### Epilogue Wrap-Up



## **Recap of the Course**

## DATA MINING & MATRICES



SVD, PCA, CUR, NMF, ICA, ...
 6 credits

Theoretical & hands-on assignments

First lecture: tomorrow at high noon in room 029, building E1 5 (MPI-SWS)

http://bit.ly/dmm15

3カ 干え メアゾび7 1時 井 14世 8 干9 ム7アメ ケの0び ロムカ26ア 7井 2 カび中 び 33世 0中 び父ケ8 の中613びJ干ア55 18 アルえ52 干 5号中 の71 キ20中 干カケえ カケび77 J び3 9メケカアえ ゾ 号9世号 6カ9 シグケ 中 6アゾメス干え3中94 干 中ゾムえ53 5 2中ム ののか たの カ542中中び86ケ1カ3世 火停ケ 3ケ ゾ397ロ ゾ干井4え又5 05世の88.97 98.77 5 37 37 7000 707 07世の88.97 98.77 5 7 98.77 7 707 07世メモクトリケル 4 5 7 × 16 7 7 00, 7 1 0 0 7 7 07世メモクトリケル 4 5 7 × 16 7 7 00, 7 1 0 0 7 0 7 1 × 7 0 7 1 0 2 7 1 7 1 0 0 7 1 7 1 0 0 7 8 1 0 7 2 3 7 8 2 8 7 7 1 0 0 2 7 1 8 1 4 1 7 6 0 5 8 1 0 7 2 7 5 9 7 9 7 1 0 2 7 1 7 1 8 2 7 1 7 1 8 2 7 1 7 1 7 1 7 2 7 7 8 1 8 4 1 7 6 0 5 8 1 0 7 7 7 7 1 0 2 7 7 1 7 0 2 7 1 9 0 5 7 7 8 1 8 4 1 7 6 0 5 8 1 7 7 7 7 1 0 2 7 7 1 7 0 3 7 8 4 6 5 4 8 4 1 7 6 0 5 7 7 7 7 7 7 1 0 3 7 8 7 7 7 8 7 8 7 8 7 8 7 8 7 8 7 8 7	1.1			V			3	3		8			E	P		Υ.		0	ŧ:					1	Ť.			0	
ムアアメ ケの00 ロムカ267 7州 2 かび中 び 33世 0中 び父ケ8 の中6)3び)モア55 18 ア州た52 〒 5号中 の71 〒20中 モカケえ カケび77) び3 9火ケカアえ ゾ 号9件号 6カ9 メケ 中 6アゾス火モえ3中94 〒 中ゾムた53 5 2中ム ののカ えの カ542中中び86ケ1カ3世 火尽ケ 3ケ ゾ3970 ゾモ世4え又5 05世の8え9〒982 79 92 7 9〒ダス78 008 7 11 11 12 ロカ2ゾ 9〒ダス78 008 7 11 12 ロカ2ゾ 00中ア84〒 ロ0 ゾ45 えいロ 号37 ス世6日から082 230 ロ797世 8日123〒828 25 51 10 51 20151 11 18えの号7 世 8日123〒828 25 51 10 51 20151 10 18えの号7 世 81 75えり州87 サケ 91 70 002 71 902 7 7 号口 8 ム 8 ロの州2ケ7 5 ケシェノシニロロ 557 7 世6日から08 937 4火ム 7 87 75 71 70 10 00 20151 10 18えの号7 世 70 7 世方30 902 7 10 10 00 00 7 世方30 902 7 7 号口 8 ム 70 7 世52えり州87 サケ 91 7 00 3 ム号大び5 4 5 7 5 0 5 0 5 0 5 0 5 0 5 0 5 0 5 0 5 0				6	Ŧ		ŧ	3		$\pm 2$	Ц		ŧL,		Γ			7	$\rho_I$	3	Х				5			t,	9
び 2 0 4 5 0 4 6 0 3 び 1 〒 7 5 5 1 8 7 兆 5 2 〒 5 5 4 0 7 1 〒 2 0 4 〒 かたえ かかび 7 7 3 び 3 9 2 ケカアえ ゾ 5 9 日号 6 か 9 2 か かち 4 2 4 4 4 5 3 5 2 4 4 0 0 か 2 0 5 日 0 8 5 9 7 9 2 5 7 4 5 3 5 7 3 4 7 7 7 5 4 4 4 2 × 5 0 5 日 0 8 5 9 7 9 2 5 7 7 5 7 7 9 2 7 7 7 5 7 4 4 4 2 × 5 0 7 7 2 7 8 7 9 2 5 7 7 7 5 7 7 7 7 7 7 7 7 7 7 7 7 7 7		0		-	3	3	3	3			t		2	-tt	7			7	9.2	÷.	$\mathbf{A}^{\dagger}$	П	ð	00	7 ()	ł.	X.	77	A
