

# Text to Matrix Generator\*

## User's Guide

Dimitrios Zeimpekis<sup>†</sup> Efstratios Gallopoulos<sup>‡</sup>

Department of Computer Engineering and Informatics,  
University of Patras, Greece

December 2008

---

\*Work conducted in the context of and supported in part by a University of Patras KARATHEODORI grant.

<sup>†</sup>e-mail: zeimpekis@gmail.com, supported in part by a Bodossaki Foundation graduate fellowship.

<sup>‡</sup>e-mail: stratis@ceid.upatras.gr.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Installation Instructions</b>	<b>1</b>
<b>3</b>	<b>Graphical User Interfaces</b>	<b>3</b>
3.1	Indexing module (tmg_gui) . . . . .	3
3.2	Dimensionality Reduction module (dr_gui) . . . . .	6
3.3	Non-Negative Factorizations module (nnmf_gui) . . . . .	9
3.4	Retrieval module (retrieval_gui) . . . . .	12
3.5	Clustering module (clustering_gui) . . . . .	14
3.6	Classification module (classification_gui) . . . . .	16
<b>A</b>	<b>Appendix: Demonstration of Use</b>	<b>22</b>
A.1	Indexing module (tmg_gui) . . . . .	22
A.2	Dimensionality Reduction module (dr_gui) . . . . .	29
A.3	Non-Negative Factorizations module (nnmf_gui) . . . . .	33
A.4	Retrieval module (retrieval_gui) . . . . .	37
A.5	Clustering module (clustering_gui) . . . . .	41
A.6	Classification module (classification_gui) . . . . .	46
<b>B</b>	<b>Appendix: Function Reference</b>	<b>50</b>
	about_tmg_gui . . . . .	50
	bisecting_nndsvd . . . . .	51
	block_diagonalize . . . . .	52
	block_nndsvd . . . . .	53
	classification_gui . . . . .	54
	clsi . . . . .	55
	clustering_gui . . . . .	56
	cm . . . . .	57
	col_normalization . . . . .	58
	col_rearrange . . . . .	59
	column_norms . . . . .	60
	compute_fro_norm . . . . .	61
	compute_scatter . . . . .	62
	create_kmeans_response . . . . .	63
	create_pddp_response . . . . .	64
	create_retrieval_response . . . . .	65
	diff_vector . . . . .	66
	dr_gui . . . . .	67
	ekmeans . . . . .	68
	entropy . . . . .	69
	get_node_scatter . . . . .	70
	gui . . . . .	71
	init_tmg . . . . .	72

knn_multi	73
knn_single	74
ks_selection	75
ks_selection1	76
llsf_multi	77
llsf_single	78
lsa	79
make_clusters_multi	80
make_clusters_single	81
make_labels	82
make_val_inds	83
merge_dictionary	84
merge_tdms	85
myperms	86
nnmf_gui	87
nnmf_mul_update	88
open_file	89
opt_2means	90
pca	91
pca_mat	92
pca_mat_afun	93
pca_propack	94
pca_propack_Atransfunc	95
pca_propack_afun	96
pca_update	97
pca_update_afun	98
pddp	99
pddp_2means	100
pddp_extract_centroids	101
pddp_optcut	102
pddp_optcut_2means	103
pddp_optcutpd	104
ps_pdf2ascii	105
retrieval_gui	106
rocchio_multi	107
rocchio_single	108
scut_knn	109
scut_llsf	110
scut_rocchio	111
sdd_tmg	112
skmeans	113
stemmer	114
strip_html	115
svd_tmg	116
svd_update	117
svd_update_afun	118

tdm_downdate . . . . .	119
tdm_update . . . . .	120
tmg . . . . .	122
tmg_gui . . . . .	124
tmg_query . . . . .	125
tmg_save_results . . . . .	127
tmg_template . . . . .	128
two_means_ld . . . . .	129
unique_elements . . . . .	130
unique_words . . . . .	131
vsm . . . . .	132

## List of Figures

1	Structure and dependencies of GUI modules of TMG. . . . .	1
2	Structure of TMG root directory. . . . .	2
3	The <code>tmg_gui</code> GUI. . . . .	3
4	The <code>dr_gui</code> GUI. . . . .	6
5	The <code>nnmf_gui</code> GUI. . . . .	9
6	The <code>retrieval_gui</code> GUI. . . . .	12
7	The <code>clustering_gui</code> GUI. . . . .	14
8	The <code>classification_gui</code> GUI. . . . .	16
9	Starting window of <code>tmg_gui</code> . . . . .	23
10	Next view of <code>tmg_gui</code> according to the user selection. . . . .	24
11	The <code>open_file</code> window. . . . .	25
12	The output “.mat” files of <code>tmg_gui</code> . . . . .	26
13	The MySQL view upon <code>tmg</code> execution. . . . .	27
14	The GUIs’ general help tab. . . . .	28
15	Starting window of <code>dr_gui</code> . . . . .	30
16	Next view of <code>dr_gui</code> according to the user selection. . . . .	31
17	The output “.mat” files of <code>dr_gui</code> . . . . .	32
18	Starting window of <code>nnmf_gui</code> . . . . .	34
19	Next view of <code>nnmf_gui</code> according to the user selection. . . . .	35
20	The output “.mat” files of <code>nnmf_gui</code> . . . . .	36
21	Starting window of <code>retrieval_gui</code> . . . . .	38
22	Next view of <code>retrieval_gui</code> according to the user selection. . . . .	39
23	The output of <code>retrieval_gui</code> . . . . .	40
24	Starting window of <code>clustering_gui</code> . . . . .	42
25	Next view of <code>clustering_gui</code> according to the user selection. . . . .	43
26	The output “.mat” files of <code>clustering_gui</code> . . . . .	44
27	The output of <code>clustering_gui</code> for PDDP. . . . .	45
28	Starting window of <code>classification_gui</code> . . . . .	47
29	Next view of <code>classification_gui</code> according to the user selection. . . . .	48

## List of Tables

1	Description of use of <code>tmg_gui</code> components. . . . .	5
2	Description of use of <code>dr_gui</code> components. . . . .	8
3	Description of use of <code>nmmf_gui</code> components. . . . .	11
4	Description of use of <code>retrieval_gui</code> components. . . . .	13
5	Description of use of <code>clustering_gui</code> components. . . . .	15
6	Description of use of <code>classification_gui</code> components. . . . .	18

## 1 Introduction

Text to Matrix Generator (TMG) is a MATLAB Toolbox that can be used for various Data Mining (DM) and Information Retrieval (IR) tasks. TMG uses the sparse matrix infrastructure of MATLAB that is especially suited for Text Mining (TM) applications where data are extremely sparse. Initially built as a preprocessing tool, TMG offers now a wide range of DM tools. In particular, TMG is composed of six Graphical User Interface (GUI) modules, presented in Figure 1 (arrows show modules dependencies).

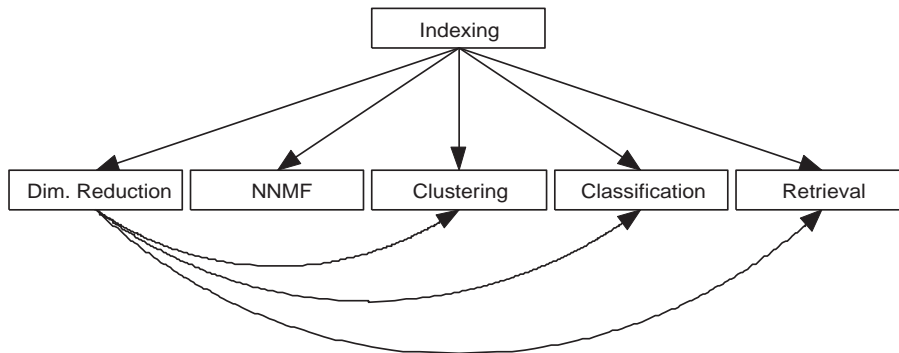


Figure 1: Structure and dependencies of GUI modules of TMG.

In the sequel, we first discuss the installation procedure of TMG and then describe in some detail the GUI's usage. In Appendix A we give a demonstration of use for all the TMG components, while Appendix B supplies a function reference.

## 2 Installation Instructions

Installation of TMG is straightforward by means of the `init_tm` script. In particular, the user has to perform the following steps:

- For MySQL functionality, install MySQL and Java Connector.
- Download TMG by filling the form from:  
[http://scgroup.hpclab.ceid.upatras.gr/scgroup/Projects/TMG/tmg\\_request.php](http://scgroup.hpclab.ceid.upatras.gr/scgroup/Projects/TMG/tmg_request.php)
- Unzip `TMG_X.XRX.zip` and start MATLAB. Figure 2 depicts the directory structure of the TMG root directory.
- Change path to the TMG root directory.
- Run `init_tm`. Give the MySQL login and password as well as the root directory of the MySQL Java Connector. The installation script creates all necessary information (including MySQL database TMG) and adds to the MATLAB path all necessary directories.

- Run gui. Alternatively, use the command line interface, type `help tmg`.

TMG requires the MySQL<sup>1</sup>, ANLS<sup>2</sup>, NNDSVD<sup>3</sup>, PROPACK<sup>4</sup>, SDDPACK<sup>5</sup> and SPQR<sup>6</sup> third party software packages. PROPACK, SDDPACK and SPQR packages are included into TMG, while the user has to download MySQL. However, we note that MySQL related software is necessary only if the user intends to use the database support implemented into TMG. Ordinary TMG will run without any problem on a Matlab 7.0 environment without any other special software.

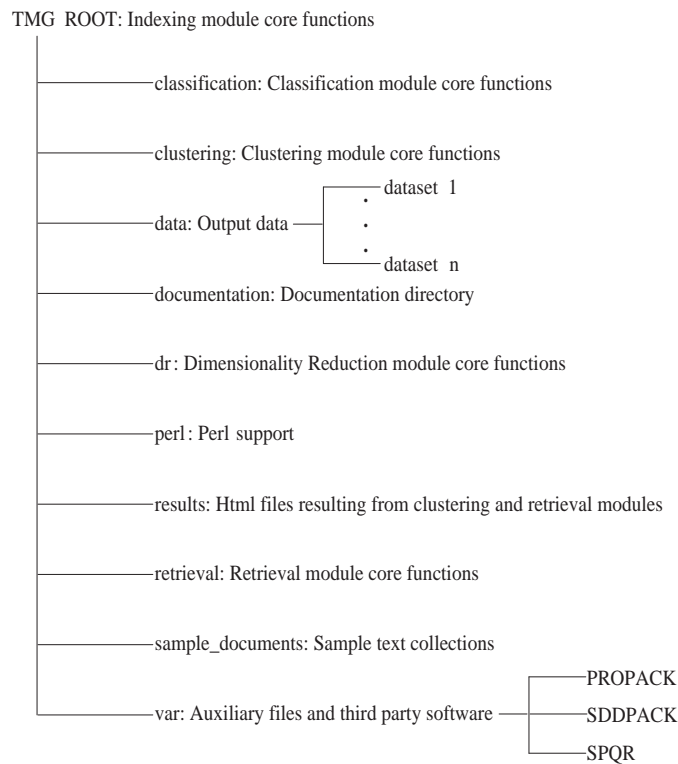


Figure 2: Structure of TMG root directory.

<sup>1</sup><http://www.mysql.com/>, <http://dev.mysql.com/downloads/connector/j/5.0.html>

<sup>2</sup><http://compbio.med.harvard.edu/hkim/nmf/index.html>

<sup>3</sup><http://www.cs.rpi.edu/~boutsc/paper1.html>

<sup>4</sup><http://soi.stanford.edu/~rmunk/PROPACK/index.html>

<sup>5</sup><http://www.cs.umd.edu/~oleary/SDDPACK/README.html>

<sup>6</sup><http://portal.acm.org/citation.cfm?id=1067972>



### 3 Graphical User Interfaces

#### 3.1 Indexing module (tmg\_gui)

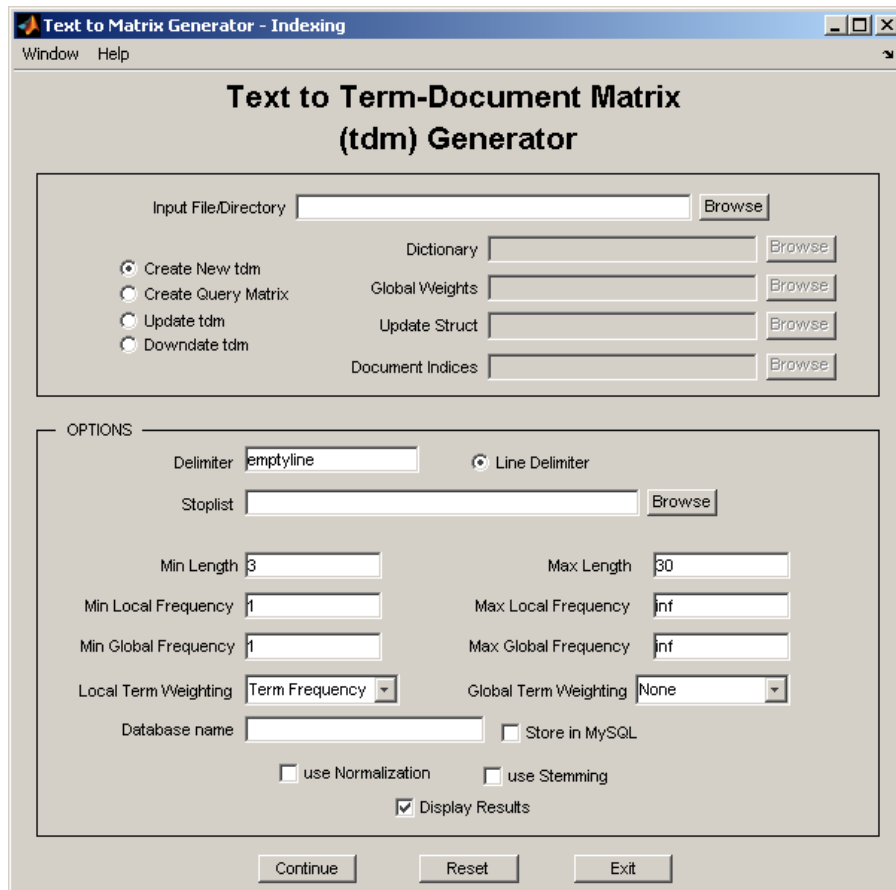


Figure 3: The tmg\_gui GUI.

