

# The Nonnegative Matrix Factorization

## a tutorial

Barbara Ball

barbaraeball@comcast.net

C. of Charleston

Mathematics Dept.

Atina Brooks

adbroom2@ncsu.edu

N.C. State U.

Statistics Dept.

Amy Langville

langvillea@cofc.edu

C. of Charleston

Mathematics Dept.

NISS NMF Workshop

February 23–24, 2007

# Outline

- Two Factorizations:
  - Singular Value Decomposition
  - Nonnegative Matrix Factorization
- Why factor anyway?
- Computing the NMF
  - Early Algorithms
  - Recent Algorithms
- Extensions of NMF

# Data Matrix

$\mathbf{A}_{m \times n}$  with rank  $r$

## Examples

term-by-document matrix

pixel intensity-by-image matrix

gene-by-DNA microarray matrix

feature-by-item matrix

user-by-purchase matrix

terrorist-by-action matrix

## SVD

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

# What is the SVD?

$$\underset{m \times n}{A} = \underset{m \times m}{U} \underset{m \times n}{\Sigma} \underset{n \times n}{V}^T$$

$$U^T U = I_m$$

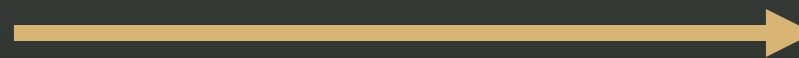
$$V^T V = I_n$$

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$$

$r$  : rank of  $A$

$$A_{m \times n} = \sum_{i=0}^{\min(m,n)} \sigma_i u_i v_i^T$$

$$A_{m \times n} = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T$$



decreasing importance

## The SVD

# Data Matrix

$\mathbf{A}_{m \times n}$  with rank  $r$

## Examples

term-by-document matrix

pixel intensity-by-image matrix

gene-by-DNA microarray matrix

feature-by-item matrix

user-by-purchase matrix

terrorist-by-action matrix

## SVD

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

## Low Rank Approximation

use  $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$  in place of  $\mathbf{A}$

$$A_{nkn} \approx \begin{bmatrix} U_{mk} & U_{mm} \end{bmatrix} \times \begin{bmatrix} S_{kk} & \\ & S_{mn} \end{bmatrix} \times \begin{bmatrix} V_{nk}^T \\ V_{nn}^T \end{bmatrix}$$

## SVD Rank Reduction

# Why use Low Rank Approximation?

- Data Compression and Storage when  $k \ll r$
- Remove noise and uncertainty
  - ⇒ improved performance on data mining task of retrieval (e.g., find similar items)
  - ⇒ improved performance on data mining task of clustering



# Properties of SVD

- basis vectors  $\mathbf{u}_i$  and  $\mathbf{v}_i$  are orthogonal

- $u_{ij}, v_{ij}$  are mixed in sign

$$\underset{\text{nonneg}}{\mathbf{A}_k} = \underset{\text{mixed}}{\mathbf{U}_k} \underset{\text{nonneg}}{\Sigma_k} \underset{\text{mixed}}{\mathbf{V}_k^T}$$

- $\mathbf{U}, \mathbf{V}$  are dense
- *uniqueness*—while there are many SVD algorithms, they all create the same (truncated) factorization
- *optimality*—of all rank- $k$  approximations,  $\mathbf{A}_k$  is optimal

$$\|\mathbf{A} - \mathbf{A}_k\|_F = \min_{\text{rank}(\mathbf{B}) \leq k} \|\mathbf{A} - \mathbf{B}\|_F$$

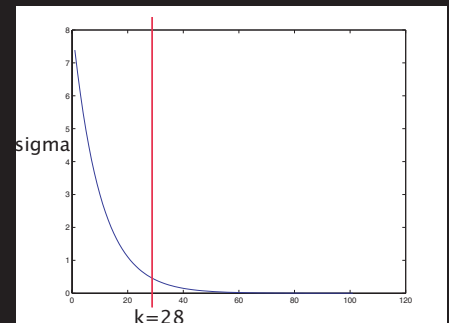
# Summary of Truncated SVD

## Strengths

- using  $\mathbf{A}_k$  in place of  $\mathbf{A}$  gives improved performance
- noise reduction isolates essential components of matrix
- best rank- $k$  approximation
- $\mathbf{A}_k$  is unique

## Weaknesses

- storage— $\mathbf{U}_k$  and  $\mathbf{V}_k$  are usually completely dense
- interpretation of basis vectors is difficult due to mixed signs
- good truncation point  $k$  is hard to determine
- orthogonality restriction



# Other Low-Rank Approximations

- **QR** decomposition
- any **URV**<sup>T</sup> factorization

- Semidiscrete decomposition (SDD)

$\mathbf{A}_k = \mathbf{X}_k \mathbf{D}_k \mathbf{Y}_k^T$ , where  $\mathbf{D}_k$  is diagonal, and elements of  $\mathbf{X}_k, \mathbf{Y}_k \in \{-1, 0, 1\}$ .

- **CUR** factorization

# Other Low-Rank Approximations

- **QR** decomposition
- any **URV**<sup>T</sup> factorization

- Semidiscrete decomposition (SDD)

$$\mathbf{A}_k = \mathbf{X}_k \mathbf{D}_k \mathbf{Y}_k^T, \text{ where } \mathbf{D}_k \text{ is diagonal, and elements of } \mathbf{X}_k, \mathbf{Y}_k \in \{-1, 0, 1\}.$$

- **CUR** factorization

**BUT**

All create basis vectors that are mixed in sign.

**Negative** elements make interpretation difficult.

# Other Low-Rank Approximations

- **QR** decomposition
- any **URV**<sup>T</sup> factorization

- Semidiscrete decomposition (SDD)

$$\mathbf{A}_k = \mathbf{X}_k \mathbf{D}_k \mathbf{Y}_k^T, \text{ where } \mathbf{D}_k \text{ is diagonal, and elements of } \mathbf{X}_k, \mathbf{Y}_k \in \{-1, 0, 1\}.$$

- **CUR** factorization

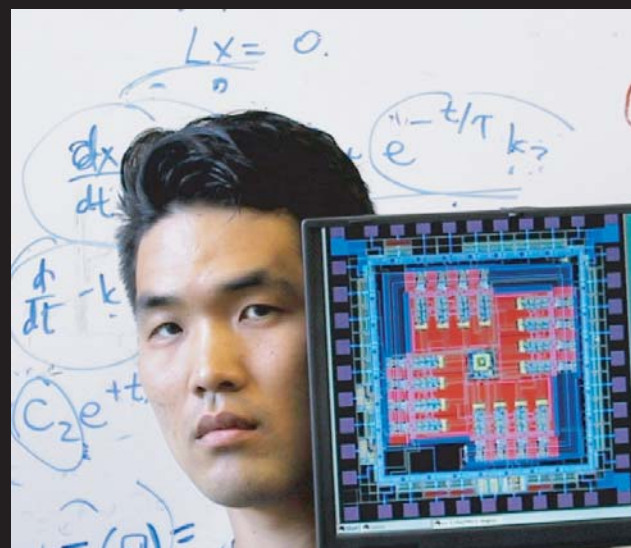
**BUT**

All create basis vectors that are mixed in sign.

**Negative** elements make interpretation difficult.

⇒ **Nonnegative Matrix Factorization**

# Nonnegative Matrix Factorization (2000)



## Daniel Lee and Sebastian Seung's Nonnegative Matrix Factorization

Idea: use low-rank approximation with nonnegative factors to improve weaknesses of trun-SVD

$$\text{SVD} \quad \mathbf{A}_k = \underset{\text{mixed}}{\mathbf{U}_k} \underset{\text{nonneg}}{\Sigma_k} \underset{\text{mixed}}{\mathbf{V}_k^T}$$

$$\text{NMF} \quad \underset{\text{nonneg}}{\mathbf{A}_k} = \underset{\text{nonneg}}{\mathbf{W}_k} \underset{\text{nonneg}}{\mathbf{H}_k}$$

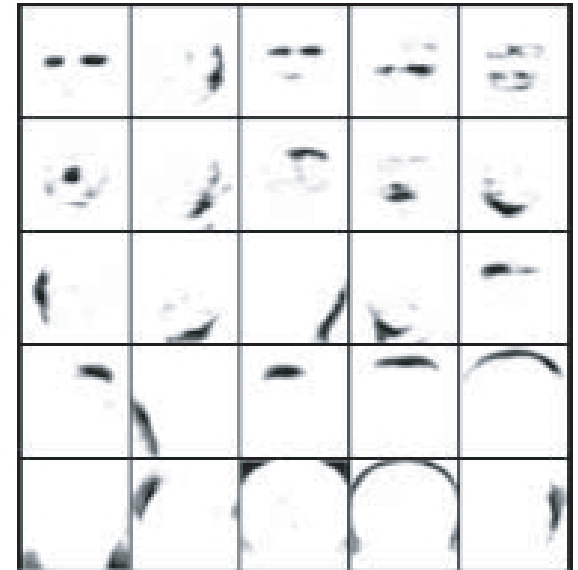
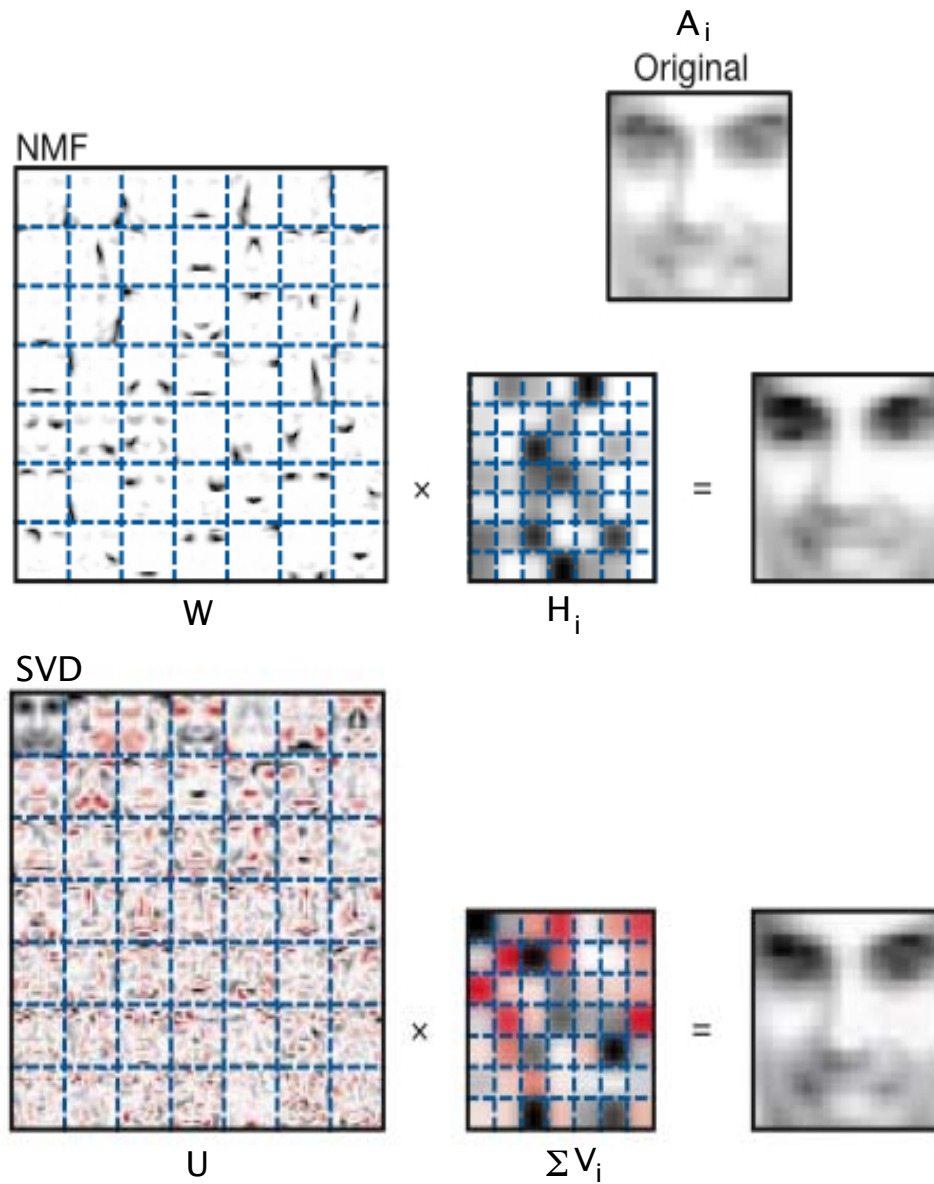
# Interpretation with NMF

