# Text to Matrix Generator Toolbox

A Brief Introduction

Eugenia Maria Kontopoulou, Dimitrios Zeimpekis
and Efstratios Gallopoulos

Department of
Computer Engineering and Informatics
University of Patras

Patras, 09/05/2014

## Documents



| Labels | Titles |
|--------|--------|
| B1 | Identifying users of social networks from their data footprint: An application of large-scale matrix factorizations |
| B2 | Data fusion based on coupled matrix and tensor factorizations |
| B3 | On incremental deterministic methods for dominant space estimation for large data sets |
| B4 | Fast projection methods for robust separable nonnegative matrix factorization |
| B5 | Experiments with randomized algorithms in the text to matrix generator toolbox |

## Term-Document Matrix (TDM)
### 33 × 5

| terms | Documents | | | | | terms | Documents | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | B1 | B2 | B3 | B4 | B5 | | B1 | B2 | B3 | B4 | B5 |
| algorithm | 0 | 0 | 0 | 0 | 2.3219 | matrix | 0.3219 | 0.3219 | 0 | 0.3219 | 0.3219 |
| applic | 2.3219 | 0 | 0 | 0 | 0 | method | 0 | 0 | 1.3219 | 1.3219 | 0 |
| base | 0 | 2.3219 | 0 | 0 | 0 | network | 2.3219 | 0 | 0 | 0 | 0 |
| coupl | 0 | 2.3219 | 0 | 0 | 0 | nonneg | 0 | 0 | 0 | 2.3219 | 0 |
| data | 0.7370 | 0.7370 | 0.7370 | 0 | 0 | project | 0 | 0 | 0 | 2.3219 | 0 |
| determinist | 0 | 0 | 2.3219 | 0 | 0 | random | 0 | 0 | 0 | 0 | 2.3219 |
| domin | 0 | 0 | 2.3219 | 0 | 0 | robust | 0 | 0 | 0 | 2.3219 | 0 |
| estim | 0 | 0 | 2.3219 | 0 | 0 | scale | 2.3219 | 0 | 0 | 0 | 0 |
| experi | 0 | 0 | 0 | 0 | 2.3219 | separ | 0 | 0 | 0 | 2.3219 | 0 |
| factor | 0.7370 | 0.7370 | 0 | 0.7370 | 0 | set | 0 | 0 | 2.3219 | 0 | 0 |
| fast | 0 | 0 | 0 | 2.3219 | 0 | social | 2.3219 | 0 | 0 | 0 | 0 |
| footprint | 2.3219 | 0 | 0 | 0 | 0 | space | 0 | 0 | 2.3219 | 0 | 0 |
| fusion | 0 | 2.3219 | 0 | 0 | 0 | tensor | 0 | 2.3219 | 0 | 0 | 0 |
| gener | 0 | 0 | 0 | 0 | 2.3219 | text | 0 | 0 | 0 | 0 | 2.3219 |
| identifi | 2.3219 | 0 | 0 | 0 | 0 | toolbox | 0 | 0 | 0 | 0 | 2.3219 |
| increment | 0 | 0 | 2.3219 | 0 | 0 | user | 2.3219 | 0 | 0 | 0 | 0 |
| larg | 1.3219 | 0 | 1.3219 | 0 | 0 | | | | | | |

✓ tf-idf          ✓ Stemming

Text Data
. . . for text mining tasks

Retrieval

## Text Data
. . . for text mining tasks

Retrieval

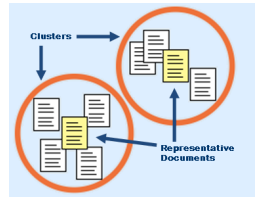Clustering

Retrieval

Clustering
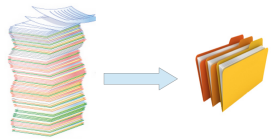
Classification

# Text to Matrix Generator

## What is TMG:

- Toolbox developed in University of Patras for text mining tasks over document collections
- Educational and Research tool

TMG: A MATLAB Toolbox for Generating Term-Document Matrices from Text Collections (ZG06b)

### Grouping Multidimensional Data

Recent Advances in Clustering
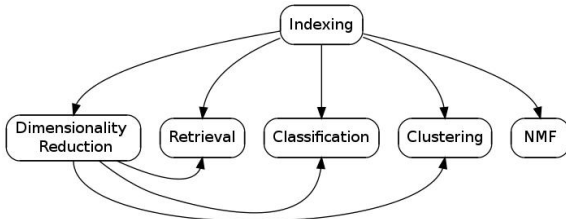Kogan, Jacob; Nicholas, Charles; Teboulle, Marc (Eds.)