TMG can be used for the construction of new and the update of existing term-document matrices (tdms) from text collections, in the form of MATLAB sparse arrays. To this end, TMG implements various steps such as:

- Removal of stopwords.
- Stemming (currently Porter stemming algorithm [11]).
- Remove of short/long terms.
- Remove of frequent/infrequent terms (locally or globally).

- Term weighting and normalization.
- Html filtering, processing of Postscript and PDF.
- Store in MySQL (optionally).

The resulting tdms can be stored as “mat” files, while text can also be stored in MySQL for further processing. TMG can also update existing tdms by efficient incremental updating or downdating operations. Finally, TMG can also construct query vectors using the existing dictionary that can be used from the retrieval and classification modules.

The indexing GUI module is depicted in Figure 3 while Table 1 describes in detail all the `tmg_gui` fields.

Field Name	Default	Description
Input File/Directory	-	Files to be parsed with resulting documents separated by “Delimiter”. Alternatively, each file in the input directory contains a single document.
Create New tdm	•	Checked if new tdm is to be created (default checked).
Create Query Matrix	-	Checked if new query matrix is to be created (default checked).
Update tdm	-	Checked if an existing tdm is to be updated with new documents. Alternatively, checked if an existing tdm is to be updated using different options (change <code>update_struct</code> ).
Downdate tdm	-	Checked if an existing tdm is to be downdated according to the “Document Indices” field.
Dictionary	-	Name of .mat file or workspace variable containing the dictionary to be used by <code>tmg_query</code> function if the “Create Query Matrix” radio button is checked.
Global Weights	-	Name of .mat file or workspace variable containing the vector of global weights to be used by <code>tmg_query</code> function if the “Create Query Matrix” radio button is checked.
Update Struct	-	Name of .mat file or workspace variable containing the structure to be updated or downdated by <code>tdm_update</code> (or <code>tdm_downdate</code> ) function if the “Update tdm” or “Downdate tdm” radio button is checked.
Document Indices	-	Name of .mat file or workspace variable containing the document indices marked for deletion when the “Downdate tdm” radio button is checked.
Field Name	Default	Description

Line Delimiter	•	Checked if the “Delimiter” takes a whole line of text.
Delimiter	emptyline	The delimiter between tmg’s view of documents. Possible values are ‘emptyline’, ‘none_delimiter’ (treats each file as single document) or any other string.
Stoplist	-	Name of file containing stopwords, i.e. common words not used in indexing.
Min Length	3	Minimum term length.
Max Length	30	Maximum term length.
Min Local Frequency	1	Minimum local term frequency.
Max Local Frequency	inf	Maximum local term frequency.
Min Global Frequency	1	Minimum global term frequency.
Max Global Frequency	inf	Maximum global term frequency.
Local Term Weighting	TF	Local term weighting function. Possible values: ‘Term Frequency’ (TF), ‘Binary’, ‘Logarithmic’, ‘Alternate Log’, ‘Augmented Normalized Term Frequency’.
Global Term Weighting	None	Global term weighting function. Possible values: ‘None’, ‘Entropy’, ‘Inverse Document Frequency (IDF)’, ‘GfIdf’, ‘Normal’, ‘Probabilistic Inverse’.
Database Name	-	The name of the folder (under ‘data’ directory) where data are to be saved (currently supported only for the “Create New tdm” module).
Store in MySQL	-	Checked if results are to be saved into MySQL (currently supported only for the “Create New tdm” module).
use Normalization	-	Indicates normalization method. Possible values: ‘None’, ‘Cosine’.
use Stemming	-	Indicates if stemming is to be applied. The algorithm currently supported is due to Porter.
Display Results	•	Display results or not to the command windows.
Continue	-	Apply the selected operation.
Reset	-	Reset window to default values.
Exit	-	Exit window.

Table 1: Description of use of tmg\_gui components.

### 3.2 Dimensionality Reduction module (dr\_gui)

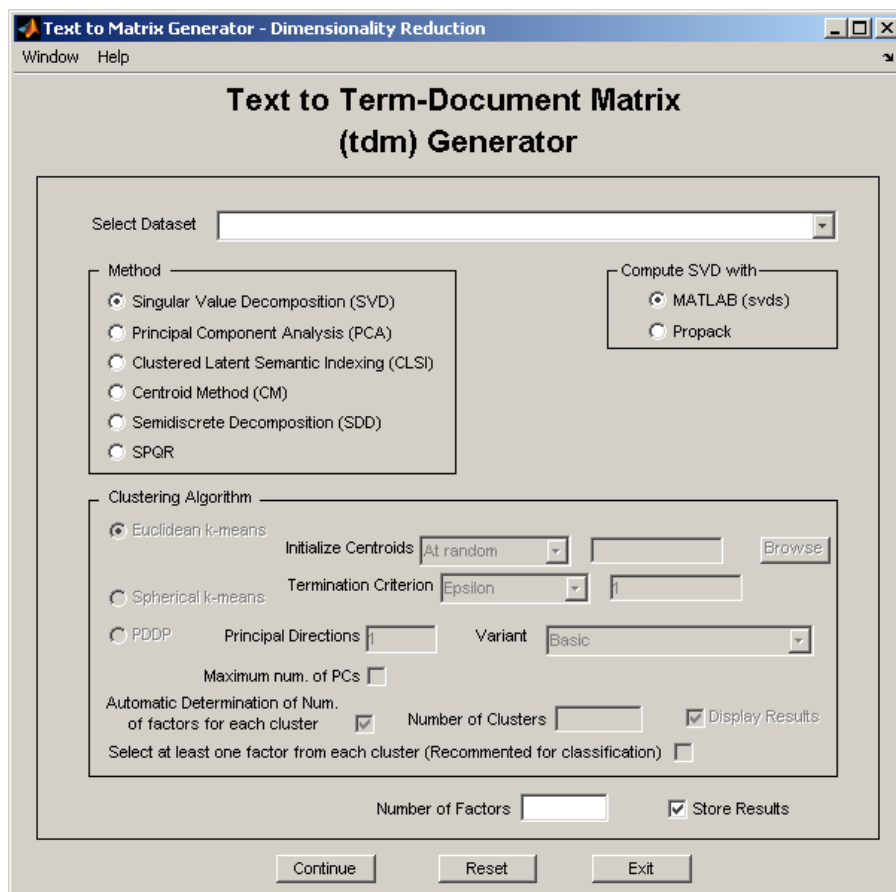


Figure 4: The dr\_gui GUI.

This module deploys a variety of powerful techniques designed to efficiently handle high dimensional data. Dimensionality Reduction (DR) is a common technique that is widely used. The target is dual: (a) more economical representation of data, and (b) better semantic representation. TMG implements six DR techniques.

- Singular Value Decomposition (SVD).
- Principal Component Analysis (PCA).
- Clustered Latent Semantic Indexing (CLSI) [16, 17].
- Centroids Method (CM) [10].
- Semidiscrete Decomposition (SDD) [8].

- SPQR Decomposition [2].

DR data can be stored as “.mat” files and used for further processing.

The dimensionality reduction GUI module is depicted in Figure 4 while Table 2 describes in detail all the `dr_gui` fields.

Field Name	Default	Description
Select Dataset	-	Select the dataset.
Singular Value Decomposition (SVD)	•	Apply the SVD method.
Principal Component Analysis (PCA)	-	Apply the PCA method.
Clustered Latent Semantic Indexing (CLSI)	-	Apply the CLSI method.
Centroid Method (CM)	-	Apply the CM method.
Semidiscrete Decomposition (SDD)	-	Apply the SDD method.
SPQR	-	Apply the SPQR method.
MATLAB (svds)	•	Check to use MATLAB function svds for the computation of the SVD or PCA.
Propack	-	Check to use PROPACK package for the computation of the SVD or PCA.
Euclidean k-means	•	Check to use the euclidean k-means clustering algorithm in the course of CLSI or CM.
Spherical k-means	-	Check to use the spherical k-means clustering algorithm in the course of CLSI or CM.
PDDP	-	Check to use the PDDP clustering algorithm in the course of CLSI or CM.
Initialize Centroids	At random	Defines the method used for the initialization of the centroid vector in the course of k-means. Possibilities are: initialize at random and supply a variable of '.mat' file with the centroids matrix.
Termination Criterion	Epsilon (1)	Defines the termination criterion used in the course of k-means. Possibilities are: use an epsilon value (default 1) and stop iteration when the objective function improvement does not exceed epsilon or perform a specific number of iterations (default 10).
Principal Directions	1	Number of principal directions used in PDDP.

Maximum num. of PCs	-	Check if the PDDP(max-l) variant is to be applied.
Variant	Basic	A set of PDDP variants. Possible values: 'Basic', 'Split with k-means', 'Optimat Split', 'Optimal Split with k-means', 'Optimal Split on Projection'.
Automatic Determination of Num. of factors for each cluster	•	Check to apply a heuristic for the determination of the number of factors computed from each cluster in the course of the CLSI algorithm.
Number of Clusters	-	Number of clusters computed in the course of the CLSI algorithm.
Display Results	•	Display results or not to the command windows.
Select at least one factor from each cluster	-	Use this option in case low-rank data are to be used in the course of classification.
Number of factors	-	Rank of approximation.
Store Results	•	Check to store results.
Continue	-	Apply the selected operation.
Reset	-	Reset window to default values.
Exit	-	Exit window.

Table 2: Description of use of `dr_gui` components.

### 3.3 Non-Negative Factorizations module (nmf\_gui)

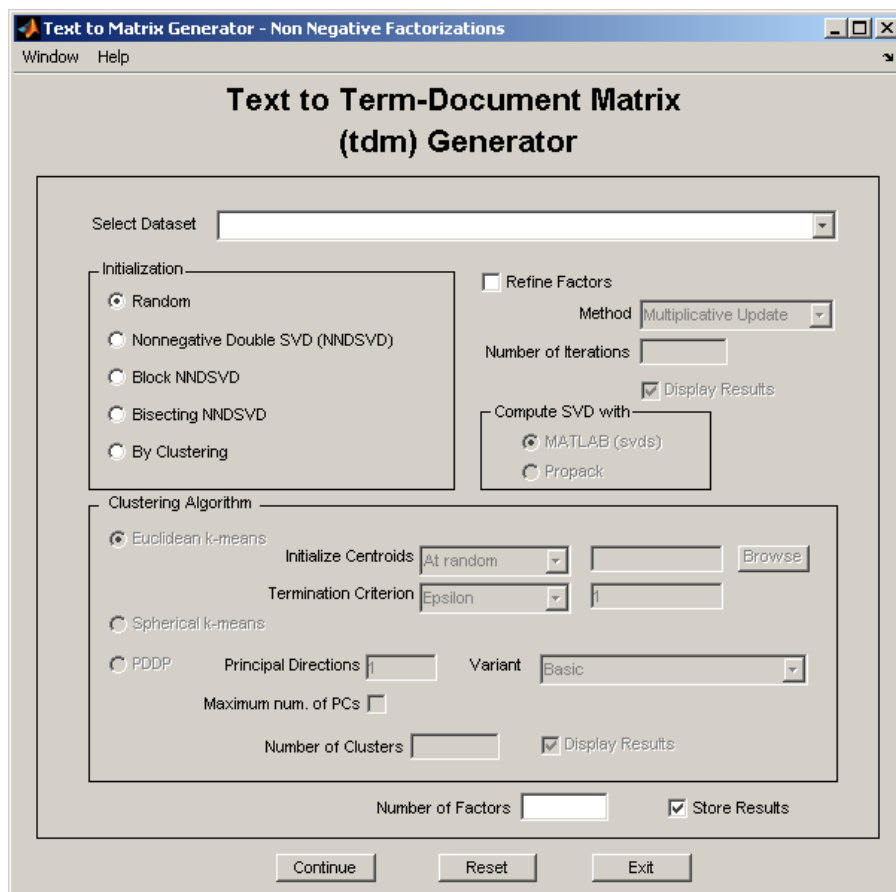


Figure 5: The nmf\_gui GUI.

This module deploys a set of Non-Negative Matrix Factorization (NNMF) techniques. Since these techniques are iterative, the final result depends on the initialization. A common approach is the random initialization of the non-negative factors, however new approaches appear to result in higher quality approximations. TMG implements four initialization techniques:

- Non-Negative Double Singular Value Decomposition (NDSVD) [4].
- Block NDSVD [20].
- Bisecting NDSVD [20].
- By clustering [13].