- columns of  $\mathbf{W}$  are the underlying basis vectors, i.e., each of the  $n$  columns of  $\mathbf{A}$  can be built from  $k$  columns of  $\mathbf{W}$ .
- columns of  $\mathbf{H}$  give the weights associated with each basis vector.

$$\mathbf{A}_k \mathbf{e}_1 = \mathbf{W}_k \mathbf{H}_{*1} = \begin{bmatrix} \vdots \\ \mathbf{w}_1 \\ \vdots \end{bmatrix} h_{11} + \begin{bmatrix} \vdots \\ \mathbf{w}_2 \\ \vdots \end{bmatrix} h_{21} + \cdots + \begin{bmatrix} \vdots \\ \mathbf{w}_k \\ \vdots \end{bmatrix} h_{k1}$$

- $\mathbf{W}_k, \mathbf{H}_k \geq 0 \Rightarrow$  immediate interpretation (additive parts-based rep.)

# Image Mining





# Image Mining Applications

- Data compression
- Find similar images
- Cluster images



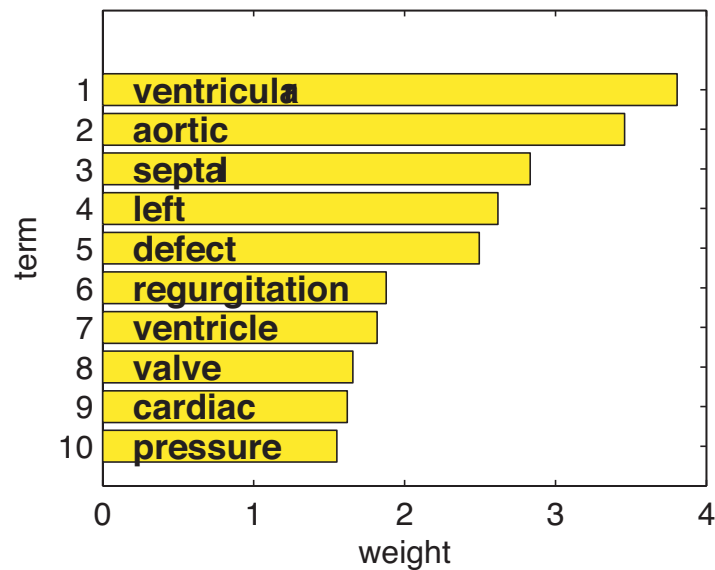
Original Image  
 $r = 400$

Reconstructed Images  
 $k = 100$

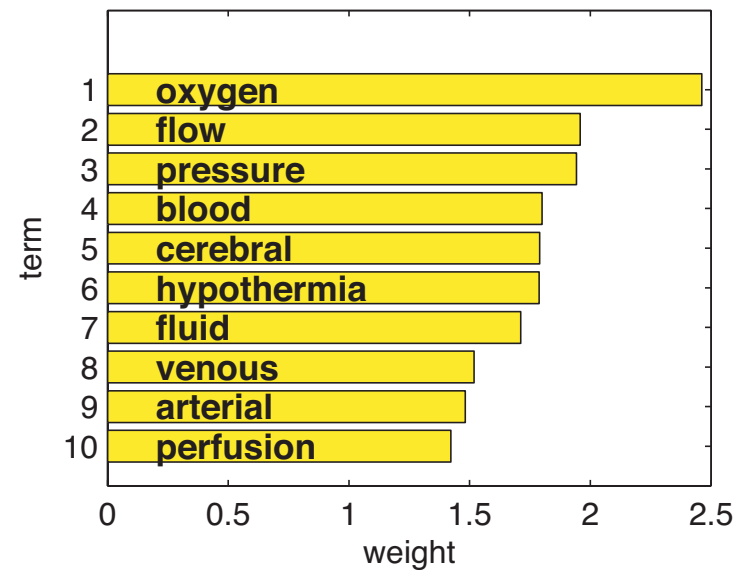
# Text Mining

## MED dataset ( $k = 10$ )

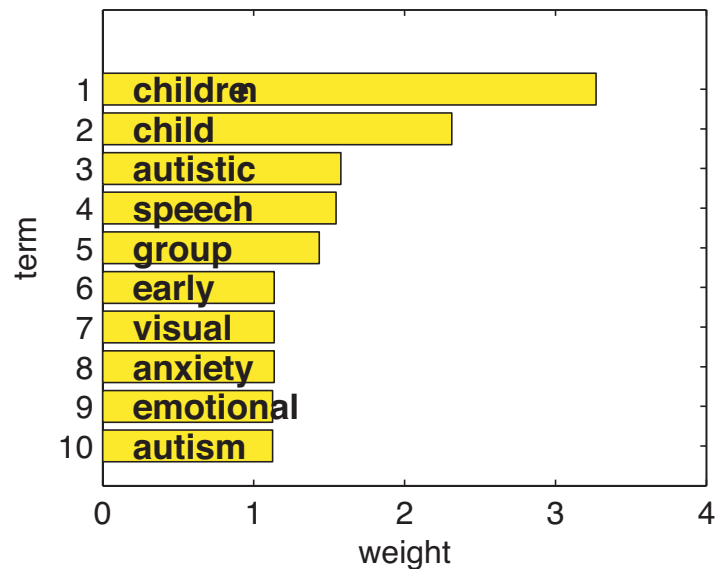
Highest Weighted Terms in Basis Vector  $W_1$



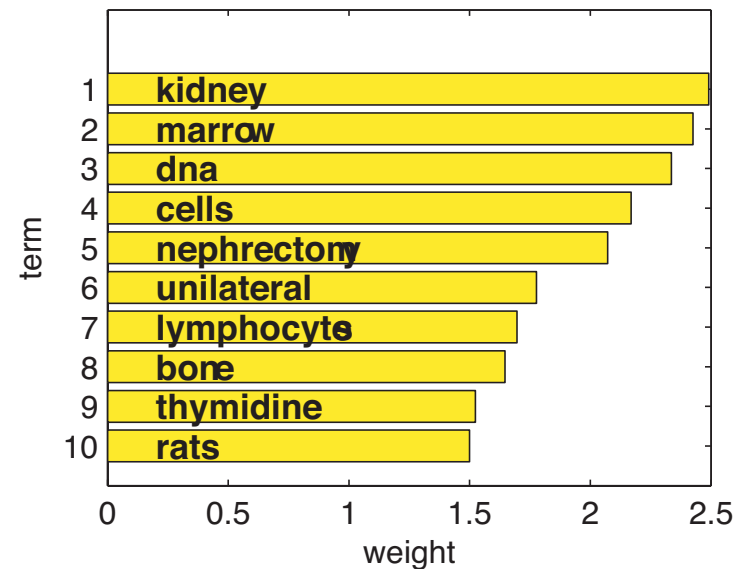
Highest Weighted Terms in Basis Vector  $W_2$



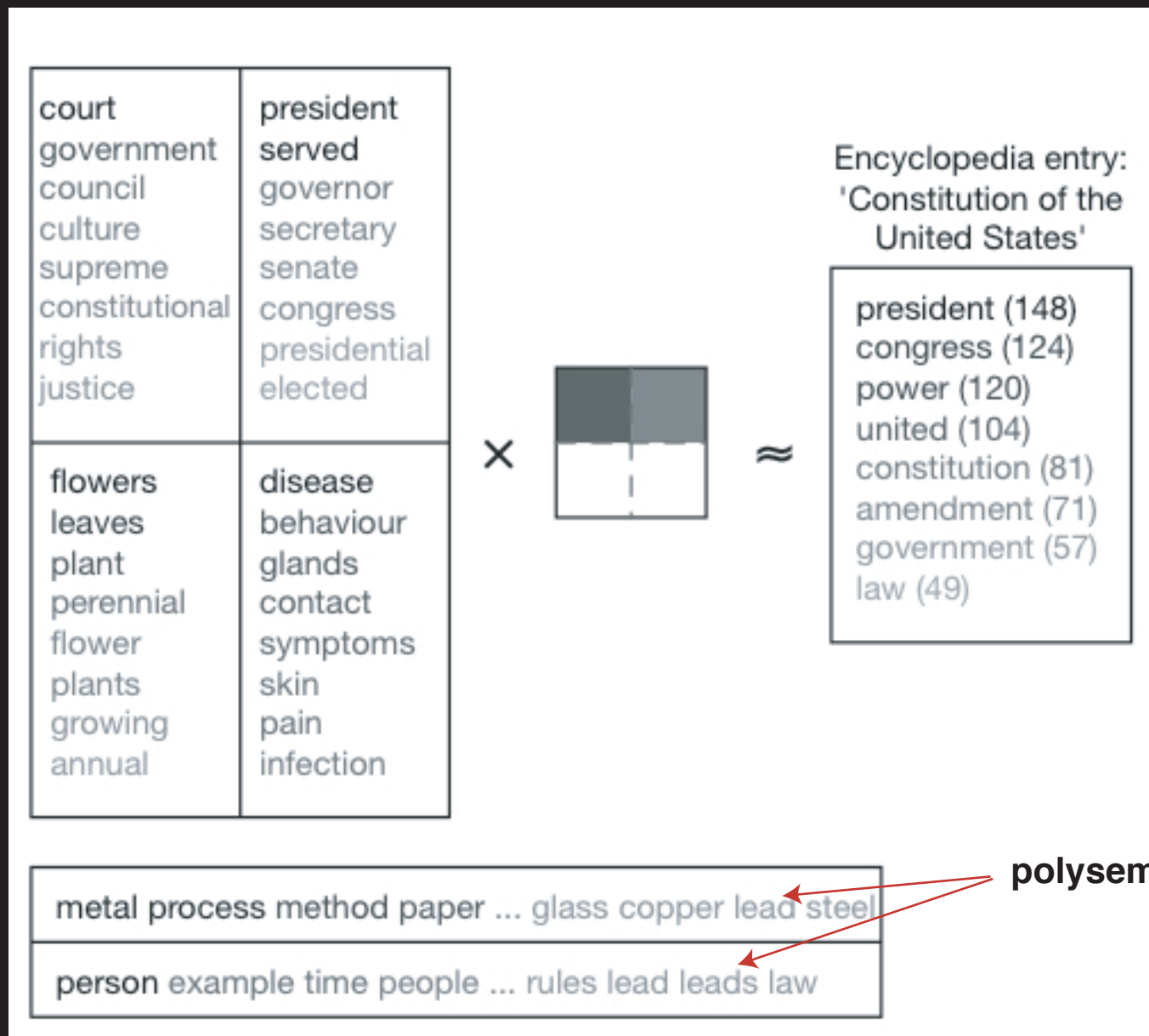
Highest Weighted Terms in Basis Vector  $W_5$