<ul> <li>〒20中 干カケミ カケび77) び3 9×ケカアミ ダ 号9件号 6カ9</li> <li>ジケ 中 6アジ××テミシキ94</li> <li>モ 中ジムミ53 5 2中ム ののカシスの カ542中中が86ケ1カ3件 ×尽ケ 3ケ ジ39アロ ジェド42×5</li> <li>0500829〒の829〒の92</li> <li>マテジ×78 008</li> <li>0779×109</li> <li>マテジ×78 008</li> <li>マテジ×78 008</li> <li>マテジ×78</li> <li>ロナジ×78</li> <li>ロナジ×100</li> <li>マテジ×100</li> <li>マテジ×100</li> <li>マテジ×100</li> <li>マテジ×100</li> <li>マテジ×100</li> <li>ロナジ×100</li> <li>ロナジ&lt;100</li> <li>ロナジ×100</li> <li>ロナジ&lt;100</li> <li>ロナジ&lt;100</li> <li>ロナジ×100</li> <li>ロナジ×100</li> <li>ロナジ×100</li> <li>ロナジ×100</li> <li>ロナジ&lt;100</li> <li>ロナジ&lt;100</li> <li>ロナジ×100</li> <li>ロナジ×100</li> <li>ロナジ&lt;100</li> <li></li></ul>	ſ	50			늰	đ		1		S d	-	ŧť	7	8	1			ĉ	d T	Ŧ	U	08	10	91	10	6	8-	$\chi$ $\phi$	Ð
ソケ 中 6アゾ火火手え3中94       干 中ソムえ533 5       2中ム ののかえんの かち42中中が86を11か3日 火陸ケ 3ケ ゾ39アロ ジナ日42×5         05008297799279927792       フレ ジラクク 1000020200020200000000000000000000000	9	ŧð			ġ.	5	1	2		τt	÷	Ŷ	e.	8	ð			ī	7	R.	÷	÷	÷	4.	÷			210	
<ul> <li>えの か542中中び8651393日 火陸ケ 35 ゾ3970 ゾ子日42×5</li> <li>05日の8297932</li> <li>07708297932</li> <li>0770847 日0 ジ45 え10 等37 ×167702ジ230 07974</li> <li>8012378286 を2 座4 ホッド3 ションクジェクシンクシングシングシングシングシングシングシングシングシングシングシングシングシングシ</li></ul>	+	00		λđ	¢				<b>c</b>		Å	Ŵ	1	Ŧ				À	ġ i	15	e	÷,	1	v	-	1		÷	Ą,
05日の8え9〒32593925939250080000000000000000000000000000000000	2	y e	A.	+	ų,		ī	í.	d,			÷	è.	ų,		¥.	+	c i		-	a	e i	$\mathbf{r}$	à	e i	E.	+	Â	e.
9子ダ火78 008 子1 07 07 07 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	di la	25	Ļ.							allen.		-1.550						0		1	Ľ,				6	10	24		F
0中ア84〒 10 945 また1 や37 文世6〒930を9230 1797世 80123〒828 そ2 84 キャド2 55年943808937 4次入 7 0〒世火于や19カ・マバ 4 10 10 10 10 10 10 10 10 10 10 10 10 10	T	20										-		1		F =	- 1				6								ň
8日)23〒828 22 月10 945 21 日 957 21 日 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 3 5 7 5 5 5 4 5 7 5 5 5 5 5 5 5 5 5 5 5 5	-	27	H									2							2				2				÷.	5	Ľ,
8日)23年828 0 〒世火干ケ1 ジカ・ジャ・シャ・シャ・シャ・シャ・シャ・シャ・シャ・シャ・ション 1 8 2 0 5 8 日の州2ケ7 3 2 シャットシューロ 1 ケ07日のム2 ブ 9 のえ7 7 号日 8 ム 8 7 干5え 1 井87 井ケ6 ブ 1 ケ07日のム2 ブ 9 のえ7 7 号日 8 ム ケサの日41 カ9 6 干ケ井8 え 号子火1 1 ムム2 びび火び 6 1 えびび 1 び モの干カ日36 5 ケ6 アム中 干7170 1 3 ム号カムびケ井 号 7 号び ラ 7 0 2 びええ 世 ジ 2 4 9 2 干85 1 火 カヤケ4 ロ 0 ム 1 ケ び ロ 世 干 干 ム 9 干 51 中5日337 火 3 ヤケ4 ロ 0 ム 1 ケ 7 0 日 世 干 干 ム メロ 井号 0 え え 7 6 4 5 0 6 井 ポ 1 ケ 7 0 2 4 2 7 世 干 1 7 0 3 4 え 中 カム 7 カ モンケセケ キャ 0 パ ヘケ パ 8 3 号 1 6 5 4 4 5 0 0 カ 1 0 × 1 0 × 1 0 × 1 0 × 1 0 × 1 0 × 1 0 × 1 0 × 1 0 × 0 ×		Λŝ			Y		X	ŝ U	2	<u>s</u> .,	0			2	2	1			2	ę		V.	. Ч	Ц.			ġ.	1	н