Resulting factors can be further refined by means of two NNMF algorithms:

- Multiplicative Update algorithm by Lee and Seung [9].
- Alternating Non-negativity-constrained Least Squares (NMF/ANLS) [7].

Field Name	Default	Description
Select Dataset	-	Select the dataset.
Random	•	Initialize at random.
Nonnegative Double SVD (NNDSVD)	-	Initialize by NNDSVD.
Block NNDSVD	-	Initialize by block NNDSVD.
Bisecting NNDSVD	-	Initialize by bisecting NNDSVD.
By Clustering	-	Initialize by clustering.
Refine factors	-	Check to run refinement algorithm.
Method	-	Refinement method (default Multiplicative Update).
Number of iterations	-	Refine by the Multiplicative Update algorithm.
Display Results	•	Display results of refinement method.
Euclidean k-means	•	Check to use the euclidean k-means clustering algorithm in the course of CLSI or CM.
Spherical k-means	-	Check to use the spherical k-means clustering algorithm in the course of CLSI or CM.
PDDP	-	Check to use the PDDP clustering algorithm in the course of CLSI or CM.
Initialize Centroids	At random	Defines the method used for the initialization of the centroid vector in the course of k-means. Possibilities are: initialize at random and supply a variable of '.mat' file with the centroids matrix.
Termination Criterion	Epsilon (1)	Defines the termination criterion used in the course of k-means. Possibilities are: use an epsilon value (default 1) and stop iteration when the objective function improvement does not exceed epsilon or perform a specific number of iterations (default 10).
Principal Directions	1	Number of principal directions used in PDDP.
Maximum num. of PCs	-	Check if the PDDP(max-1) variant is to be applied.
Variant	Basic	A set of PDDP variants. Possible values: 'Basic', 'Split with k-means', 'Optimat Split', 'Optimal Split with k-means', 'Optimal Split on Projection'.



MATLAB (svds)	•	Check to use MATLAB function svds for the computation of the SVD or PCA.
Propack	-	Check to use PROPACK package for the computation of the SVD or PCA.
Number of Clusters	-	Number of clusters computed.
Display Results	•	Display results or not to the command windows.
Store Results	•	Check to store results.
Continue	-	Apply the selected operation.
Reset	-	Reset window to default values.
Exit	-	Exit window.

Table 3: Description of use of `nmmf_gui` components.

### 3.4 Retrieval module (retrieval\_gui)

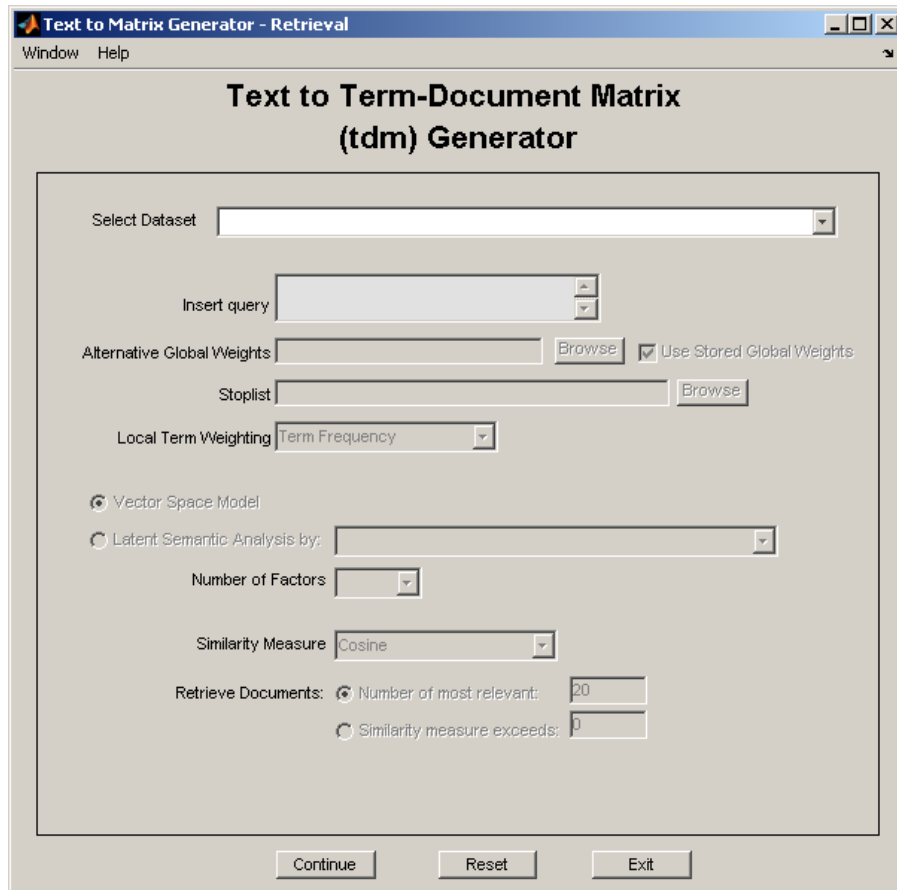


Figure 6: The retrieval\_gui GUI.

TMG offers two alternatives for Text Mining.

- Vector Space Model (VSM) [12].
- Latent Semantic Analysis (LSA) [1, 5],

using a combination of any DR technique and Latent Semantic Indexing (LSI). Using the corresponding GUI, the user can apply a question to an existing dataset using any of the aforementioned techniques and get HTML response.

The retrieval GUI module is depicted in Figure 6 while Table 4 describes in detail all the retrieval\_gui fields.

Field Name	Default	Description
Select Dataset	-	Select the dataset.
Insert Query	•	The query to be executed.
Alternative Global Weights	-	Global weights vector used for the construction of the query vector.
Use Stored Global Weights	•	Use the global weights vector found on the container directory of the dataset.
Stoplist	-	Use a stoplist.
Local Term Weighting	TF	The local term weighting to be used.
Vector Space Model	•	Apply the Vector space Model retrieval method.
Latent Semantic Analysis	-	The method used in the course of the Latent Semantic Analysis technique. Possible values: 'Singular Value Decomposition', 'Principal Component Analysis', 'Clustered Latent Semantic Analysis', 'Centroid Method', 'Semidiscrete Decomposition', 'SPQR'.
Number of Factors	-	Select the number of factors used during the retrieval process.
Similarity Measure	Cosine	Similarity measure used during the retrieval process.
Number of most relevant	•	Defines the number of most relevant documents returned for a query.
Similarity measure exceeds	-	Defines the minimum similarity measure value for which a document is treated as relevant to the query.
Continue	-	Apply the selected operation.
Reset	-	Reset window to default values.
Exit	-	Exit window.

Table 4: Description of use of `retrieval_gui` components.

### 3.5 Clustering module (clustering\_gui)

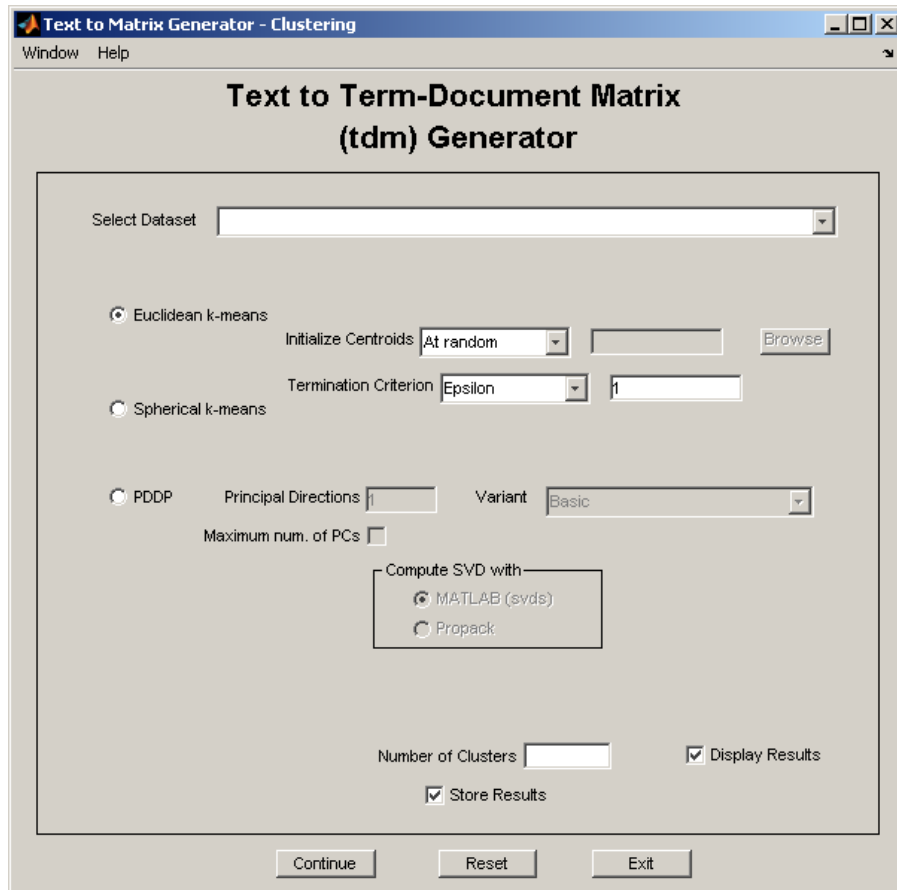


Figure 7: The clustering\_gui GUI.

TMG implements three clustering algorithms.

- k-means.
- Spherical k-means [6].
- Principal Direction Divisive Partitioning (PDDP) [3, 15].

Regarding PDDP, TMG implements the basic algorithm as well as the PDDP(l) [15] along with some recent hybrid variants of PDDP and kmeans [19].

The clustering GUI module is depicted in Figure 7 while Table 5 describes in detail all the clustering\_gui fields.

Field Name	Default	Description
Select Dataset	-	Select the dataset.
Euclidean k-means	•	Check to use the euclidean k-means clustering algorithm.
Spherical k-means	-	Check to use the spherical k-means clustering algorithm.
PDDP	-	Check to use the PDDP clustering algorithm.
Initialize Centroids	At random	Defines the method used for the initialization of the centroid vector in the course of k-means. Possibilities are: initialize at random and supply a variable of '.mat' file with the centroids matrix.
Termination Criterion	Epsilon (1)	Defines the termination criterion used in the course of k-means. Possibilities are: use an epsilon value (default 1) and stop iteration when the objective function improvement does not exceed epsilon or perform a specific number of iterations (default 10).
Principal Directions	1	Number of principal directions used in PDDP.
Maximum num. of PCs	-	Check if the PDDP(max-l) variant is to be applied.
Variant	Basic	A set of PDDP variants. Possible values: 'Basic', 'Split with k-means', 'Optimat Split', 'Optimal Split with k-means', 'Optimal Split on Projection'.
MATLAB (svds)	•	Check to use MATLAB function svds for the computation of the SVD in the course of PDDP.
Propack	-	Check to use PROPACK package for the computation of the SVD in the course of PDDP.
Number of Clusters	-	Number of clusters computed.
Display Results	•	Display results or not to the command windows.
Store Results	•	Check to store results.
Continue	-	Apply the selected operation.
Reset	-	Reset window to default values.
Exit	-	Exit window.

Table 5: Description of use of clustering\_gui components.

### 3.6 Classification module (classification\_gui)

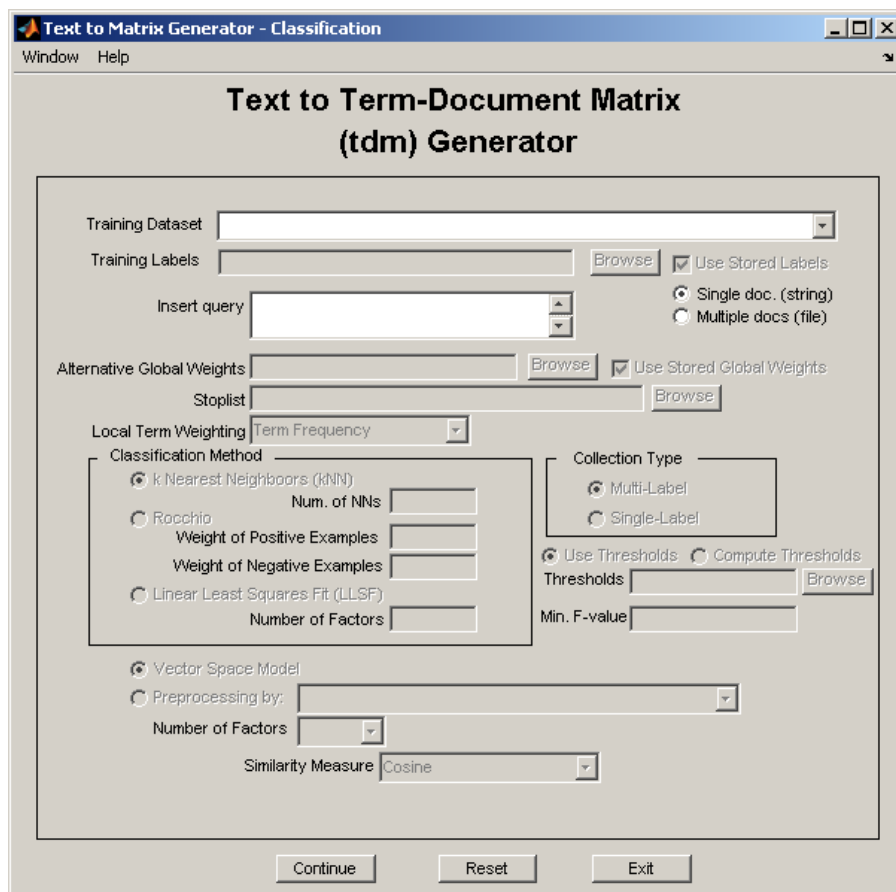


Figure 8: The classification\_gui GUI.

