

Text Mining Approaches for Email Surveillance

Massive Data Sets Workshop, Stanford/Yahoo!

Michael W. Berry and Murray Browne

Department of Computer Science, UTK

June 22, 2006

Collaborators

- ▶ Pau'l Pauca, Bob Plemmons (Wake Forest)
- ▶ Amy Langville (College of Charleston)
- ▶ David Skillicorn, Parambir Keila (Queens U.)
- ▶ Stratis Gallopoulos, Ioannis Antonellis (U. Patras)

Enron Background

Non-Negative Matrix Factorization (NMF)

Electronic Mail Surveillance

Surveillance Tool Prototype

Conclusions and References

Email Collection

- ▶ By-product of the FERC investigation of Enron (originally contained 15 million email messages).
- ▶ This study used the improved corpus known as the Enron Email set, which was edited by Dr. William Cohen at CMU.
- ▶ This set had over 500,000 email messages. The majority were sent in the 1999 to 2001 timeframe.

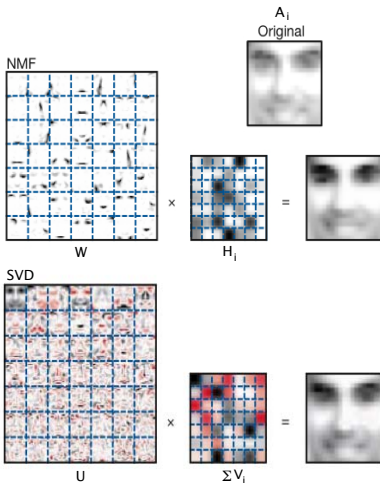
Enron Historical 1999-2001

- ▶ Ongoing, problematic, development of the Dabhol Power Company (DPC) in the Indian state of Maharashtra.
- ▶ Deregulation of the Calif. energy industry, which led to rolling electricity blackouts in the summer of 2000 (and subsequent investigations).
- ▶ Revelation of Enron's deceptive business and accounting practices that led to an abrupt collapse of the energy colossus in October, 2001; Enron filed for bankruptcy in December, 2001.

NMF Origins

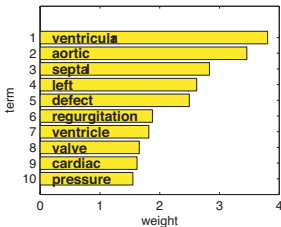
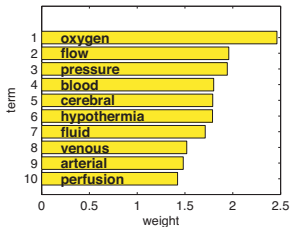
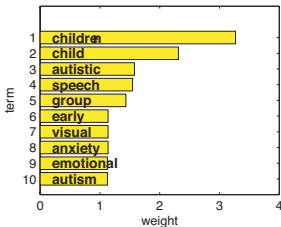
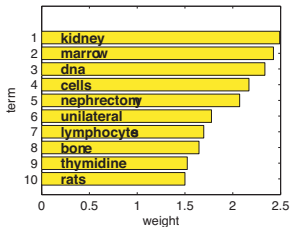
- ▶ NMF (Nonnegative Matrix Factorization) can be used to approximate high-dimensional data having nonnegative components.
- ▶ Lee and Seung (1999) demonstrated its use as a *sum-by-parts* representation of image data in order to both identify and classify image *features*.
- ▶ [Xu et al., 2003] demonstrated how NMF-based indexing could outperform SVD-based Latent Semantic Indexing (LSI) for some information retrieval tasks.

NMF for Image Processing



Sparse NMF verses Dense SVD Bases; Lee and Seung (1999)

NMF for Text Mining (Medlars)

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 

Interpretable NMF feature vectors; Langville et al. (2006)

Derivation

- ▶ Given an $m \times n$ term-by-message (sparse) matrix X .
- ▶ Compute two reduced-dim. matrices W, H so that $X \simeq WH$; W is $m \times r$ and H is $r \times n$, with $r \ll n$.
- ▶ **Optimization problem:**

$$\min_{W, H} \|X - WH\|_F^2,$$

subject to $W_{ij} \geq 0$ and $H_{ij} \geq 0$, $\forall i, j$.

- ▶ **General approach:** construct initial estimates for W and H and then improve them via alternating iterations.