2006, XII, 268 p.

Grouping Multidimensional Data
2006, pp 187-210

TMG: A MATLAB Toolbox for Generating Term-Document Matrices from Text Collections
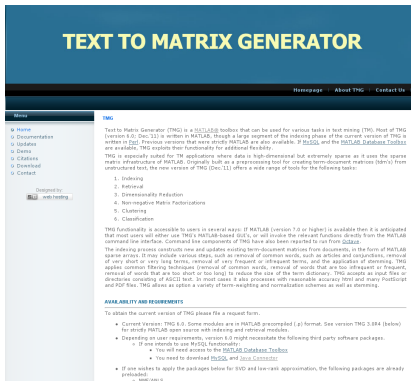
D. Zeimpekis, E. Gallopoulos

# Text to Matrix Generator

## What is TMG:

- Toolbox developed in University of Patras for text mining tasks over document collections
- Educational and Research tool

## Implementation:

- over 17.000 lines of `matlab` and `perl`
- takes advantage from sparse technology provided by MATLAB
- first version by Zeimpekis (´06)

TMG: A MATLAB Toolbox for Generating Term-Document Matrices from Text Collections (ZG06b)

### Grouping Multidimensional Data



Recent Advances in Clustering
Kogan, Jacob; Nicholas, Charles; Teboulle, Marc (Eds.)
2006, XII, 268 p.

Grouping Multidimensional Data
2006, pp 187-210

TMG: A MATLAB Toolbox for Generating Term-Document Matrices from Text Collections

D. Zeimpekis, E. Gallopoulos

Six basic modules:

1. Indexing
2. Dimensionality Reduction
3. Non-Negative Matrix Factorizations
4. Retrieval
5. Clustering
6. Classification

# How can I find TMG?

## Free under request from:

http://scgroup20.ceid.upatras.gr:8000/tmg/



## More than 4000 requests worldwide . . .

Caltech, Maryland, Purdue, Carnegie Mellon, Tennessee, Berkeley, Texas, Minnesota, Stanford, MIT, Columbia Renault, Leuven, Max-Planck, Michigan, Oxford, Philips, Princeton, Illinois, ETH, RPI, Los Alamos, Toronto, Queen Mary, St Andrews, Colorado, Texas, Livermore, Mathworks, Yahoo!, . . .

Part I

Introduction in version 6.0R7

# Generate, Update and Downdate Term-by-Document Matrices I

## Graphical User Interface



## Purpose

Document Collection

⇓

**Term-by-Document Matrix**

## Procedure

# Generate, Update and Downdate Term-by-Document Matrices III

## Supported non-ASCII formats

| | ver.5.0R6 | Filter ver.5.0R6 | ver. 6.0R7 | Filter ver. 6.0R7 |
|---|---|---|---|---|
| **doc** | $\times$ | $\times$ | $\sqrt{}$ | TIKA |
| **docx** | $\times$ | $\times$ | $\sqrt{}$ | TIKA |
| **htm** | $\sqrt{}$ | `strip_html` | $\sqrt{}$ | `strip_html` |
| **html** | $\sqrt{}$ | `strip_html` | $\sqrt{}$ | TIKA |
| **odt** | $\times$ | $\times$ | $\sqrt{}$ | TIKA |
| **pdf** | $\sqrt{}$ | `ps2ascii` | $\sqrt{}$ | `ps2ascii` |
| **ps** | $\sqrt{}$ | `ps2ascii` | $\sqrt{}$ | `ps2ascii` |
| **rtf** | $\times$ | $\times$ | $\sqrt{}$ | TIKA |
| **tex** | $\times$ | $\times$ | $\sqrt{}$ | Untex |