Highest Weighted Terms in Basis Vector  $W_6$



# Text Mining



- polysems broken across several basis vectors  $\mathbf{w}_i$

# Text Mining Applications

- Data compression  $\mathbf{W}_k \mathbf{H}_k$
- Find similar terms  $0 \leq \cos(\theta) = \mathbf{W}_k \mathbf{H}_k \mathbf{q} \leq 1$
- Find similar documents  $0 \leq \cos(\theta) = \mathbf{q}^T \mathbf{W}_k \mathbf{H}_k \leq 1$
- Cluster documents

# Clustering with the NMF

## Clustering Terms

- use rows of  $\mathbf{W}_{m \times k}$

$$= \begin{matrix} & \begin{matrix} cl.1 & cl.2 & \dots & cl.k \end{matrix} \\ \begin{matrix} term1 \\ term2 \\ \vdots \end{matrix} & \begin{pmatrix} .9 & 0 & \dots & .3 \\ .1 & .8 & \dots & .2 \\ \vdots & \vdots & \ddots & \vdots \end{pmatrix} \end{matrix}$$

## Clustering Documents

- use cols of  $\mathbf{H}_{k \times n}$

$$= \begin{matrix} & \begin{matrix} doc1 & doc2 & \dots & docn \end{matrix} \\ \begin{matrix} cl.1 \\ \vdots \\ cl.k \end{matrix} & \begin{pmatrix} .4 & 0 & \dots & .5 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & .8 & \dots & .2 \end{pmatrix} \end{matrix}$$

soft clustering is very natural

# The Enron Email Dataset

(SAS)

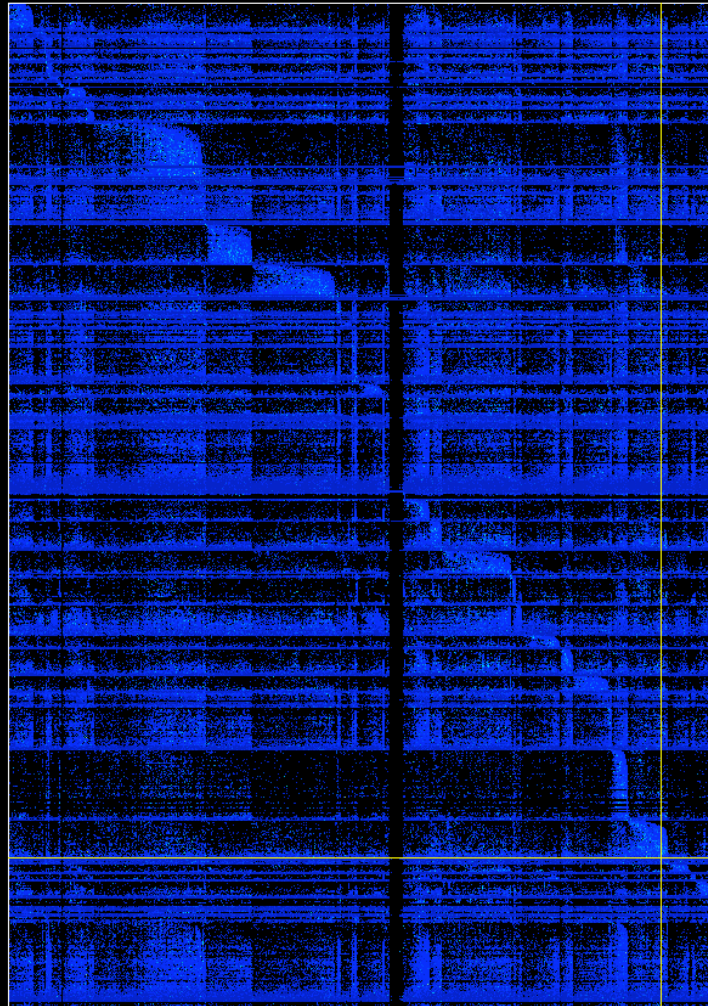
- PRIVATE email collection of 150 Enron employees during 2001
- 92,000 terms and 65,000 messages
- Term-by-Message Matrix

$$\begin{array}{c} \vdots \\ subpoena \\ dynegey \\ \vdots \end{array} \begin{pmatrix} fastow1 & fastow2 & skilling1 & \dots \\ \vdots & \vdots & \vdots & \dots \\ 2 & 0 & 1 & \dots \\ 0 & 3 & 0 & \dots \\ \vdots & \vdots & \vdots & \dots \end{pmatrix}$$

Year	1950	1951	1952	1953	1954	1955	1956	1957	1958	1959	1960	1961	1962	1963	1964	1965	1966	1967	1968	1969	1970	1971	1972	1973	1974	1975	1976	1977	1978	1979	1980	1981	1982	1983	1984	1985	1986	1987	1988	1989	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100																														
Population	1,000,000	1,050,000	1,100,000	1,150,000	1,200,000	1,250,000	1,300,000	1,350,000	1,400,000	1,450,000	1,500,000	1,550,000	1,600,000	1,650,000	1,700,000	1,750,000	1,800,000	1,850,000	1,900,000	1,950,000	2,000,000	2,050,000	2,100,000	2,150,000	2,200,000	2,250,000	2,300,000	2,350,000	2,400,000	2,450,000	2,500,000	2,550,000	2,600,000	2,650,000	2,700,000	2,750,000	2,800,000	2,850,000	2,900,000	2,950,000	3,000,000	3,050,000	3,100,000	3,150,000	3,200,000	3,250,000	3,300,000	3,350,000	3,400,000	3,450,000	3,500,000	3,550,000	3,600,000	3,650,000	3,700,000	3,750,000	3,800,000	3,850,000	3,900,000	3,950,000	4,000,000	4,050,000	4,100,000	4,150,000	4,200,000	4,250,000	4,300,000	4,350,000	4,400,000	4,450,000	4,500,000	4,550,000	4,600,000	4,650,000	4,700,000	4,750,000	4,800,000	4,850,000	4,900,000	4,950,000	5,000,000	5,050,000	5,100,000	5,150,000	5,200,000	5,250,000	5,300,000	5,350,000	5,400,000	5,450,000	5,500,000	5,550,000	5,600,000	5,650,000	5,700,000	5,750,000	5,800,000	5,850,000	5,900,000	5,950,000	6,000,000	6,050,000	6,100,000	6,150,000	6,200,000	6,250,000	6,300,000	6,350,000	6,400,000	6,450,000	6,500,000	6,550,000	6,600,000	6,650,000	6,700,000	6,750,000	6,800,000	6,850,000	6,900,000	6,950,000	7,000,000	7,050,000	7,100,000	7,150,000	7,200,000	7,250,000	7,300,000	7,350,000	7,400,000	7,450,000	7,500,000	7,550,000	7,600,000	7,650,000	7,700,000	7,750,000	7,800,000	7,850,000	7,900,000	7,950,000	8,000,000	8,050,000	8,100,000	8,150,000	8,200,000	8,250,000	8,300,000	8,350,000	8,400,000	8,450,000	8,500,000	8,550,000	8,600,000	8,650,000	8,700,000	8,750,000	8,800,000	8,850,000	8,900,000	8,950,000	9,000,000	9,050,000	9,100,000	9,150,000	9,200,000	9,250,000	9,300,000	9,350,000	9,400,000	9,450,000	9,500,000	9,550,000	9,600,000	9,650,000	9,700,000	9,750,000	9,800,000	9,850,000	9,900,000	9,950,000	10,000,000



```
row [78198]:  destroy
col [59746]:  sanders-rsempra8
val [78198][59746]:  0.459000
```





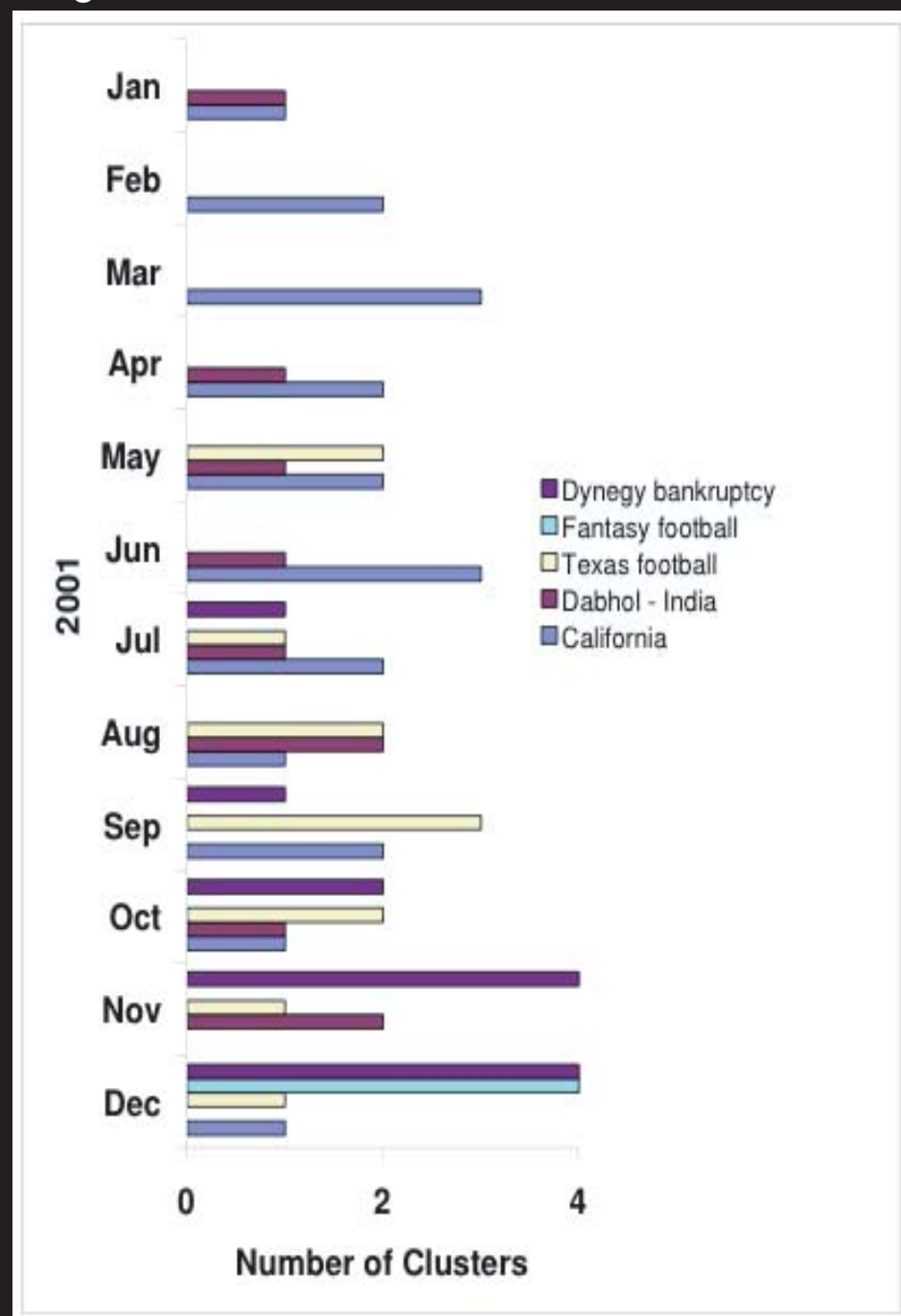
# Text Mining Applications

- Data compression  $\mathbf{W}_k \mathbf{H}_k$
- Find similar terms  $0 \leq \cos(\theta) = \mathbf{W}_k \mathbf{H}_k \mathbf{q} \leq 1$
- Find similar documents  $0 \leq \cos(\theta) = \mathbf{q}^T \mathbf{W}_k \mathbf{H}_k \leq 1$
- Cluster documents
- Topic detection and tracking