0+#火キケ1火カ(ソ) (1000	ĩ	À	X		Σ				b.			2			1		×				r Á			Ŋ.			$\mathbf{x}_{i}$	ιI	Ċ,
8 ロのサ2ケ73スケスケンション100000000000000000000000000000000000	¢	Q	9							~				4				t,		K	X	K	R.	V.	7	t+	$\chi_1$		Û
2 )」等等1于20等升 の)1 ク07ロのム2ア」 9のえ7 7 等口 8 4 87于5え」升87升ケ6 ア」 ケ0ア7世カ3の グ申46 え 甘8廿口 ケ ケ州の口41カ9 6于ケ州8 え等子×1」ムム2びびなび 6 」えびびしび すの于カロ36号5 ケ6アム中 干7170」3ム号カムびケ州 号 7号び号グ0 2びええ甘グ2492于85」× カ中ケ4口0ム」ケ× 中の ケ び口甘干干ム 9干 51 申5口337×8ダ 487世4び6え口び グ州 ケ0びカ」0×世 ×ロ카号0ええ7645の6州州 グケ 2ア甘子甘71 アの 34え中カム7カ	+	V	구 (	0.5	8							Π	A C		20	0 r		4			4		$\mathbf{V}$	4	Ϋ́,	ιo			8
87〒5えり外87分ケ6 アリ ケ0ア7世カ3の ゾ中46 え 世8世口 ケ ケ州の口41カ9 6〒ケ州8 え <b>号子×1</b> リムム2びび×び 6 リえびびりび 干の干カ口36号5 ケ6アム中 干7170 <b>)</b> 3ム号カムびケ州 号 7号び号グ0 2びええせグ2492〒85リメ カ中ケ4口0 <b>ム</b> リケメ 中の ケ び口世干干ム 9干 51 中5口337×8グ 487世4び6え口び グ州 ケ0びカリ0×世 ×ロ州号0ええ7645の6州州 グケ 27世子世71 アの 34え中カム7カ	А	8		미성		7	۲.	ŝ (e	6		7	Ω,	A Ç		7	0 7	e.	1	0	9	t	붙이	$\mathbf{Z}_{1}$	Ŧ	I ÷	12	Ų.	L.	3
ケ州の口41カ9 6干ケ州8 え号子火1ノムム2びび火び 6 リえびびりび 干の干カロ36号5 ケ6アム中 干7170J3ム号カムびケ州 号 7号び号ゾ0 2びえたせゾ2492干85リメ カ中ケ4口0ムリケメ 中の ケ び口甘干干ム 9干 51 中5口337×8ダ 487世4び6え口び ゾ州 ケ0びカリ0×世 ×ロ州号0ええ7645の6州州 ゾケ 2ア甘子甘71 アの 34え中カム7カ	ネ	Π	H	84		J.m.	91	Т	N	Q	34	Ċ,	+.	Y	0	d.		U	7	9	4	ŧΝ	18	t	( .	56	Ŧ	78	
〒の干カロ36号5 ケ6アム中 干717013ム号カムびケ州 号 7号び号グ0 2びええせゾ2492干851火 カ中ケ4ロ0ムリケ火 中の ケ びロ世干干ム 9干 51 中5ロ337火8ダ 487世4び6えロび ゾ州 ケ0びカリ0火世 火口州号0ええ7645の6州州 ゾケ 27世子世71 アの 34え中カム7カ サゾケサケキモロびょのケィズト 83 屋166キ4号 07ロロサア14号かん	£0	10	RU.	ŝ١		9	ίU	C U		$S_A$	A	ι.	13	ŧ		Jan.		8	it r	÷Ŧ	9	Ð	Ť	1	ŀΓ	10	Ht.	4	
20% えたせび 2492〒85 JX カキケ4口0ムJケX キの ケ び口世干干ム 9〒 51 キ5口337X8ダ 487世406え口び ダ州 ケ00カ J0X世 X口外号0ええ7645の6州州 ダケ 27世子世71 アの 34えキカム7カ サゾケサケキケキチロバムのケムパー 83 8166キ4号 07円ロサア14号カの	0	V.		2			it-	10	$\mathbf{A}$	÷Ę	A	8	U O	7	Г	Ŧ₹	1	ф,	47	63	Ŕ	đ	-	ð	Ēr	1t	Ŧ	λŦ	
9〒 51 中5ロ337×87 487世406とロび グ升 ケ0びカJ0×世 ×ロ外号0とと7645の6升升 グケ 27世子世71 アの 34と中カム7カ サゾケサケキモのダイのケイダー 83 祭166キ4祭 07ロロサア1 中陸キの	4	ŦŦ	H	ΠÈ		t.	01			χ <del>γ</del> .	1	Δ	Ô E	14	7	15	÷	X	ιċ	8	Ŧ	9 P		12	6-	t s		2.0	
X田井琴02.2764506升升 ダケ 27世子世71 アの 342中カム7カ サゾケサケナキのダイのケイダー 83 8166キ48 07円サア1車巻キの	+	¥0	Π	÷'n	0	t.	413	è		U F	E	a	O P	H	7	8 1			V.	Ŷ	7	c c	П	d		Г	5		ρ
a + a + b + b + b + b + b + b + b + b +	+	53	÷	1	N	5	6.3	ř		1		Ŧ	-	- 0		43				h	Ċ.	e n	a	7	ŧ	ŧ	년.		Y
	à	$\pm 9$		r ÷	4	nn	Ť	ì		91	+	à	a r	5		e e	i.			4	à	Ň	ň			4	4.		

