedback for j^+ but not for j^-
 - Can be large, but does not really need to be constructed
- Under certain assumptions, maximum likelihood estimate for L and R is

$$\arg\max_{\boldsymbol{L}\in\mathbb{R}^{m\times r},\boldsymbol{R}\in\mathbb{R}^{r\times n}}\prod_{(i,j^+,j^-)\in D}p(j^+>_ij^-)$$
$$=\arg\max_{\boldsymbol{L}\in\mathbb{R}^{m\times r},\boldsymbol{R}\in\mathbb{R}^{r\times n}}\sum_{(i,j^+,j^-)\in D}\log p(j^+>_ij^-)$$

- We can use stochastic gradient ascent to find the estimates
 - Gradient estimate obtained by sampling a single triple
 - In practice, add regularization

BPR in action



Online shopping: Rossmann

Video Rental: Netflix

Figure 6: Area under the ROC curve (AUC) prediction quality for the Rossmann dataset and a Netflix subsample. Our BPR optimizations for matrix factorization BPR-MF and k-nearest neighbor BPR-kNN are compared against weighted regularized matrix factorization (WR-MF) [5, 10], singular value decomposition (SVD-MF), k-nearest neighbor (Cosine-kNN) [2] and the most-popular model. For the factorization methods BPR-MF, WR-MF and SVD-MF, the model dimensions are increased from 8 to 128 dimensions. Finally, np_{max} is the theoretical upper bound for any non-personalized ranking method.

Contextual information

- We often have additional information; e.g.,
 - Demographics of users (age, sex, city, ...)
 - ▶ Information about items (e.g., genre, actors, directors, ...)
- How to exploit this additional information with LFM models?
 - 1. Build a separate model and combine predictions (stacking)
 - 2. As above, but but learn model parameters jointly
 - 3. Extend LFMs to directly incorporate context
- A popular approach for (3) are factorization machines
 - Open source implementation: libfm
 - Can be combined with BPR
 - Advertisement: also useful for <u>extracting relations from</u> <u>natural-language text</u>
Factorization machines (intuition)

	Feature vector x									Та	irge	ət y											
X ⁽¹⁾	1	0	0		1	0	0	0		0.3	0.3	0.3	0		13	0	0	0	0		5	5	y ⁽¹⁾
X ⁽²⁾	1	0	0		0	1	0	0		0.3	0.3	0.3	0		14	1	0	0	0		3	3	y ⁽²⁾
X ⁽³⁾	1	0	0		0	0	1	0		0.3	0.3	0.3	0		16	0	1	0	0		1		y ⁽²⁾
X ⁽⁴⁾	0	1	0		0	0	1	0		0	0	0.5	0.5		5	0	0	0	0		4	ŀ	y ⁽³⁾
X ⁽⁵⁾	0	1	0		0	0	0	1		0	0	0.5	0.5		8	0	0	1	0		Ę	5	y ⁽⁴⁾
X ⁽⁶⁾	0	0	1		1	0	0	0		0.5	0	0.5	0		9	0	0	0	0		1		y ⁽⁵⁾
X ⁽⁷⁾	0	0	1		0	0	1	0		0.5	0	0.5	0		12	1	0	0	0		Ę	5	y ⁽⁶⁾
	A	B Us	C er		П	NH	SW Movie	ST Э		TI Otl	NH her N	SW lovie	ST s rate	ed	Time	"	NH _ast I	SW Novie	ST e rate	ed			

- Rows = observed entries
- Columns = information about observed entry
 - Which user rated which item? In what context?
 - Must be constructable at prediction time for unobserved entries (!)
- Associate a latent feature vector with each column
 - Prediction \approx sum of pairwise inner products of weighted feature vectors
 - As before, learn latent features using SGD to fit observed data

Outline

- 1. Collaborative Filtering
- 2. Matrix Completion
- 3. Algorithms
- 4. Variants
- 5. Summary

Lessons learned

- Collaborative filtering methods learn from past user behavior
 - Latent factor models are best-performing single approach
 - Many variants exist
 - In practice, often combined with other methods
- Users and items are represented in common latent factor space
 - Holistic matrix-factorization approach
 - Similar users/item placed at similar positions
 - Low-rank assumption = few "factors" influence user preferences
- Close relationship to matrix completion problem
 - Reconstruct a partially observed low-rank matrix
 - Many applications
- SGD is simple and practical algorithm to solve LFMs in summation form

Suggested reading

- Y. Koren, R. Bell, C. Volinsky Matrix factorization techniques for recommender systems IEEE Computer, 42(8), p. 30–37, 2009 http://ieeexplore.ieee.org/xpl/articleDetails.jsp?arnumber=5197422
- E. Candès, B. Recht
 - Exact matrix completion via convex optimization Communications of the ACM, 55(6), p. 111–119, 2012 http://doi.acm.org/10.1145/2184319.2184343
- https://en.wikipedia.org/wiki/Matrix_calculus
- And references in the above articles and slides

Data Mining and Matrices 05 – Non-Negative Matrix Factorization

Prof. Dr. Rainer Gemulla

Universität Mannheim

FSS 2018

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 may involve negative values in factors and components
- Negative values describe the absence of something
- Often no natural interpretation

Goal: Find a decomposition that is more natural to non-negative data.

Example (SVD)

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



Here are the corresponding components:



Negative values can make interpretation unnatural or difficult.

Outline

1. Non-Negative Matrix Factorization

- 2. Algorithms
- 3. Probabilistic Latent Semantic Analysis
- 4. Summary

Non-negative matrix factorization (NMF)

Definition (Non-negative matrix factorization, basic form)

Given a non-negative matrix $D \in \mathbb{R}^{m \times n}_+$, a non-negative matrix factorization of rank k is

 $D \approx LR$,

where $\boldsymbol{L} \in \mathbb{R}^{m \times r}_+$ and $\boldsymbol{R} \in \mathbb{R}^{r \times n}_+$ are both non-negative.

- Additive decomposition: factors and components non-negative \rightarrow No cancellation effects
- Rows of **R** can be thought as "parts"
- Row of *D* obtained by mixing (or "assembling") parts as described in *L*
- Smallest *r* such that *D* = *LR* exists is called **non-negative** rank of *D*

 $\mathsf{rank}(\mathbf{D}) \le \mathsf{rank}_+(\mathbf{D}) \le \min\{m, n\}$

Example (NMF)

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



Here are the corresponding components:



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

Decomposing faces (SVD)



Lee and Seung, 1999.

Decomposing faces (NMF)



Decomposing digits (NMF)



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

Cichocki et al., 2009.

Some applications

- Text mining (more later)
- Bioinformatics
- Microarray analysis
- Mineral exploration
- Neuroscience
- Image understanding
- Air pollution research
- Chemometrics
- Spectral data analysis
- Linear sparse coding
- Image classification
- Clustering
- Neural learning process
- Sound recognition
- Remote sensing
- Object characterization



Figure 4.8 Illustration for (a) benchmark used in large-scale experiments with 10 nonnegative sources; (b) Typical 1000 mixtures; (c) Ten estimated components by using FAST HALS NMF from the observations mutrix V of dimension 1000 × 1000. (d) Performance expressed via the PSNR using the Beta HALS NMF algorithm for β = 0.5, 1, 1.5, 2 and 3.

Cichocki et al., 2009

Gaussian NMF

• Gaussian NMF is the most basic form of non-negative factorizations:

minimize
$$\|\boldsymbol{D} - \boldsymbol{L}\boldsymbol{R}\|_F^2$$

s.t. $\boldsymbol{L} \in \mathbb{R}^{m \times r}_+$
 $\boldsymbol{R} \in \mathbb{R}^{r \times n}_+$

- Truncated SVD minimizes the same objective (but without non-negativity constraints)
- Many other variants exist
 - Different objective functions (e.g., KL-divergence)
 - ► Additional regularizations (e.g., *L*₁-regularization)
 - ► Different constraints (e.g., orthogonality of **R**)
 - Different compositions (e.g., 3 matrices)
 - multi-layer NMF, semi-NMF, sparse NMF, tri-NMF, symmetric NMF, orthogonal NMF, non-smooth NMF (nsNMF), overlapping NMF, convolutive NMF (CNMF), k-Means, ...

k-Means can be seen as a variant of NMF



k-Means factors correspond to prototypical faces.

NMF is not unique

• Factors are not "ordered"



- One way of ordering: decreasing Frobenius norm of components (i.e., order by ||*I_kr^T_k*||_F = ||*I_k*|| ||*r_k*||)
- Factors/components are not unique



Additional constraints or regularization can encourage uniqueness.

NMF is not hierarchical

Rank-2 NMF



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

Outline

1. Non-Negative Matrix Factorization

2. Algorithms

- 3. Probabilistic Latent Semantic Analysis
- 4. Summary

NMF is difficult

We focus on minimizing $L(\boldsymbol{L}, \boldsymbol{R}) = \|\boldsymbol{D} - \boldsymbol{L}\boldsymbol{R}\|_{\boldsymbol{F}}^2$.

- For varying *m*, *n*, and *r*, problem is NP-hard
- When rank(D) = 1 (or r = 1), can be solved in polynomial time
 - 1. Take first non-zero column of **D** as $L_{m \times 1}$
 - 2. Determine $\mathbf{R}_{1 \times n}$ entry by entry (using the fact that $\mathbf{d}_i = \mathbf{L} \mathbf{r}_{1i}$)
- Problem is not convex
 - Local optimum may not correspond to global optimum
 - Generally little hope to find global optimum
- But: Problem is biconvex
 - For fixed **R**, $f(\mathbf{L}) = \|\mathbf{D} \mathbf{L}\mathbf{R}\|_{F}^{2}$ is convex

 $f(\boldsymbol{L}) = \sum_{i} \|\boldsymbol{d}_{i}^{T} - \boldsymbol{I}_{i}^{T} \boldsymbol{R}\|_{F}^{2}$ (chain rule)

$$\nabla_{l_k} f(\boldsymbol{L}) = -2(\boldsymbol{d}_i^T - \boldsymbol{l}_i^T \boldsymbol{R}) \boldsymbol{r}_k \qquad \text{(product rule)}$$
$$\nabla_{l_k}^2 f(\boldsymbol{L}) = 2\boldsymbol{r}_k^T \boldsymbol{r}_k \ge 0 \qquad \text{(convex in } l_{ik}; \text{ can show: also in } \boldsymbol{L})$$

(convex in I_{ik} ; can show: also in \boldsymbol{L})

- For fixed **L**, $f(\mathbf{R}) = \|\mathbf{D} \mathbf{L}\mathbf{R}\|_{F}^{2}$ is convex
- Allows for efficient algorithms

General framework

- Gradient descent generally slow
- Stochastic gradient descent often also slow
- Key approach: alternating minimization
 - 1: Pick starting point \boldsymbol{L}_0 and \boldsymbol{R}_0
 - 2: while not converged
 - 3: Keep **R** fixed, optimize **L**
 - 4: Keep *L* fixed, optimize *R*
- Update steps 3 and 4 easier than full problem
- Also called alternating projections or (block) coordinate descent
- Starting point
 - Random
 - Multi-start initialization: try multiple random starting points, run a few epochs, continue with best
 - Based on SVD
 - ▶ ...