# Text Mining Applications

Enron email messages 2001

Feature Index ( $k$ )	Cluster Size	Topic Description	Dominant Terms
10	497	California	ca, <b>cpuc</b> , gov, <b>socalgas</b> , sempra, org, sce, gmssr, aelaw, ci
23	43	Louise Kitchen named top woman by Fortune	evp, <b>fortune</b> , britain, woman, <b>ceo</b> , avon, fiorinai, cfo, hewlett, packard
26	231	Fantasy football	game, wr, qb, play, rb, season, injury, updated, fantasy, image
33	233	Texas longhorn football newsletter	UT, orange, longhorn[s], texas, true, truorange, recruiting, oklahoma defensive
34	65	Enron collapse	<b>partnership[s]</b> , <b>fastow</b> , shares, <b>sec</b> , stock, shareholder, investors, equity, <b>lay</b>
39	235	Emails about India	<b>dahhol</b> , <b>dpc</b> , <b>india</b> , <b>mseb</b> , <b>maharashtra</b> , indian, lenders, delhi, foreign, minister
46	127	Enron collapse	dow, debt, reserved, wall, copyright jones, cents, analysts, reuters, spokesman



# Recommendation Systems

purchase  
history  
matrix

$$\mathbf{A} = \begin{matrix} & \begin{matrix} \text{User 1} & \text{User 2} & \dots & \text{User n} \end{matrix} \\ \begin{matrix} \text{Item 1} \\ \text{Item 2} \\ \vdots \\ \text{Item m} \end{matrix} & \begin{pmatrix} 1 & 5 & \dots & 0 \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & \dots & 2 \end{pmatrix} \end{matrix}$$

- Create profiles for classes of users from basis vectors  $\mathbf{w}_i$
- Find similar users
- Find similar items

# Microarray Study

*Kim and Tidor, 2003*

- 300 experiments with 5436 *S. cerevisiae* genes
- expression for a gene described by the expression in experiment divided by control experiment of wild type under typical conditions
- basis vector represented by an experiment, containing a relative expression for each gene and its related feature

# Functional Relationships

**Table 3.** The 58 Predictions That Could Be Validated by YPD of the 100 Strongest Functional Relationships Detected by NMF

## Coregulated

dir1	ecm34
gyp1	yap7
adel6	sir1
hpt1	sir1
rml2	ymr293c
cbp2	mmp133
mmp133	rml2
cnb1	yor072w
adel6	ymr041c
gid1	utr4
cla4 (haploid)	KAR2 (tet promoter)
yel001c	ymr141c
ckb2	gon4
arg5,6	rpl8a
mrt4	rpl12a
clb6	whi2
erp2	ymr141c
erp2	yel001c
erp2	yor015w
rpl12a	yel033w
ckb2	rtg1
eca39	ras1

## Identical Genes

lsw1	lsw1, lsw2
dig1, dig2	dig1, dig2 (haploid)
fks1 (haploid)	FKS1 (tet promoter)
bub3	bub3 (haploid)

## Binding

cla4 (haploid)	CDC42 (tet promoter)
qcr2 (haploid)	rip1
far1 (haploid)	ste4 (haploid)
bub1 (haploid)	bub3
bub1 (haploid)	bub3 (haploid)

## Cell Wall

fks1 (haploid)	2-deoxy-D-glucose
2-deoxy-D-glucose	Glucosamine
gis1	Tunicamycin
fks1 (haploid)	Glucosamine
yer083c	Tunicamycin
ste12 (haploid)	ste5 (haploid)

## Mating

ste5 (haploid)	ste7 (haploid)
fus3, kss1 (haploid)	ste5 (haploid)
ste18 (haploid)	ste5 (haploid)
ste12 (haploid)	ste18 (haploid)
ste18 (haploid)	ste7 (haploid)
fus3, kss1 (haploid)	ste18 (haploid)
fus3, kss1 (haploid)	ste7 (haploid)
fus3, kss1 (haploid)	ste12 (haploid)
ste12 (haploid)	ste7 (haploid)

## Ergosterol Pathway

erg3 (haploid)	Itraconazole
erg2	Itraconazole
yer044c (haploid)	ERG11 (tet promoter)
ERG11 (tet promoter)	Itraconazole
erg3 (haploid)	ERG11 (tet promoter)
erg3 (haploid)	yer044c (haploid)
erg2	erg3 (haploid)
erg2	yer044c (haploid)
erg2	ERG11 (tet promoter)

## Vacuolar ATPase

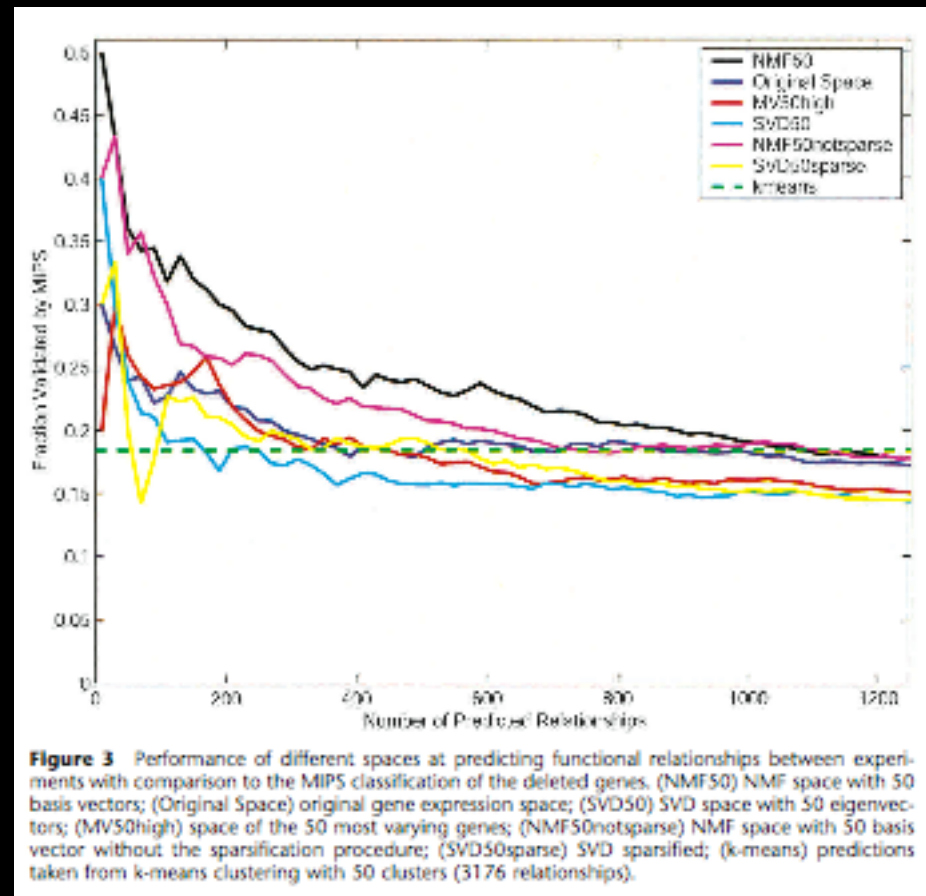
cup5	mac1
mac1	vma8
cup5	vma8

# More Functional Relationships

**Table 4.** The 42 Predictions of Functional Relationships That Could Not Be Verified on YPD From the 100 Strongest Relationships Detected

rtg1	vps8
are1, are2 (haploid)	yor015w
pex12	yea4
ckb2	yel008w
yer002w	ymr034c
mrt4	yel033w
ckb2	rts1
mrp133	ymr293c
imp2	yer050c
cbp2	pet111
cyt1	pet111
yer034w	ynd1
rps24a	ymr014w
yel001c	yor015w
ymr014w	yer006c
aep2	rml2
aep2	mrp133
ymr014w	yor078w
rml2	yer050c
mrp133	yer050c
aep2	imp2
sr1	ymr041c
ymr034c	yor015w
pfd2	yor051c
ymr025w	ymr029c
ckb2	vps8
msu1	ymr293c
sbh2	yer084w
mrp133	msu1
imp2	ymr293c
rtg1	rts1
msu1	yer050c
msu1	rml2
ymr003w	ymr034c
aep2	msu1
CDC42 (tet promoter)	KAR2 (tet promoter)
rps24a	yor078w
pfd2	yel044w
gcn4	yel008w
yer050c	ymr293c
aep2	yer050c
aep2	ymr293c