Multiplicative Method (MM)

- ▶ Multiplicative update rules for W and H (Lee and Seung, 1999):
 1. Initialize W and H with non-negative values, and scale the columns of W to unit norm.
 2. Iterate for each c, j , and i until convergence or after k iterations:
 - 2.1 $H_{cj} \leftarrow H_{cj} \frac{(W^T X)_{cj}}{(W^T WH)_{cj} + \epsilon}$
 - 2.2 $W_{ic} \leftarrow W_{ic} \frac{(XH^T)_{ic}}{(WHH^T)_{ic} + \epsilon}$
 - 2.3 Scale the columns of W to unit norm.
- ▶ Setting $\epsilon = 10^{-9}$ will suffice [Shahnaz et al., 2006].

Normalization, Complexity, and Convergence

- ▶ Important to normalize X initially and the basis matrix W at each iteration.
- ▶ When optimizing on a unit hypersphere, the column (or feature) vectors of W , denoted by W_k , are effectively mapped to the surface of the hypersphere by repeated normalization.
- ▶ MM implementation of NMF requires $\mathcal{O}(rmn)$ operations per iteration; Lee and Seung (1999) proved that $\|X - WH\|_F^2$ is monotonically non-increasing with MM.
- ▶ From a nonlinear optimization perspective, MM/NMF can be considered a **diagonally-scaled gradient descent method**.

Hoyer's Method

- ▶ From neural network applications, Hoyer (2002) enforced statistical sparsity for the weight matrix H in order to enhance the parts-based data representations in the matrix W .
- ▶ Mu et al. (2003) suggested a regularization approach to achieve statistical sparsity in the matrix H : **point count regularization**; penalize the *number* of nonzeros in H rather than $\sum_{ij} H_{ij}$.
- ▶ Goal of increased sparsity – better representation of *parts* or *features* spanned by the corpus (X) [Shahnaz et al., 2006].

GD-CLS – Hybrid Approach

- ▶ First use MM to compute an approximation to W for each iteration – a gradient descent (**GD**) optimization step.
- ▶ Then, compute the weight matrix H using a constrained least squares (**CLS**) model to penalize non-smoothness (i.e., non-sparsity) in H – common Tikhonov regularization technique used in image processing (Prasad et al., 2003).
- ▶ Convergence to a non-stationary point evidenced (but no formal proof given to date).

GD-CLS Algorithm

1. Initialize W and H with non-negative values, and scale the columns of W to unit norm.
2. Iterate until convergence or after k iterations:
 - 2.1 $W_{ic} \leftarrow W_{ic} \frac{(XH^T)_{ic}}{(WHH^T)_{ic} + \epsilon}$, for c and i
 - 2.2 Rescale the columns of W to unit norm.
 - 2.3 Solve the constrained least squares problem:

$$\min_{H_j} \{ \|X_j - WH_j\|_2^2 + \lambda \|H_j\|_2^2 \},$$

where the subscript j denotes the j^{th} column, for $j = 1, \dots, m$.

- ▶ Any negative values in H_j are set to zero. The parameter λ is a regularization value that is used to balance the reduction of the metric $\|X_j - WH_j\|_2^2$ with enforcement of smoothness and sparsity in H [Shahnaz et al., 2006].

INBOX Collection

- ▶ Parsed *inbox* folder of all 150 accounts (users) via **GTP** (General Text Parser); 495-term stoplist used and extracted terms must appear in more than 1 email and more than once globally.

PRIVATE Collection

- ▶ Parsed all mail directories (of all 150 accounts) with the exception of `all_documents`, `calendar`, `contacts`, `deleted_items`, `discussion_threads`, `inbox`, `notes_inbox`, `sent`, `sent_items`, and `_sent_mail`; 495-term stoplist used and extracted terms must appear in more than 1 email and more than once globally.
- ▶ Distribution of messages sent in the year 2001:

Month	Msgs	Terms	Month	Msgs	Terms
Jan	3,621	17,888	Jul	3,077	17,617
Feb	2,804	16,958	Aug	2,828	16,417
Mar	3,525	20,305	Sep	2,330	15,405
Apr	4,273	24,010	Oct	2,821	20,995
May	4,261	24,335	Nov	2,204	18,693
Jun	4,324	18,599	Dec	1,489	8,097

Term Weighting Schemes

- ▶ For $m \times n$ term-by-message matrix $X = [x_{ij}]$, define

$$x_{ij} = l_{ij} g_i d_j,$$

where l_{ij} is the local weight for term i occurring in message j , g_i is the global weight for term i in the subcollection, and d_j is a document normalization factor (set $d_j = 1$).