Example

Ignore non-negativity for now. Consider the regularized least-square error:

$$L(L, R) = \|D - LR\|_F^2 + \frac{\lambda}{2}(\|L\|_F^2 + \|R\|_F^2)$$

By setting $m = n = r = 1$, $D = (1)$ and $\lambda = 0.1$, we obtain
 $L(l, r) = (1 - lr)^2 + 0.05(l^2 + r^2)$



$$\begin{aligned} \nabla_{l}L &= -2r(1-lr) + 0.1l \\ \nabla_{r}L &= -2l(1-lr) + 0.1r \\ \text{Local optima:} \\ & \left(\sqrt{\frac{19}{20}}, \sqrt{\frac{19}{20}}\right), \ \left(-\sqrt{\frac{19}{20}}, -\sqrt{\frac{19}{20}}\right) \\ \text{Stationary point:} \ (0,0) \end{aligned}$$

Example (Alternating Least Squares, ALS)



Example (Alternating Least Squares, ALS)



Alternating non-negative least squares (ANLS)

• Uses non-negative least squares approximation of *L* and *R*:

 $\underset{\boldsymbol{L} \in \mathbb{R}^{m \times r}_{+}}{\operatorname{argmin}} \|\boldsymbol{D} - \boldsymbol{L}\boldsymbol{R}\|_{F}^{2} \quad \text{and} \quad \underset{\boldsymbol{R} \in \mathbb{R}^{r \times n}_{+}}{\operatorname{argmin}} \|\boldsymbol{D} - \boldsymbol{L}\boldsymbol{R}\|_{F}^{2}$

- Equivalently: find non-negative least squares solution to LR = D
- Common (bad!) hack: Solve unconstrained least squares problems and "remove" negative values. E.g., when columns (rows) of *L* (*R*) are linearly independent, set

$$m{L} = [m{D}m{R}^+]_\epsilon$$
 and $m{R} = [m{L}^+m{D}]_\epsilon$

where

•
$$\mathbf{R}^+ = \mathbf{R}^{\mathsf{T}}_{\mathsf{T}}(\mathbf{R}\mathbf{R}^{\mathsf{T}})^{-1}$$
 is the pseudo-inverse of \mathbf{R}

• $\boldsymbol{L}^+ = (\boldsymbol{L}^T \boldsymbol{L})^{-1} \boldsymbol{L}^T$ is the pseudo-inverse of \boldsymbol{L}

• $[a]_{\epsilon} = \max\{\epsilon, a\}$ for $\epsilon = 0$ or some small constant (e.g., $\epsilon = 10^{-9}$)

- Difficult to analyze due to non-linear update steps
- Often slow convergence to a "bad" local minimum (better when regularized or good NLS solver used)

Example (ANLS, common bad hack)



Example (ANLS, common bad hack)



Hierarchical alternating least squares (HALS)

- Optimize just one component, then next component, and so on
- Let **D**^(k) be the residual matrix (error) when k-th component is removed:

$$\boldsymbol{D}^{(k)} = \boldsymbol{D} - \boldsymbol{L}\boldsymbol{R} + \boldsymbol{I}_k \boldsymbol{r}_k^T = \boldsymbol{D} - \sum_{k' \neq k} \boldsymbol{I}_{k'} \boldsymbol{r}_{k'}^T$$

- HALS optimizes $\|\boldsymbol{D}^{(k)} \boldsymbol{I}_k \boldsymbol{r}_k^T\|_F^2$ for k = 1, 2, ..., r, 1, ... (thus: optimize only *k*-th component)
- E.g., in each iteration, set (once or multiple times):

$$\boldsymbol{I}_{k}^{T} \leftarrow \frac{1}{\|\boldsymbol{r}_{k}\|_{F}^{2}} \left[\boldsymbol{D}^{(k)} \boldsymbol{r}_{k} \right]_{\epsilon} \quad \text{and} \quad \boldsymbol{r}_{k}^{T} \leftarrow \frac{1}{\|\boldsymbol{I}_{k}\|_{F}^{2}} \left[\boldsymbol{I}_{k}^{T} \boldsymbol{D}^{(k)} \right]_{\epsilon}$$

- $D^{(k)}$ can be incrementally maintained \rightarrow fast implementation $D^{(k+1)} = D^{(k)} + I_k r_k^T - I_{k+1} r_{k+1}^T$
- Good performance in practice
- Converges to stationary point when initialized with positive matrix and sufficiently small ϵ

Multiplicative updates

• Gradient descent step with step size η_{kj}

$$r_{kj} \leftarrow r_{kj} + \eta_{kj} ([\boldsymbol{L}^T \boldsymbol{D}]_{kj} - [\boldsymbol{L}^T \boldsymbol{L} \boldsymbol{R}]_{kj})$$

• Setting $\eta_{kj} = \frac{r_{kj}}{[L^T L R]_{kj}}$, we obtain the multiplicative update rules

$$m{L} \leftarrow m{L} \circ rac{m{D}m{R}^T}{m{L}m{R}m{R}^T}$$
 and $m{R} \leftarrow m{R} \circ rac{m{L}^Tm{D}}{m{L}^Tm{L}m{R}}$

where multiplication (\circ) and division are element-wise

- Does not necessarily find optimum *L* (or *R*), but can be shown to never increase loss
- Faster than ANLS (no NLS problems solved), easy to implement and parallelize
- Zeros in factors are problematic (never changed, division by 0)

$$m{L} \leftarrow m{L} \circ rac{[m{D}m{R}^T]_\epsilon}{m{L}m{R}m{R}^T + \epsilon}$$
 and $m{R} \leftarrow m{R} \circ rac{[m{L}^Tm{D}]_\epsilon}{m{L}^Tm{L}m{R} + \epsilon}$

Lee and Seung, 2001. Cichocki et al., 2006.

Example (multiplicative updates)



Outline

- 1. Non-Negative Matrix Factorization
- 2. Algorithms
- 3. Probabilistic Latent Semantic Analysis
- 4. Summary

Topic modeling

• Consider a document-word matrix constructed from some corpus

		air	water	pollution	power	democrat	republican
$ ilde{m{P}}=$	doc 1/	3	2	8	4	0	0)
	doc 2	1	4	12	2	0	0
	doc 3	0	0	0	4	10	11
	doc 4	0	0	0	3	8	5
	doc 5^{\prime}	1	1	1	1	1	1 /

- Documents seem to talk about two "topics"
 - 1. Environment (with words air, water, pollution, power)
 - 2. Congress (with words democrat, republican, power)
- Note: "power" is polysemous (electrical power, political power)

Can we automatically detect topics in documents?

Let's try SVD

• Rank-2 truncated SVD of example matrix



- Can see clear topic-driven differences for documents in $ilde{u}_1$ and $ilde{u}_2$
- Similarly for words in $\tilde{\pmb{v}}_1$ and $\tilde{\pmb{v}}_2$
- Hard to interpret
- Called latent semantic analysis

Let's try NMF

• Rank-2 NMF of example matrix



- Can see clear topic-driven differences in \tilde{L} and \tilde{R}
- Easier to interpret: large value \rightarrow relevant for topic
- Decomposition is sparse
- Related to probabilistic latent semantic analysis (pLSA)
 - An even more interpretable factorization
 - Assumptions are well-understood
 - Example of a topic model
- Polysemy of "power" hinted at by its use in multiple topics
- Topic distributions often much more complex

A probabilistic viewpoint

• Let's normalize \tilde{P} such that the entries sum to unity

air	wat	pol	pow	dem	rep		air	wat	pol j	pow	dem	rep
3	2	8	4	0	0		0.04	0.02	0.1	0.05	0	0
1	4	12	2	0	0		0.01	0.05	0.14	0.02	0	0
0	0	0	4	10	11		0	0	0	0.05	0.12	0.13
0	0	0	3	8	5		0	0	0	0.04	0.1	0.06
1	1	1	1	1	1		0.01	0.01	0.01	0.01	0.01	0.01
		j	Ď			-				כ		

• Put all words in an urn, draw, and denote by D and W the result. The probability to draw word w from document d is given by

$$P(D=d,W=w)=p_{dw}$$

- Matrix **P** can represent any such probability distribution
- pLSA tries to find a distribution that is "close" to **P** but exposes information about topics

Probabilities: Shortcut notation

Let X and Y be discrete random variables with possible values Val(X) and Val(Y). Let $x \in Val(X)$ and $y \in Val(Y)$.

Expression	Shortcut notation
P(X = x)	P(x)
P(X = x, Y = y)	P(x,y)
$P(X = x \mid Y = y)$	$P(x \mid y)$
$\forall x. P(X = x) = f(x)$	P(X) = f(X)
$\forall x.\forall y.P(X = x \mid Y = y) = f(x, y)$	$P(X \mid Y) = f(X, Y)$

- P(x), P(x, y), P(x | y) are numbers (= probabilities)
- P(X) and P(X | Y) are functions (= probability distributions)
- Can be thought of as functions from $Val(X) \rightarrow [0,1]$ or $Val(X) \times Val(Y) \rightarrow [0,1]$, respectively
- f_y(X) = P(X | y) is often referred to as conditional probability distribution (CPD)
Probabilities: Important properties $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ (inclusion-exclusion) $P(\bar{A}) = 1 - P(A)$ If $B \supset A$, $P(B) = P(A) + P(B \setminus A) > P(A)$ $P(X, Y) = P(Y \mid X)P(X)$ (product rule) $P(X) = \sum_{y} P(X, y)$ (sum rule, discrete) $P(X) = \int_{\mathcal{X}} P(X, y)$ (sum rule, continuous) $P(X \mid Y) = \frac{P(Y \mid X)P(X)}{P(Y)}$ (Bayes theorem) E[aX+b] = aE[X]+b(linearity of expectation) E[X + Y] = E[X] + E[Y] $E_Y[E_X[X \mid Y]] = E[X]$ (law of total expectation)

Probabilistic latent semantic analysis (pLSA)

Definition (pLSA, NMF formulation)

Given a rank r, find matrices L, Σ , and R such that

$P \approx L\Sigma R$

where

- $L_{m \times r}$ is a column-stochastic matrix,
- $\Sigma_{r \times r}$ is a non-negative, diagonal matrix that sums to unity, and
- $R_{r \times n}$ is a row-stochastic matrix.
- **Column-stochastic** = each column is probability vector = each column is non-negative and sums to 1
- $\bullet~\approx$ is usually taken to be the (generalized) KL divergence
- Regularization or "tempering" may be necessary to avoid overfitting

Example

• pLSA factorization of example matrix



- Rank r corresponds to number of topics
- σ_z corresponds to overall relative frequency of topic z
- I_{dz} corresponds to contribution of document d to topic z
- r_{zw} corresponds to frequency of word w in topic z
- pLSA constraints allow for probabilistic interpretation

$$P(d, w) \approx [L\Sigma R]_{dw} = \sum_{z} \sigma_{z} I_{dz} r_{zw} = \sum_{z} P(z) P(d \mid z) P(w \mid z)$$

• pLSA model assumes conditional independence, i.e., it assumes that words and documents are conditionally independent given a topic \rightarrow restricted space of distributions

Another example

Concepts (10 of 128) extracted from Science Magazine articles (12K)

P(w z)	universe	0.0439	drug	0.0672	cells	0.0675	sequence	0.0818	years	0.156
	galaxies	0.0375	patients	0.0493	stem	0.0478	sequences	0.0493	million	0.0556
	clusters	0.0279	drugs	0.0444	human	0.0421	genome	0.033	ago	0.045
	matter	0.0233	clinical	0.0346	cell	0.0309	dna	0.0257	time	0.0317
	galaxy	0.0232	treatment	0.028	gene	0.025	sequencing	0.0172	age	0.0243
	cluster	0.0214	trials	0.0277	tissue	0.0185	map	0.0123	year	0.024
	cosmic	0.0137	therapy	0.0213	cloning	0.0169	genes	0.0122	record	0.0238
	dark	0.0131	trial	0.0164	transfer	0.0155	chromosome	0.0119	early	0.0233
	light	0.0109	disease	0.0157	blood	0.0113	regions	0.0119	billion	0.0177
	density	0.01	medical	0.00997	embryos	0.0111	human	0.0111	history	0.0148
			-							
	bacteria	0.0983	male	0.0558	theory	0.0811	immune	0.0909	stars	0.0524
	bacteria bacterial	0.0983 0.0561	male females	0.0558 0.0541	theory physics	0.0811	immune response	0.0909 0.0375	stars star	0.0524 0.0458
	bacteria bacterial resistance	0.0983 0.0561 0.0431	male females female	0.0558 0.0541 0.0529	theory physics physicists	0.0811 0.0782 0.0146	immune response system	0.0909 0.0375 0.0358	stars star astrophys	0.0524 0.0458 0.0237
(z)	bacteria bacterial resistance coli	0.0983 0.0561 0.0431 0.0381	male females female males	0.0558 0.0541 0.0529 0.0477	theory physics physicists einstein	0.0811 0.0782 0.0146 0.0142	immune response system responses	0.0909 0.0375 0.0358 0.0322	stars star astrophys mass	0.0524 0.0458 0.0237 0.021
(x x)	bacteria bacterial resistance coli strains	0.0983 0.0561 0.0431 0.0381 0.025	male females female males sex	0.0558 0.0541 0.0529 0.0477 0.0339	theory physics physicists einstein university	0.0811 0.0782 0.0146 0.0142 0.013	immune response system responses antigen	0.0909 0.0375 0.0358 0.0322 0.0263	stars star astrophys mass disk	0.0524 0.0458 0.0237 0.021 0.0173
P(w z)	bacteria bacterial resistance coli strains microbiol	0.0983 0.0561 0.0431 0.0381 0.025 0.0214	male females female males sex reproductive	0.0558 0.0541 0.0529 0.0477 0.0339 0.0172	theory physics physicists einstein university gravity	0.0811 0.0782 0.0146 0.0142 0.013 0.013	immune response system responses antigen antigens	0.0909 0.0375 0.0358 0.0322 0.0263 0.0184	stars star astrophys mass disk black	0.0524 0.0458 0.0237 0.021 0.0173 0.0161
P(w z)	bacteria bacterial resistance coli strains microbiol microbiol	0.0983 0.0561 0.0431 0.0381 0.025 0.0214 0.0196	male females female males sex reproductive offspring	0.0558 0.0541 0.0529 0.0477 0.0339 0.0172 0.0168	theory physics physicists einstein university gravity black	0.0811 0.0782 0.0146 0.0142 0.013 0.013 0.013	immune response system responses antigen antigens immunity	0.0909 0.0375 0.0358 0.0322 0.0263 0.0184 0.0176	stars star astrophys mass disk black. gas	0.0524 0.0458 0.0237 0.021 0.0173 0.0161 0.0149
P(w z)	bacteria bacterial resistance coli strains microbiol microbial strain	0.0983 0.0561 0.0431 0.0381 0.025 0.0214 0.0196 0.0165	male females female males sex reproductive offspring sexual	0.0558 0.0541 0.0529 0.0477 0.0339 0.0172 0.0168 0.0166	theory physics physicists einstein university gravity black theories	0.0811 0.0782 0.0146 0.0142 0.013 0.013 0.013 0.0127 0.01	immune response system responses antigen antigens immunity immunology	0.0909 0.0375 0.0358 0.0322 0.0263 0.0184 0.0176 0.0145	stars star astrophys mass disk black gas stellar	0.0524 0.0458 0.0237 0.021 0.0173 0.0161 0.0149 0.0127
P(w z)	bacteria bacterial resistance coli strains microbiol microbial strain salmonella	0.0983 0.0561 0.0431 0.025 0.0214 0.0196 0.0165 0.0163	male females female males sex reproductive offspring sexual reproduction	0.0558 0.0541 0.0529 0.0477 0.0339 0.0172 0.0168 0.0166 0.0143	theory physics physicists einstein university gravity black theories aps	0.0811 0.0782 0.0146 0.0142 0.013 0.013 0.013 0.0127 0.01 0.00987	immune response system responses antigen antigens immunity immunology antibody	0.0909 0.0375 0.0358 0.0322 0.0263 0.0184 0.0176 0.0145 0.014	stars star astrophys mass disk black gas stellar astron	0.0524 0.0458 0.0237 0.021 0.0173 0.0161 0.0149 0.0127 0.0125