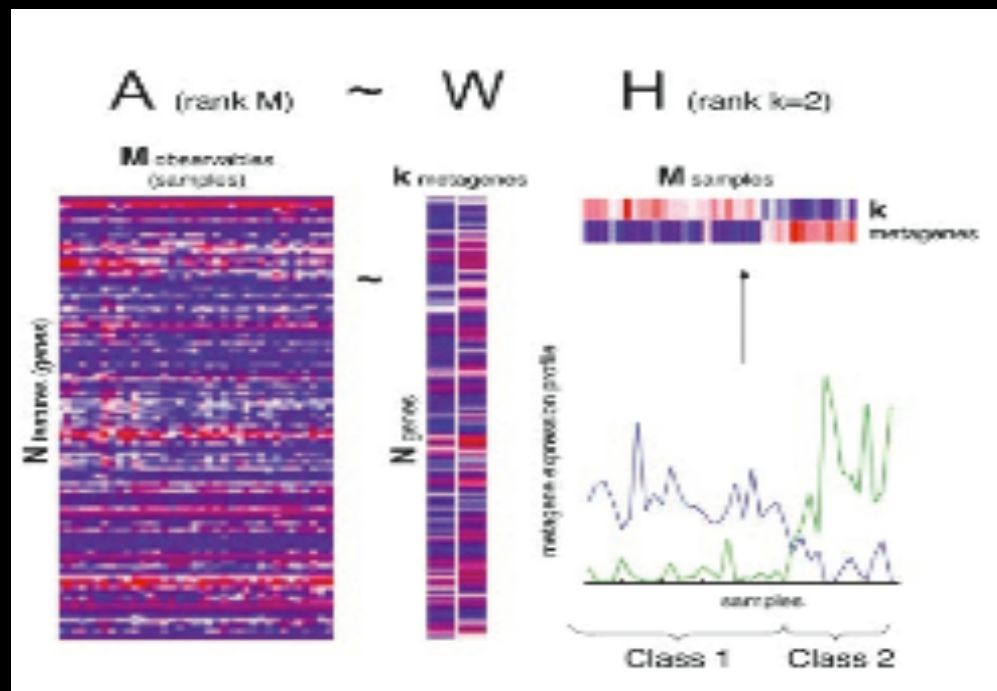
# Comparative Study



# Metagenes Study

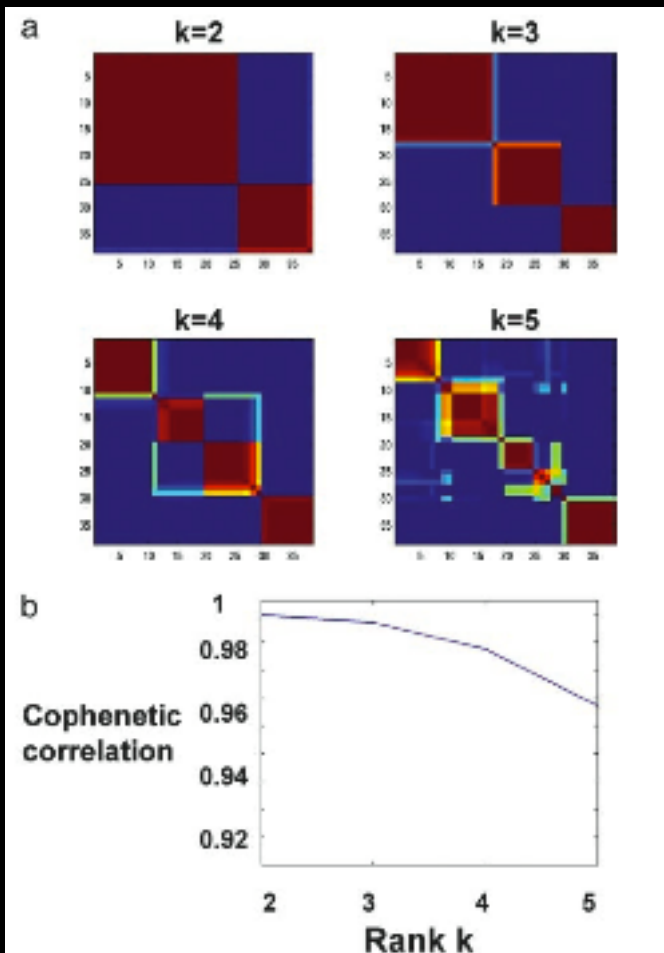
*Brunet et al 2004*

Data: gene expressions x samples





# Leukemia Samples



**Fig. 4.** (a) Reordered consensus matrices averaging 50 connectivity matrices computed at  $k = 2-5$  for the leukemia data set with the 5,000 most highly varying genes according to their coefficient of variation. Samples are hierarchically clustered by using distances derived from consensus clustering matrix entries, colored from 0 (deep blue, samples are never in the same cluster) to 1 (dark red, samples are always in the same cluster). Compositions of the leukemia clusters determined by HC of consensus matrices are as follows: for  $k = 2$ : {(25 ALL), (11 AML and 2 ALL)},  $k = 3$ : {(17 ALL-B), (8 ALL-T and 1 ALL-B), (11 AML and 1 ALL-B)},  $k = 4$ : {(11 ALL-B), (7 ALL-B and 1 AML), (8 ALL-T and 1 ALL-B), (10 AML)}, (b) Cophenetic correlation coefficients for hierarchically clustered matrices in a.

# Samples from Medulloblastoma Tumors

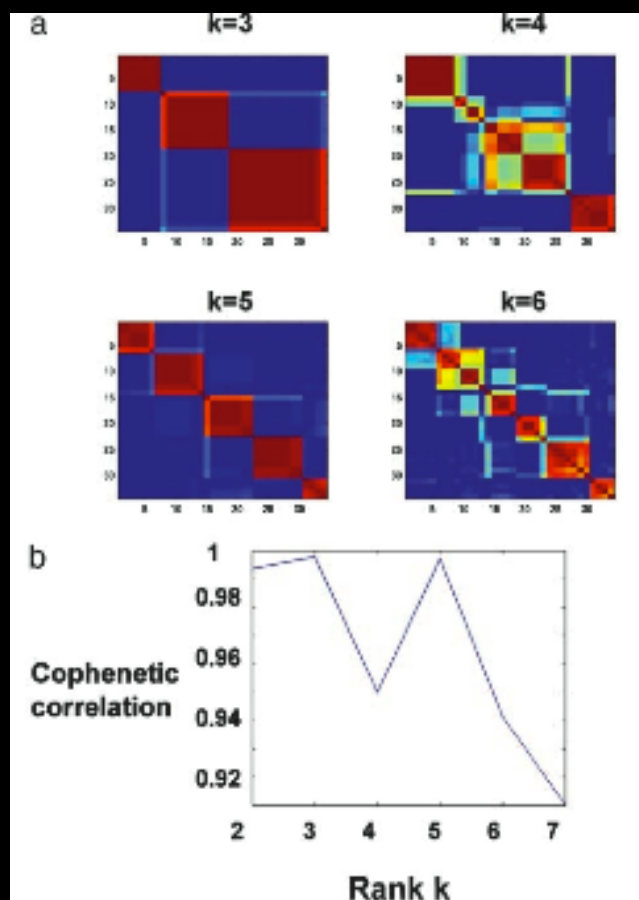
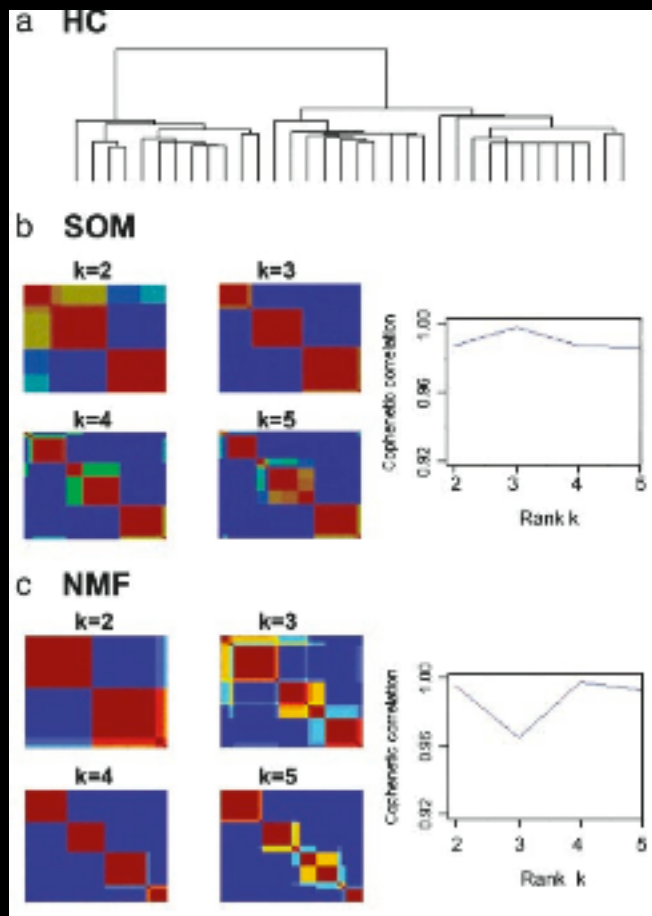


Fig. 6. (a) NMF model selection for a data set of 25 classic and 9 desmoplastic medulloblastoma tumors [ $n = 5,893$ ;  $M = 34$  (14)]. At each rank  $k$ , a consensus matrix, averaging 50 connectivity matrices, is reordered by using HC (color map as Fig. 4). In addition to a robust two-class partition (not shown), the consensus is strong for  $k = 3, 5$ , indicating reproducible partitioning of samples into two, three, and five classes but not four or six. (b) Cophenetic correlation coefficients corresponding to the HC of consensus matrices for  $k = 2-7$  shows a dip at  $k = 4$ , where reproducibility is poor, and suggests  $k = 5$  as the largest number of classes recognized by NMF for this data set.

# CNS Embryonal Tumors



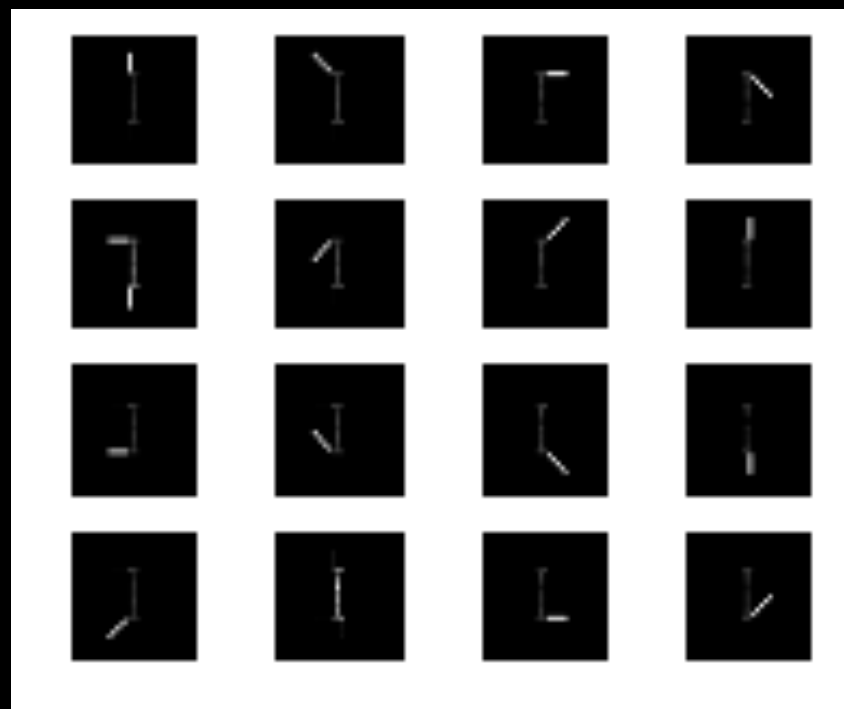
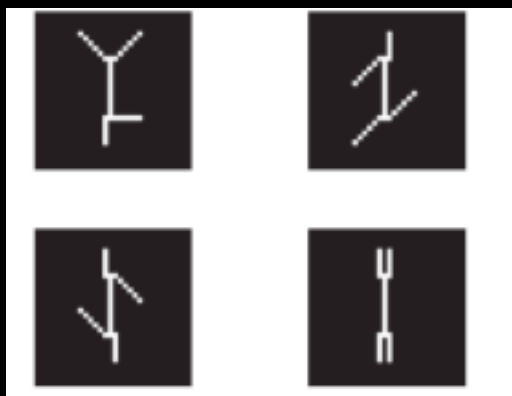
**Fig. 7.** Analysis of central nervous system embryonal tumors using 5,560 genes. The data set consists of 34 samples, including 10 classic medulloblastomas, 10 malignant gliomas, 10 rhabdoids, and 4 normals. (a) The dendrogram from HC indicates two or three major subclasses but gives no clear indication of a four-class split. (b) Reordered consensus matrices for  $k = 2-5$  centroid SOM clusterings from 20 initial conditions. Cophenetic correlation argues for a three-class decomposition. (c) Reordered consensus matrices for 20 NMF initial conditions (50 NMF iterations each), for  $k = 2-5$  (color scale same as Fig. 2). Cophenetic correlation coefficient suggests the existence of at most four robust classes.