### A womb?

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





### Matrices in data mining



	Data	Matrix	Mining	
Book 1	/ 5	0	3	
Book 2	0	0	7	
Book 3	4	6	5	$\mathcal{I}$
	Docum	ent-term	matrix	

	Jan	Jun	Sep	
Saarbrücken	(1)	11	10 \	
Helsinki	6.5	10.9	8.7	
Cape Town	$\setminus 15.7$	7.8	8.7	
			_	

Cities and monthly temperatures

# Matrix decompositions in data mining

- A common goal in data mining is to find regularities (or patterns) in the data
  - Often, to summarize the data
- A matrix decomposition presents the data as a sum of "simple" elements, i.e. patterns
  - but there's also other uses... stay tuned!

#### Intuition for Matrix Multiplication

Matrix **AB** is a sum of k matrices **a**<sub>l</sub>**b**<sub>l</sub><sup>T</sup>
 obtained by multiplying the *l*-th column of **A** with the *l*-th row of **B**



#### "The SVD is the Swiss Army knife of matrix decompositions"

– Diane O'Leary, 2006



## Why is SVD important?

- It gives us the dimensions of the fundamental subspaces
- It lets us compute various norms
- It tells about sensitivity of linear systems
- It gives us optimal solutions to least-squares linear systems
- It gives us the least-error rank-k decomposition
- Every matrix has one

## The z-scores

- The z-scores are attributes whose values are transformed by
  - centering them to 0 by removing the (column) mean from each value
  - normalizing the magnitudes by dividing every value with the (column) standard deviation

$$X' = \frac{X - \mu}{\sigma}$$

#### Issues with the rankselection methods

 Require computing the full SVD first or otherwise computationally heavy

Guttman–Kaiser

Guttman–Kaiser

scree entro

entropy-based rando

random flips

random flips

Require subjective evaluation



scree

entropy-based

## Example use of SVD

- Data: people's ratings on different wines
- Scatterplot of first two LSV
  - SVD doesn't know what the data is
- Conclusion: winelovers like red and white alike, others are more biased



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

#### SVD gives directions of largest variances

- The singular vectors give the directions of the largest variances
  - First singular vector points to the direction of the largest variance
  - Second to the second-largest
    - Spans a hyperplane with the first
- The projection distance to these hyperplanes is minimal over all hyperplanes (Eckart–Young)



#### Very general idea for solving SVD

- SVD is unique
  - If **U** and **V** are orthogonal s.t.  $U^T A V = \Sigma$ , then  $U \Sigma V^T$  is the SVD of **A**
- Idea: find orthogonal **U** and **V** s.t.  $U^T A V$  is as desired
  - Iterative process: find orthogonal  $U_1$ ,  $U_2$ , ... and set  $U = U_1U_2U_3$ ...
    - Still orthogonal

#### Householder reflections

• A Householder reflection is *n*-by-*n* matrix

$$\boldsymbol{P} = \boldsymbol{I} - \beta \boldsymbol{v} \boldsymbol{v}^T$$
 where  $\beta = \frac{2}{\boldsymbol{v}^T \boldsymbol{v}}$ 

• If we set  $v = x - ||x||_2 e_1$ , then  $Px = ||x||_2 e_1$ 

• 
$$\boldsymbol{e}_1 = (1, 0, 0, ..., 0)^T$$

- Note:  $\mathbf{PA} = \mathbf{A} (\beta \mathbf{v})(\mathbf{v}^T \mathbf{A})$  where  $\beta = 2/(\mathbf{v}^T \mathbf{v})$ 
  - We never have to compute matrix *P*

## **Givens rotations**

- Householder is too crude to give identity
- Givens rotations are rank-2 corrections to the identity of form

$$\boldsymbol{G}(i,k,\theta) = \begin{pmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & \cos(\theta) & \cdots & \sin(\theta) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & -\sin(\theta) & \cdots & \cos(\theta) & \cdots & 0 \\ \vdots & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix} k$$

Pauli Miettinen

## Non-negative matrix factorization



Forcing the factors to be non-negative can, and often will, improve the interpretability of the factorization

#### Some comments on NMF

- NMF is **not** unique
  - If X is nonnegative and with nonnegative inverse, then WXX<sup>-1</sup>H is equivalent valid decomposition
- Computing NMF (and non-negative rank) is NP-hard
  - This was open until 2008

## NMF example: faces



Row of original



=







#### General idea for solving NMF

- NMF is not convex, but it is **biconvex** 
  - If **W** is fixed,  $\frac{1}{2} \|\mathbf{A} \mathbf{W}\mathbf{H}\|_F^2$  is convex
- Start from random W and repeat
  - Fix **W** and update **H**
  - Fix *H* and update *W*
- until the error doesn't decrease anymore

## The NMF-ALS algorithm

- 1.  $W \leftarrow random(n, k)$
- 2. repeat
  - 2.1.  $H \leftarrow [W^+A]_+$
  - 2.2.  $W \leftarrow [AH^+]_+$
- 3. until convergence