## Update

Update the `TDM` by inserting new documents

## Downdate

Downdate the `TDM` by extracting useless documents

# Dimensionality Reduction I

## Graphical User Interface

## Purpose

| Handling High Dimensional Data | Reducing noise |
|---|---|

$\Downarrow$        $\Downarrow$

| **Economical representation** | **Better semantic representation** |
|---|---|

# Dimensionality Reduction II

## Available Methods

1. Singular Value Decomposition (SVD)
   - ✓ MATLAB svds
   - ✓ PROPACK svd (Larsen (Lar))

2. Centroids Method (CM) (Park, Jeon & Rosen (PJR03))

3. Semidiscrete Decomposition (SDD) (Kolda & O'Leary (KO00))

4. **Clustered LSI (CLSI)** (Zeimpekis & Gallopoulos (ZG05; ZG06a))

5. Sparse Pivoted QR Decomposition (SPQR) (Berry, Pulatova & Stewart (BPS05))

6. Principal Component Analysis (PCA)

---

SDD and SPQR call routines available from Netlib(TOMS)

# Nonnegative Matrix Factorizations (NMF) I

## Purpose

Graphical User Interface

Factorizations on Nonnegative Matrices
$\Downarrow$
**Nonegative Factors**

$\Downarrow$

Preserving non-negativity
$\Downarrow$
**Better semantic representation**



✓ Final results depend on initialization          ✓ Resulting factors can be refined

# Nonnegative Matrix Factorizations (NMF) II

## Initialization Techniques

1. Random Initialization
2. **Nonnegative Double SVD NNDSVD** (Boutsidis & Gallopoulos (BG08))
3. **Block Nonnegative Double SVD** (Zeimpekis & Gallopoulos (ZG08))
4. **Bisecting Nonnegative Double SVD** (Zeimpekis & Gallopoulos (ZG08))
5. By Clustering (Wild, Curry, Dougherty (WCD04))

NNDSVD uses prepared implementation

## Factors Refinement

1. Multiplicative Update Algorithm (Lee & Seung (LS01))
2. Alternating Non-Negative-Constrained Least Squares (NMF/ANLS) (Kim & Park (KH08))

NMF/ANLS uses prepared implementation

## Graphical User Interface

## Purpose

Queries over a dataset

⇓

**Retrieve all relevant documents via a HTML response**

## Available Methods

1. Vector Space Model (`VSM`) (Salton, Wong, & Yang (SWY75))
2. Latent Semantic Analysis (`LSA`) (Berry et al. (BDJ99; Dee+90))

> `LSA` can be combined with any `DR` or `NMF` technique

## Graphical User Interface

## Purpose

Collection of documents as a TDM

⇩

**Clusters of related documents**

## Available Methods

1. Euclidean k-means
2. Spherical k-means(DM01)
3. Principal Direction Divisive Partitioning (PDDP) (Boley (Bol97))
4. **PDDP(l)** (Zeimpekis & Gallopoulos (ZG03))
5. **PDDP(l)** with some hybrid variants of PDDP and kmeans (Zeimpekis & Gallopoulos (ZG03))

## PDDP(l) Variants

✓ Split with k-means
✓ Optimal Split
✓ Optimal Split with k-means
✓ Optimal Split on Projections

Graphical User Interface

Purpose

Collection of documents as training
TDM
+
List of training labels

⇓

**Assign new documents to related
classes (labels)**

## Available Methods

1. k Nearest Neighboors (`knn`)
2. Rocchio
3. Linear Least Squares Fit (`LLSF`) (Yang & Chute (YC92))

> ✓ Combination with `CLSI`, `CM` and `SVD` DR techniques
> ✓ Implementations for multilabel and singlelabel collections

Goal:

☺ Make TMG more user friendly

**Goal:**

☺ Make TMG more user friendly

**Work in Progress:**

- Smarter parsing $\rightarrow$ boost parsing time
- Increase the degrees of freedom during parsing phase (e.g. stoplist, incorporation of new filters)
- Manual writing using MATLAB `publish`
- New stemming algorithms (e.g. greek stemmer)
- GUIs makeover
- Incorporation of new capabilities (e.g. WordNet, Wordle)

Thank you!

(BDJ99)    M.W. Berry, Z. Drmač, and E. R. Jessup. ``Matrices, vector spaces, and Information Retrieval''. In: *SIAM Rev.* 41 (1999), pp. 335-362.