# When Does Non-Negative Matrix Factorization Give a Correct Decomposition into Parts?

*Donoho and Stodden, 2003*

- Set of weighted generators, non-negative
- Each combination separable from other combinations
- All combinations represented in dataset

# Example of a Separable Factorial Articulation Family



# Properties of NMF

- basis vectors  $\mathbf{w}_i$  are not  $\perp \Rightarrow$  can have overlap of topics
- can restrict  $\mathbf{W}$ ,  $\mathbf{H}$  to be sparse
- $\mathbf{W}_k, \mathbf{H}_k \geq 0 \Rightarrow$  immediate interpretation (additive parts-based rep.)

**EX:** large  $w_{ij}$ 's  $\Rightarrow$  basis vector  $\mathbf{w}_i$  is mostly about terms  $j$

**EX:**  $h_{i1}$  how much  $doc_1$  is pointing in the “direction” of topic vector  $\mathbf{w}_i$

$$\mathbf{A}_k \mathbf{e}_1 = \mathbf{W}_k \mathbf{H}_{*1} = \begin{bmatrix} \vdots \\ \mathbf{w}_1 \\ \vdots \end{bmatrix} h_{11} + \begin{bmatrix} \vdots \\ \mathbf{w}_2 \\ \vdots \end{bmatrix} h_{21} + \cdots + \begin{bmatrix} \vdots \\ \mathbf{w}_k \\ \vdots \end{bmatrix} h_{k1}$$

- NMF is algorithm-dependent:  $\mathbf{W}$ ,  $\mathbf{H}$  not unique

# Report Card for SVD and NMF

Subject	SVD	NMF
low rank approximation improves performance on data mining tasks	A	A
noise reduction isolates essential components of matrix	A	A
quality of low rank approximation, $\ \mathbf{A} - \mathbf{A}_k\ $	A <sup>+</sup>	B <sup>+</sup>
uniqueness of low rank approximation	A <sup>+</sup>	B
storage of low rank factors	D	A
interpretation of vectors in low rank factors	C	A
choosing truncation point $k$	C	C
orthogonality restriction on vectors in low rank factors	B	A
updating the factors in the factorization	B	C
downdating the factors in the factorization	A	C

# Computation of NMF

(Lee and Seung 2000)

MEAN SQUARED ERROR OBJECTIVE FUNCTION

$$\begin{aligned} \min \quad & \| \mathbf{A} - \mathbf{WH} \|_F^2 \\ \text{s.t.} \quad & \mathbf{W}, \mathbf{H} \geq 0 \end{aligned}$$

## Nonlinear Optimization Problem

- convex in  $\mathbf{W}$  or  $\mathbf{H}$ , but not both  $\Rightarrow$  tough to get global min
- huge # unknowns:  $mk$  for  $\mathbf{W}$  and  $kn$  for  $\mathbf{H}$   
(EX:  $\mathbf{A}_{70K \times 10K}$  and  $k=10$  topics  $\Rightarrow$  800K unknowns)
- above objective is one of many possible



# Other Objective Functions

DIVERGENCE OBJECTIVE FUNCTION

$$\min \sum_{i,j} (\mathbf{A}_{ij} \log \frac{\mathbf{A}_{ij}}{[\mathbf{WH}]_{ij}} - \mathbf{A}_{ij} + [\mathbf{WH}]_{ij})$$

WEIGHTED MEAN SQUARED ERROR OBJECTIVE FUNCTION

$$\min \|\mathbf{B}_{\cdot} * (\mathbf{A} - \mathbf{WH})\|_F^2$$

WEIGHTED DIVERGENCE OBJECTIVE FUNCTION

$$\min \sum_{i,j} \mathbf{B}_{ij\cdot} * (\mathbf{A}_{ij} \log \frac{\mathbf{A}_{ij}}{[\mathbf{WH}]_{ij}} - \mathbf{A}_{ij} + [\mathbf{WH}]_{ij})$$

BREGMAN DIVERGENCE CLASS OF OBJECTIVE FUNCTIONS

(coming tomorrow—Inderjit Dhillon)

SUITE OF OTHER DIVERGENCE OBJECTIVE FUNCTIONS

(NMFLAB—Cichocki)

**Table 1.** Amari Alpha-NMF algorithms

<b>Amari alpha divergence:</b> $D_A^{(\alpha)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{y_{ik}^\alpha z_{ik}^{1-\alpha} - \alpha y_{ik} + (\alpha - 1)z_{ik}}{\alpha(\alpha - 1)}$	
Algorithm:	$x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^\alpha \right)^{\frac{\omega_X}{\alpha}} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^\alpha \right)^{\frac{\omega_A}{\alpha}} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$
<i>Pearson distance:</i> ( $\alpha = 2$ ): $D_A^{(\alpha=2)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{(y_{ik} - [\mathbf{A} \mathbf{X}]_{ik})^2}{[\mathbf{A} \mathbf{X}]_{ik}},$	
Algorithm:	$x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^2 \right)^{\frac{\omega_X}{2}} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^2 \right)^{\frac{\omega_A}{2}} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$
<i>Hellinger distance:</i> ( $\alpha = \frac{1}{2}$ ): $D_A^{(\alpha=0.5)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{(y_{ik} - [\mathbf{A} \mathbf{X}]_{ik})^2}{[\mathbf{A} \mathbf{X}]_{ik}},$	
Algorithm:	$x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \sqrt{\frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}}} \right)^{2\omega_X} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \sqrt{\frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}}} \right)^{2\omega_A} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$

**Table 2.** Amari Alpha-NMF algorithms (continued)

<p><i>Kullback-Leibler divergence: (<math>\alpha \rightarrow 1</math>):</i></p> $\lim_{\alpha \rightarrow 1} D_A^{(\alpha)}(y_{ik}    z_{ik}) = \sum_{ik} y_{ik} \log \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} - y_{ik} + [\mathbf{A}\mathbf{X}]_{ik},$	
<p>Algorithm:</p>	$x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\omega_X} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\omega_A} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj}$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$
<p><i>Dual Kullback-Leibler divergence: (<math>\alpha \rightarrow 0</math>):</i></p> $\lim_{\alpha \rightarrow 0} D_A^{(\alpha)}(y_{ik}    z_{ik}) = \sum_{ik} [\mathbf{A}\mathbf{X}]_{ik} \log \frac{[\mathbf{A}\mathbf{X}]_{ik}}{y_{ik}} + y_{ik} - [\mathbf{A}\mathbf{X}]_{ik}$	
<p>Algorithm:</p>	$x_{jk} \leftarrow \left( x_{jk} \prod_{i=1}^m \left( \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\omega_X a_{ij}} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \prod_{k=1}^N \left( \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\tilde{\eta}_j x_{jk}} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj}$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$

**Table 3.** Other generalized NMF algorithms

<p><i>Beta generalized divergence:</i></p> $D_K^{(\beta)}(y_{ik}  z_{ik}) = \sum_{ik} y_{ik} \frac{y_{ik}^{\beta-1} - [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1}}{\beta(\beta-1)} + [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1} \frac{[\mathbf{A}\mathbf{X}]_{ik} - y_{ik}}{\beta}$ <p>Kompass algorithm:</p> $x_{jk} \leftarrow x_{jk} \frac{\sum_{i=1}^m a_{ij} (y_{ik}/[\mathbf{A}\mathbf{X}]_{ik}^{2-\beta})}{\sum_{i=1}^m a_{ij} [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1} + \varepsilon}$ $a_{ij} \leftarrow \left( a_{ij} \frac{\sum_{k=1}^N x_{jk} (y_{ik}/[\mathbf{A}\mathbf{X}]_{ik}^{2-\beta})}{\sum_{k=1}^N x_{jk} [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1} + \varepsilon} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$
<p><i>Triangular discrimination:</i></p> $D_T^{(\beta)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{y_{ik}^\beta z_{ik}^{1-\beta} - \beta y_{ik} + (\beta-1)z_{ik}}{\beta(\beta-1)}$ <p>Algorithm:</p> $x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \left( \frac{2y_{ik}}{y_{ik} + [\mathbf{A}\mathbf{X}]_{ik}} \right)^2 \right)^{\omega_X} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \left( \frac{2y_{ik}}{y_{ik} + [\mathbf{A}\mathbf{X}]_{ik}} \right)^2 \right)^{\omega_A} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj}, \quad 0 < \omega_X < 2, \quad 0 < \omega_A < 2$
<p><i>Itakura-Saito distance:</i></p> $D_{IS}(y_{ik}  z_{ik}) = \sum_{ik} \frac{y_{ik}}{z_{ik}} - \log \left( \frac{y_{ik}}{z_{ik}} \right) - 1$ <p>Algorithm:</p> $\mathbf{X} \leftarrow \mathbf{X} \odot [(\mathbf{A}^T \mathbf{P}) \oslash (\mathbf{A}^T \mathbf{Q} + \varepsilon)].^\beta$ $\mathbf{A} \leftarrow \mathbf{A} \odot [(\mathbf{P}\mathbf{X}^T) \oslash (\mathbf{Q}\mathbf{X}^T + \varepsilon)].^\beta$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj}, \quad \beta = [0.5, 1]$ $\mathbf{P} = \mathbf{Y} \oslash (\mathbf{A}\mathbf{X} + \varepsilon).^2, \quad \mathbf{Q} = \mathbf{1} \oslash (\mathbf{A}\mathbf{X} + \varepsilon)$

**Table 4.** Generalized SMART NMF adaptive algorithms and corresponding loss functions - part I.

Generalized SMART algorithms	
$a_{ij} \leftarrow a_{ij} \exp \left( \sum_{k=1}^N \tilde{\eta}_j x_{jk} \rho(y_{ik}, z_{ik}) \right), \quad x_{jk} \leftarrow x_{jk} \exp \left( \sum_{i=1}^m \eta_j a_{ij} \rho(y_{ik}, z_{ik}) \right),$ $a_j = \sum_{i=1}^m a_{ij} = 1, \quad \forall j, \quad a_{ij} \geq 0, \quad y_{ik} > 0, \quad z_{ik} = [\mathbf{AX}]_{ik} > 0, \quad x_{jk} \geq 0$	
Divergence: $D(\mathbf{Y} \parallel \mathbf{AX})$	Error function: $\rho(y_{ik}, z_{ik})$
Dual Kullback-Leibler I-divergence: $D_{KL2}(\mathbf{AX} \parallel \mathbf{Y})$	
$\sum_{ik} \left( z_{ik} \ln \frac{z_{ik}}{y_{ik}} + y_{ik} - z_{ik} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{y_{ik}}{z_{ik}} \right),$
Relative Arithmetic-Geometric divergence: $D_{RAG}(\mathbf{Y} \parallel \mathbf{AX})$	
$\sum_{ik} \left( (y_{ik} + z_{ik}) \ln \left( \frac{y_{ik} + z_{ik}}{2y_{ik}} \right) + y_{ik} - z_{ik} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right),$
Symmetric Arithmetic-Geometric divergence: $D_{SAG}(\mathbf{Y} \parallel \mathbf{AX})$	
$2 \sum_{ik} \left( \frac{y_{ik} + z_{ik}}{2} \ln \left( \frac{y_{ik} + z_{ik}}{2\sqrt{y_{ik}z_{ik}}} \right) \right),$	$\rho(y_{ik}, z_{ik}) = \frac{y_{ik} - z_{ik}}{2z_{ik}} + \ln \left( \frac{2\sqrt{y_{ik}z_{ik}}}{y_{ik} + z_{ik}} \right),$
J-divergence: $D_J(\mathbf{Y} \parallel \mathbf{AX})$	
$\sum_{ik} \left( \frac{y_{ik} - z_{ik}}{2} \ln \left( \frac{y_{ik}}{z_{ik}} \right) \right),$	$\rho(y_{ik}, z_{ik}) = \frac{1}{2} \ln \left( \frac{y_{ik}}{z_{ik}} \right) + \frac{y_{ik} - z_{ik}}{2z_{ik}},$