- ▶ Schemes used in parsing INBOX and PRIVATE subcollections:

Name	Local	Global
txx	Term Frequency $l_{ij} = f_{ij}$	None $g_i = 1$
lex	Logarithmic $l_{ij} = \log(1 + f_{ij})$	Entropy (Define: $p_{ij} = f_{ij} / \sum_j f_{ij}$) $g_i = 1 + (\sum_j p_{ij} \log(p_{ij})) / \log n$

Computational Complexity

- Rank-50 NMF ($X \simeq WH$) computed on a 450MHz (Dual) UltraSPARC-II processor using 100 iterations:

Collection	Mail Messages	Dictionary Terms	λ	Time (sec.)
INBOX	44,872	80,683	0.1	1,471
			0.01	1,451
			0.001	1,521
PRIVATE	65,031	92,133	0.1	51,489
			0.01	51,393
			0.001	51,562

PRIVATE with Log-Entropy Weighting

- Identify rows of H from $X \simeq WH$ or H^k with $\lambda = 0.1$; $r = 50$ feature vectors (W_k) generated by GD-CLS:

Feature Index (k)	Cluster Size	Topic Description	Dominant Terms
10	497	California	ca, cpuc , gov , socalgas , sempra, org, sce, gmssr, aelaw, ci
23	43	Louise Kitchen named top woman by Fortune	evp, fortune , britain, woman, ceo , avon, fiorina, cfo, hewlett, packard
26	231	Fantasy football	game, wr, qb, play, rb, season, injury, updated, fantasy, image

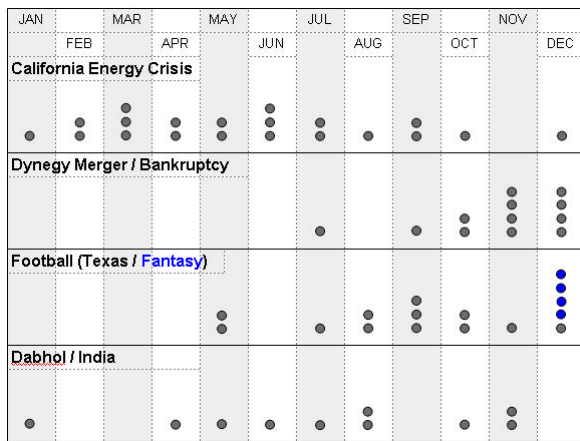
(Cluster size \equiv no. of H^k elements $>$ $row_{max}/10$)

PRIVATE with Log-Entropy Weighting

- ▶ Additional topic clusters of significant size:

Feature Index (k)	Cluster Size	Topic Description	Dominant Terms
33	233	Texas longhorn football newsletter	UT, orange, longhorn[s], texas, true, truorange, recruiting, oklahoma, defensive
34	65	Enron collapse	partnership[s] , fastow , shares, sec , stock, shareholder, investors, equity, lay
39	235	Emails about India	dabhol , dpc , india , mseb , maharashtra , indian, lenders, delhi, foreign, minister

2001 Topics Tracked by GD-CLS



$r = 50$ features, **lex** term weighting, $\lambda = 0.1$

Two Penalty Term Formulation

- ▶ Introduce smoothing on W_k (feature vectors) in addition to H^k :

$$\min_{W,H} \{ \|X - WH\|_F^2 + \alpha \|W\|_F^2 + \beta \|H\|_F^2 \},$$

where $\|\cdot\|_F$ is the Frobenius norm.

- ▶ Constrained NMF (CNMF) iteration [Piper et al., 2004]:

$$H_{cj} \leftarrow H_{cj} \frac{(W^T X)_{cj} - \beta H_{cj}}{(W^T WH)_{cj} + \epsilon}$$

$$W_{ic} \leftarrow W_{ic} \frac{(XH^T)_{ic} - \alpha W_{ic}}{(WHH^T)_{ic} + \epsilon}$$

Term Distribution in Feature Vectors

Terms	Wt	Lambda			Alpha			Topics
		0.1	0.01	0.001	0.1	0.01	0.001	
Blackouts	0.508				4	6	4	Cal
Stocks	0.511				2			Collapse
UT	0.517				2			Texasfoot
Chronicle	0.523				3	2	3	
Indian	0.527				2			India
Fastow	0.531				5	3	4	Collapse
Gas	0.531					2	2	
CFO	0.556				2		2	Kitchen
Californians	0.557					3		Cal
Solar	0.570				2			
Partnerships	0.576				6	2	5	Collapse
Workers	0.577					3	2	
Maharashtra	0.591				2		2	India
Mseb	0.605				2			India
Beach	0.611			2				
Ljm	0.621				3		3	Collapse
Tues	0.626		2	2				
IPPS	0.644			2		2		Cal
Rebates	0.647					2		
Ljm2	0.688				2		2	Collapse