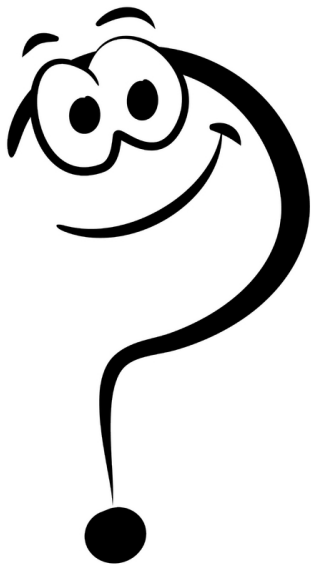(BG08)     C. Boutsidis and E. Gallopoulos. ``SVD based initialization: A head start for nonnegative matrix factorization''. In: *Pattern Recognition* 41 (Apr. 2008), pp. 1350-1362.

(Bol97)    D. Boley. ``Principal Direction Divisive Partitioning''. In: *Data Mining and Knowledge Discovery* 2 (1997), pp. 325-344.

(BPS05)    M. W. Berry, S. A. Pulatova, and G. W. Stewart. ``Algorithm 844: Computing sparse reduced-rank approximations to sparse matrices''. In: *ACM TOMS* 31.2 (June 2005), 252—269.

(Dee+90)   S. Deerwester et al. ``Indexing by Latent Semantic Analysis''. In: *Journal of the American Society for Information Science* 41.6 (1990), pp. 391-407.

(DM01)     I. S. Dhillon and D. S. Modha. ``Concept decompositions for large sparse text data using clustering''. In: *Machine Learning* 42.1 (2001), pp. 143-175.

(KH08)     H. Kim and H.Park. In: *SIAM J. M. Anal. and Appl.* (2008).

(KO00)    T. G. Kolda and D. P. O'Leary. ``Algorithm 805: computation and uses of the semidiscrete matrix decomposition''. In: *ACM TOMS* 26.3 (2000), pp. 415-435.

(Lar)     R.M. Larsen. *PROPACK: A software package for the symmetric eigenvalue problem and singular value problems on Lanczos and Lanczos bidiagonalization with partial reorthogonalization.*

(LS01)    D. D. Lee and H. S. Seung. ``Algorithms for Non-negative Matrix Factorization''. In: *NIPS*. MIT Press, 2001, pp. 556-562.

(PJR03)   H. Park, M. Jeon, and J.B. Rosen. ``Lower Dimensional Representation of Text Data Based on Centroids and Least Squares''. In: *BIT Numerical Mathematics* 43.2 (2003), pp. 427-448.

(SWY75)   G. Salton, A. Wong, and C. S. Yang. ``A vector space model for automatic indexing''. In: *Comm. ACM* 18.11 (1975), pp. 613-620.

(WCD04)   S. Wild, J. Curry, and A. Dougherty. ``Improving non-negative matrix factorizations through structured initialization.'' In: *Pattern Recognition* 37.11 (2004), pp. 2217-2232.

(YC92)    Y.Yang and C. G. Chute. ``A linear least squares fit mapping method for information retrieval from natural language texts''. In: (1992), pp. 447-453.

(ZG03)     D. Zeimpekis and E. Gallopoulos. ``PDDP(*l*): Towards a flexible principal direction divisive partitioning clustering algorithm''. In: *Proc. Workshop on Clustering Large Data Sets (held in conjunction with the Third IEEE Int'l. Conf. Data Min.)* Ed. by D. Boley et al. Melbourne, FL, 2003, pp. 26-35.

(ZG05)     D. Zeimpekis and E. Gallopoulos. ``CLSI: A Flexible Approximation Scheme from Clustered Term-Document Matrices''. In: *Proc. 5th SIAM Int'l Conf. Data Mining*. Ed. by H. Kargupta *et al.* SIAM. Philadelphia, 2005, pp. 631-635.

(ZG06a)    D. Zeimpekis and E. Gallopoulos. ``Linear and Non-Linear Dimensional Reduction via Class Representatives for Text Classification''. In: *6th Int'l. Conf. Data Mining (ICDM'06)*. Los Alamitos, CA, USA: IEEE Computer Society, 2006, pp. 1172-1177.

(ZG06b)    D. Zeimpekis and E. Gallopoulos. ``TMG: A MATLAB toolbox for generating term document matrices from text collections''. In: *Grouping Multidimensional Data: Recent Advances in Clustering*. Ed. by J. Kogan, C. Nicholas, and M. Teboulle. Berlin: Springer, 2006, 187—210.

(ZG08)     D. Zeimpekis and E. Gallopoulos. ``Document clustering using nmf based on spectral information''. In: *Text Mining Workshop (Atlanta)*. 2008.