TMG implements three classification algorithms.

- $k$  Nearest Neighbors (kNN).
- Rocchio.
- Linear Least Squares Fit (LLSF) [14].

All these algorithms can be combined with CLSI, CM and SVD DR techniques.

The classification GUI module is depicted in Figure 8 while Table 6 describes in detail all the classification\_gui fields.

Field Name	Default	Description
Training Dataset	-	The training dataset.
Training Labels	-	The labels of the training dataset.
Use Stored Labels	•	Check to use the stored vector of labels of training documents in the container folder.
Insert query	-	The test document(s).
Single doc. (string)	•	Check if a single test document is to be inserted.
Multiple docs (file)	-	Check if multiple test document are to be inserted.
Filename	-	In 'Multiple docs (file)' is checked, insert the filename containing the test documents.
Delimiter	-	In 'Multiple docs (file)' is checked, insert the delimiter to be used for the test documents.
Line Delimiter	•	In 'Multiple docs (file)' is checked, check if delimiter of test documents' file takes a whole l of text.
Alternative Global Weights	-	Global weights vector used for the construction of the test documents' vectors.
Use Stored Global Weights	•	Use the global weights vector found on the container directory of the training dataset.
Stoplist	-	Use a stoplist.
Local Term Weighting	TF	The local term weighting to be used.
k Nearest Neighbors (kNN)	•	Check if the kNN classifier is to be applied.
Num. of NNs	-	Number of Nearest Neighbors in kNN classifier.
Rocchio	-	Check if Rocchio classifier is to be applied.
Weight of Positive Examples	-	The weight of the positive examples in the formation of the centroids vectors in Rocchio.
Weight of Negative Examples	-	The weight of the negative examples in the formation of the centroids vectors in Rocchio.
Linear Least Squares Fit (LLSF)	-	Check if LLSF classifier is to be applied.
Number of Factors	-	Number of factors used in the course of LLSF.
Multi-Label	•	Check if classifier is to be applied for a multi-label collection.
Single-Label	-	Check is classifier is to be applied for a single-label collection.
Use Thresholds	•	If 'Multi-Label' radio button is checked, use a stored vector of thresholds.
Compute Thresholds	-	If 'Multi-Label' radio button is checked, compute thresholds.
Thresholds	-	If 'Multi-Label' and 'Use Thresholds' radio buttons are checked, supply a stored vector of thresholds.

Min. F-value	-	If 'Multi-Label' and 'Compute Thresholds' radio buttons are checked, supply minimum F1 value used in the thresholding algorithm.
Vector Space Model	•	Use the basic Vector Space Model.
Preprocessing by	-	Use preprocessed training data with: 'Singular Value Decomposition', 'Principal Component Analysis', 'Clustered Latent Semantic Analysis', 'Centroid Method', 'Semidiscrete Decomposition', 'SPQR'.
Number of Factors	-	Number of factors for preprocessed training data.
Similarity Measure	Cosine	The similarity measure to be used.
Continue	-	Apply the selected operation.
Reset	-	Reset window to default values.
Exit	-	Exit window.

Table 6: Description of use of `classification_gui` components.



## Acknowledgments

TMG was conceived after a motivating discussion with Andrew Knyazev regarding a collection of MATLAB tools we had put together to aid in our clustering experiments. We thank our colleagues Ioannis Antonellis, Anastasios Zouzias, Efi Kokiopoulou and Constantine Bekas for many helpful suggestions, Jacob Kogan and Charles Nicholas for inviting us to contribute to [18], Elias Houstis for his help in the initial phases of this research and Michael Berry, Tamara Kolda, Rasmus Munk Larsen, Christos Boutsidis and Haesun Park for letting us use and distribute SPQR, SDDPACK, PROPACK, NNDSVD and ANLS software respectively. Special thanks are due to many of the users for their constructive comments regarding TMG. This research was supported in part by a University of Patras “Karatheodori” grant. The first author was also supported by a Bodossaki Foundation graduate fellowship.

## References

- [1] M. Berry, Z. Drmac, and E. Jessup, *Matrices, vector spaces, and information retrieval*, SIAM Review **41** (1998), 335–362.
- [2] M. W. Berry, S. A. Pulatova, and G. W. Stewart, *Computing sparse reduced-rank approximations to sparse matrices*, ACM TOMS **31** (2005), no. 2.
- [3] D. Boley, *Principal direction divisive partitioning*, Data Mining and Knowledge Discovery **2** (1998), no. 4, 325–344.
- [4] C. Boutsidis and E. Gallopoulos, *Svd-based initialization: A head start on non-negative matrix factorization*, Pattern Recognition **41** (2008), no. 4, 1350–1362.
- [5] S. Deerwester, S. Dumais, G. Furnas, T. Landauer, and Harshman R., *Indexing by Latent Semantic Analysis*, Journal of the American Society for Information Science **41** (1990), no. 6, 391–407.
- [6] I. S. Dhillon and D. S. Modha, *Concept decompositions for large sparse text data using clustering*, Machine Learning **42** (2001), no. 1, 143–175.
- [7] H. Kim and H. Park, *Non-negative matrix factorization based on alternating non-negativity-constrained least squares and the active set method*, SIAM Journal of Matrix Analysis and Applications (2008), to appear.
- [8] T. Kolda and D. O’Leary, *Algorithm 805: computation and uses of the semidiscrete matrix decomposition*, ACM TOMS **26** (2000), no. 3.
- [9] D. D. Lee and H. S. Seung, *Algorithms for Non-Negative Matrix Factorizations*, Advances in Neural Information Processing Systems **13** (2001), 556–562.
- [10] H. Park, M. Jeon, and J. Rosen, *Lower dimensional representation of text data based on centroids and least squares*, BIT **43** (2003).
- [11] M.F. Porter, *An algorithm for suffix stripping*, Program (1980), no. 3, 130–137.

- [12] G. Salton, C. Yang, and A. Wong, *A Vector-Space Model for Automatic Indexing*, Communications of the ACM **18** (1975), no. 11, 613–620.
- [13] S. Wild, J. Curry, and A. Dougherty, *Improving non-negative matrix factorizations through structured initialization*, Pattern Recognition **37** (2004), 2217–2232.
- [14] Y. Yang and C. Chute, *A linear least squares fit mapping method for information retrieval from natural language texts*, In 14th Conf. Comp. Linguistics, 1992.
- [15] D. Zeimpekis and E. Gallopoulos, *PDDP(l): Towards a Flexing Principal Direction Divisive Partitioning Clustering Algorithms*, Proc. IEEE ICDM '03 Workshop on Clustering Large Data Sets (Melbourne, Florida) (D. Boley, I. Dhillon, J. Ghosh, and J. Kogan, eds.), 2003, pp. 26–35.
- [16] D. Zeimpekis and E. Gallopoulos, *CLSI: A flexible approximation scheme from clustered term-document matrices*, In Proc. SIAM 2005 Data Mining Conf. (Newport Beach, California) (H. Kargupta, J. Srivastava, C. Kamath, and A. Goodman, eds.), April 2005, pp. 631–635.
- [17] D. Zeimpekis and E. Gallopoulos, *Linear and non-linear dimensional reduction via class representatives for text classification*, In Proc. of the 2006 IEEE International Conference on Data Mining (Hong Kong), December 2006, pp. 1172–1177.
- [18] D. Zeimpekis and E. Gallopoulos, *TMG: A MATLAB toolbox for generating term-document matrices from text collections*, Grouping Multidimensional Data: Recent Advances in Clustering (J. Kogan, C. Nicholas, and M. Teboulle, eds.), Springer, Berlin, 2006, pp. 187–210.
- [19] D. Zeimpekis and E. Gallopoulos, *k-means steering of spectral divisive clustering algorithms*, In Proc. of Text Mining Workshop (Minneapolis), 2007.
- [20] D. Zeimpekis and E. Gallopoulos, *Document clustering using nmf based on spectral information*, In Proc. of Text Mining Workshop (Atlanta), 2008.



## A Appendix: Demonstration of Use

### A.1 Indexing module (`tmg_gui`)

Assume we want to run `tmg.m` for the following input:

- filename: `sample_documents/sample1`
- delimiter: `emptyline`
- line\_delimiter: `yes`
- stoplist: `common_words`
- minimum length: `3`
- maximum length: `30`
- minimum local frequency: `1`
- maximum local frequency: `inf`
- minimum global frequency: `1`
- maximum global frequency: `inf`
- local term weighting: `logarithmic`
- global term weighting: `IDF`
- normalization: `cosine`
- stemming: `-`

and store results to directory “`sample1`” and MySQL.

1. Initially select the operation you want to perform, by pressing the corresponding radio button at the upper frame.
2. The selection of a radio button activates the required fields in the GUI, while deactivating the rest fields.

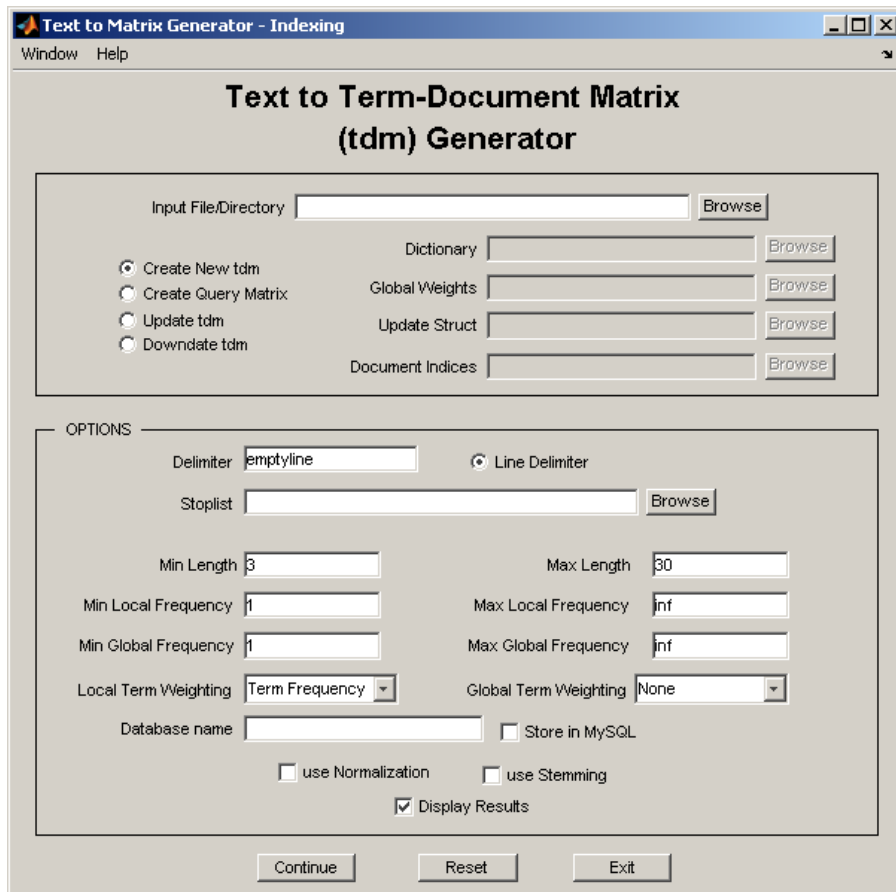


Figure 9: Starting window of tmg\_gui.

3. Fill in the required fields, by pressing the check buttons, editing the edit boxes or selecting the appropriate files/variables by pressing the “Browse” button.

Text to Matrix Generator - Indexing

Window Help

### Text to Term-Document Matrix (tdm) Generator

Input File/Directory:

Create New tdm  
 Create Query Matrix  
 Update tdm  
 Downdate tdm

Dictionary:    
Global Weights:    
Update Struct:    
Document Indices:

OPTIONS

Delimiter:   Line Delimiter

Stoplist:

Min Length:  Max Length:

Min Local Frequency:  Max Local Frequency:

Min Global Frequency:  Max Global Frequency:

Local Term Weighting:  Global Term Weighting:

Database name:   Store in MySQL

use Normalization  use Stemming

Display Results

Figure 10: Next view of `tmg_gui` according to the user selection.

4. The user can select a file or a variable by pressing the corresponding browse button.

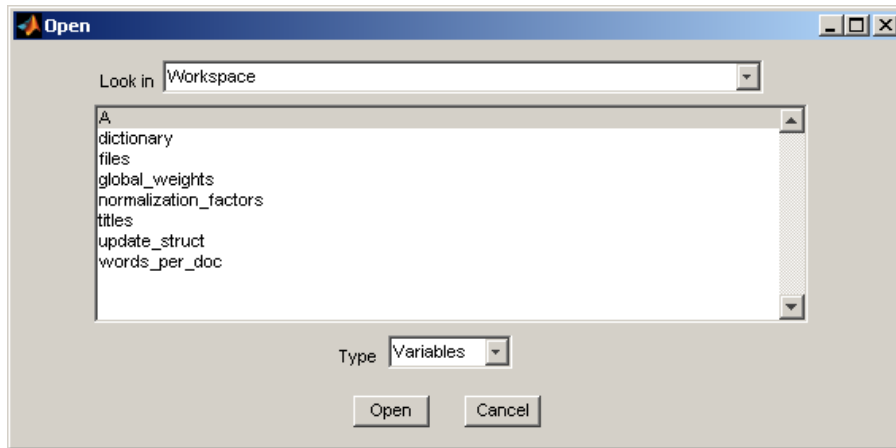


Figure 11: The open\_file window.