"Topics" (also: concepts, aspects) described by distributions of words.

pLSA geometry

• Rewrite probabilistic formulation

$$P(d, w) = \sum_{z} P(z)P(d \mid z)P(w \mid z)$$
$$P(w \mid d) = \sum_{z} P(w \mid z)P(z \mid d)$$

- Generative process
 - 1. Pick a document according to P(d)
 - 2. Select a topic acc. to $P(z \mid d)$
 - 3. Select a word acc. to $P(w \mid z)$



Figure 2. Sketch of the probability simplex and a convex region spanned by class-conditional probabilities in the aspect model.

Example

• pLSA factorization of example matrix, rewritten formulation



- This formulation more directly exposes topic distribution of documents (more in assignment)
- Next: Applications & a suitable definition of " \approx "

Applications of pLSA

- Topic modeling
- Clustering documents
- Clustering terms
- Information retrieval
 - Treat query q as a "new" document
 - Using the second formulation of the previous slide, determine $P(z \mid q)$, keeping $P(w \mid z)$ fixed (called "fold in" the query)
 - ▶ In more detail, find P(z | q) such that $P(w | q) \approx \sum_{z} P(w | z) P(z | q)$
 - Retrieve documents with similar topic mixture (= P(z | d))
 - Can deal with synonymy and polysemy
- Better generalization performance than LSA (=SVD), esp. with tempering
- Full Bayesian treatment: Latent Dirichlet Allocation (LDA)



Kullback-Leibler divergence (1)

- Let $\tilde{\pmb{P}}$ be the unnormalized word-count data and denote by N total number of words
- Probability of seeing \tilde{P} when drawing N words with replacement

$$P(ilde{oldsymbol{P}}) \propto \prod_{d=1}^m \prod_{w=1}^n P(d,w)^{ ilde{
ho}_{dw}}$$

• Conditional probability $P(\tilde{P} \mid L, \Sigma, R)$ of seeing the data given a pLSA model is as above with P(d, w) given by

$$P(d, w \mid \boldsymbol{L}, \boldsymbol{\Sigma}, \boldsymbol{R}) = [\boldsymbol{L} \boldsymbol{\Sigma} \boldsymbol{R}]_{dw}$$

- pLSA maximizes this probability = likelihood $\mathcal{L}(\tilde{\textbf{P}} \mid \textbf{L}, \pmb{\Sigma}, \textbf{R})$ of data given model
- Equivalent to maximizing the log-likelihood log $\mathcal{L}(\tilde{P} \mid L, \Sigma, R)$

Kullback-Leibler divergence (2)

Let
$$\hat{P} = L\Sigma R$$
.
 $\log \mathcal{L}(\tilde{P} \mid \hat{P}) \propto \sum_{d=1}^{m} \sum_{w=1}^{n} \tilde{p}_{dw} \log \hat{p}_{dw}$
 $\propto \sum_{d=1}^{m} \sum_{w=1}^{n} p_{dw} \log \hat{p}_{dw}$
 $= -\sum_{d=1}^{m} \sum_{w=1}^{n} p_{dw} \log \frac{1}{\hat{p}_{dw}}$
Cross entropy $H(P, \hat{P})$
 $= -\sum_{d=1}^{m} \sum_{w=1}^{n} p_{dw} \log \frac{p_{dw}}{\hat{p}_{dw}} + c_P$
Kullback-Leibler divergence $D_{KL}(P \parallel \hat{P})$

• Takeaway: maximizing the log-likelihood \equiv minimize cross entropy \equiv minimizing the Kullback-Leibler divergence Gaussier and Goutte, 2005.

Background: Entropy

- Consider any discrete contribution over values $\{a_1, \ldots, a_n\}$
- We can represent this distribution using a probability vector $\boldsymbol{p} \in \mathbb{R}^n$, where value a_i is associated with probability p_i
- Suppose Anna selects a random value A according to p, i.e., $P(A = a_i) = p_i$
- She wants to tell Bob which value she selected
 - ► Anna and Bob agree upfront on a codeword (bitstring) for each value *a_i*
 - On average, how many bits does Anna have to send to Bob so that he can determine the value of A?
- Answer: at best $H(\mathbf{p})$ bits, where $H(\mathbf{p})$ is the Shannon entropy

$$H(\boldsymbol{p}) = \sum_{i=1}^{n} p_i \log \frac{1}{p_i}$$

- p_i = probability that Anna selects value a_i
- $\log \frac{1}{p_i} =$ number of bits of codeword for a_i

Background: Cross entropy and KL divergence

- Let's modify the game a bit
 - Suppose Anna cheats and gives Bob the wrong distribution \boldsymbol{q}
 - They use an optimal code for \boldsymbol{q} , i.e., codeword for a_i has $\log \frac{1}{a_i}$ bits
 - But Anna later selects values according to p, not q
 - How many bits are sent on average?
- Answer: cross entropy of p and q

$$H(\boldsymbol{p},\boldsymbol{q}) = \sum_{i=1}^{n} p_i \log \frac{1}{q_i}$$

- p_i = probability that Anna selects value a_i
- $\log \frac{1}{q_i}$ = number of bits of codeword for a_i
- On average, how many additional bits are now sent?
- Answer: Kullback-Leibler divergence

$$D_{\mathsf{KL}}(\boldsymbol{p} \| \boldsymbol{q}) \stackrel{ ext{def}}{=} H(\boldsymbol{p}, \boldsymbol{q}) - H(\boldsymbol{p}) = \sum_{i=1}^{n} p_i \log rac{p_i}{q_i}$$

Kullback-Leibler divergence (3)

- KL divergence is a measure of "difference" between distributions
- In our setting, we have

$$D_{\mathsf{KL}}(oldsymbol{P} \| \hat{oldsymbol{P}}) = \sum_{d,w} p_{dw} \log rac{p_{dw}}{\hat{
ho}_{dw}}$$

- Interpretation: expected number of extra bits for encoding a value drawn from P using an optimum code for distribution \hat{P}
- $D_{\mathsf{KL}}(\boldsymbol{P} \| \hat{\boldsymbol{P}}) \geq 0$
- $D_{\mathsf{KL}}(\boldsymbol{P} \| \boldsymbol{P}) = 0$
- In general, $D_{\mathsf{KL}}(\boldsymbol{P} \| \hat{\boldsymbol{P}}) \neq D_{\mathsf{KL}}(\hat{\boldsymbol{P}} \| \boldsymbol{P})$
- NMF-based pLSA algorithms often minimize the generalized KL divergence (exercise)

$$D_{\mathsf{GKL}}(ilde{oldsymbol{P}}\| ilde{oldsymbol{\hat{P}}}) = \sum_{d,w} \Bigl(ilde{
ho}_{dw}\lograc{ ilde{
ho}_{dw}}{ ilde{oldsymbol{
ho}}_{dw}} - ilde{
ho}_{dw} + ilde{oldsymbol{\hat{p}}}_{dw} \Bigr),$$

where $ilde{m{ eta}} = ilde{m{L}} ilde{m{R}}$

Multiplicative updates for GKL (w/o tempering)

- We first find a decomposition $\tilde{P} \approx \tilde{L}\tilde{R}$, where \tilde{L} and \tilde{R} are non-negative matrices
 - ► E.g., using multiplicative update rules

$$\begin{split} \tilde{\boldsymbol{L}} &\leftarrow \tilde{\boldsymbol{L}} \circ \frac{\tilde{\boldsymbol{P}}}{\tilde{\boldsymbol{L}}\tilde{\boldsymbol{R}}} \tilde{\boldsymbol{R}}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{1}_r/(\tilde{\boldsymbol{R}}\boldsymbol{1}_n)) \\ \tilde{\boldsymbol{R}} &\leftarrow \tilde{\boldsymbol{R}} \circ \operatorname{diag}(\boldsymbol{1}_r/(\boldsymbol{1}_m^{\mathsf{T}}\tilde{\boldsymbol{L}})^{\mathsf{T}}) \tilde{\boldsymbol{L}}^{\mathsf{T}} \frac{\tilde{\boldsymbol{P}}}{\tilde{\boldsymbol{L}}\tilde{\boldsymbol{R}}} \end{split}$$

- GKL is non-increasing under these update rules
- Normalize by rescaling columns of \tilde{L} and rows of \tilde{R} to obtain

$$\begin{split} \boldsymbol{L} &= \tilde{\boldsymbol{L}} \operatorname{diag}(\boldsymbol{1}_r / (\boldsymbol{1}_m^T \tilde{\boldsymbol{L}})^T) & (\operatorname{colums of } \boldsymbol{L} \operatorname{sum to } 1) \\ \boldsymbol{R} &= \operatorname{diag}(\boldsymbol{1}_r / (\tilde{\boldsymbol{R}} \boldsymbol{1}_n)) \tilde{\boldsymbol{R}} & (\operatorname{rows of } \boldsymbol{R} \operatorname{sum to } 1) \\ \tilde{\boldsymbol{\Sigma}} &= \operatorname{diag}((\boldsymbol{1}_m^T \tilde{\boldsymbol{L}})^T \circ (\tilde{\boldsymbol{R}} \boldsymbol{1}_n)) & (\operatorname{scaling factors}) \\ \boldsymbol{\Sigma} &= \tilde{\boldsymbol{\Sigma}} / \sum_z \tilde{\boldsymbol{\Sigma}}_{zz} & (\boldsymbol{\Sigma} \operatorname{sums to } 1) \end{split}$$

Outline

- 1. Non-Negative Matrix Factorization
- 2. Algorithms
- 3. Probabilistic Latent Semantic Analysis
- 4. Summary

Lessons learned

- Non-negative matrix factorization (NMF) appears natural for non-negative data
- NMF encourages parts-based decomposition, interpretability, and (sometimes) sparseness
- Many variants, many applications
- Usually solved via alternating minimization algorithms
 - Alternating non-negative least squares (ANLS)
 - Projected gradient local hierarchical ALS (HALS)
 - Multiplicative updates
- pLSA is a well-known approach to topic modeling
 - Can be seen as an NMF

Literature

• David Skillicorn

Understanding Complex Datasets: Data Mining with Matrix Decompositions (Chapter 8) Chapman and Hall, 2007

• Andrzej Cichocki, Rafal Zdunek, Anh Huy Phan, and Shun-ichi Amari

Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation Wiley, 2009

- Yifeng Li and Alioune Ngom *The NMF MATLAB Toolbox* <u>http://cs.uwindsor.ca/ li11112c/nmf</u>
- Renaud Gaujoux and Cathal Seoighe *NMF R package*

http://cran.r-project.org/web/packages/NMF/index.html

• References given at bottom of slides

Data Mining and Matrices 06 – Spectral Clustering

Prof. Dr. Rainer Gemulla

Universität Mannheim

FSS 2018

Graph mining

- Graphs everywhere, e.g.
 - Internet
 - World wide web
 - Social networks
 - Protein-protein interactions
 - Similarity graphs
- Goals of graph mining
 - ► As data mining: classification, clustering, outliers, patterns
 - Output often also one or more graphs
 - Interesting subgraphs (e.g., communities, near-cliques, clusters)
 - Interesting vertices (e.g., influential bloggers, PageRank, outliers)
 - Web mining (e.g., topic predicition, classification)
 - Web usage mining (e.g., frequent subgraphs, patterns)
 - Recommender systems (e.g., movie recommendation)
 - Knowledge bases (e.g., link prediction)

۰...