British National Corpus (BNC) Noun Conservation

- ▶ In collaboration with P. Keila and D. Skillicorn (Queens Univ.)
- ▶ 289,695 email subset (all mail folders - not just private)
- ▶ Smoothing solely applied to NMF W matrix
($\alpha = 0.001, 0.01, 0.1, 0.25, 0.50, 0.75, 1.00$ with $\beta = 0$)
- ▶ Log-entropy term weighting applied to the term-by-message matrix X
- ▶ Monitor top ten nouns for each feature vector (ranked by descending component values) and extract those appearing in two or more features; topics assigned manually.

BNC Noun Distribution in Feature Vectors

Noun	GF	Entropy	Alpha							Topic
			0.001	0.01	0.1	0.25	0.50	0.75	1.00	
Waxman	680	0.424	2		2	2	2	2		Downfall
Lieberman	915	0.426	2	2	2	2			2	Downfall
Scandal	679	0.428	2				2		2	Downfall
Nominee(s)	544	0.436		4	3	2		2	2	
Barone	470	0.437	2	2	2				2	Downfall
MEADE	456	0.437							2	Downfall
Fichera	558	0.438	2			2				California blackout
Prabhu	824	0.445	2	2	2	2		2	2	India-strong
Tata	778	0.448							2	India-weak
Rupee(s)	323	0.452	3	4	4	4	3	4	2	India-strong
Soybean(s)	499	0.455	2	2	2	2	2	2	2	
Rushing	891	0.486	2	2	2					Football - college
Dlrs	596	0.487							2	
Janus	580	0.488	2	3				2	3	India-weak
BSES	451	0.498	2	2					2	India-weak
Caracas	698	0.498						2		
Escondido	326	0.504	2			2				California/Blackout
Promoters	180	0.509	2							Energy/Scottish
Aramco	188	0.550	2							India-weak
DOORMAN	231	0.598		2						Bawdy/Real Estate

Hoyer Sparsity Constraint

- ▶ $\text{sparseness}(\mathbf{x}) = \frac{\sqrt{n} - \|\mathbf{x}\|_1 / \|\mathbf{x}\|_2}{\sqrt{n-1}}$, [Hoyer, 2004]
- ▶ Imposed as a penalty term of the form

$$J_2(\mathbf{W}) = (\omega \|\text{vec}(\mathbf{W})\|_2 - \|\text{vec}(\mathbf{W})\|_1)^2,$$

where $\omega = \sqrt{mk} - (\sqrt{mk} - 1)\gamma$ and $\text{vec}(\cdot)$ transforms a matrix into a vector by column stacking.

- ▶ Desired sparseness in \mathbf{W} is specified by setting $\gamma \in [0, 1]$; *sparseness* is zero iff all vector components are equal (up to signs) and is one iff the vector has a single nonzero.

Sample Benchmarks

- ▶ Elapsed CPU times for CNMF on a 3.2GHz Intel Xeon 3.2GHz (1024KB cache, 4.1GB RAM)
- ▶ $k = 50$ feature vectors generated, log-entropy noun-weighting used on $7,424 \times 289,695$ noun-by-message matrix, random $\mathbf{W}_0, \mathbf{H}_0$

W-Constraint	Iterations	Parameters	CPU time
L_2 norm	100	$\alpha = 0.1, \beta = 0$	19.6m
L_2 norm	100	$\alpha = 0.01, \beta = 0$	20.1m
L_2 norm	100	$\alpha = 0.001, \beta = 0$	19.6m
Hoyer	30	$\alpha = 0.01, \beta = 0, \gamma = 0.8$	2.8m
Hoyer	30	$\alpha = 0.001, \beta = 0, \gamma = 0.8$	2.9m