**Table 5.** Generalized SMART NMF adaptive algorithms and corresponding loss functions - part II.

Relative Jensen-Shannon divergence: $D_{RJS}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} \left( 2y_{ik} \ln \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right) + z_{ik} - y_{ik} \right),$	$\rho(y_{ik}, z_{ik}) = \frac{y_{ik} - z_{ik}}{2z_{ik}} + \ln \left( \frac{2\sqrt{y_{ik}z_{ik}}}{y_{ik} + z_{ik}} \right),$
Dual Jensen-Shannon divergence: $D_{DJS}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} y_{ik} \ln \left( \frac{2z_{ik}}{z_{ik} + y_{ik}} \right) + y_{ik} \ln \left( \frac{2y_{ik}}{z_{ik} + y_{ik}} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{z_{ik} + y_{ik}}{2y_{ik}} \right),$
Symmetric Jensen-Shannon divergence: $D_{SJS}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} y_{ik} \ln \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right) + z_{ik} \ln \left( \frac{2z_{ik}}{y_{ik} + z_{ik}} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{y_{ik} + z_{ik}}{2z_{ik}} \right),$
Triangular discrimination: $D_T(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} \left\{ \frac{(y_{ik} - z_{ik})^2}{y_{ik} + z_{ik}} \right\},$	$\rho(y_{ik}, z_{ik}) = \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right)^2 - 1,$
Bose-Einstein divergence: $D_{BE}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} y_{ik} \ln \left( \frac{(1 + \alpha)y_{ik}}{y_{ik} + \alpha z_{ik}} \right) + \alpha z_{ik} \ln \left( \frac{(1 + \alpha)z_{ik}}{y_{ik} + \alpha z_{ik}} \right),$	$\rho(y_{ik}, z_{ik}) = \alpha \ln \left( \frac{y_{ik} + \alpha z_{ik}}{(1 + \alpha)z_{ik}} \right),$

# Early NMF Algorithms

- Alternating Least Squares
  - Paatero 1994
  - ALS algorithms that incorporate sparsity
- Multiplicative update rules
  - Lee-Seung 2000
  - Hoyer 2002
- Gradient Descent
  - Hoyer 2004
  - Berry-Plemmons 2004

# PMF Algorithm: Paatero & Tapper 1994

MEAN SQUARED ERROR—ALTERNATING LEAST SQUARES

$$\begin{aligned} \min \quad & \| \mathbf{A} - \mathbf{WH} \|_F^2 \\ \text{s.t.} \quad & \mathbf{W}, \mathbf{H} \geq 0 \end{aligned}$$

---

$\mathbf{W} = \text{abs}(\text{randn}(m,k));$

for  $i = 1 : \text{maxiter}$

$\text{LS}$  for  $j = 1 : n = \#docs$ , solve

$$\begin{aligned} \min_{\mathbf{H}_{*j}} \quad & \| \mathbf{A}_{*j} - \mathbf{WH}_{*j} \|_2^2 \\ \text{s.t.} \quad & \mathbf{H}_{*j} \geq 0 \end{aligned}$$

$\text{LS}$  for  $j = 1 : m = \#terms$ , solve

$$\begin{aligned} \min_{\mathbf{W}_{j*}} \quad & \| \mathbf{A}_{j*} - \mathbf{W}_{j*} \mathbf{H} \|_2^2 \\ \text{s.t.} \quad & \mathbf{W}_{j*} \geq 0 \end{aligned}$$

end

---



# ALS Algorithm

---

**W** = abs(randn(m,k));

for i = 1 : maxiter

LS      solve matrix equation  $\mathbf{W}^T \mathbf{W} \mathbf{H} = \mathbf{W}^T \mathbf{A}$  for **H**

NONNEG **H** = **H**. \* (**H** >= 0)

LS      solve matrix equation  $\mathbf{H} \mathbf{H}^T \mathbf{W}^T = \mathbf{H} \mathbf{A}^T$  for **W**

NONNEG **W** = **W**. \* (**W** >= 0)

end

---

# ALS Summary

## Pros

- + fast
- + works well in practice
- + speedy convergence
- + only need to initialize  $\mathbf{W}^{(0)}$
- + 0 elements not *locked*

## Cons

- no sparsity of  $\mathbf{W}$  and  $\mathbf{H}$  incorporated into mathematical setup
- ad hoc nonnegativity: negative elements are set to 0
- ad hoc sparsity: negative elements are set to 0
- no convergence theory

# Early NMF Algorithms

- Alternating Least Squares
  - Paatero 1994
  - ALS algorithms that incorporate sparsity
- Multiplicative update rules
  - Lee-Seung 2000
  - Hoyer 2002
- Gradient Descent
  - Hoyer 2004
  - Berry-Plemmons 2004

# NMF Algorithm: Lee and Seung 2000

MEAN SQUARED ERROR OBJECTIVE FUNCTION

$$\begin{aligned} \min \quad & \| \mathbf{A} - \mathbf{WH} \|_F^2 \\ \text{s.t.} \quad & \mathbf{W}, \mathbf{H} \geq 0 \end{aligned}$$

---

```
W = abs(randn(m,k));  
H = abs(randn(k,n));  
for i = 1 : maxiter  
    H = H .* (WTA) ./ (WTWH + 10-9);  
    W = W .* (AHT) ./ (WHHT + 10-9);  
end
```

---

Many parameters affect performance (k, obj. function, sparsity constraints, algorithm, etc.).

— NMF is not unique!

(proof of convergence to fixed point based on E-M convergence proof)

# NMF Algorithm: Lee and Seung 2000

DIVERGENCE OBJECTIVE FUNCTION

$$\min \sum_{i,j} (\mathbf{A}_{ij} \log \frac{\mathbf{A}_{ij}}{[\mathbf{WH}]_{ij}} - \mathbf{A}_{ij} + [\mathbf{WH}]_{ij})$$

*s.t.*    $\mathbf{W}, \mathbf{H} \geq 0$

---

$\mathbf{W} = \text{abs}(\text{randn}(m,k));$

$\mathbf{H} = \text{abs}(\text{randn}(k,n));$

for  $i = 1 : \text{maxiter}$

$\mathbf{H} = \mathbf{H} .* (\mathbf{W}^T (\mathbf{A} ./ (\mathbf{WH} + 10^{-9}))) ./ \mathbf{W}^T \mathbf{ee}^T;$

$\mathbf{W} = \mathbf{W} .* ((\mathbf{A} ./ (\mathbf{WH} + 10^{-9})) \mathbf{H}^T) ./ \mathbf{ee}^T \mathbf{H}^T;$

end

---

(proof of convergence to fixed point based on E-M convergence proof)

(objective function tails off after 50-100 iterations)

# Multiplicative Update Summary

## Pros

- + convergence theory: guaranteed to converge to fixed point
- + good initialization  $\mathbf{W}^{(0)}, \mathbf{H}^{(0)}$  speeds convergence and gets to better fixed point

## Cons

- fixed point may be local min or saddle point
- good initialization  $\mathbf{W}^{(0)}, \mathbf{H}^{(0)}$  speeds convergence and gets to better fixed point
- slow: many M-M multiplications at each iteration
- hundreds/thousands of iterations until convergence
- no sparsity of  $\mathbf{W}$  and  $\mathbf{H}$  incorporated into mathematical setup
- 0 elements *locked*

**Table 1.** Amari Alpha-NMF algorithms

<b>Amari alpha divergence:</b> $D_A^{(\alpha)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{y_{ik}^\alpha z_{ik}^{1-\alpha} - \alpha y_{ik} + (\alpha - 1)z_{ik}}{\alpha(\alpha - 1)}$	
Algorithm:	$x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^\alpha \right)^{\frac{\omega_X}{\alpha}} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^\alpha \right)^{\frac{\omega_A}{\alpha}} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$
<i>Pearson distance:</i> ( $\alpha = 2$ ): $D_A^{(\alpha=2)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{(y_{ik} - [\mathbf{A} \mathbf{X}]_{ik})^2}{[\mathbf{A} \mathbf{X}]_{ik}},$	
Algorithm:	$x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^2 \right)^{\frac{\omega_X}{2}} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \left( \frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}} \right)^2 \right)^{\frac{\omega_A}{2}} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$
<i>Hellinger distance:</i> ( $\alpha = \frac{1}{2}$ ): $D_A^{(\alpha=0.5)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{(y_{ik} - [\mathbf{A} \mathbf{X}]_{ik})^2}{[\mathbf{A} \mathbf{X}]_{ik}},$	
Algorithm:	$x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \sqrt{\frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}}} \right)^{2\omega_X} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \sqrt{\frac{y_{ik}}{[\mathbf{A} \mathbf{X}]_{ik}}} \right)^{2\omega_A} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$ $0 < \omega_X < 2, \quad 0 < \omega_A < 2$

**Table 2.** Amari Alpha-NMF algorithms (continued)

*Kullback-Leibler divergence: ( $\alpha \rightarrow 1$ ):*

$$\lim_{\alpha \rightarrow 1} D_A^{(\alpha)}(y_{ik} || z_{ik}) = \sum_{ik} y_{ik} \log \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} - y_{ik} + [\mathbf{A}\mathbf{X}]_{ik},$$

Algorithm:

$$\begin{aligned} x_{jk} &\leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\omega_X} \right)^{1+\alpha_{sX}} \\ a_{ij} &\leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\omega_A} \right)^{1+\alpha_{sA}} \\ a_{ij} &\leftarrow a_{ij} / \sum_p a_{pj} \\ 0 &< \omega_X < 2, \quad 0 < \omega_A < 2 \end{aligned}$$

*Dual Kullback-Leibler divergence: ( $\alpha \rightarrow 0$ ):*

$$\lim_{\alpha \rightarrow 0} D_A^{(\alpha)}(y_{ik} || z_{ik}) = \sum_{ik} [\mathbf{A}\mathbf{X}]_{ik} \log \frac{[\mathbf{A}\mathbf{X}]_{ik}}{y_{ik}} + y_{ik} - [\mathbf{A}\mathbf{X}]_{ik}$$

Algorithm:

$$\begin{aligned} x_{jk} &\leftarrow \left( x_{jk} \prod_{i=1}^m \left( \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\omega_X a_{ij}} \right)^{1+\alpha_{sX}} \\ a_{ij} &\leftarrow \left( a_{ij} \prod_{k=1}^N \left( \frac{y_{ik}}{[\mathbf{A}\mathbf{X}]_{ik}} \right)^{\tilde{\eta}_j x_{jk}} \right)^{1+\alpha_{sA}} \\ a_{ij} &\leftarrow a_{ij} / \sum_p a_{pj} \\ 0 &< \omega_X < 2, \quad 0 < \omega_A < 2 \end{aligned}$$