Spectral analysis of matrices associated with graphs is an important tool in graph mining. Our focus: spectral clustering and link analysis.

A graph is a matrix is a graph

- Let G = (V, E) be a (weighted) graph
- Vertices $V = \{v_1, \ldots, v_n\}$
- Edge (i, j) ∈ E has positive weight w_{ij} (or 1 if graph is unweighted)
- Convention: absent edges $(i, j) \notin E$ have weight $w_{ij} = 0$
- Adjacency matrix W is $n \times n$ matrix with entries w_{ij}
- Undirected graph \implies **W** symmetric (**W** = **W**^T)
- (Out-)degree of vertex *i* given by $d_i = \sum_j w_{ij} = \boldsymbol{w}_i^T \boldsymbol{1}$
- Degree matrix **D** is $n \times n$ diagonal matrix with $d_{ii} = d_i$



Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition
- 4. Graph Laplacian
- 5. Spectral Clustering
- 6. Summary

k-Means example (1)



k-Means cannot detect non-convex clusters well.

k-Means example (2)



k-Means is sensitive to skew in cluster sizes.

A better clustering



In this clustering, points within a cluster are close to their neighbors, but not necessarily to all the points in the cluster.

Graph-based clustering

- 1. Given a dataset, construct a *similarity graph* modeling local neighborhood relationships
- 2. Partition the similarity graph using suitable graph cuts



Similarity graph

Clustering

Discussion

- Clustering
 - 1. Points within a cluster should be similar
 - 2. Points in different clusters should be dissimlar
- *k*-Means is global
 - 1. All points within a cluster should be similar (close)
 - 2. Points in different clusters should be dissimilar (far apart)
- Graph-based clustering is local
 - 1. Neighboring points within a cluster should be similar (close)
 - 2. Points in different clusters should be dissimilar (far apart)

Another example (1)



Global distances can be misleading.

Another example (2)



Distances to neighbors may capture structure better.

Which cut? (1)

- G = (V, E): Undirected, weighted similarity graph
- Cut = partitioning of vertices into two partitions $A \subset V$ and $\bar{A} = V \setminus A$
- For us: A and \overline{A} correspond to clusters
- Minimum cut is cut that minimizes weight

$$\operatorname{cut}(A, \overline{A}) = \sum_{i \in A, j \in \overline{A}} w_{ij}$$

- Can be solved efficiently (in P)
- Often not useful in practice, e.g., may separate a single vertex \rightarrow Need to balance cut weight and cluster sizes

Which cut? (2)

• Minimum ratio cut (penalize different sizes w.r.t. vertices)

$$\mathsf{RatioCut}(A,ar{A}) = \left(rac{1}{|A|} + rac{1}{|ar{A}|}
ight) \sum_{i \in A, j \in ar{A}} w_{ij}$$

• Minimum normalized cut (penalize different sizes w.r.t. degrees)

$$\operatorname{Ncut}(A, \overline{A}) = \left(\frac{1}{\operatorname{vol}(A)} + \frac{1}{\operatorname{vol}(\overline{A})}\right) \sum_{i \in A, j \in \overline{A}} w_{ij},$$

where vol(A) = $\sum_{i \in A} d_i$

• Both problems are NP-hard

Spectral clustering is a relaxation of RatioCut or Ncut, is simple to implement, and can be solved efficiently.

Which cut? (3)

- Recall clustering objectives
 - 1. Points in same cluster should be similar (maximize within-cluster similarity)
 - 2. Points in different clusters should be dissimilar (minimize between-cluster similarity)
- $(1) = \operatorname{vol}(A)$ and $\operatorname{vol}(ar{A})$ are both large
- (2) = minimize $cut(A, \overline{A})$
- cut, RatioCut, and Ncut all implement (2)
- Only Ncut additionally implements (1)
- Ncut captures both goals \rightarrow usually good choice

Outline

1. Graph-Based Clustering

2. Similarity Graphs

- 3. Background: Eigendecomposition
- 4. Graph Laplacian
- 5. Spectral Clustering
- 6. Summary

Similarity graph

- In graph-based clustering, we first need to construct a similarity graph
 - Need notion of similarity and a way to construct an appropriate graph based on that similarity
 - How to do this depends on domain
 - Generally difficult
 - Can significantly affect results
- E.g., similarity for documents
 - Similarity of words in documents
 - Similarity of words in documents with tf/idf weighting
 - Similarity of topic distribution of documents
 - ▶ ...
- Here:
 - How to obtain a similarites from Euclidean data
 - How to construct a similarity graph from similarities

From distances to similarities

- Sometimes: Need to "convert" distances to similarities
- Large distance $\delta_{ij} \iff$ small similarity w_{ij} (and vice versa)
- Simplest choice: reciprocal (problematic, unbounded)

$$w_{ij} = rac{1}{\delta_{ij}}$$

• Common choice for Euclidian data: Gaussian kernel (in [0,1])

$$w_{ij} = \exp(-\delta_{ij}^2/(2\sigma^2)),$$

where δ_{ij} is Euclidean distance $(||\mathbf{x}_i - \mathbf{x}_j||)$

 Parameter σ controls what is considered local (large σ = large neighborhood)



From distances to similarities (examples)





 $\sigma = 0.01$

(too small)



 $\sigma = 0.4$

(good)



 $\sigma = 3$ (too large)

Full graph

- Connect all pairs of vertices
- Weigh edges by similarity
- Generally expensive, not feasible for large datasets





ϵ -Neighborhood graph

- Pick neighborhood size ϵ
- Connect vertices of distance $\leq \epsilon$
- Unweighted or weighted by similarity





 ϵ too small



 $\epsilon \, \operatorname{good}$



Skewed clusters: ϵ too large for red, too small for black
Nearest neighbor graphs

- Pick number k of neighbors
- Directed k-nearest neighbor graph
 - ▶ Add directed edge (i,j) if j is among k closest neighbors of i
 - But: need undirected graph for well-defined similarities
- (Symmetric) k-nearest neighbor graph
 - Connect (i, j) if (i, j) or (j, i) in directed k-NN graph (OR)
 - ► Each node has at least k, but potentially more than k "neighbors"
- Mutual k-nearest neighbor graph
 - Connect (i, j) if (i, j) and (j, i) in directed k-NN graph (AND)
 - Each node has at most k, but potentially less than k "neighbors"
- Weigh edges by similarity



directed 2-NN





symmetric 2-NN

mutual 2-NN

k-Nearest neighbor graph (examples)

Symmetric k-NN graph



k = 1 (too small)



 $k = 10 \pmod{4}$



Skewed, $k = 10 \pmod{2}$





Discussion (1)

- Construction of similarity graph non-trivial and not well-understood
- Clustering results sensitive to choice of graph
- Which similarity function?
 - Should capture similarity of most-similar objects well (other edges pruned by neighborhood graphs anyway)
 - Gaussian kernel common choice for Euclidean data
 - Generally application-dependent
- Which graph?
 - Fully connected graph often too large + requires suitable similarity function, dense graph
 - \blacktriangleright $\epsilon\text{-neighborhood graph cannot deal well with clusters of different densities$
 - ► k-NN graph can connect points in regions with different densities → Generally recommended choice, sparse graph
 - Mutual k-NN graph is somewhere in between

Discussion (2)

- Which parameters? (ϵ or k, σ)
 - ϵ or k should be small so that graph is sparse
 - But large enough to ensure that similarity graph is connected (or at least has fewer components than desired clusters)
 - Otherwise: clustering sizes arbitrarily unbalanced, sensitive to outliers
 - ▶ k-NN: try various values (start with, e.g., $k = O(\log(n))$
 - Mutual k-NN: no good heuristics known
 - ► e-N: around length of longest edge in minimal spanning tree (problematic with outliers or clusters that are far apart)
 - σ: no. neighbors with similarity significantly larger than 0 "neither too small nor too large" (e.g., mean distance to k-th nearest neighbor, or using minimal spanning tree)

Skilled data miners do not run out of jobs.

Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition
- 4. Graph Laplacian
- 5. Spectral Clustering
- 6. Summary

Eigenvectors and eigenvalues

• A non-zero vector $\boldsymbol{v} \in \mathbb{R}^n$ is an eigenvector of $\boldsymbol{A} \in \mathbb{R}^{n \times n}$ if

 $Av = \lambda v$

- λ is the corresponding **eigenvalue**
- The eigenvalues are the roots of the characteristic polynomial

 $p_{\boldsymbol{A}}(\lambda) = \det(\boldsymbol{A} - \lambda \boldsymbol{I})$

• We can factor $p_{\mathbf{A}}(\lambda)$ as

$$p_{\mathbf{A}}(\lambda) = (\lambda - \lambda_1)^{n_1} \cdots (\lambda - \lambda_N)^{n_N},$$

where $\sum_{i} n_i = n$

- *n_i* is called the algebraic multiplicity of λ_i
- ► There are 1 ≤ m_i ≤ n_i linearly independent eigenvectors associated with eigenvalue λ_i
- *m_i* is called the geometric multiplicity of λ_i
- Note: Some eigenvectors can be complex
- Collection of eigenvalues is called spectrum of A

Eigendecomposition

• The eigendecomposition of $\mathbf{A} \in \mathbb{R}^{n \times n}$ is given by

 $\boldsymbol{A} = \boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^{-1},$

where

- **Q** is square and has eigenvectors as its columns
- \blacktriangleright **\Lambda** is diagonal and has eigenvalues on its diagonal
- Does not always exist; if it does, **A** is called diagonalizable
- Some properties
 - ▶ When A is symmetric, there exists an eigendecomposition where Q is real and orthogonal
 - $\bullet \mathbf{A} = \mathbf{B}\mathbf{B}^{\mathsf{T}} = (\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{\mathsf{T}})(\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{\mathsf{T}})^{\mathsf{T}} = \mathbf{U}\boldsymbol{\Sigma}^{2}\mathbf{U}^{\mathsf{T}}$
 - E.g., $\mathbf{A} = \mathbf{B}^T \mathbf{B} = (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^T (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T) = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^T$
 - E.g., $\boldsymbol{A} = \boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^{-1}$, then $\operatorname{tr}(\boldsymbol{A}) = \operatorname{tr}(\boldsymbol{\Lambda}) = \sum_{i} \lambda_{i}$
 - ▶ rank (**A**) = number of non-zero eigenvalues (counting multiplicity)

Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition
- 4. Graph Laplacian
- 5. Spectral Clustering
- 6. Summary

Graph Laplacian

Definition

Let G be an undirected graph with positive edge weights. Denote by \boldsymbol{W} the (weighted) adjacency matrix of G, and by \boldsymbol{D} the degree matrix of G. Then

L = D - W

is called the (unnormalized) graph Laplacian of G.

Note that self edges $(w_{ii} > 0)$ do not affect the graph Laplacian.

Graph Laplacians are the main tool for spectral clustering, but they have many other uses too (e.g., label propagation, graph drawing).

Properties of the graph Laplacian (1)

Theorem

For every vector
$$\mathbf{x} \in \mathbb{R}^n$$
, $f(\mathbf{x}) = \mathbf{x}^T \mathbf{L} \mathbf{x} = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_{ij} (x_i - x_j)^2$.

- x assigns a real value to each vertex
- f(x) is a quadratic form and small when "similar" vertices (which are connected with high-weight edges) take similar values

Proof.

$$\begin{aligned} \mathbf{x}^{T} \mathbf{L} \mathbf{x} &= \mathbf{x}^{T} \mathbf{D} \mathbf{x} - \mathbf{x}^{T} \mathbf{W} \mathbf{x} = \sum_{i=1}^{n} d_{i} x_{i}^{2} - \sum_{i=1}^{n} \sum_{j=1}^{n} w_{ij} x_{i} x_{j} \\ &= \frac{1}{2} \left(\sum_{i=1}^{n} d_{i} x_{i}^{2} - 2 \sum_{i=1}^{n} \sum_{j=1}^{n} w_{ij} x_{i} x_{j} + \sum_{j=1}^{n} d_{j} x_{j}^{2} \right) \\ &= \frac{1}{2} \left(\sum_{i=1}^{n} \sum_{j=1}^{n} w_{ij} (x_{i}^{2} - 2x_{i} x_{j} + x_{j}^{2}) \right) = \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} w_{ij} (x_{i} - x_{j})^{2} \end{aligned}$$

Properties of the graph Laplacian (2)

$$\mathbf{x}^T \mathbf{L} \mathbf{x} = \frac{1}{2} \sum_{i,j} w_{ij} (x_i - x_j)^2$$







Properties of the graph Laplacian (3)

A matrix $\mathbf{A}^{n \times n}$ is called **positive semi-definite** if $\mathbf{x}^T \mathbf{A} \mathbf{x} \ge 0$ for any $\mathbf{x} \in \mathbb{R}^n$.