5. Press the “Continue” button in order to perform the selected operation.
6. Results have been saved to the workspace. Furthermore, directory “sample1” has been created under “TMG\_HOME/data” with each output variable stored to a single “.mat” file.

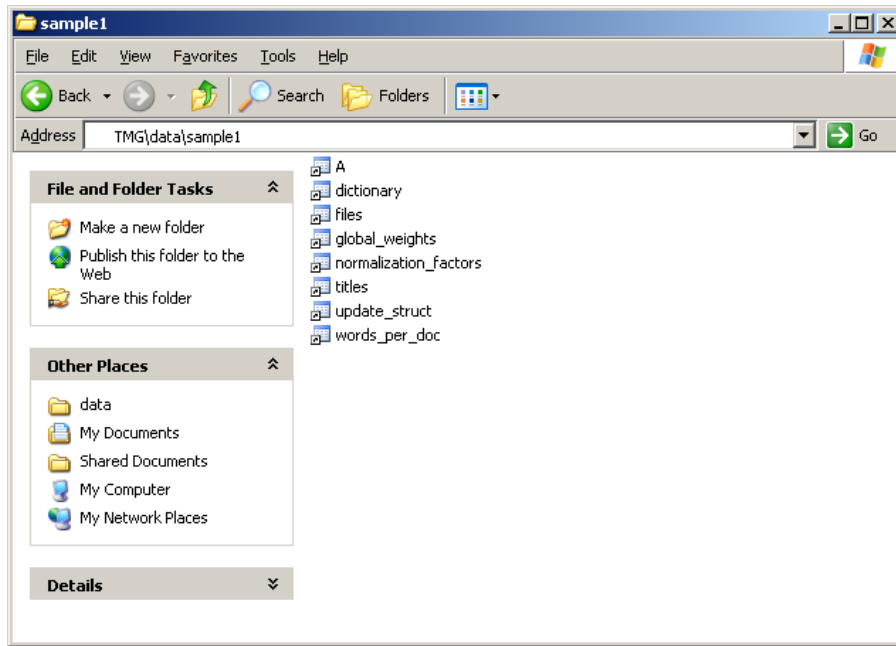
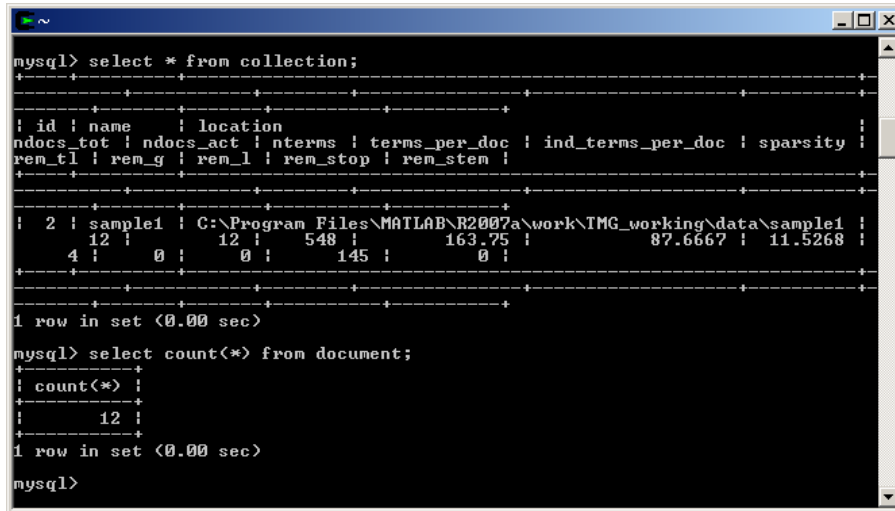


Figure 12: The output “.mat” files of `tmg_gui`.



7. Results have also been saved in MySQL (used for further processing, e.g. retrieval\_gui).



```
mysql> select * from collection;
+-----+-----+-----+-----+-----+-----+-----+
| id | name | location | ndocs_tot | ndocs_act | nterms | terms_per_doc | ind_terms_per_doc | sparsity | rem_tl | rem_g | rem_l | rem_stop | rem_sten |
+-----+-----+-----+-----+-----+-----+-----+
| 2 | sample1 | C:\Program Files\MATLAB\R2007a\work\IMG_working\data\sample1 | 12 | 12 | 548 | 163.75 | 87.6667 | 11.5268 | 4 | 0 | 0 | 145 | 0 |
+-----+-----+-----+-----+-----+-----+
1 row in set (0.00 sec)

mysql> select count(*) from document;
+-----+
| count(*) |
+-----+
| 12 |
+-----+
1 row in set (0.00 sec)

mysql>
```

Figure 13: The MySQL view upon tmg execution.

8. Press the “Reset” button in order to change the input.

- For further documentation type “help tmg\_gui” at the MATLAB command window, or select the “Documentation” tab from the “Help” menu.

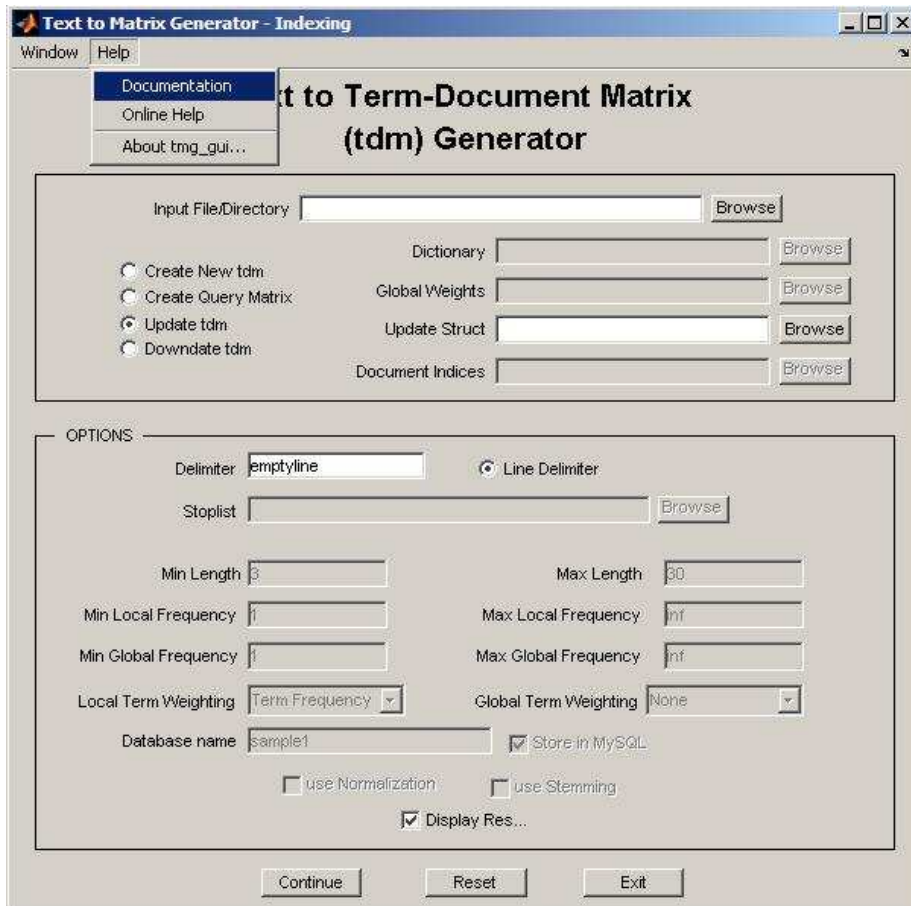


Figure 14: The GUIs' general help tab.

- In order to update a tdm, give the “input file/directory” and the update\_struct corresponding to the initial collection. In case you just want to alter some options, give a blank “input file/directory” and change the corresponding fields of update\_struct.
- In order to downdate a tdm, give the update\_struct corresponding to the initial collection and the document indices vector you want to remove.
- In order to construct a term-query matrix, give the dictionary char array of the initial collection and the corresponding vector of global weights (optional).

## A.2 Dimensionality Reduction module (`dr_gui`)

Suppose we have processed a collection with `tmg_gui`, construct a `tdm` with 1,033 documents and 12,184 terms (corresponding to the well-known MEDLINE collection) and store the results to “`TMG_HOME/data/medline`”. Assume then, we want to construct a low-rank approximation of the TDM, using the Clustered Latent Semantic Indexing (CLSI) technique for the following input:

- compute SVD with: Propack
- clustering algorithm: PDDP
- principal directions: 1
- maximum number of PCs: -
- variant: basic
- automatic determination of num. of factors from each cluster: yes
- number of clusters: 10
- number of factors: 100

and you want to store results to directory “`medline`”.

1. Initially select the operation you want to perform, by pressing the corresponding radio button at the upper left frame.
2. The selection of a radio button activates the required fields in the GUI, while deactivating the rest fields.

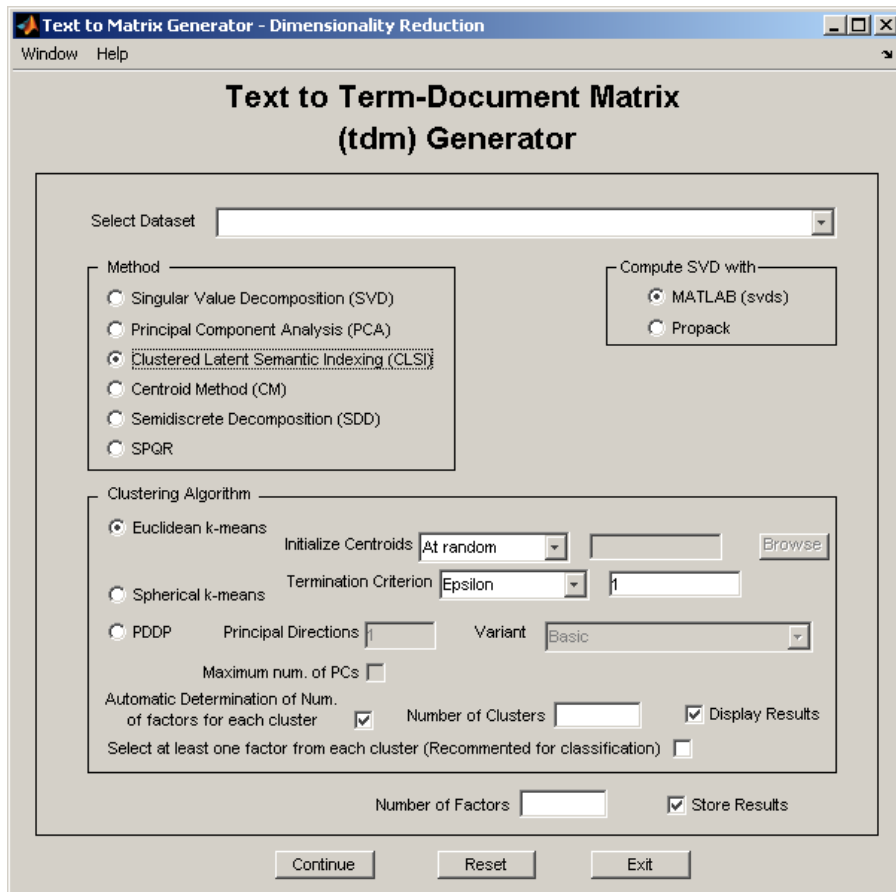


Figure 15: Starting window of dr\_gui.

3. Fill in the required fields, by pressing the check buttons, editing the edit boxes or selecting the appropriate files/variables by pressing the “Browse” button.

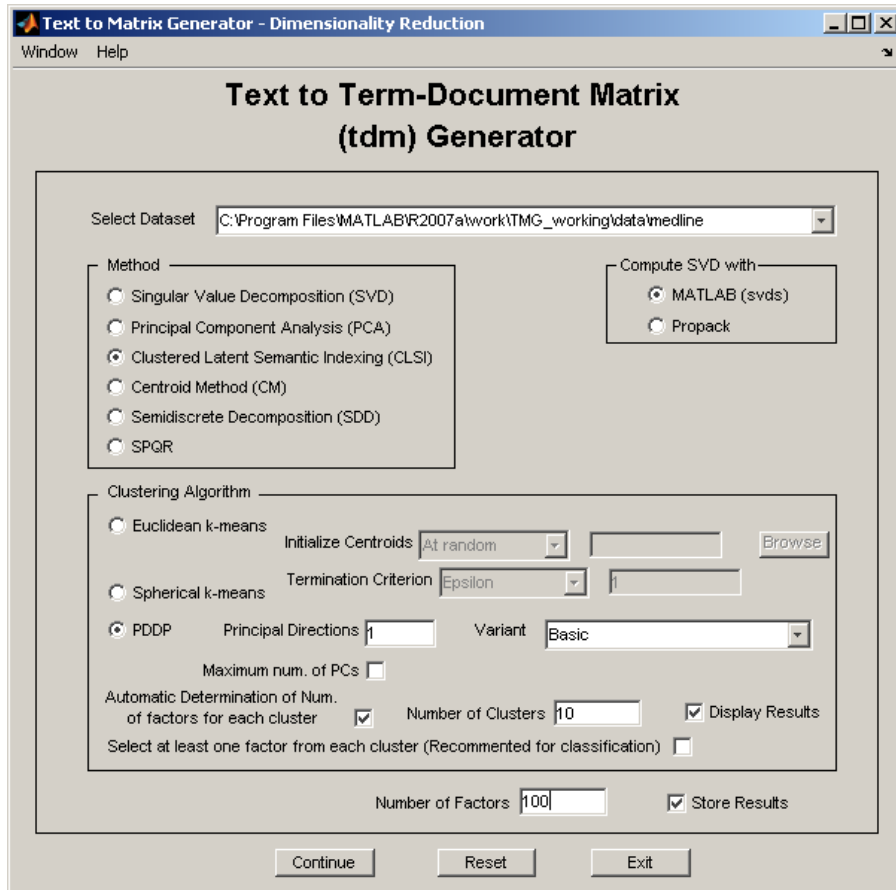


Figure 16: Next view of dr\_gui according to the user selection.