# The NMF gradient descent algorithm

- 1.  $W \leftarrow random(n, k)$
- 2.  $H \leftarrow random(k, m)$
- 3. repeat

3.1. 
$$H \leftarrow H - \varepsilon_H \frac{\partial f}{\partial H}$$
  
3.2.  $W \leftarrow W - \varepsilon_W \frac{\partial f}{\partial W}$ 

#### 4. until convergence

# The NMF multiplicative updates algorithm

- 1.  $W \leftarrow random(n, k)$
- 2.  $H \leftarrow random(k, m)$
- 3. repeat

3.1. 
$$h_{ij} \leftarrow h_{ij} \frac{(\boldsymbol{W}^T \boldsymbol{A})_{ij}}{(\boldsymbol{W}^T \boldsymbol{W} \boldsymbol{H})_{ij} + \varepsilon}$$
  
3.2.  $w_{ij} \leftarrow w_{ij} \frac{(\boldsymbol{A} \boldsymbol{H}^T)_{ij}}{(\boldsymbol{W} \boldsymbol{H} \boldsymbol{H}^T)_{ij} + \varepsilon}$ 

#### 4. until convergence

## **Geometry of NMF**

NMF factors Data points Convex cone Projections



## Hoyer's sparse NMF

 Hoyer (2004) considers the following sparsity function for *n*-dimensional vector *x*

sparsity(**x**) = 
$$\frac{\sqrt{n} - \left(\sum_{i} |x_{i}|\right) / \sqrt{\sum_{i} x_{i}^{2}}}{\sqrt{n} - 1}$$

- sparsity( $\mathbf{x}$ ) = 1 iff nnz( $\mathbf{x}$ ) = 1
- sparsity( $\mathbf{x}$ ) = 0 iff  $|\mathbf{x}_i| = |\mathbf{x}_j|$  for all i, j

Hoyer 2004

DMM, summer 2015

## Semi-orthogonal NMF

 In semi-orthogonal NMF we restrict *H* to roworthogonal:

minimize  $||\mathbf{A} - \mathbf{W}\mathbf{H}||_F$  s.t.  $\mathbf{H}\mathbf{H}^T = \mathbf{I}$  and  $\mathbf{W}$  and  $\mathbf{H}$  are nonnegative

- Solutions are unique (up to permutations)
- The problem is "equivalent" to k-means
  - In the sense that the optimal solutions have the same value

## **NMF and clustering**

• In *k*-means, we minimize

$$\sum_{j=1}^{k} \sum_{i \in C_j} \|\boldsymbol{a}_i - \boldsymbol{\mu}_j\|_2^2 = \sum_{j=1}^{k} \sum_{i=1}^{n} \boldsymbol{G}_{ij} \|\boldsymbol{a}_i - \boldsymbol{\mu}_j\|_2^2$$

- $\boldsymbol{\mu}_j$  is the centroid of the *j*th cluster  $C_j$
- G is n-by-k cluster assignment matrix
  - $G_{ij} = 1$  if  $i \in C_j$  and 0 otherwise
- Equivalently:  $\|\mathbf{A} \mathbf{GM}\|_F^2$  Type of NMF if A is nonnegative!
  - **M** is *k*-by-*m* containing the centroids as its rows

#### **Orthogonal tri-factor NMF**

- We can find NMF where both W and H are (column/row) orthogonal
  - Often too restrictive; cannot handle different scales
- In orthogonal nonnegative tri-factorization we add third non-negative matrix S: minimize ||A – WSH||<sub>F</sub> s.t. W<sup>T</sup>W = I, HH<sup>T</sup> = I, and all matrices are non-negative

#### Generalized KL-divergence and matrix factorizations

- The standard KL-divergence requires P and Q be probability distributions (e.g.  $\sum_i P(i) = 1$ )
  - The generalized KL-divergence (or I-divergence) removes this requirement:  $D_{GKL}(P||Q) = \sum_{i} \left( P(i) \ln \frac{P(i)}{Q(i)} - P(i) + Q(i) \right)$
- In NMF,  $P = \mathbf{A}$  and  $Q = \mathbf{WH}$ :  $D_{GKL}(\mathbf{A} || \mathbf{WH}) = \sum_{i,j} \left( \mathbf{A}_{ij} \ln \frac{\mathbf{A}_{ij}}{(\mathbf{WH})_{ii}} - \mathbf{A}_{ij} + (\mathbf{WH})_{ij} \right)$

## **NMF for GKL**

• The update rules for multiplicative GKL NMF algorithm are

$$\boldsymbol{H}_{kj} \leftarrow \boldsymbol{H}_{kj} \frac{\sum_{i=1}^{n} \boldsymbol{W}_{ik} (\boldsymbol{A}_{ij} / (\boldsymbol{W} \boldsymbol{H})_{ij})}{\sum_{i=1}^{n} \boldsymbol{W}_{ik}}$$

$$\boldsymbol{W}_{ik} \leftarrow \boldsymbol{W}_{ik} \frac{\sum_{j=1}^{m} (\boldsymbol{A}_{ij} / (\boldsymbol{W} \boldsymbol{H})_{ij}) \boldsymbol{H}_{kj}}{\sum_{j=1}^{m} \boldsymbol{H}_{kj}}$$