Theorem

L is symmetric and positive semi-definite.

- Implies that $f(\mathbf{x}) = \mathbf{x}^T \mathbf{L} \mathbf{x}$ is a convex function
- Implies that $\mathbf{L} = \mathbf{P}\mathbf{P}^T$ for some \mathbf{P} (oriented incidence matrix)

Proof. Since **D** and **W** are symmetric, so is **L**. Since $\mathbf{x}^T \mathbf{L} \mathbf{x} \ge 0$ (see slide 30) for all $\mathbf{x} \in \mathbb{R}^n$, **L** is positive semi-definite.

Properties of the graph Laplacian (4)

Theorem

The smallest eigenvalue of L is zero, the corresponding eigenvector is the constant one vector 1.

$$\begin{array}{c|c} \hline 1 & 1 & 1 \\ \hline & & 1 \\ \hline & & 1 \\ \hline & & 1 \\ \lambda_3 = 0 \end{array}$$

Proof. The row sums satisfy L1 = 0 = 01 by construction. For smallest, see next slide.

Properties of the graph Laplacian (5)

Theorem

All eigenvalues are non-negative and real-valued, i.e., $\lambda_1 \ge \ldots \ge \lambda_{n-1} \ge \lambda_n = 0.$

Proof. All eigenvalues of a symmetric matrix are real. If $\boldsymbol{L}\boldsymbol{v} = \lambda \boldsymbol{v}$, then $\boldsymbol{v}^T \boldsymbol{L} \boldsymbol{v} = \lambda \|\boldsymbol{v}\|^2 \ge 0$ and thus $\lambda \ge 0$.

Connected graphs

Theorem

If G is connected, then eigenvalue 0 has multiplicity 1, i.e., $\lambda_{n-1} > 0$.



Proof. Recall that **1** is an eigenvector of **L** with eigenvalue 0. Suppose that $\mathbf{0} \neq \mathbf{v} \neq c\mathbf{1}$ is an eigenvector of **L** with eigenvalue λ . Since *G* is connected, this implies that there are two neighboring vertices i' and j' such that $v_{i'} \neq v_{j'}$. Now

$$\lambda \|\mathbf{v}\|^{2} = \mathbf{v}^{T} \mathbf{L} \mathbf{v} = \frac{1}{2} \sum_{i,j} w_{ij} (v_{i} - v_{j})^{2} \ge w_{i'j'} (v_{i'} - v_{j'})^{2} > 0$$

so that $\lambda > 0$.

Connected components

Theorem

The multiplicity k of eigenvalue 0 is equal to the number of connected components G_1, \ldots, G_k of G. The corresponding eigenspace is spanned by the indicator vectors $\mathbf{1}_{G_i}$ (value 1 for vertices in G_i , value 0 otherwise).

Proof. Let L_1, \ldots, L_k be the graph Laplacian of the connected components. Order w.l.o.g. the vertices by their component so that

$$\boldsymbol{L} = \begin{pmatrix} \boldsymbol{L}_1 & & & \\ & \boldsymbol{L}_2 & & \\ & & \ddots & \\ & & & \boldsymbol{L}_k \end{pmatrix}$$

Since L is block-diagonal, the spectrum of L is given by the union of the spectra of the L_i . The corresponding eigenvectors are the eigenvectors of L_i , filled with 0 at positions of other blocks.

Connected components (example)







Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition

4. Graph Laplacian

- 5. Spectral Clustering
- 5.1 Unnormalized Spectral Clustering
- 5.2 Normalized Spectral Clustering
- 5.3 Variants

6. Summary

Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition
- 4. Graph Laplacian
- 5. Spectral Clustering5.1 Unnormalized Spectral Clustering5.2 Normalized Spectral Clustering5.3 Variants

6. Summary

Algorithm

- Algorithm to construct k clusters
 - 1. Construct similarity graph \boldsymbol{W}
 - 2. Compute its (unnormalized) graph Laplacian L
 - Compute the last k eigenvectors u_n,..., u_{n-k+1} of L (i.e., having k smallest eigenvalues)
 - 4. Construct $n \times k$ matrix $\boldsymbol{U} = (\boldsymbol{u}_n \quad \boldsymbol{u}_{n-1} \quad \cdots \quad \boldsymbol{u}_{n-k+1})$
 - 5. Cluster the rows of U using k-means
- Simple, easy to implement
- Main trick: represent (or "embed") each vertex into R^k (= rows of U)
- Change of representation enhances clustering properties in the data

Why does this work? Why are we interested in the smallest eigenvalues?

Unnormalized spectral clustering (example, 1)



Similarity graph



Spectral clustering (k = 2)





Unnormalized spectral clustering (example, 2)



Similarity graph



Spectral clustering (k = 2)





Unnormalized spectral clustering (example)



Spectral clustering (k = 2)



Spectral clustering (k = 4)



Spectral clustering (k = 3)



Spectral clustering (k = 5)

Why does spectral clustering work? (1)

• Consider the minimum ratio cut problem (k = 2)

$$\min_{A \subset V} \mathsf{RatioCut}(A, \bar{A}) = \min_{A \subset V} \sum_{i \in A, j \in \bar{A}} w_{ij} \left(\frac{1}{|A|} + \frac{1}{|\bar{A}|} \right)$$

• Given A, set $\boldsymbol{x} \in \mathbb{R}^n$ such that

$$x_i = \begin{cases} \sqrt{|\bar{A}|/|A|} & \text{if } v_i \in A \\ -\sqrt{|A|/|\bar{A}|} & \text{if } v_i \in \bar{A} \end{cases}$$

- Easy to show
 - 1. $\mathbf{x}_{n}^{T} \mathbf{L} \mathbf{x} = n \cdot \operatorname{RatioCut}(A, \bar{A})$

2.
$$\sum_{i=1}^{n} x_i = 0 \text{ so that } \mathbf{x} \perp \mathbf{1}$$

3.
$$\|\mathbf{x}\|^2 = n$$

Why does spectral clustering work? (2)

- Minimum ratio cut can be rewritten as minimize $\mathbf{x}^T \mathbf{L} \mathbf{x}$ subject to $\mathbf{x} \perp \mathbf{1}$ $\|\mathbf{x}\| = \sqrt{n}$ \mathbf{x} takes form defined in previous slide
- Still NP-hard; relax by dropping discreteness constraint minimize $\mathbf{x}^T \mathbf{L} \mathbf{x}$ subject to $\mathbf{x} \perp \mathbf{1}$ $\|\mathbf{x}\| = \sqrt{n}$
- By *Rayleigh-Ritz theorem*: solution is eigenvector corresponding to second-smallest eigenvalue (appropriately normalized)

•
$$\boldsymbol{u}_{n-1}^T \boldsymbol{L} \boldsymbol{u}_{n-1} = \boldsymbol{u}_{n-1}^T \lambda_{n-1} \boldsymbol{u}_{n-1} = n \lambda_{n-1}$$

- Thus: $\lambda_{n-1} \leq \min_{A \subset V} \operatorname{RatioCut}(A, \overline{A})$
- Similar arguments for k > 2: solutions of relaxation
 eigenvectors u_{n-1},..., u_{n-k+1} (see exercise)

Why does spectral clustering work? (3)

Recall

$$x_i = \begin{cases} \sqrt{|\bar{A}|/|A|} & \text{if } v_i \in A \\ -\sqrt{|A|/|\bar{A}|} & \text{if } v_i \in \bar{A} \end{cases}$$

in the constrained solution

• Need to obtain clustering from \boldsymbol{u}_{n-1}

- Generally does not satisfy the discreteness constraints
- Simple heuristic: use sign as cluster indicator (threshold=0)
- Optimal: pick threshold that minimizes RatioCut (*optimal* thresholding)
- k-means often used in practice (also works well for k > 2)
- In general, no guarantees to obtain good solution
- But: popular because simple, standard linear algebra problem that tends to work well in practice
- Approximation of balanced graph cuts (up to constant factor) still hard

Cockroach graph

- Example where spectral clustering performs particularly bad
- Minimum ratio cut 8/n
- Spectral clustering ratio cut: 1
- Spectral clustering is O(n) times worse



Minimum ratio cut



Ratio cut with spectral clustering and sign heuristic

Discussion (1)

- Computation of eigenvectors
 - Graph can be very large
 - But Laplacian is sparse
 - Efficient algorithms for finding the eigendecomposition of such matrices exist
- Number k of clusters
 - Difficult problem
 - Standard approaches can be used
 - ► Eigengap heuristic: choose k such that eigenvalue λ₁,..., λ_{n-k} large, eigenvalues λ_{n-k+1},..., λ_n small



Discussion (2)



Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition
- 4. Graph Laplacian
- 5. Spectral Clustering5.1 Unnormalized Spectral Clustering5.2 Normalized Spectral Clustering5.3 Variants

6. Summary

Normalized graph Laplacians

Definition

There are two common normalizations of the graph Laplacian:

$$\begin{split} {\pmb L}_{sym} &= {\pmb D}^{-1/2} {\pmb L} {\pmb D}^{-1/2} = {\pmb I} - {\pmb D}^{-1/2} {\pmb W} {\pmb D}^{-1/2} \\ {\pmb L}_{rw} &= {\pmb D}^{-1} {\pmb L} = {\pmb I} - {\pmb D}^{-1} {\pmb W} \end{split}$$

- Normalization is performed w.r.t. degree
- **L**_{sym} is symmetric, **L**_{rw} is not

Normalized spectral clustering

- Normalized graph Laplacians have similar spectral properties
- Normalized spectral clustering (using L_{rw})
 - 1. Construct similarity graph \boldsymbol{W}
 - 2. Compute its normalized graph Laplacian L_{rw}
 - Compute the last k eigenvectors u_n,..., u_{n-k+1} of L_{rw} (i.e., having k smallest eigenvalues)
 - 4. Construct $n \times k$ matrix $\boldsymbol{U} = (\boldsymbol{u}_n \quad \boldsymbol{u}_{n-1} \quad \cdots \quad \boldsymbol{u}_{n-k+1})$
 - 5. Cluster the rows of U using k-means
- Normalized spectral clustering is a relaxation of Ncut
- Better behaved from statistical point of view

The normalized spectral clustering algorithm above is often method of choice.

Some examples



Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition
- 4. Graph Laplacian

5. Spectral Clustering

- 5.1 Unnormalized Spectral Clustering
- 5.2 Normalized Spectral Clustering
- 5.3 Variants

6. Summary

Some variants of spectral clustering

- p-spectral clustering
 - Use $\sum_{i,j} w_{ij} |x_i x_j|^p$
 - Standard spectral clustering is obtained for p = 2
 - As $p \to 1$, we obtain provably better (Cheeger) cuts
- Constrained spectral clustering
 - Way to incorporate domain knowledge and side information
 - Intuitively, obtained clustering must satisfy some application-defined constraints
 - E.g., must-link and cannot-link constraints (i.e., provide vertex pairs that must or must not end up in same cluster)
 - ► E.g., partial labels (i.e., provide labels for some vertexes) → graph-based semi-supervised learning (tutorial)
- Active area of research

Outline

- 1. Graph-Based Clustering
- 2. Similarity Graphs
- 3. Background: Eigendecomposition
- 4. Graph Laplacian
- 5. Spectral Clustering
- 6. Summary
Lessons learned

- Graphs can be represented by matrices (and vice versa)
 - Adjacency matrix
 - Degree matrix
 - Walk matrix
 - Graph Laplacian
- Spectral properties of these matrices relate to properties of the graph
- Spectral clustering
 - Find non-convex clusters in similarity graphs
 - Good clustering \approx good graph cut (RatioCut or Ncut)
 - Related to smallest eigenvectors of graph Laplacian
 - Many useful variants

Literature

• David Skillicorn

Understanding Complex Datasets: Data Mining with Matrix Decompositions (Chapter 4) Chapman and Hall, 2007

Ulrike von Luxburg
 A Tutorial on Spectral Clustering
 Statistics and Computing, 17(4), 2007
 http://www.kyb.mpg.de/publications/attachments/
 Luxburg07_tutorial_4488%5B0%5D.pdf

Data Mining and Matrices 08 – Tensors

Prof. Dr. Rainer Gemulla

Universität Mannheim

FSS 2018

- 1. What Is a Tensor?
- 2. Tensor Basics
- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 5. Wrap-Up

What is a tensor?

A tensor is a multi-way extension of a matrix.

- A multi-dimensional array
- A multi-linear map



Why tensors?