BNC Noun Distribution in Sparsified Feature Vectors

Noun	GF	Entropy	Alpha							Topic
			0.001	0.01	0.01	0.25	0.50	0.75	1.00	
Fleischer	903	0.409			3					Downfall
Coale	836	0.414			2					Downfall
Waxman	680	0.424	2		2	2	2	2		Downfall
Businessweek	485	0.424			2					
Lieberman	915	0.426	2	2		2			2	Downfall
Scandal	679	0.428	2		2		2		2	Downfall
Nominee(s)	544	0.436		4	2	2		2	2	
Barone	470	0.437	2	2	2				2	Downfall
MEADE	456	0.437							2	Downfall
Fichera	558	0.438	2		2	2				California blackout
Prabhu	824	0.445	2	2	3	2		2	2	India-strong
Tata	778	0.448							2	India-weak
Rupee(s)	323	0.452	3	4	3	4	3	4	2	India-strong
Soybean(s)	499	0.455	2	2	3	2	2	2	2	
Rushing	891	0.486	2	2						Football - college
Dlrs	596	0.487			2				2	
Janus	580	0.488	2	3	2			2	3	India-weak
BSES	451	0.498	2	2	2				2	India-weak
Caracas	698	0.498						2		
Escondido	326	0.504	2			2				California/Blackout
Promoters	180	0.509	2							Energy/Scottish
Aramco	188	0.550	2							India-weak
DOORMAN	231	0.598		2						Bawdy/Real Estate

MailMiner - CS365/Spring 2005

- ▶ Course Project in CS365/Programming Languages (Spring 2005)
- ▶ Student authors: C. Mollenhour, J. Russell, R. Warren
- ▶ Recent modifications by H. Gonzales and K. Rankin
- ▶ Specifications:
 1. Rank emails by keyword frequencies (OR based) and highlight keywords in hyperlinked email files; ranked results can be saved to file.
 2. Provide additional boolean operators for query (AND/NOT)
 3. Accept NMF-generated email clusters and delineate members within the ranked results display.

Query: california AND blackout; sort by term frequencies

The screenshot shows the mailMiner application window. The search query is "california AND blackout". The results are displayed in a table with columns: Rank, Subject, Sender, Date, and Cluster. The top results are duplicates of "Energy Issues" and "PowerMarketers.com Daily Power Report for 12 April 2001". The results for "California Sing-along" are highlighted in blue, showing a cluster of "Humor_calif".

Rank	Subject	Sender	Date	Cluster
108	Energy Issues	miyung.buster@enron.com	04/02/2001 6:05 AM	duplicates
69	Energy Issues	miyung.buster@enron.com	04/06/2001 5:45 AM	energy_newsfeed
69	Energy Issues	miyung.buster@enron.com	04/06/2001 5:45 AM	duplicates
27	PowerMarketers.com Daily Power Report for 12 April 2001	pmadpr@worldnet.att.net	04/11/2001 7:58 PM	energy_newsfeed
20	PowerMarketers.com Daily Power Report for 9 April 2001	pmadpr@worldnet.att.net	04/08/2001 7:49 PM	energy_newsfeed
13	PowerMarketers.com Daily Power Report for 13 April 2001	pmadpr@worldnet.att.net	04/12/2001 7:45 PM	energy_newsfeed
4	State-by-state forecast for summer	Karen.denne@enron.com	04/10/2001 9:04 AM	energy_newsfeed
4	State-by-state forecast for summer	Karen.denne@enron.com	04/10/2001 9:04 AM	duplicates
3	California Sing-along	jeffery.fawcett@enron.com	04/30/2001 11:17 AM	Humor_calif
3	California Sing-along	jeffery.fawcett@enron.com	04/30/2001 11:17 AM	duplicates

File: lokay-m/personal/158
 To: lorna.brannan@enron.com, kevin.hyatt@enron.com, tk.lohman@enron.com, ...
 From: jeffery.fawcett@enron.com
 Date: Mon, 30 Apr 2001 11:17 AM
 Subject: California Sing-along

And to start off the week....

California's new State Song! The Rolling **Blackout** Theme Song
 (To the theme music from the TV western "Rawhide")
 Rollin', rollin', rollin',
 Though the state is golden,
 Keep them blackouts rollin', statewide.
 A little colder weather,
 And we all freeze together.

Search successful. Results: 21

Query: godbole AND report; sort by timestamp

The screenshot shows the mailMiner application window. The search bar contains the query "godbole AND report". The interface includes a sidebar with a hierarchical tree view of clusters, a main table of search results, and a bottom status bar.