 The columns of *W* are normalized to sum to unity after every iteration

## pLSI generative process

- Pick a document according to P(d)
- Select a topic according to
   P(z | d)
- Select a word according to
   P(w | z)



## The CX decomposition

- In the CX decomposition we are given a matrix A and a rank k, and we need to select k columns of A into matrix C and build matrix X s.t. we minimize ||A CX||<sub>ξ</sub>
  - ξ is either *F* or 2
  - A.k.a. column subset selection problem (CSSP)

## Related idea: RRQR

 The rank-revealing QR (RRQR) factorization *k*-by-*k* upper-triangular w/ positive diagonal of matrix **A** is *k*-by-(*m*-*k*) *n*-by-*n* orthogonal  $\mathbf{A} \Pi = \mathbf{Q} \mathbf{R} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ 0 & \mathbf{R}_{22} \end{pmatrix} (n-k)-by-(m-k)$ Permutation matrix that satisfies n-by-m kth singular value of A  $\frac{\sigma_k(\boldsymbol{A})}{p_1(k,m)} \leq \sigma_{\min}(\boldsymbol{R}_{11}) \leq \sigma'_k(\boldsymbol{A})$  $\sigma_{k+1}(\boldsymbol{A}) \leq \sigma_{\max}(\boldsymbol{R}_{22}) \leq p_2(k,m)\sigma_{k+1}(\boldsymbol{A})$ Some polynomial on k and m

## **Geometry of NNCX**

Columns in **C** 0.9 0.8 0.7 Columns not in C<sup>0.7</sup> 0.5 Convex cone 0.4 0.3 Projections 0.2 0.1 1.5 0 1.5 0.5 1 0.5 0 0
## The CUR decomposition

 In the CUR decomposition we are given matrix A and integers c and r, and our task is to select c columns of A to matrix C and r rows to matrix R, and build c-by-r matrix U minimizing ||A – CUR||<sub>F</sub>

• Often 
$$c = r = k$$

## ICA definition

- Setting. Let  $x_j \in \mathbb{R}$ , j=1,...,n be observed random variables. Assume there exists nlatent random variables  $s_i \in \mathbb{R}$  and latent coefficients  $a_{ij}$  such that  $x_j = \sum_i a_{ij}s_i$  for all j.
  - x = sA and for T observations, X = SA
     where X and S have T rows
- Problem. Find A and s given x

# ICA assumptions (important slide!)

- Original signals s<sub>i</sub> are mutually statistically independent
- At most one original signal signa
- The mixing matrix A is square and invertible
  - This is not necessary but simplifies the theory



### **Computing ICA: Central limit theorem**

 Average of i.i.d. variables converges to normal distribution

• 
$$\sqrt{n}\left(\left(\frac{1}{n}\sum_{i=1}^{n}X_{i}\right)-\mu\right)\overset{d}{\rightarrow}N(0,\sigma^{2})$$
 as  $n\rightarrow\infty$ 

- Hence  $(X_1 + X_2)/2$  is "more Gaussian" than  $X_1$  or  $X_2$ alone
  - For i.i.d. zero-centered non-Gaussian X<sub>1</sub> and X<sub>2</sub>
- Hence, we can try to find components s that are "maximally non-Gaussian"

## FastICA for one IC

• Noticing that  $||\mathbf{w}||^2 = 1$  by constraint and taking infinite step update, we get

$$\mathbf{w} \leftarrow E[(\mathbf{z}\mathbf{w}')^{3}\mathbf{z}] - 3\mathbf{w}$$

- Again set  $\mathbf{w} \leftarrow \mathbf{w}/||\mathbf{w}||$  after every update
- Expectation has naturally to be estimated
- No theoretical guarantees but works in practice

# Multiple components

- So far we have found only one component
  - To find more, remember that vectors *w<sub>i</sub>* are orthogonal (whitening!)
- General idea:
  - Find one vector w
  - Find second that is orthogonal to the first one
  - Find third that is orthogonal to the two previous ones, etc.

# Putting it all together

- Start with random **B** and  $\gamma$ , choose learning rates  $\delta$  and  $\delta_{\gamma}$
- Iterate until convergence
  - *y* ← *Bx* and normalize *y* to unit variance
  - $\gamma_i \leftarrow (1 \delta_\gamma)\gamma_{i-1} + \delta_\gamma E[-tanh(y_i)y_i + (1 tanh(y_i)^2)]$ 
    - if  $\gamma_i > 0$ , use super-Gaussian g; o/w sub-Gaussian g
  - $\boldsymbol{B} \leftarrow \boldsymbol{B} + \delta(\boldsymbol{I} + \sum_{t} \boldsymbol{g}(\boldsymbol{y}_{t})^{T} \boldsymbol{y}_{t}) \boldsymbol{B}$

## **Graph Laplacians**

• The Laplacian matrix *L* of a graph is the adjacency matrix subtracted from the degree matrix

$$\boldsymbol{L} = \boldsymbol{\Delta} - \boldsymbol{A} = \begin{pmatrix} \sum_{j \neq 1} a_{1,j} & -a_{1,2} & \cdots & -a_{1,n} \\ -a_{2,1} & \sum_{j \neq 2} a_{2,j} & \cdots & -a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ -a_{n,1} & -a_{n,2} & \cdots & \sum_{j \neq n} a_{n,j} \end{pmatrix}$$