- Tensors can be used when matrices are not enough
- A matrix can represent a binary relation (or function)
 - Entities and attributes
 - Vertices and edges
- A tensor can represent an *n*-ary relation (or function)
 - A relationship between n sets of "things"
 - E.g. Subject-predicate-object data
- A tensor can represent a set of binary relations (or functions)
 - Vertices and typed edges ("colored")
 - E.g., images (rows and columns) for a set of persons
- A tensor can represent a sequence of binary (or *n*-ary) relations (or functions)
 - Objects and their attributes over time
 - But: using tensors for time series should be approached with care

Two examples

CT scans

Knowledge bases







1. What Is a Tensor?

- 2. Tensor Basics
- 2.1 Indexing
- 2.2 Matricization
- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 5. Wrap-Up

Terminology

- We say a tensor is an *N*-way array
 - E.g., a matrix is a 2-way array
- In the literature, one also finds
 - N-dimensional tensor (confusing: is a 3-dimensional vector a 1-dimensional tensor?)
 - rank-N tensor
 (but: we have a different use for the word rank)
- A 3-way tensor
 - Has three modes: rows, columns, and tubes
 - ► Can be *I*-by-*J*-by-*K*-dimensional
 - ► Also used: *M*-by-*N*-by-*K* or *I*₁-by-*I*₂-by-*I*₃
- Notation
 - Tensors: $\mathcal{X}, \mathcal{Y}, \mathcal{Z}, \ldots$
 - Sizes: I, J, K, ...
 - ▶ Indexes: *i*, *j*, *k*, ...

1. What Is a Tensor?

Tensor Basics
 Indexing
 Matricization

- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 5. Wrap-Up

Fibers

- We can extract subarrays by suitable indexing operations
- If we fix all dimensions, we obtain a scalar (e.g., x_{ijk})
- If we fix all but one dimension, we obtain a fiber
 - Convention: fibers are extracted as column vectors



Slices

- If we fix all but two dimensions, we obtain a slice
 - Convention: slices are extracted as matrices
 - Rows correspond to first omitted dimension, columns to second



Example

$$\boldsymbol{\mathcal{X}} = \begin{pmatrix} 5 & 7\\ 1 & 3\\ 2 & 4 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 3\\ 2 & 4 \end{pmatrix} \begin{pmatrix} 5 & 7\\ 6 & 8 \end{pmatrix}$$

• Column
$$\boldsymbol{x}_{:11} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$
, row $\boldsymbol{x}_{1:1} = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$, tube $\boldsymbol{x}_{11:} = \begin{pmatrix} 1 \\ 5 \end{pmatrix}$

• Frontal slice
$$\boldsymbol{X}_{::1} = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$$

• Lateral slice
$$\boldsymbol{X}_{:1:} = \begin{pmatrix} 1 & 5 \\ 2 & 6 \end{pmatrix}$$

• Horizontal slice
$$\boldsymbol{X}_{1::} = \begin{pmatrix} 1 & 5 \\ 3 & 7 \end{pmatrix}$$

1. What Is a Tensor?

Tensor Basics
 Indexing
 Matricization

- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 5. Wrap-Up

Matricization

- We can often express tensor operations using suitable matrices
 - Allows use of standard linear algebra packages
 - We'll repeatedly need this, so we discuss it first
 - In particular: matricization, Kroenecker product, Khatri-Rao product, Hadamard product
- Tensor matricization (*flattening*, *unfolding*) unfolds an *N*-way tensor into a matrix
 - ► Mode-*n* matricization X_(n) arranges the mode-*n* fibers as columns of a matrix
 - As many rows as size of *n*-th mode
 - As many columns as the product of the sizes of the other modes
 - ► Cumbersome to express formally: If X is an I₁ × I₂ × ... × I_n tensor, then X_(n) contains x_{i₁i₂...i_N} at position (i_n, j), where

$$j = 1 + \sum_{k=1}^{N} (i_k - 1) J_k[k \neq n]$$
 with $J_k = \prod_{m=1}^{k-1} I_m[m \neq n]$

Different authors may use different orders (!)

Example

$$\boldsymbol{\mathcal{X}} = \begin{pmatrix} 5 & 7\\ 16 & 3\\ 2 & 4 \end{pmatrix} \rightarrow \underbrace{\begin{pmatrix} 1 & 3\\ 2 & 4 \end{pmatrix}}_{\boldsymbol{X}_1} \underbrace{\begin{pmatrix} 5 & 7\\ 6 & 8 \end{pmatrix}}_{\boldsymbol{X}_2}$$

• Column
$$\mathbf{x}_{:11} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$
, row $\mathbf{x}_{1:1} = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$, tube $\mathbf{x}_{:11} = \begin{pmatrix} 1 \\ 5 \end{pmatrix}$
• $\mathbf{X}_{(1)} = \begin{pmatrix} 1 & 3 & 5 & 7 \\ 2 & 4 & 6 & 8 \end{pmatrix} = (\mathbf{X}_1 \quad \mathbf{X}_2)$
• $\mathbf{X}_{(2)} = \begin{pmatrix} 1 & 2 & 5 & 6 \\ 3 & 4 & 7 & 8 \end{pmatrix}$
• $\mathbf{X}_{(3)} = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{pmatrix}$

Kronecker matrix product

- Rescales one matrix with each element of another matrix
- $m_1 imes n_1$ and $m_2 imes n_2$ matrices give $m_1 m_2 imes n_1 n_2$ matrix

$$\boldsymbol{A} \otimes \boldsymbol{B} = \begin{pmatrix} a_{11}\boldsymbol{B} & a_{12}\boldsymbol{B} & \cdots & a_{1n_1}\boldsymbol{B} \\ a_{21}\boldsymbol{B} & a_{22}\boldsymbol{B} & \cdots & a_{2n_1}\boldsymbol{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m_11}\boldsymbol{B} & a_{m_12}\boldsymbol{B} & \cdots & a_{m_1n_1}\boldsymbol{B} \end{pmatrix}$$

• For example,

$$\begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \otimes \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 3 & 3 & 3 \\ 1 & 1 & 1 & 3 & 3 & 3 \\ 2 & 2 & 2 & 4 & 4 & 4 \\ 2 & 2 & 2 & 4 & 4 & 4 \end{pmatrix}$$

Khatri-Rao matrix product

- Column-wise Kroenecker product
 - Rescales the columns of one matrix with each element in the respective column of another matrix
 - Number of columns must match
- $m_1 \times n$ and $m_2 \times n$ matrices give $m_1 m_2 \times n$ matrix

$$\mathbf{A} \odot \mathbf{B} = \begin{pmatrix} \mathbf{a}_1 \otimes \mathbf{b}_1 & \mathbf{a}_2 \otimes \mathbf{b}_2 & \cdots & \mathbf{a}_n \otimes \mathbf{b}_n \end{pmatrix}$$
$$= \begin{pmatrix} a_{11}\mathbf{b}_1 & a_{12}\mathbf{b}_2 & \cdots & a_{1n}\mathbf{b}_n \\ a_{21}\mathbf{b}_1 & a_{22}\mathbf{b}_2 & \cdots & a_{2n}\mathbf{b}_n \\ \vdots & \vdots & \ddots & \vdots \\ a_{m_11}\mathbf{b}_1 & a_{m_12}\mathbf{b}_2 & \cdots & a_{m_1n}\mathbf{b}_n \end{pmatrix}$$

• For example,

$$\begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \odot \begin{pmatrix} 1 & 100 \\ 10 & 1000 \end{pmatrix} = \begin{pmatrix} 1 & 300 \\ 10 & 3000 \\ 2 & 400 \\ 20 & 4000 \end{pmatrix}$$

Hadamard matrix product

- Element-wise product
 - Rescales every entry of one matrix with the corresponding entry of another matrix
 - Number of rows and columns must match
- Two $m \times n$ matrices give $m \times n$ matrix

$$\boldsymbol{A} * \boldsymbol{B} = \begin{pmatrix} a_{11}b_{11} & a_{12}b_{12} & \cdots & a_{1n}b_{1n} \\ a_{21}b_{21} & a_{22}b_{22} & \cdots & a_{2n}b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}b_{m1} & a_{m2}b_{m2} & \cdots & a_{mn}b_{mn} \end{pmatrix}$$

• For example,

$$\begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} * \begin{pmatrix} 1 & 100 \\ 10 & 1000 \end{pmatrix} = \begin{pmatrix} 1 & 300 \\ 20 & 4000 \end{pmatrix}$$

Some identities

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{A}\mathbf{C} \otimes \mathbf{B}\mathbf{D}$$
$$(\mathbf{A} \otimes \mathbf{B})^{\dagger} = \mathbf{A}^{\dagger} \otimes \mathbf{B}^{\dagger}$$
$$\mathbf{A} \odot \mathbf{B} \odot \mathbf{C} = (\mathbf{A} \odot \mathbf{B}) \odot \mathbf{C} = \mathbf{A} \odot (\mathbf{B} \odot \mathbf{C})$$
$$(\mathbf{A} \odot \mathbf{B})^{\mathsf{T}} (\mathbf{A} \odot \mathbf{B}) = \mathbf{A}^{\mathsf{T}} \mathbf{A} * \mathbf{B}^{\mathsf{T}} \mathbf{B}$$
$$(\mathbf{A} \odot \mathbf{B})^{\dagger} = (\mathbf{A}^{\mathsf{T}} \mathbf{A} * \mathbf{B}^{\mathsf{T}} \mathbf{B})^{\dagger} (\mathbf{A} \odot \mathbf{B})^{\mathsf{T}}$$

- Here $\textbf{\textit{A}}^{\dagger}$ denotes the pseudo-inverse of $\textbf{\textit{A}}$
- Using these equations may save considerable space and time; e.g.
 - Let $\boldsymbol{A}, \boldsymbol{B}$ be $n \times k$ matrices with $n \gg k$
 - $\boldsymbol{A} \odot \boldsymbol{B}$ is an $n^2 \times k$ matrix
 - $(\mathbf{A} \odot \mathbf{B})^T (\mathbf{A} \odot \mathbf{B}), \ \mathbf{A}^T \mathbf{A}, \ \mathbf{B}^T \mathbf{B} \text{ and } \mathbf{A}^T \mathbf{A} * \mathbf{B}^T \mathbf{B} \text{ are } k \times k$ matrices

- 1. What Is a Tensor?
- 2. Tensor Basics
- 3. The CP Decomposition
- 3.1 Definition
- 3.2 Tensor Rank
- 3.3 Computing the CP Decomposition
- 3.4 Variants
- 4. The Tucker Decomposition
- 5. Wrap-Up

- 1. What Is a Tensor?
- 2. Tensor Basics

3. The CP Decomposition 3.1 Definition

- 3.2 Tensor Rank
- 3.3 Computing the CP Decomposition
- 3.4 Variants
- 4. The Tucker Decomposition
- 5. Wrap-Up

Outer product

- Tensors require extensions to the standard linear algebra operations for matrices
- We'll discuss those in turn whenever needed
- A multi-way vector outer product is a tensor where each element is the product of corresponding elements in the vectors

$$\mathcal{X} = \mathbf{a} \circ \mathbf{b} \circ \mathbf{c} \qquad \iff \qquad x_{ijk} = a_i b_j c_k$$



• Observe: $\boldsymbol{a} \circ \boldsymbol{b} = \boldsymbol{a} \boldsymbol{b}^T$

Example

- $\boldsymbol{a} = \begin{pmatrix} 1 & 2 & 3 \end{pmatrix}^{T}$ • $\boldsymbol{b} = \begin{pmatrix} 1 & 2 & 4 \end{pmatrix}^{T}$ • $\boldsymbol{c} = \begin{pmatrix} 1 & 10 \end{pmatrix}^{T}$
- $\mathcal{X} = \mathbf{a} \circ \mathbf{b} \circ \mathbf{c}$ is a $3 \times 3 \times 2$ tensor
- The frontal slices are then

$$\boldsymbol{X}_{1} = (\boldsymbol{a} \circ \boldsymbol{b})c_{1} = (\boldsymbol{a}\boldsymbol{b}^{T})c_{1} = \begin{pmatrix} 1 & 2 & 4 \\ 2 & 4 & 8 \\ 3 & 6 & 12 \end{pmatrix}$$
$$\boldsymbol{X}_{2} = (\boldsymbol{a} \circ \boldsymbol{b})c_{2} = (\boldsymbol{a}\boldsymbol{b}^{T})c_{2} = \begin{pmatrix} 10 & 20 & 40 \\ 20 & 40 & 80 \\ 30 & 60 & 120 \end{pmatrix}$$

Rank-1 tensors

- A matrix decomposition represents the given matrix as the product of two (or more) factor matrices
- The *rank* of a matrix **M** is the
 - Number of linearly independent rows (row rank)
 - Number of linearly independent columns (column rank)
 - Minimum number of rank-1 matrices needed to be summed up to get *M* (*Schein rank*)
 - All definitions are equivalent
- Let's first look at rank-1 tensors
 - A rank-1 matrix is an outer product of two vectors
 - ► We define: An *N*-way tensor is a rank-1 tensor if it can be written as an outer product of *N* vectors
 - E.g., the tensor on the previous slide