**Table 3.** Other generalized NMF algorithms

<p><i>Beta generalized divergence:</i></p> $D_K^{(\beta)}(y_{ik}  z_{ik}) = \sum_{ik} y_{ik} \frac{y_{ik}^{\beta-1} - [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1}}{\beta(\beta-1)} + [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1} \frac{[\mathbf{A}\mathbf{X}]_{ik} - y_{ik}}{\beta}$ <p>Kompass algorithm:</p> $x_{jk} \leftarrow x_{jk} \frac{\sum_{i=1}^m a_{ij} (y_{ik}/[\mathbf{A}\mathbf{X}]_{ik}^{2-\beta})}{\sum_{i=1}^m a_{ij} [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1} + \varepsilon}$ $a_{ij} \leftarrow \left( a_{ij} \frac{\sum_{k=1}^N x_{jk} (y_{ik}/[\mathbf{A}\mathbf{X}]_{ik}^{2-\beta})}{\sum_{k=1}^N x_{jk} [\mathbf{A}\mathbf{X}]_{ik}^{\beta-1} + \varepsilon} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj},$
<p><i>Triangular discrimination:</i></p> $D_T^{(\beta)}(y_{ik}  z_{ik}) = \sum_{ik} \frac{y_{ik}^\beta z_{ik}^{1-\beta} - \beta y_{ik} + (\beta-1)z_{ik}}{\beta(\beta-1)}$ <p>Algorithm:</p> $x_{jk} \leftarrow \left( x_{jk} \left( \sum_{i=1}^m a_{ij} \left( \frac{2y_{ik}}{y_{ik} + [\mathbf{A}\mathbf{X}]_{ik}} \right)^2 \right)^{\omega_X} \right)^{1+\alpha_{sX}}$ $a_{ij} \leftarrow \left( a_{ij} \left( \sum_{k=1}^N x_{jk} \left( \frac{2y_{ik}}{y_{ik} + [\mathbf{A}\mathbf{X}]_{ik}} \right)^2 \right)^{\omega_A} \right)^{1+\alpha_{sA}}$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj}, \quad 0 < \omega_X < 2, \quad 0 < \omega_A < 2$
<p><i>Itakura-Saito distance:</i></p> $D_{IS}(y_{ik}  z_{ik}) = \sum_{ik} \frac{y_{ik}}{z_{ik}} - \log \left( \frac{y_{ik}}{z_{ik}} \right) - 1$ <p>Algorithm:</p> $\mathbf{X} \leftarrow \mathbf{X} \odot [(\mathbf{A}^T \mathbf{P}) \oslash (\mathbf{A}^T \mathbf{Q} + \varepsilon)]^\beta$ $\mathbf{A} \leftarrow \mathbf{A} \odot [(\mathbf{P}\mathbf{X}^T) \oslash (\mathbf{Q}\mathbf{X}^T + \varepsilon)]^\beta$ $a_{ij} \leftarrow a_{ij} / \sum_p a_{pj}, \quad \beta = [0.5, 1]$ $\mathbf{P} = \mathbf{Y} \oslash (\mathbf{A}\mathbf{X} + \varepsilon)^2, \quad \mathbf{Q} = \mathbf{1} \oslash (\mathbf{A}\mathbf{X} + \varepsilon)$

**Table 4.** Generalized SMART NMF adaptive algorithms and corresponding loss functions - part I.

Generalized SMART algorithms	
$a_{ij} \leftarrow a_{ij} \exp \left( \sum_{k=1}^N \tilde{\eta}_j x_{jk} \rho(y_{ik}, z_{ik}) \right), \quad x_{jk} \leftarrow x_{jk} \exp \left( \sum_{i=1}^m \eta_j a_{ij} \rho(y_{ik}, z_{ik}) \right),$ $a_j = \sum_{i=1}^m a_{ij} = 1, \quad \forall j, \quad a_{ij} \geq 0, \quad y_{ik} > 0, \quad z_{ik} = [\mathbf{AX}]_{ik} > 0, \quad x_{jk} \geq 0$	
Divergence: $D(\mathbf{Y} \parallel \mathbf{AX})$	Error function: $\rho(y_{ik}, z_{ik})$
Dual Kullback-Leibler I-divergence: $D_{KL2}(\mathbf{AX} \parallel \mathbf{Y})$	
$\sum_{ik} \left( z_{ik} \ln \frac{z_{ik}}{y_{ik}} + y_{ik} - z_{ik} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{y_{ik}}{z_{ik}} \right),$
Relative Arithmetic-Geometric divergence: $D_{RAG}(\mathbf{Y} \parallel \mathbf{AX})$	
$\sum_{ik} \left( (y_{ik} + z_{ik}) \ln \left( \frac{y_{ik} + z_{ik}}{2y_{ik}} \right) + y_{ik} - z_{ik} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right),$
Symmetric Arithmetic-Geometric divergence: $D_{SAG}(\mathbf{Y} \parallel \mathbf{AX})$	
$2 \sum_{ik} \left( \frac{y_{ik} + z_{ik}}{2} \ln \left( \frac{y_{ik} + z_{ik}}{2\sqrt{y_{ik}z_{ik}}} \right) \right),$	$\rho(y_{ik}, z_{ik}) = \frac{y_{ik} - z_{ik}}{2z_{ik}} + \ln \left( \frac{2\sqrt{y_{ik}z_{ik}}}{y_{ik} + z_{ik}} \right),$
J-divergence: $D_J(\mathbf{Y} \parallel \mathbf{AX})$	
$\sum_{ik} \left( \frac{y_{ik} - z_{ik}}{2} \ln \left( \frac{y_{ik}}{z_{ik}} \right) \right),$	$\rho(y_{ik}, z_{ik}) = \frac{1}{2} \ln \left( \frac{y_{ik}}{z_{ik}} \right) + \frac{y_{ik} - z_{ik}}{2z_{ik}},$

**Table 5.** Generalized SMART NMF adaptive algorithms and corresponding loss functions - part II.

Relative Jensen-Shannon divergence: $D_{RJS}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} \left( 2y_{ik} \ln \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right) + z_{ik} - y_{ik} \right),$	$\rho(y_{ik}, z_{ik}) = \frac{y_{ik} - z_{ik}}{2z_{ik}} + \ln \left( \frac{2\sqrt{y_{ik}z_{ik}}}{y_{ik} + z_{ik}} \right),$
Dual Jensen-Shannon divergence: $D_{DJS}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} y_{ik} \ln \left( \frac{2z_{ik}}{z_{ik} + y_{ik}} \right) + y_{ik} \ln \left( \frac{2y_{ik}}{z_{ik} + y_{ik}} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{z_{ik} + y_{ik}}{2y_{ik}} \right),$
Symmetric Jensen-Shannon divergence: $D_{SJS}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} y_{ik} \ln \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right) + z_{ik} \ln \left( \frac{2z_{ik}}{y_{ik} + z_{ik}} \right),$	$\rho(y_{ik}, z_{ik}) = \ln \left( \frac{y_{ik} + z_{ik}}{2z_{ik}} \right),$
Triangular discrimination: $D_T(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} \left\{ \frac{(y_{ik} - z_{ik})^2}{y_{ik} + z_{ik}} \right\},$	$\rho(y_{ik}, z_{ik}) = \left( \frac{2y_{ik}}{y_{ik} + z_{ik}} \right)^2 - 1,$
Bose-Einstein divergence: $D_{BE}(\mathbf{Y}  \mathbf{AX})$	
$\sum_{ik} y_{ik} \ln \left( \frac{(1 + \alpha)y_{ik}}{y_{ik} + \alpha z_{ik}} \right) + \alpha z_{ik} \ln \left( \frac{(1 + \alpha)z_{ik}}{y_{ik} + \alpha z_{ik}} \right),$	$\rho(y_{ik}, z_{ik}) = \alpha \ln \left( \frac{y_{ik} + \alpha z_{ik}}{(1 + \alpha)z_{ik}} \right),$

# Early NMF Algorithms

- Alternating Least Squares
  - Paatero 1994
  - ALS algorithms that incorporate sparsity
- Multiplicative update rules
  - Lee-Seung 2000
  - Hoyer 2002
- Gradient Descent
  - Hoyer 2004
  - Berry-Plemmons 2004

# Sparsity Measures

- Berry et al.  $\|\mathbf{x}\|_2^2$
- Hoyer  $spar(\mathbf{x}_{n \times 1}) = \frac{\sqrt{n} - \|\mathbf{x}\|_1 / \|\mathbf{x}\|_2}{\sqrt{n} - 1}$
- Diversity measure  $E^{(p)}(\mathbf{x}) = \sum_{i=1}^n |x_i|^p, \quad 0 \leq p \leq 1$   
 $E^{(p)}(\mathbf{x}) = - \sum_{i=1}^n |x_i|^p, \quad p < 0$

Rao and Kreutz-Delgado: algorithms for minimizing  $E^{(p)}(\mathbf{x})$   
s.t.  $\mathbf{Ax} = \mathbf{b}$ , but expensive iterative procedure

- Ideal  $nnz(\mathbf{x})$  not continuous, NP-hard to use this in optim.

# NMF Algorithm: Berry et al. 2004

GRADIENT DESCENT-CONSTRAINED LEAST SQUARES

---

**W** = abs(randn(m,k)); (scale cols of **W** to unit norm)

**H** = zeros(k,n);

for i = 1 : maxiter

    CLS for j = 1 : #docs, solve

$$\min_{\mathbf{H}_{*j}} \|\mathbf{A}_{*j} - \mathbf{W}\mathbf{H}_{*j}\|_2^2 + \lambda \|\mathbf{H}_{*j}\|_2^2$$

$$\text{s.t. } \mathbf{H}_{*j} \geq 0$$

    GD **W** = **W** .\* (**AH**<sup>T</sup>) ./ (**WHH**<sup>T</sup> + 10<sup>-9</sup>); (scale cols of **W**)

end

---

# NMF Algorithm: Berry et al. 2004

GRADIENT DESCENT-CONSTRAINED LEAST SQUARES

---

```
W = abs(randn(m,k));                                (scale cols of W to unit norm)
H = zeros(k,n);
for i = 1 : maxiter
    CLS for j = 1 : #docs, solve
        
$$\min_{\mathbf{H}_{*j}} \|\mathbf{A}_{*j} - \mathbf{W}\mathbf{H}_{*j}\|_2^2 + \lambda \|\mathbf{H}_{*j}\|_2^2$$

        s.t.  $\mathbf{H}_{*j} \geq 0$ 
        solve for H:  $(\mathbf{W}^T \mathbf{W} + \lambda \mathbf{I}) \mathbf{H} = \mathbf{W}^T \mathbf{A};$  (small matrix solve)
    GD W = W .* (AHT) ./ (WHHT + 10-9);      (scale cols of W)
end
```

---

(objective function tails off after 15-30 iterations)

# Berry et al. 2004 Summary

## Pros

- + fast: less work per iteration than most other NMF algorithms
- + fast: small # of iterations until convergence
- + sparsity parameter for  $\mathbf{H}$

## Cons

- 0 elements in  $\mathbf{W}$  are *locked*
- no sparsity parameter for  $\mathbf{W}$
- ad hoc nonnegativity: negative elements in  $\mathbf{H}$  are set to 0, could run `lsqnonneg` or `snnls` instead
- no convergence theory