- The Laplacian is symmetric and positive semi-definite
  - Undirected graphs
  - Has n real, non-negative, orthogonal eigenvalues

 $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \ldots \geq \lambda_n \geq 0$ 

### **Graph cuts using matrices**

$$\text{RatioCut} = \sum_{i=1}^{k} \frac{W(C_i, V \setminus C_i)}{|C_i|} = \sum_{i=1}^{k} \frac{\boldsymbol{c}_i^T \boldsymbol{L} \boldsymbol{c}_i}{\|\boldsymbol{c}_i\|^2}$$

NormalizedCut = 
$$\sum_{i=1}^{k} \frac{W(C_i, V \setminus C_i)}{vol(C_i)} = \sum_{i=1}^{k} \frac{c_i^T L c_i}{c_i^T \Delta c_i}$$

# Spectral clustering pseudo-code

#### Assume connected graph



#### Definition of the BMF

#### Boolean Matrix Factorization (BMF)

The (exact) **Boolean matrix factorization** of a binary matrix  $A \in \{0,1\}^{m \times n}$  expresses it as a Boolean product of two factor matrices,  $B \in \{0,1\}^{m \times k}$  and  $C \in \{0,1\}^{k \times n}$ . That is  $A = B \boxtimes C$ .

- Typically (in data mining), k is given, and we try to find B and C to get as close to A as possible
- Normally the optimization function is the squared Frobenius norm of the residual,  $\|\mathbf{A} (\mathbf{B} \boxtimes \mathbf{C})\|_F^2$ 
  - Equivalently,  $| \boldsymbol{A} \oplus (\boldsymbol{B} \boxtimes \boldsymbol{C}) |$  where
    - $\star$  |**A**| is the sum of values of **A** (number of 1s for binary matrices)
    - ★  $\oplus$  is the element-wise exclusive-or (1+1=0)
  - The alternative definition is more "combinatorial" in flavour

#### BMF and (quasi-)biclique covers



- A **biclique** is a complete bipartite graph
  - Each left-hand-side verted is connected to each right-hand-side vertex
- Each rank-1 binary matrix defines a biclique (subgraph)
  - If  $\boldsymbol{v} \in \{0,1\}^m$  and  $\boldsymbol{u} \in \{0,1\}^n$ , then  $\boldsymbol{v}\boldsymbol{u}^T$  is a biclique between  $v_i \in V$  and  $u_i \in U$  for which  $\boldsymbol{v}_i = \boldsymbol{u}_i = 1$
- Exact BMF corresponds to covering each edge of the graph with at least one biclique
  - In approximate BMF, quasi-bicliques cover most edges

#### BMF and the Set Basis problem



- In the Set Basis problem, we are given a set system (U, S), and our task is to find collection C ⊆ 2<sup>U</sup> such that we can cover each set S ∈ S with a union of some sets of C
  - For each  $S \in S$ , there is  $C_S \subseteq C$  such that  $S = \bigcup_{C \in C_S} C$
- A set basis corresponds to exact BMF
  - The size of the smallest set basis is the Boolean rank
- N.B.: this is the same problem as covering with bicliques

#### Example of $\pm \text{PSC}$ and Basis Usage

defines the sets



# Can we find the original pattern?



Pauli Miettinen

## Planted partition model

- Let  $G(\varphi, \mathbf{P})$  be a random graph distribution where  $\varphi: V \rightarrow \{1,...,k\}$  partition the vertices to k classes and  $\mathbf{P}=(p_{ij})$  is a k-by-k matrix with  $p_{ij} \in [0,1]$ . Include edge (i, j) with probability  $p_{\varphi(i)\varphi(j)}$ .
- **Example**: planted clique. Let  $\varphi(v) = 1$  iff v is in the clique. Set  $p_{11} = 1$  and  $p_{ij} = p$  elsewhere
- **Problem**. Given a sample **G**' from  $G(\varphi, \mathbf{P})$ , find a partition  $\varphi'$  s.t.  $\varphi'(\mathbf{v}) = \varphi'(\mathbf{u})$  iff  $\varphi(\mathbf{v}) = \varphi(\mathbf{u})$

### **Planted partition results**

- Now if  $||\mathbf{g}_u \mathbf{g}_v||$  is large enough when  $\varphi(v) \neq \varphi(u)$ , we can find  $\varphi$ 
  - Depends on the above error bounds
- With more complicated error bounds we get:
  - If *s* is the size of a planted clique, then there is a constant *c* s.t. for sufficiently large *n* we can recover  $\varphi$  with probability  $1 \delta$  if

$$\frac{1-p}{p} > C\left(\frac{n}{s^2} + \frac{\log(n/\delta)}{s}\right)$$

### Maximum clique as rank minimization

• Maximum *n*-vertex clique in graph G = (V, E)can be found with the following program



### Nuclear norm relaxation

- The rank minimization problem is NP-hard
- We can relax it to nuclear norm minimization:

min 
$$\|X\|_*$$
  
s.t.  $\sum_{i \in V} \sum_{j \in V} x_{ij} \ge n^2 \leftarrow \text{can be replaced with 1}$   
 $x_{ij} = 0$  if  $\{i, j\} \notin E$  and  $i \neq j$ 