The CP decomposition



- The CP decomposition approximates \mathcal{X} by summing up R rank-1 tensors
- R is called the size of the CP decomposition
- Can also be written using N factor matrices: $X = \llbracket A, B, C \rrbracket$
 - For each mode, create a matrix with the corresponding vectors of each rank-1 component as columns
 - $\blacktriangleright \text{ E.g., } \boldsymbol{A} = \begin{pmatrix} \boldsymbol{a}_1 & \boldsymbol{a}_2 & \cdots & \boldsymbol{a}_R \end{pmatrix}$
 - We have $x_{ijk} \approx \sum_{r=1}^{R} a_{ir} b_{jr} c_{kr}$

The many names of the CP decomposition

Name	Proposed by
Polyadic Form of a Tensor	Hitchcock, 1927
PARAFAC (Parallel Factors)	Harshman, 1970
CANDECOMP or CAND (Canonical decomposition)	Carroll and Chang, 1970
Topographic Components Model	Möcks, 1988
CP (CANDECOMP/PARAFAC)	Kiers, 2000

• Chemometrics, psychometrics, phonetics, sensor data processing, telecommunications, neuroscience, medical data, text mining, social network analysis, image compression, computer vision,

Example: TOPHITS

- Builds a web pages-by-web pages-by-anchor text tensor to study the link structure and link topics of web pages
 - w_{ijk} is the number of times page *i* links to page *j* using term *k*
- The CP decomposition of this tensor behaves akin to HITS
 - In size-1 CP, *a* gives hub scores, *b* gives authority scores, and *c* gives term weights
 - ► In size-*R* CP, the data is split up into multiple topics, each with its own hubs, authorities, and term weights
 - Helps to avoid topic drift, not restricted to focused subgraphs
- Used a higher-order power method to compute the CP

	Topics	Authorities			
SCORE	Term	SCORE	Ноѕт		
1st Principal Factor					
0.23	java	0.86	java.sun.com		
0.18	sun	0.38	developers.sun.com		
0.17	platform	0.16	docs.sun.com		
0.16	solaris	0.14	see.sun.com		
0.16	developer	0.14	www.sun.com		
0.15	edition	0.09	www.samag.com		
0.15	download	0.07	developer.sun.com		
0.14	info	0.06	sunsolve.sun.com		
0.12	software	0.05	access1.sun.com		
		0.05	iforce sup com		

- 1. What Is a Tensor?
- 2. Tensor Basics

3. The CP Decomposition

3.1 Definition

3.2 Tensor Rank

- 3.3 Computing the CP Decomposition
- 3.4 Variants

4. The Tucker Decomposition

5. Wrap-Up

Tensor rank

- The rank of a tensor is the minimum number of rank-1 tensors needed to represent the tensor exactly
 - ► Given by the exact CP decomposition of size R = rank (X), also called rank decomposition
 - Generalizes the notion of Schein rank for matrices



- Tensors behave differently than matrices
 - Much more complicated on the one hand
 - (Some) nicer properties on the other hand

Tensor rank oddities (1)

• The rank of a (real-valued) tensor can be different over reals and over complex numbers.

Tensor	Over $\mathbb R$	Over $\mathbb C$
$oldsymbol{X}_1 = egin{pmatrix} 1 & 0 \ 0 & 1 \end{pmatrix}$	$\boldsymbol{A} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{pmatrix}$	$\mathbf{A} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -i & i \end{pmatrix}$
$oldsymbol{X}_2 = egin{pmatrix} 0 & 1 \ -1 & 0 \end{pmatrix}$	$oldsymbol{B} = egin{pmatrix} 1 & 0 & 1 \ 0 & 1 & 1 \end{pmatrix}$	$oldsymbol{B} = rac{1}{\sqrt{2}} egin{pmatrix} 1 & 1 \ i & -i \end{pmatrix}$
	$oldsymbol{\mathcal{C}}=\left(egin{array}{cccc} 1&1&0\-1&1&1 \end{array} ight)$	$\boldsymbol{C} = \begin{pmatrix} 1 & 1 \\ -i & i \end{pmatrix}$

- Determining the rank is NP-hard
 - No "straightforward" algorithm known
- Over reals, the rank can be larger than the largest dimension
 - rank $(\mathcal{X}) \leq \min \{ IJ, IK, JK \}$ ofr $I \times J \times K$ tensor
 - Better bounds known only for few special cases

Tensor rank oddities (2)

- For matrices, the Eckart-Young theorem tells us that
 - We can determine a best low-rank approximation
 - Factors of best rank-k approximation are part of best higher-rank approximations
- There are tensors of rank R that can be approximated arbitrarily well with tensors of rank $R^\prime < R$
 - Such tensors are called degenerate
 - There are no *best* low-rank approximations for degenerate tensors
 - ► The smallest such *R*′ is called the **border rank** of the tensor
- There are tensors for which the factors of the best rank-one approximation are not part of the best rank-2 approximation
 - Cannot find factors sequentially (but must find them simultaneously)

Tensor rank oddities (3)

- The rank decomposition is often essentially unique
 - Finally some good news!
 - Not true for matrices (SVD is essentially unique only due to its additional constraints)
- Essentially unique = only scaling and permutations allowed
- Holds under "mild conditions"

- 1. What Is a Tensor?
- 2. Tensor Basics

3. The CP Decomposition

- 3.1 Definition
- 3.2 Tensor Rank

3.3 Computing the CP Decomposition

- 3.4 Variants
- 4. The Tucker Decomposition
- 5. Wrap-Up

How to compute the CP decomposition?

- Here we discuss the alternating least squares (ALS) approach
 - Popular method
 - Simple to understand and implement
 - Not guaranteed to converge to minimum or stationary point
 - Objective merely decreases over iterations (until converged)
 - Can take many iterations until convergence
 - Final solution heavily dependent on starting point
- Many alternative methods have been proposed (and are being proposed)

ALS for CP

- Our goal is to optimize $\min_{\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}} \| \boldsymbol{\mathcal{X}} [\![\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}]\!] \|_{F}$
- Method is as before
 - 1. Fix \boldsymbol{B} and \boldsymbol{C} and solve for \boldsymbol{A}
 - 2. Solve for **B**, then solve for **C**
 - 3. Repeat until convergence
- To implement this method, we rewrite the CP decomposition

$$\begin{aligned} \boldsymbol{X}_{(1)} &\approx \boldsymbol{A}(\boldsymbol{C} \odot \boldsymbol{B})^{\mathsf{T}} \\ \boldsymbol{X}_{(2)} &\approx \boldsymbol{B}(\boldsymbol{C} \odot \boldsymbol{A})^{\mathsf{T}} \\ \boldsymbol{X}_{(3)} &\approx \boldsymbol{C}(\boldsymbol{B} \odot \boldsymbol{A})^{\mathsf{T}} \end{aligned}$$

• Now we can "read off" the solution; e.g.

$$\begin{split} \min_{\boldsymbol{A}} \|\boldsymbol{X}_{(1)} - \boldsymbol{A}(\boldsymbol{C} \odot \boldsymbol{B})^{T}\|_{F} \\ \boldsymbol{A} &= \boldsymbol{X}_{(1)} (\underbrace{(\boldsymbol{C} \odot \boldsymbol{B})^{T}}_{R \times JK})^{\dagger} \\ \boldsymbol{A} &= \boldsymbol{X}_{(1)} (\boldsymbol{C} \odot \boldsymbol{B}) (\underbrace{\boldsymbol{C}^{T} \boldsymbol{C} * \boldsymbol{B}^{T} \boldsymbol{B}}_{R \times R})^{\dagger} \end{split}$$
Outline

- 1. What Is a Tensor?
- 2. Tensor Basics

3. The CP Decomposition

- 3.1 Definition
- 3.2 Tensor Rank
- 3.3 Computing the CP Decomposition
- 3.4 Variants

4. The Tucker Decomposition

5. Wrap-Up

Variants of the CP decomposition

Many variants of the CP decomposition have been proposed; e.g.,

• Explicit modeling of component weights; e.g.,

$$\mathcal{X} \approx \llbracket \boldsymbol{\lambda}; \boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C} \rrbracket = \sum_{r=1}^{R} \lambda_r \, \boldsymbol{a}_r \circ \boldsymbol{b}_r \circ \boldsymbol{c}_r,$$

where $\boldsymbol{\lambda} \in \mathbb{R}^R$ and the columns of \boldsymbol{A} , \boldsymbol{B} , and \boldsymbol{C} have unit norm

- Non-negative CP decomposition (NNCP)
- Tensor completion (to handle missing entries)
- PARAFAC2 decomposition: jointly factor K matrices such that *X_k* ≈ *U_kS_kV^T*, where *V^T* is shared for all matrices (*shared* subspace learning)
- The INDSCAL decomposition

. . .



The INDSCAL decomposition (definition)

- Recall: CP decomposition decomposes a 3-way tensor X using three factor matrices A, B, and C
 - $\boldsymbol{\mathcal{X}} \approx \llbracket \boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C} \rrbracket$
 - Or element-wise: $x_{ijk} \approx \sum_{r=1}^{R} a_{ir} b_{jr} c_{kr}$
- The INDSCAL decomposition decomposes a 3-way tensor \mathcal{X} into two factor matrices \boldsymbol{A} and \boldsymbol{C}
 - Matrix A is used for first two modes
 - $\blacktriangleright \ \mathcal{X} \approx \llbracket \mathbf{A}, \mathbf{A}, \mathbf{C} \rrbracket$
 - Or element-wise: $x_{ijk} \approx \sum_{r=1}^{R} a_{ir} a_{jr} c_{kr}$
 - ➤ X must be an I × I × K tensor, i.e., the first two modes must have the same dimensionality
- Observe that reconstructed frontal slices are symmetric (i.e., [[A, A, C]]_{ijk} = [[A, A, C]]_{jik})
 - ► If we know that the frontal slices of X are symmetric, INDSCAL won't destroy this structure

The INDSCAL decomposition (why?)

- INDSCAL stands for "Individual Differences in Scaling"
 - Useful to analyze multiple symmetric matrices referring to the same objects
 - Distance, similarity, covariance, Laplacian matrices, ...
 - Recently also applied to knowledge bases (under the name DistMult)
- Assume K subjects each ranked the similarity of I objects
 - This gives us K similarity matrices, each $I \times I$
 - INDSCAL assumes that similarity decisions for each subject depends on a set of latent factors (A), and that those factors are weighted differently by different subjects (C)
- Best method for computing INDSCAL still open; common hack:
 - ► Compute normal CP and hope that **A** and **B** converge
 - ▶ Then force **A** and **B** equal and update **C** one more time

Psychological data: "How similar are these countries?"



Psychological data: "How similar are these countries?"



Psychological data: "How similar are these countries?"



FIGURE 13

Outline

- 1. What Is a Tensor?
- 2. Tensor Basics
- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 4.1 *n*-Mode Product
- 4.2 The Tucker Decomposition
- 4.3 The Tucker2, RESCAL, and DEDICOM Decompositions

5. Wrap-Up

Outline

- 1. What Is a Tensor?
- 2. Tensor Basics
- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 4.1 *n*-Mode Product
- 4.2 The Tucker Decomposition
- 4.3 The Tucker2, RESCAL, and DEDICOM Decompositions
- 5. Wrap-Up

Tensor *n*-mode product

- Let X be an N-way tensor of size I₁ × I₂ × ··· × I_N, U a matrix of size J × I_n, and v vector of size I_n
- The *n*-mode vector product $\mathcal{X} \times_n \mathbf{v}$ of \mathcal{X} with \mathbf{v} computes the inner product of each mode-*n* fiber with \mathbf{v}
 - ▶ Result is tensor of order N-1; size $I_1 \times \cdots \times I_{n-1} \times I_{n+1} \times \cdots \times I_N$

• Element-wise:
$$[\boldsymbol{\mathcal{X}} \times_{n} \boldsymbol{v}]_{i_{1}\cdots i_{n-1}i_{n+1}\cdots i_{N}} = \sum_{i_{n}=1}^{I_{n}} x_{i_{1}i_{2}\cdots i_{N}} v_{i_{n}}$$

- The *n*-mode (matrix) product $\mathcal{X} \times_n U$ of \mathcal{X} with U multiplies each mode-*n* fiber with U (from the left)
 - ▶ Result is tensor of order *N*; size $I_1 \times \cdots \times I_{n-1} \times J \times I_{n+1} \times \cdots \times I_N$
 - ▶ In terms of unfolding: $\mathcal{Y} = \mathcal{X} \times_n \mathcal{U} \iff \mathcal{Y}_{(n)} = \mathcal{UX}_{(n)}$
 - Or element-wise: $[\boldsymbol{\mathcal{X}} \times_n \boldsymbol{U}]_{i_1 \cdots i_{n-1} j i_{n+1} \cdots i_N} = \sum_{i_n=1}^{I_n} x_{i_1 i_2 \cdots i_N} u_{j i_n}$

