Data Mining and Matrices 01 – Introduction

Rainer Gemulla, Pauli Miettinen

April 18, 2013

Outline

What is data mining?

2 What is a matrix?

3 Why data mining and matrices?



"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 data
 - ▶ 1M customer transactions are performed at Walmart per hour
 - 25M Netflix customers view and rate hundreds of thousands of movies
 - 40B photos have been uploaded to Facebook
 - 200M active Twitter users write 400M tweets per day
 - ► 4.6B mobile-phone subscriptions worldwide
- Government, health care, news, stocks, books, web search, ...

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



Prediction





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 pipeline

b ai				V.			40		÷	1			ε			V			01	t8						t				0	+
				64			8		t	12	а		-tt						23	7.0	7)	X				14	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.	ċ,		Ě.	n
	U.	4	n											1			Ŧ					2					2	÷.	ñ	H	Ď
+	÷		п							<u>P</u>												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		F	0	ę.	J,	X	X	A	Q I		L	4			6	٤.	٩L	륃	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	1	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 J.J. Sylvester in 1850



System of linear equations

• Systems of linear equations can be written as matrices

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

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

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

• 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



X	У
/ -3.84	-2.21
-3.33	-2.19
-2.55	-1.47
-2.46	-1.25
-1.49	-0.76
-1.67	-0.39
-1.3	-0.59
:	÷
1.59	0.78
1.53	1.02
1.45	1.26
1.86	1.18
2.04	0.96
2.42	1.24
2.32	2.03
2.9	1.35 /

Linear maps

 \bullet 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 f_1 written as a matrix

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

 $\bullet\,$ 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

	Bread	Butter	Beer			Data	Ma	trix	Mining	5
Anna	/ 1	1	0		Book 1 /	5	()	3	
Bob	1	1	1		Book 2	0	()	7	
Charlie	0	1	1)	Book 3 \	4	6	5	5	
	Custome	er transac	tions		Ľ	Docur	nent-t	erm i	matrix	
	Avatar	The Ma	trix	Up			Jan	Jun	Sep	
Alice	(4		2	Saarbrück	en /	1	11	10 `	\
Bob	3	2			Helsinki		6.5	10.9	8.7	
Charlie	5			3 /	Cape Tow	n 🔪	15.7	7.8	8.7)
Incomplete rating matrix Cities and monthly temperatures										
N A	1.00		C	1		A 44 14 14				

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

What is a matrix?

- A means to describe computation
 - Rotation
 - Rescaling
 - Permutation > Linear operators
 - Projection
 - • •
- A means to describe data

Rows	Columns	Entries					Attr	
<i>Objects</i> Equations Data points	<i>Attributes</i> Variables Axes	<i>Values</i> Coefficients Coordinates		$ \begin{pmatrix} A_{11} \\ A_{21} \\ \vdots \end{pmatrix} $	$A_{12} \\ A_{22} \\ \vdots$	 	$egin{array}{c} A_{1j} \ A_{2j} \ dots \end{array}$)
Vertices	Vertices	Edges	Object i	A _{i1}	A _{i2}	•••	A _{ij}	••••
:	•	•			÷	·	:	·)

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

bute j

Outline

What is data mining?

2 What is a matrix?





17 / 27

Key tool: Matrix decompositions

A matrix decomposition of a data matrix ${\bf D}$ is given by three matrices ${\bf L},$ ${\bf M},$ ${\bf R}$ such that

 $\mathbf{D}=\mathbf{LMR},$

where

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

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

$$\mathbf{D}_{ij} = \sum_{k,k'} \mathbf{L}_{ik} \mathbf{M}_{kk'} \mathbf{R}_{k'j}$$



Example: Singular value decomposition



 $L_{50\times 2}$

 $M_{2\times 2}$









Example: Non-negative matrix factorization



Lee and Seung. Learning the parts of objects by non-negative matrix factorization. Nature, 1999.

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	R
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	

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.

(L)

Other matrix decompositions

- Singular value decomposition (SVD)
- k-means
- Non-negative matrix factorization (NMF)
- Semi-discrete decomposition (SDD)
- Boolean matrix decomposition (BMF)
- Independent component analysis (ICA)
- Matrix completion
- Probabilistic matrix factorization

• . . .

What can we do with matrix decompositions?

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

Matrix decompositions can make data mining algorithms more effective. They may also provide insight into the data by themselves.

Factor interpretion of matrix decompositions

Assume that \mathbf{M} is diagonal. Consider object *i*.

- Row of **R** = part (or piece), called *latent factor* ("latent object")
- \square Entry of \mathbf{M} = weight of corresponding part
- Row of **MR** = weighted part
- Row of L = "view" of corresponding row of D in terms of the weighted parts (r pieces of information)
- r forces "compactness" (often r < n)

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

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



Other interpretions

- Geometric interpretation
 - ► Transformation of *n*-dimensional space in *r*-dimensional space
 - Row of R = axis
 - Row of C = coordinates
- Component interpretation
 - **D** is viewed as consisting of *r* layers (of same shape as **D**)
 - k-th layer described by $\mathbf{L}_{*k}\mathbf{M}_{kk}\mathbf{R}_{k*}$
 - $\mathbf{P} = \sum_{k} \mathbf{L}_{*k} \mathbf{M}_{kk} \mathbf{R}_{k*}$
- 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

What is data mining?

2 What is a matrix?

3 Why data mining and matrices?



Lessons learned

- Data mining = from data to knowledge \rightarrow Prediction, clustering, outlier detection, local patterns
- Many different data types can be represented with a matrix \rightarrow Linear equations, data points, maps, graphs, relational data, ...
- Common interpretation: rows = objects, columns = attributes
- Matrix decompositions reveal structure in the data \rightarrow $\mathbf{D}=\mathbf{LMR}$
- Many different decompositions with different applications exist \rightarrow SVD, *k*-means, NMF, SDD, BMF, ICA, completion, ...
- Factor interpretation: objects described by "latent attributes"

Suggested reading

David Skillicorn

Understanding Complex Datasets: Data Mining with Matrix Decompositions (Chapters 1–2) Chapman and Hall, 2007