4. Press the “Continue” button in order to perform selected operation.
5. Results have been saved to the workspace. Furthermore, directory “clsi/k\_100” has been created under “TMG\_HOME/data/medline” with each output variable stored to a single “.mat” file.

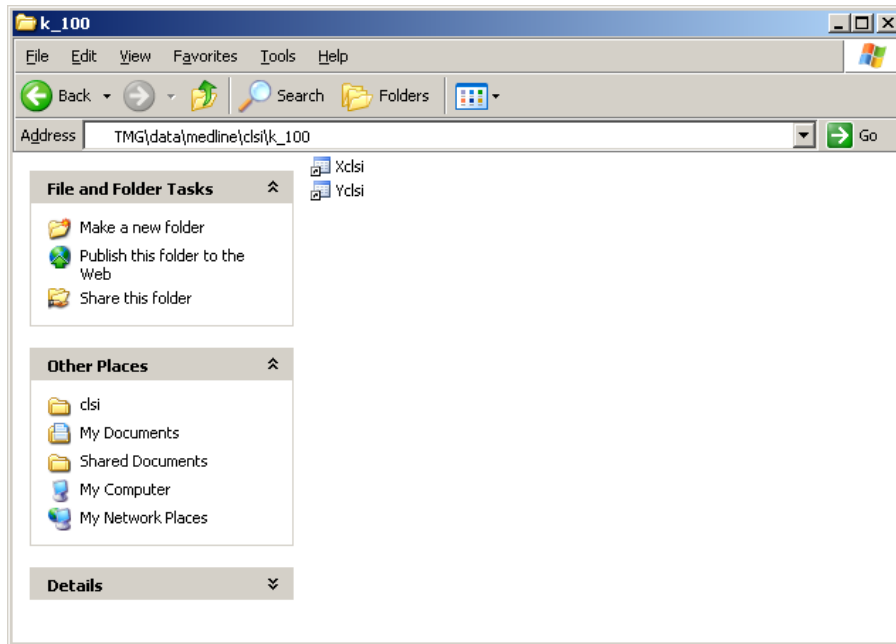


Figure 17: The output “.mat” files of `dr_gui`.

6. Press the “Reset” button in order to change the input.

### A.3 Non-Negative Factorizations module (nnmf\_gui)

Assume we have processed a collection with `tmg_gui`, construct a `tdm` with 1,033 documents and 12,184 terms (corresponding to the well-known MEDLINE collection) and store the results to “`TMG_HOME/data/medline`”. Assume then, we want to construct a non-negative factorization of the TDM, using the Multiplicative Update algorithm initializing by the block NNDSVD technique for the following input:

- initialization: Block NNDSVD
- refine factors: yes
- method: Multiplicative update
- number of iterations: 10
- compute SVD with: Propack
- clustering algorithm: PDDP
- principal directions: 1
- maximum number of PCs: -
- variant: basic
- number of clusters: 10
- number of factors: 10

and you want to store results to directory “`medline`”.

1. Initially select the operation you want to perform, by pressing the corresponding radio button at the upper left frame.
2. The selection of a radio button activates the required fields in the GUI, while deactivating the rest fields.

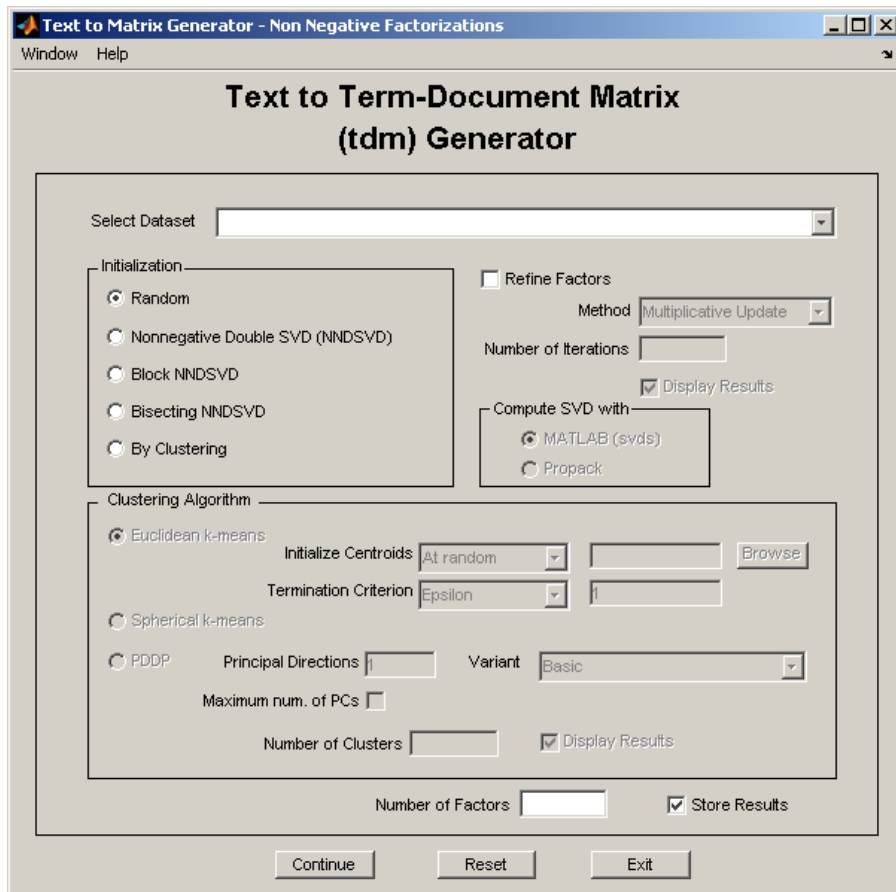


Figure 18: Starting window of nnmf\_gui.



3. Fill in the required fields, by pressing the check buttons, editing the edit boxes or selecting the appropriate files/variables by pressing the “Browse” button.

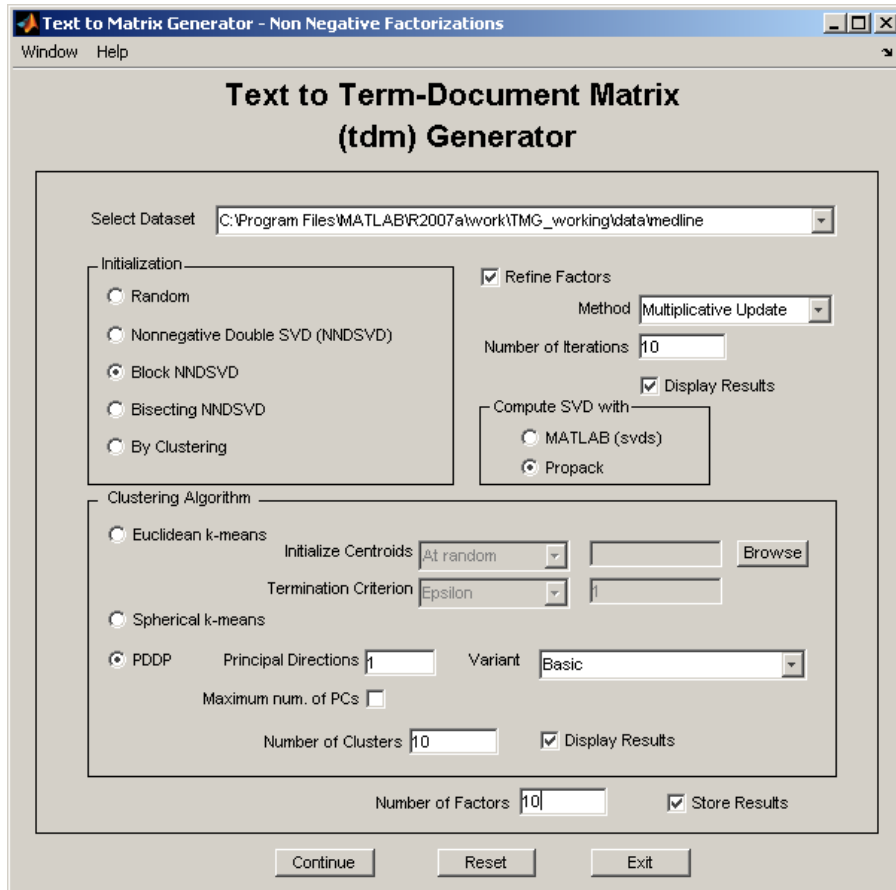


Figure 19: Next view of nmmf\_gui according to the user selection.

4. Press the “Continue” button in order to perform selected operation.
5. Results have been saved to the workspace. Furthermore, directory “nmmf/k\_10/mlup” has been created under “TMG\_HOME/data/medline” with each output variable stored to a single “.mat” file.

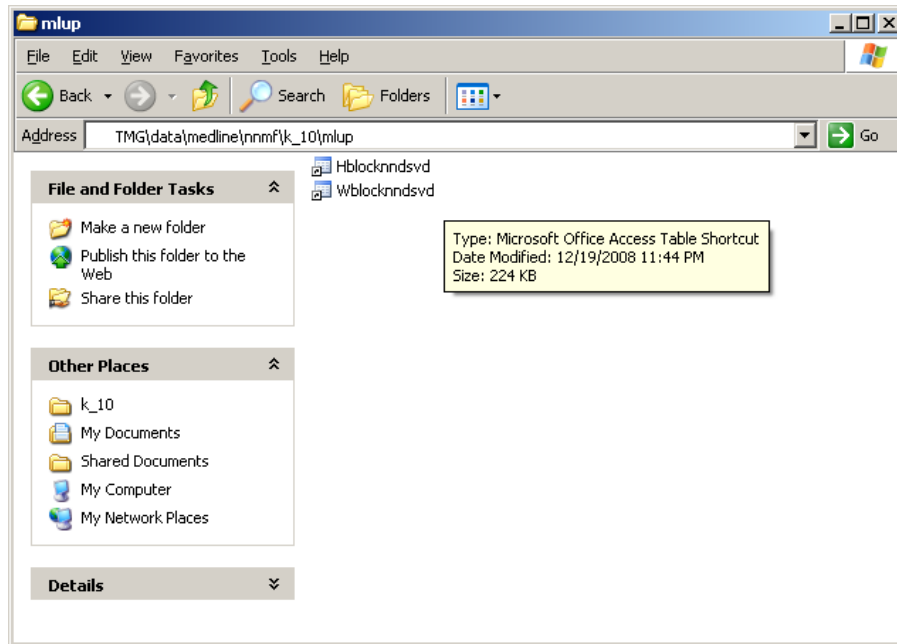


Figure 20: The output “.mat” files of nmmf\_gui.

6. Press the “Reset” button in order to change the input.

#### A.4 Retrieval module (`retrieval_gui`)

Suppose we have processed a collection with `tmg_gui`, construct a `tdm` with 1,033 documents and 12,184 terms (corresponding to the well-known MEDLINE collection) and store the results to “`TMG_HOME/data/medline`”. Assume then, we want to retrieve the relevant documents to a specific query for the following input:

- insert query: ‘the crystalline lens in vertebrates, including humans’
- use stored global weights: yes
- stoplist: `common_words`
- local term weighting: Term Frequency
- latent semantic analysis by: Clustered Latent Semantic Indexing
- number of factors: 100
- similarity measure: Cosine
- number of most relevant documents: 5

1. Initially select the retrieval method you want to apply, by pressing the corresponding radio button.
2. The selection of a radio button activates the required fields in the GUI, while deactivating the rest fields.