Example



•
$$\mathcal{X} \times_3 \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 6 & 10 \\ 8 & 12 \end{pmatrix}$$
 $\mathcal{X} \times_3 \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 5 \begin{pmatrix} 6 \\ 0 \\ 8 \end{pmatrix} \begin{pmatrix} 10 \\ 2 \\ 4 \end{pmatrix} \begin{pmatrix} 6 \\ 12 \end{pmatrix}$
• $\mathcal{X} \times_1 \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 3 & 11 \\ 7 & 15 \end{pmatrix}$ $\mathcal{X} \times_1 \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} \begin{pmatrix} 1 & 5 \\ 6 \\ 3 \\ 2 \end{pmatrix} \begin{pmatrix} 8 \\ 12 \\ 12 \end{pmatrix} \begin{pmatrix} 8 \\ 12 \\ 12 \end{pmatrix}$

Rewriting the CP decomposition

- Useful identities
 - $\blacktriangleright \mathcal{X} \times_n \boldsymbol{I}_{I_n} = \mathcal{X}$
 - $\blacktriangleright \ \mathcal{X} \times_m \mathbf{A} \times_n \mathbf{B} = \mathcal{X} \times_n \mathbf{B} \times_m \mathbf{A} \text{ for } m \neq n$
 - $\blacktriangleright \quad \mathcal{X} \times_n \mathbf{A} \times_n \mathbf{B} = \mathcal{X} \times_n (\mathbf{B}\mathbf{A})$
- We can rewrite a size-R CP decomposition as

$$\llbracket \boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C} \rrbracket = \boldsymbol{\mathcal{I}}_R \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C}$$

- ➤ *I_R* is the *R* × *R* × *R* identity tensor (superdiagonal, 1s on the superdiagonal)
- What happens if we use another tensor instead of \mathcal{I}_R ?

Outline

- 1. What Is a Tensor?
- 2. Tensor Basics
- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 4.1 *n*-Mode Product
- 4.2 The Tucker Decomposition
- 4.3 The Tucker2, RESCAL, and DEDICOM Decompositions
- 5. Wrap-Up

The Tucker decomposition

• The Tucker3 decomposition decomposes a 3-way tensor as

 $\boldsymbol{\mathcal{X}} pprox \boldsymbol{\mathcal{G}} imes_1 \boldsymbol{A} imes_2 \boldsymbol{B} imes_3 \boldsymbol{C}$

- \mathcal{X} is $I \times J \times K$
- \mathcal{G} is a $P \times Q \times R$ tensor, called the core tensor
- **A**, **B**, **C** are $I \times P$, $J \times Q$, and $K \times R$, resp.
- We write $[\![\boldsymbol{\mathcal{G}} ; \boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}]\!]$ for short and obtain



Tucker, 1966

Some remarks

- Core tensor can be thought of a compressed version of ${\mathcal X}$
- Many degrees of freedom
 - Parameters P, Q, and R
 - ▶ Properties of **A**, **B**, and **C** (often taken to be orthogonal)
 - CP decomposition is special case (\mathcal{G} superdiagonal, P = Q = R)
- Not unique

$$\llbracket \mathcal{G}; \mathcal{A}, \mathcal{B}, \mathcal{C}
rbracket = \llbracket \mathcal{G} imes_1 \mathcal{U} imes_2 \mathcal{V} imes_3 \mathcal{W}; \mathcal{A} \mathcal{U}^{-1}, \mathcal{B} \mathcal{V}^{-1}, \mathcal{C} \mathcal{W}^{-1}
rbracket$$

- ▶ Used to "simplify" core (e.g., many zeros, all-orthogonal, ...)
- Can be generalized to N-way tensors
- Can be computed using ALS with matricized forms

$$\begin{split} \boldsymbol{X}_{(1)} &\approx \boldsymbol{A}\boldsymbol{G}_{(1)}(\boldsymbol{C}\otimes\boldsymbol{B})^{T} \\ \boldsymbol{X}_{(2)} &\approx \boldsymbol{B}\boldsymbol{G}_{(2)}(\boldsymbol{C}\otimes\boldsymbol{A})^{T} \\ \boldsymbol{X}_{(3)} &\approx \boldsymbol{C}\boldsymbol{G}_{(3)}(\boldsymbol{B}\otimes\boldsymbol{B})^{T} \end{split}$$

► If factor matrices are column-orthogonal, update $\mathcal{G} \leftarrow \mathcal{X} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T$

Higher-order SVD

- Tucker's "Method I" for computing the decomposition
 - Idea is to first find components that capture variation in each mode (independently of the other modes)
 - Set $\mathbf{A} = \mathbf{U}_1$ = leading P left singular vectors of $\mathbf{X}_{(1)}$
 - Set $B = U_2$ = leading Q left singular vectors of $X_{(2)}$
 - Set $C = U_3$ = leading R left singular vectors of $X_{(3)}$
 - Set core tensor $\boldsymbol{\mathcal{G}} = \boldsymbol{\mathcal{X}} \times_1 \boldsymbol{U}_1^T \times_2 \boldsymbol{U}_2^T \times_3 \boldsymbol{U}_3^T$
- Today known as <u>higher-order SVD</u> (HOSVD)
- If *P*, *Q*, and *R* are smaller than rank of their corresponding matricized forms, we obtain a truncated HOSVD



Discussion

- SVD of **X**_(n) gives *n*-mode singular vectors (**U**_n) and *n*-mode singular values
- Indeed, many useful properties of the SVD transfer over
 - Factors and slices of core tensor are ordered (Frobenius norm of slice = corresponding *n*-mode singular value)
 - Essentially unique
 - Core tensor is all-orthogonal (inner product between pairs of slices = 0)
- Core tensor cannot be made superdiagonal in general
- Truncated HOSVD is not optimal w.r.t. reconstruction error
 - But it holds

 $\| \boldsymbol{\mathcal{X}} - \hat{\boldsymbol{\mathcal{X}}} \|^2 \leq$ sum of squares of truncated singular values

Good starting point for ALS (or other methods)

- Data: 7943-pixel B&W photographs of 28 people in 5 poses under 3 illumination setups performing 3 different expressions
 - $28 \times 5 \times 3 \times 3 \times 7943$ tensor (10M elements)



- Data: 7943-pixel B&W photographs of 28 people in 5 poses under 3 illumination setups performing 3 different expressions
 - $28 \times 5 \times 3 \times 3 \times 7943$ tensor (10M elements)



 $\boldsymbol{U}_{\text{pix}} (= \boldsymbol{U}_5)$ contains the eigenfaces (SVD of pixels-by-picture matrix)

- Data: 7943-pixel B&W photographs of 28 people in 5 poses under 3 illumination setups performing 3 different expressions
 - $28 \times 5 \times 3 \times 3 \times 7943$ tensor (10M elements)



Some visualizations of $\mathcal{G} \times_5 \mathcal{U}_{pix}$ (decreasing "importance" from top to bottom / left to right)

- Data: 7943-pixel B&W photographs of 28 people in 5 poses under 3 illumination setups performing 3 different expressions
 - $28 \times 5 \times 3 \times 3 \times 7943$ tensor (10M elements)



Some visualizations of $\mathcal{G} \times_2 \mathcal{U}_{\text{pos}} \times_3 \mathcal{U}_{\text{ill}} \times_4 \mathcal{U}_{\text{exp}} \times_5 \mathcal{U}_{\text{pix}}$. The rows show the componets w.r.t. people. The columns refer to different viewpoints, illuminations, and expressions.

Outline

- 1. What Is a Tensor?
- 2. Tensor Basics
- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 4.1 *n*-Mode Product
- 4.2 The Tucker Decomposition
- 4.3 The Tucker2, RESCAL, and DEDICOM Decompositions
- 5. Wrap-Up

The Tucker2 decomposition

- Recall: Tucker3 decomposition decomposes 3-way tensor \mathcal{X} into three factor matrices A, B, and C, and a core tensor \mathcal{G}
- The Tucker2 decomposition decomposes a 3-way tensor into core and *two* factor matrices
 - Equivalently: third factor matrix is taken as the identity matrix
 - If data tensor is $I \times J \times K$, then core is $P \times Q \times K$



Why Tucker2?

- Use Tucker2 if you don't want to "compress" one of the modes
 - E.g., too small dimension (e.g., $500 \times 300 \times 3$)
 - E.g., want to handle this mode separately
 - ► For example, if third mode is time, we might first do Tucker2 and then analyze the slices { G_k } over time
- Tucker2 is slightly simpler than Tucker3
 - We have $\boldsymbol{X}_k = \boldsymbol{A}\boldsymbol{G}_k\boldsymbol{B}^T$ for each frontal slice k
 - ► When using ALS, we can update each frontal slice separately during update of *G*
- Forms basis for RESCAL and DEDICOM (up next)

The RESCAL decomposition

- The RESCAL decomposition combines Tucker2 and INDSCAL
- Given an I-by-I-by-K tensor X and a rank R, find an I × R factor matrix A and an R × R × K core tensor R such that

 $\boldsymbol{\mathcal{X}} pprox \boldsymbol{\mathcal{R}} imes_1 \boldsymbol{A} imes_2 \boldsymbol{A}$

• I.e., minimize

$$\sum_{k=1}^{K} \|\boldsymbol{X}_k - \boldsymbol{A}\boldsymbol{R}_k \boldsymbol{A}^T\|_F^2.$$

• In practice, add regularization to avoid overfitting

RESCAL and subject-object-predicate data (1)

- Proposed with subject-object-predicate data in mind
 - Designed for settings with few relations
 - ► E.g., <u>YAGO knowledge base</u>: < 100 relations, millions of entities
- Used for
 - Link prediction (by looking at $\hat{\mathcal{X}}$)
 - Assessment of entity similarity (by looking at A)
 - Assessment of relation similarity (by looking at \mathcal{R})
 - Currently a hot research area; many methods can be seen as constrained variants of RESCAL



RESCAL and subject-object-predicate data (2)



Why RESCAL?

- RESCAL computes an *embedding* for every entity (*A_i*) and for every relation (*R_k*)
- With Tucker2, entity embeddings were different depending on whether the entity appears as subject or object
 - Having one embedding per entity ensures "information flow"
 - Facts with entity as subject help predicting facts with entity as obj.
- With INDSCAL, reconstructed relations were symmetric
 - That's not the case with RESCAL because *R_k* tells us how to mix embeddings to reconstruct the *k*-th relation
 - For symmetric relations, R_k should be symmetric
 - For asymmetric relations, \boldsymbol{R}_k should be asymmetric
- Computation using ALS-style methods
 - ► To update **A**, we use mode-1 matricization of RESCAL (as before)
 - $\boldsymbol{X}_{(1)} \approx \boldsymbol{A} \boldsymbol{R}_{(1)} (\boldsymbol{I} \otimes \boldsymbol{A})^T$
 - ► Hard because **A** appears on left- and right-hand side
 - Original paper uses hack (treat right-hand side as fixed + cleverness)

The DEDICOM decomposition

- The DEDICOM decomposition is a precursor to RESCAL
 - Mixing matrix is the same for all relations
 - But entity factors are weighed differently
 - $X_k = AD_k RD_k A^T$, where D_k is diagonal and contains the weights of each factor for the k-th relation
- Many of the ideas in RESCAL are based on DEDICOM
 - But DEDICOM is more involved to scale to large data
 - And DEDICOM generally needs larger R due to shared relation matrix



Outline

- 1. What Is a Tensor?
- 2. Tensor Basics
- 3. The CP Decomposition
- 4. The Tucker Decomposition
- 5. Wrap-Up

Lessons learned

- Tensors generalize matrices
 - Many matrix concepts generalize well
 - But some don't, and some behave very differently
- Compared to matrix decomposition methods, tensor algorithms are in their youth
- Tensor decompositions are used in many different fields of science
 - Sometimes the wheel gets re-invented multiple times
 - Traditionally tensor problems were dense
 → Fewer algorithms for decomposing sparse tensors
- There are many tensor decompositions related to CP and Tucker
 - Need care to select the one that's best suited for the task at hand
 - Computational complexity can be an issue

Suggested reading

- Skillicorn, Ch. 9
- Kolda & Bader, <u>Tensor Decompositions and Applications</u>, SIAM Rew. 51(3), 2009
 - A great survey on tensor decompositions, includes many variants and applications
- Acar & Yener, <u>Unsupervised Multiway Data Analysis</u>: A <u>Literature Survey</u>, IEEE Trans. Knowl. Data Eng. 21(1), 2009
 - Another survey, shorter and more focused on applications
- All the papers linked at the bottom parts of the slides