- The maximum clique is a valid solution and the unique optimizer under certain conditions
  - When this is the case, we can find the clique

### Destructive noise models for bicliques

- So far we've added each edge independently with probability p
  - Erdős–Rényi random graph model
- We can also follow the preferential attachment model
  - Barabási–Albert random graph model
  - Some vertices have big changes on neighbors, others less
    - If the noise follows the B–A model, it can't have large bicliques ⇒ easy

### Results

- Erdős–Rényi: The minimum size of the original biclique  $\zeta = \log(NM)$
- Barabási–Albert: log  $N \ll \zeta \ll \sqrt{N}$

### On exam

## Format & basic info

- Written exam
- 28 July 2015 from 12:00–14:00
  - Times are sharp!
- Lecture hall 001, building E1 3

# What you can and cannot bring

- You can (must) bring
  - writing equipments & student ID
  - one (1) A4-sized "cheat sheet" paper
- You cannot bring (use)
  - electronic devices (incl. phones and pocket calculators and electric pencil sharpeners)
  - any other notes than the cheat sheet (incl. lecture slides, assignments, etc.)

## **Cheat sheet**

- Must contain your name!
- A4-sized paper, text can be on both sides
- Any content is OK (as long as its legal)
  - Use your discretion what you think is important or consider hard
- Can be made with computer or be handwritten (or with typewriter)

# What is covered in the exam?

- All lectures between 21 April to 21 July
  - Lecture on 21 July is a wrap-up, no new contents
- All pen-and-paper and first two programming assignments
- The chapters of books and articles cited in the lecture slides

# What kind of questions are there in the exam?

- Simple mathematical proofs
  - Similar to those in homework assignments
- Developing variations of presented algorithms
  - "Explain how would you compute ABC decomposition with the following constraints"
- Short texts or longer essays comparing different decomposition methods and/or explaining their use cases and interpretations
  - "What are the main differences between ABC and XYZ?" "Given this-and-that kind of data, how would you interpret its ABC decomposition?"
- Short questions about features and properties of decompositions and methods
  - "Explain briefly the main idea behind algorithms computing ABC." "True or false: computing the optimal XYZ decomposition (w.r.t. the Frobenius norm) is NP-hard."

#### Pauli Miettinen

# Exam checking day

- 31 July from 12:15 to 14:00
  - D5 rotunda (E1 4, 4th floor, left from elevator)
  - your only chance

### **Re-Exam**

- Open for everybody
  - Bonus points count, better of two exams
- You must register via email by 7 August
   12 noon
- Date & place TBD (late September)

### Advertisements

# Follow-up seminar

- I'm planning for a follow-up seminar in next semester
  - Topics involve deeper dive into new matrix (and maybe tensor) factorization methods for data analysis
  - Current plan (subject to change): Block seminar with two days of presentations in January/early February
- Limited attendance with first-come-first-served basis
  - Send me e-mail if you want to get notified as soon as registration to the seminar is possible

# HiWi & MSc student positions

- I have positions for HiWis & MSc students on data mining
  - Matrices, tensors, and other stuff tailored to taste
- This course is important evaluation point
  - I will not consider any application before I know the results of this course
    - ⇒ If interested, apply in late August/early
      September

# Ask Me Anything

#### Johnson–Lindenstrauss lemma

- Finding the decomposition can be expensive
- Decompositions give only *global* guarantees
  Any pair of points can have very different distances
- Can we guarantee *local* similarity?

**Johnson–Lindenstrauss lemma**. Given  $\varepsilon > 0$  and an integer *n*, let *k* be a positive integer such that  $k \ge k_0 = O(\varepsilon^{-2}\log n)$ . For every set *X* of *n* points in  $\mathbb{R}^d$  there exists  $F: \mathbb{R}^d \to \mathbb{R}^k$  such that for all  $\mathbf{x}_i, \mathbf{x}_j \in X$  $(1 - \varepsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2 \le \|\mathbf{F}(\mathbf{x}_i) - \mathbf{F}(\mathbf{x}_j)\|^2 \le (1 + \varepsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2$ 

#### How to find the projections?

- We need to find an *k*-by-*d* matrix  $\mathbf{R} = (r_{ij})$  such that function  $\mathbf{x} \mapsto \mathbf{R}\mathbf{x}$  satisfies JL
- Remarkably, if we select  $r_{ij} \sim N(0,1)$ , *R* satisfies JL with high probability

– That is, JL holds for *all* points of X with high probability

• Achlioptas has show that we can also select  $\Pr[r_{ij} = 1] = 1/2$  and  $\Pr[r_{ij} = -1] = 1/2$  or  $\Pr[r_{ij} = 1] = 1/6$ ,  $\Pr[r_{ij} = 0] = 2/3$ ,  $\Pr[r_{ij} = -1] = 1/6$ - Sparse matrix
## **Spurious correlations**

## **Age of Miss America**

correlates with

## Murders by steam, hot vapours and hot objects



http://www.tylervigen.com/spurious-correlations

Pauli Miettinen