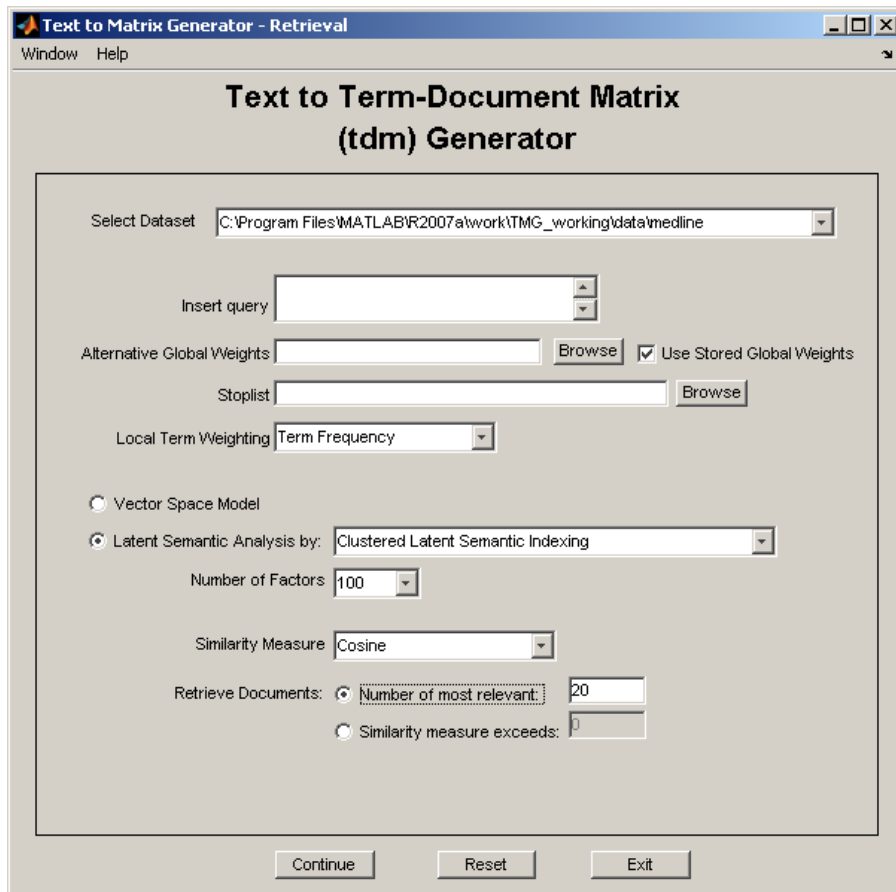


Figure 21: Starting window of retrieval\_gui.

3. Fill in the required fields, by pressing the check buttons, editing the edit boxes or selecting the appropriate files/variables by pressing the “Browse” button.

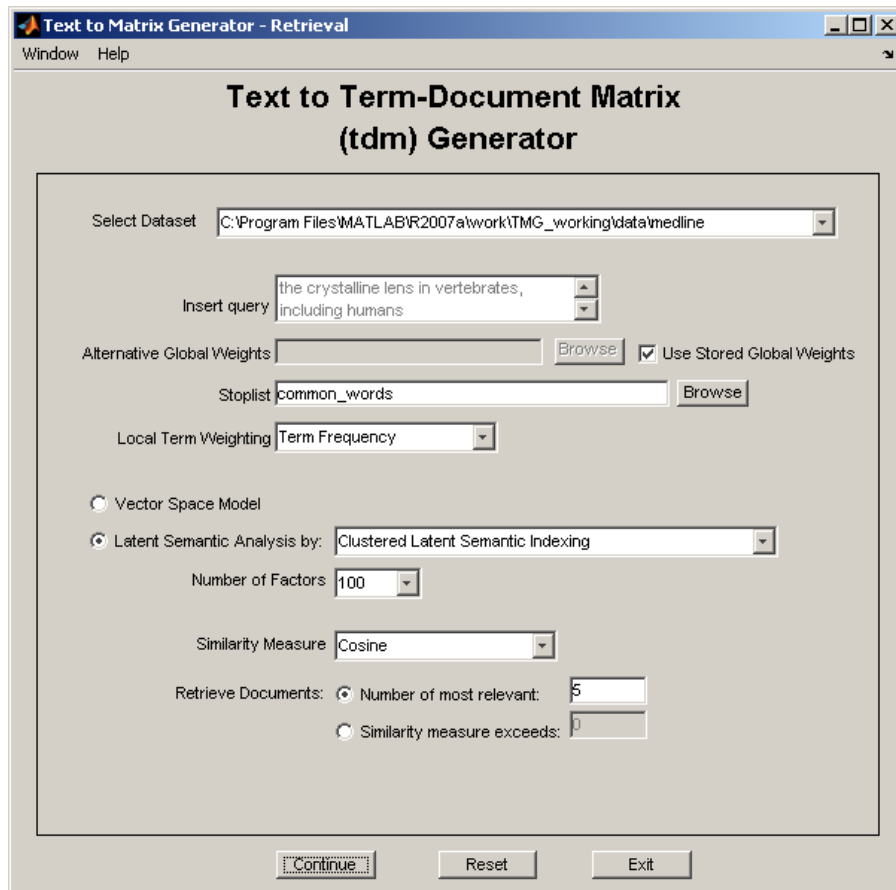


Figure 22: Next view of retrieval\_gui according to the user selection.

4. Press the “Continue” button in order to perform selected operation.
5. Results have been saved to the workspace.
6. Furthermore, in case data have been stored to MySQL, the user gets an html response.

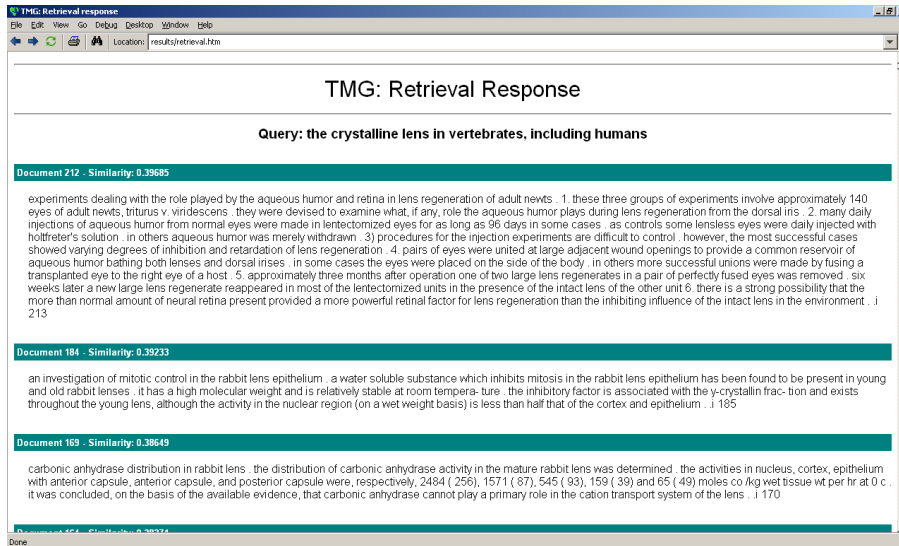


Figure 23: The output of `retrieval_gui`.

7. Press the “Reset” button in order to change the input.

## A.5 Clustering module (clustering\_gui)

Suppose we have processed a collection with `tmg_gui`, construct a `tdm` with 1,033 documents and 12,184 terms (corresponding to the well-known MEDLINE collection) and store the results to “`TMG_HOME/data/medline`”. Assume then, we want to cluster the TDM, using the PDDP clustering algorithm with the following input:

- principal directions: 1
- maximum number of PCs: -
- variant: basic
- number of clusters: 5

and you want to store results to directory “`medline`”.

1. Initially select the clustering algorithm you want to apply, by pressing the corresponding radio button.
2. The selection of a radio button activates the required fields in the GUI, while deactivating the rest fields.

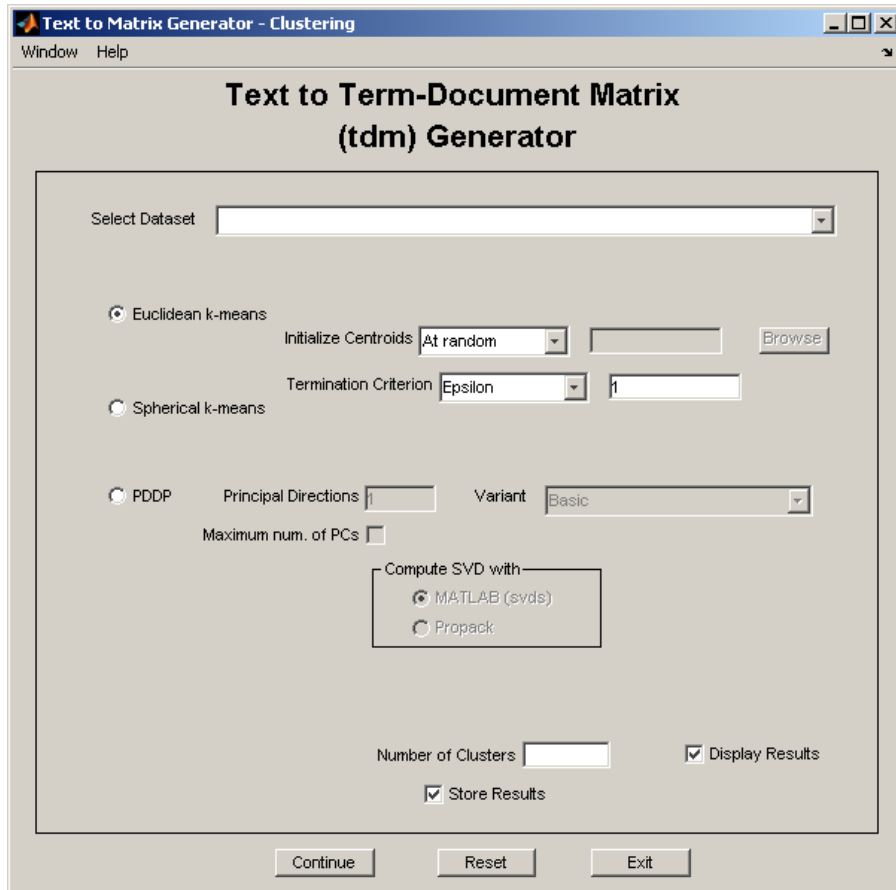


Figure 24: Starting window of clustering\_gui.



3. Fill in the required fields, by pressing the check buttons, editing the edit boxes or selecting the appropriate files/variables by pressing the “Browse” button.

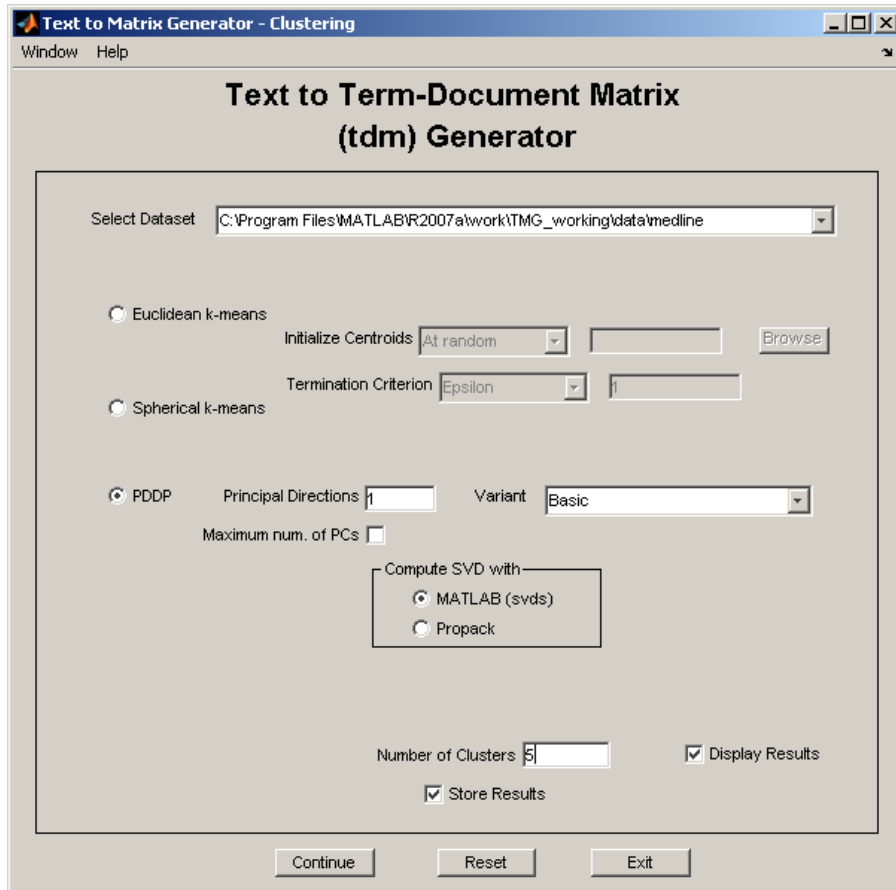


Figure 25: Next view of clustering\_gui according to the user selection.

4. Press the “Continue” button in order to perform selected operation.
5. Results have been saved to the workspace. Furthermore, directory “kmeans/k\_10” has been created under “TMG\_HOME/data/medline” with each output variable stored to a single “.mat” file.

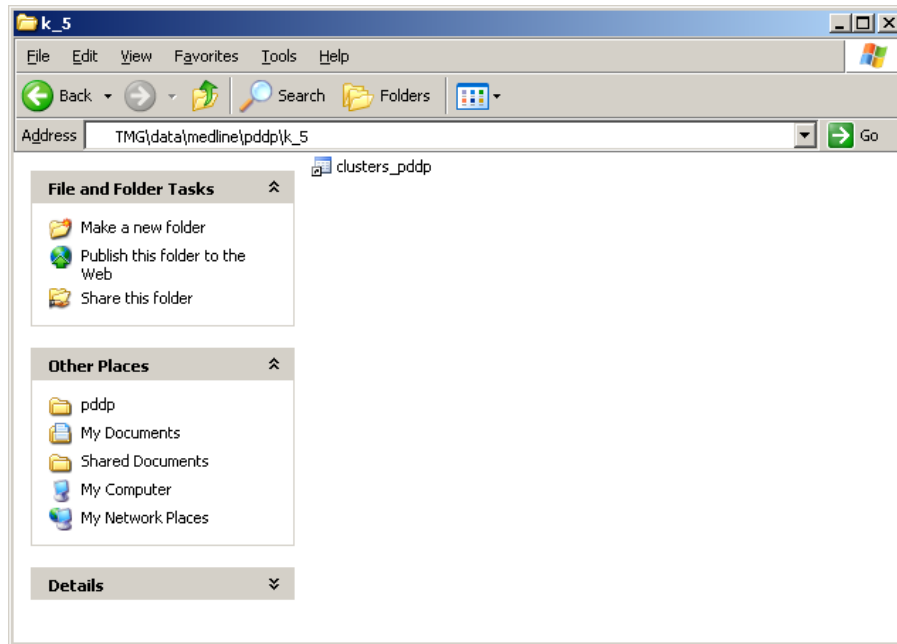


Figure 26: The output “.mat” files of clustering\_gui.