Search Bar: godbole AND report

Buttons: Search, Open search in new tab, New Search Tab, Close Search Tab, Clear Results, View Full Text, Export Results List, Load New Dataset

Cluster Tree (Left):

- duplicates
- Dyriegy
- education
- energy_newsfeed
- Enron_online
- EnronOnline
- FERC_DOE
- finance_conference
- government_state
- government_us
- Houston_FireSelect
- Humor
- Humor_calif
- India
- India_Dabhol**
 - /snozz/homes/cs365/enron-private/2001/Apr/kaminski-vje/15.
 - /snozz/homes/cs365/enron-private/2001/Apr/kean-sj/india/59.
 - /snozz/homes/cs365/enron-private/2001/Apr/kean-sj/india/58.
 - /snozz/homes/cs365/enron-private/2001/Apr/kean-sj/india/57.
- India_Dabhol_Anshuman
- India_energy
- kitchen_daily

Search Results Table:

Rank	Subject	Sender	Date	Cluster
9	Godbole Report	sandeep.kohli@enron.com	04/16/2001 5:08 AM	India_Dabhol
11	RE: Final Draft Media Statement: DPC reaction to G...	loretta.brelsford@enron.com	04/17/2001 12:04 PM	India_Dabhol
11	RE: Final Draft Media Statement: DPC reaction to G...	loretta.brelsford@enron.com	04/17/2001 12:04 PM	duplicates
29	From The Enron India Newsdesk - April 23rd news...	sandeep.kohli@enron.com	04/23/2001 4:01 AM	duplicates
28	From The Enron India Newsdesk - April 23rd news...	sandeep.kohli@enron.com	04/23/2001 2:01 PM	energy_new...

Status Bar: Search successful. Results: 5

Annotation Project

- ▶ Subset of 2001 PRIVATE collection:

Month	Total	Classified	Usable
Jan, Sep	5591	1100	699
Feb	2804	900	460
Mar	3525	1200	533
Apr	4273	1500	705
May	4261	1800	894
June	4324	1025	538
Total	24778	7525	3829

- ▶ Approx. 40 topics identified (after NMF initial clustering with $k = 50$ features) by two MailMiner users.

Annotation Project, contd.

- ▶ Human classifiers: M. Browne (extensive background reading on Enron collapse) and B. Singer (junior Economics major).
- ▶ Classify email content versus type (see UC Berkeley Enron Email Analysis Group
http://bailando.sims.berkeley.edu/enron_email.html)
- ▶ Potential U. Penn LDC (Linguistic Data Consortium) submission (see <http://www.ldc.upenn.edu>)

Conclusions

- ▶ GD-CLS Algorithm can effectively produce a *parts-based* approximation $X \simeq WH$ of a sparse term-by-message matrix X .
- ▶ Smoothing on the features matrix (W) as opposed to the weight matrix H forces more reuse of higher weighted terms.
- ▶ Surveillance systems based on GD-CLS could be used to monitor discussions without the need to isolate or perhaps incriminate individual employees.
- ▶ Potential applications include the monitoring/tracking of company morale, employee feedback to policy decisions, and extracurricular activities

Future Work

- ▶ Further work needed in determining effects of alternative term weighting schemes (for X) and choices of control parameters (α, β, λ) on quality of the basis vectors W_k .
- ▶ How does document (or message) clustering change with different column ranks (r) in the matrix W ?
- ▶ Use MailMiner and similar text mining software to produce a topic **annotated** Enron email subset for the public domain.
- ▶ Explore use of NMF for automated gene classification;
Semantic Gene Organizer (K. Heinrich, PhD Thesis 2006)

For Further Reading

- ▶ F. Shahnaz, M.W. Berry, V.P. Pauca, and R.J. Plemmons.
Document Clustering Using Nonnegative Matrix Factorization.
Info. Proc. & Management 42(2), 2006, pp. 373-386.
- ▶ J. Piper, P. Pauca, R. Plemmons, and M. Giffin.
Object Characterization from Spectral Data using ICA ...
Proc. AMOS Technical Conference, Maui, HI, September 2004.
- ▶ P. Hoyer.
Non-negative Matrix Factorization with Sparseness Constraints.
J. Machine Learning Research 5, 2004, pp. 1457-1469.
- ▶ W. Xu, X. Liu, and Y. Gong.
Document-Clustering based on Non-neg. Matrix Factorization.
Proceedings of SIGIR'03, Toronto, CA, 2003, pp. 267-273.

SMD07 Text Mining Workshop



- ▶ Fifth workshop held since 2001; M.W. Berry and M. Castellanos (Organizers)
- ▶ Regular papers and poster papers accepted
- ▶ Text mining topics include: information retrieval and extraction, machine learning, natural language processing, mathematical models, software environments, case studies
- ▶ URL: <http://www.siam.org/sdm07>