6. The user gets an html response that summarizes the clustering result.

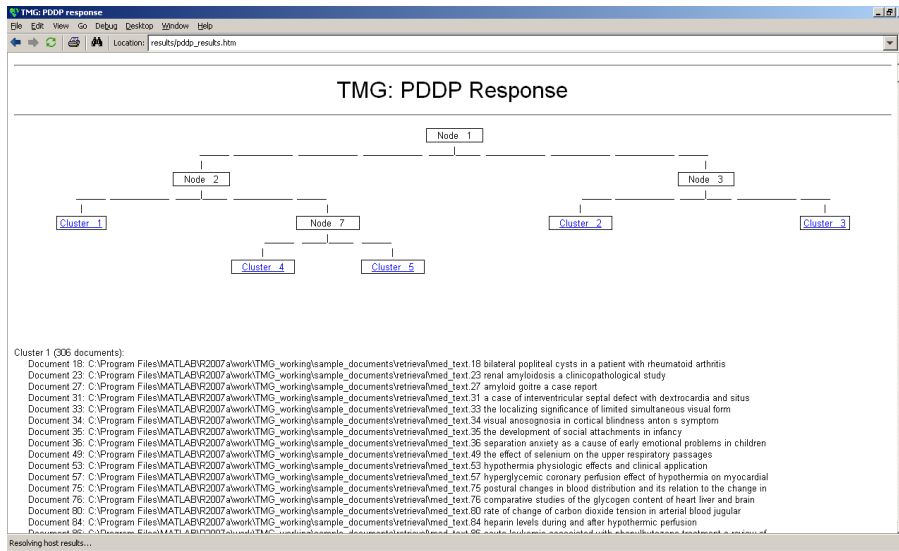


Figure 27: The output of clustering\_gui for PDDP.

7. Press the "Reset" button in order to change the input.

## A.6 Classification module (classification\_gui)

Suppose we have processed a collection with `tmg_gui`, construct a `tdm` with 6,495 documents and 21,764 terms (a single label dataset corresponding to the well-known modapte split of the Reuters-21578 collection) and store the results to “`TMG_HOME/data/reuters`”. Assume then, we want to classify the test part of the modapte split, using the k-Nearest Neighbors classifier for the following input:

- Multiple docs (file): yes
- filename: `sample_document/reuters.test`
- delimiter: `</reuters>`
- line delimiter: yes
- use stored global weights: yes
- stoplist: `common_words`
- local term weighting: Term Frequency
- classification method: k Nearest Neighbors (kNN)
- num. of NNs: 10
- collection type: Single-Label
- preprocessed by: Clustered Latent Semantic Indexing
- number of factors: 100
- similarity measure: Cosine

1. Initially select the classification algorithm you want to apply, by pressing the corresponding radio button at the left frame.
2. The selection of a radio button activates the required fields in the GUI, while deactivating the rest fields.

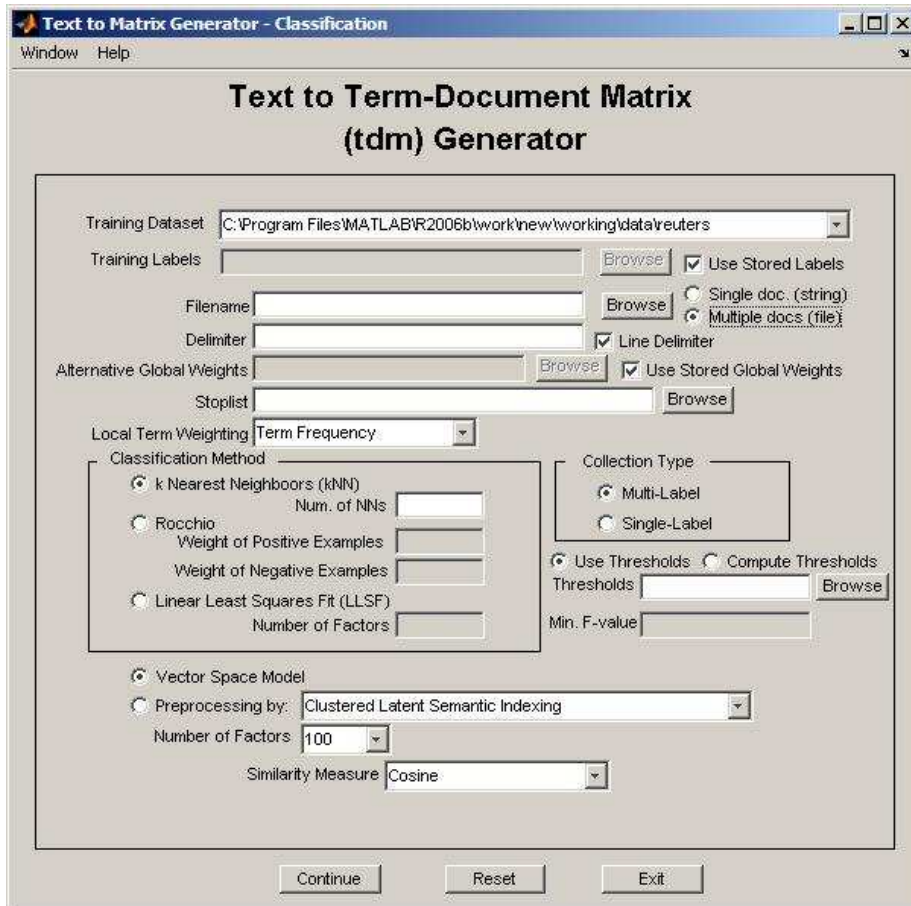


Figure 28: Starting window of classification\_gui.

3. Fill in the required fields, by pressing the check buttons, editing the edit boxes or selecting the appropriate files/variables by pressing the "Browse" button.

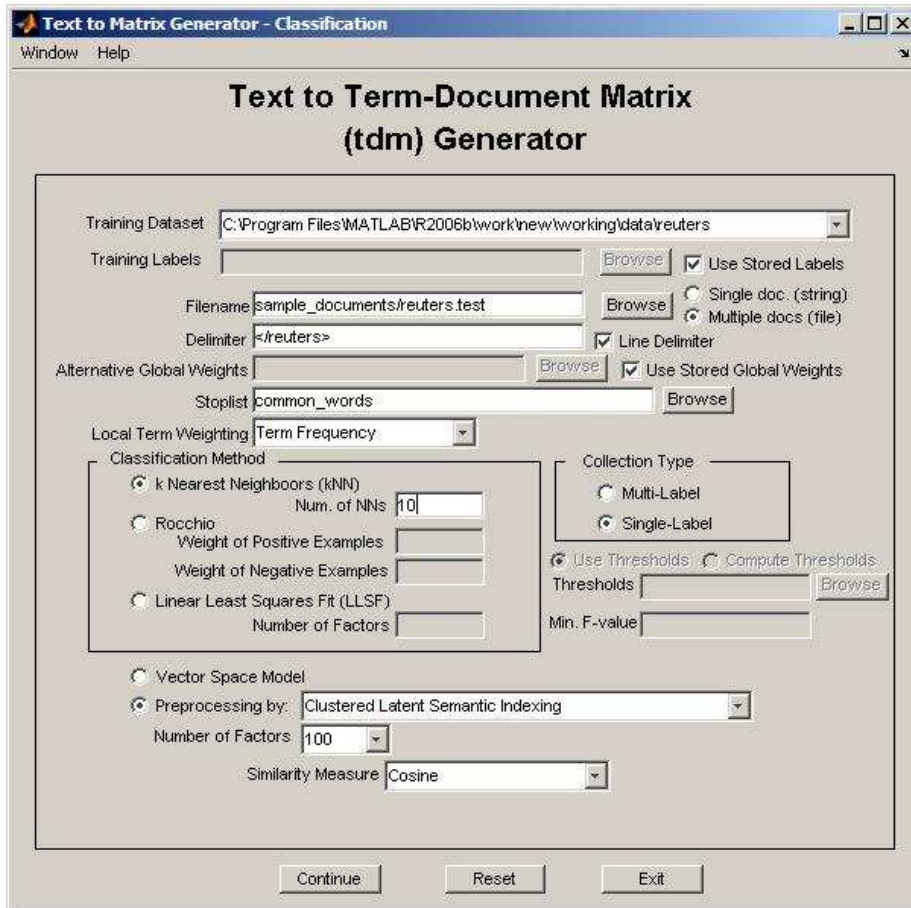


Figure 29: Next view of classification\_gui according to the user selection.

4. Press the “Continue” button in order to perform selected operation.
5. Results have been saved to the workspace.
6. Press the “Reset” button in order to change the input.

## B Appendix: Function Reference

about_tmg_gui
ABOUT_TMG_GUI ABOUT_TMG_GUI displays information for TMG.



bisecting\_nndsvd

**BISECTING\_NNDSVD** - a bisecting form of the the Non-Negative Double Singular Value Decomposition Method [2].

**BISECTING\_NNDSVD** applies a bisecting form of the the Non-Negative Double Singular Value Decomposition Method [2] using a PDDP-like [2] clustering Method.

[W, H]=**BISECTING\_NNDSVD**(A, k, svd\_method) returns a non-negative rank-k approximation of the input matrix A using svd\_method for the SVD (possible values svds, propack).

**REFERENCES:**

[1] D.Boley, Principal Direction Divisive Partitioning, Data Mining and Knowledge Discovery 2 (1998), no. 4, 325-344.

[2] C. Boutsidis and E. Gallopoulos. SVD-based initialization: A head start on nonnegative matrix factorization. Pattern Recognition, Volume 41, Issue 4, Pages 1350-1362, April 2008.

block.diagonalize

**BLOCK\_DIAGONALIZE** - reorders a matrix heuristically using a clustering result

[A, N\_ROWS, N\_COLS, ROW\_INDS, COL\_INDS]=**BLOCK\_DIAGONALIZE**(A, CLUSTERS) reorders matrix A using the clustering result represented by the structure CLUSTERS. N\_ROWS and N\_COLS store the last row and column index for each row and column block respectively, while ROW\_INDS and COL\_INDS contain the permuted row and column indices.

## block\_nndsvd

**BLOCK\_NNDSVD** - computes a non-negative rank-L approximation of the input matrix using the Clustered Latent Semantic Indexing Method [2] and the Non-Negative Double Singular Value Decomposition Method [1].

$[X, Y]=\text{BLOCK\_NNDSVD}(A, \text{CLUSTERS}, L, \text{FUNC}, \text{ALPHA\_VAL}, \text{SVD\_METHOD})$  computes a non-negative rank-L approximation  $X*Y$  of the input matrix  $A$  with the Clustered Latent Semantic Indexing Method [2], and the Non-Negative Double Singular Value Decomposition Method [1], using the cluster structure information from **CLUSTERS** [3]. **FUNC** denotes the method used for the selection of the number of factors from each cluster. Possible values for

**FUNC**:

- 'f': Selection using a heuristic method from [2] (see **KS\_SELECTION**).
- 'f1': Same as 'f' but use at least one factor from each cluster.
- 'equal': Use the same number of factors from each cluster.

**ALPHA\_VAL** is a value in  $[0, 1]$  used in the number of factors selection heuristic [2]. Finally, **SVD\_METHOD** defines the method used for the computation of the SVD (svds or propack).

### REFERENCES:

- [1] C. Boutsidis and E. Gallopoulos. SVD-based initialization: A head start on nonnegative matrix factorization. *Pattern Recognition*, Volume 41, Issue 4, Pages 1350-1362, April 2008.
- [2] D. Zeimpekis and E. Gallopoulos. CLSI: A Flexible Approximation Scheme from Clustered Term-Document Matrices. In *Proc. 5th SIAM International Conference on Data Mining*, pages 631635, Newport Beach, California, 2005.
- [3] D. Zeimpekis and E. Gallopoulos. Document Clustering using NMF based on Spectral Information. In *Proc. Text Mining Workshop 2008 held in conjunction with the 8th SIAM International Conference on Data Mining*, Atlanta, 2008.

`classification_gui`

**CLASSIFICATION\_GUI**

CLASSIFICATION\_GUI is a graphical user interface for all classification functions of the Text to Matrix Generator (TMG) Toolbox.

## clsi

CLSI - computes a rank-L approximation of the input matrix using the Clustered Latent Semantic Indexing Method [1]

$[X, Y]=\text{CLSI}(A, \text{CLUSTERS}, L, \text{FUNC}, \text{ALPHA\_VAL}, \text{SVD\_METHOD})$

computes the rank-L approximation  $X*Y$  of the input matrix  $A$  with the Clustered Latent Semantic Indexing Method [1], using the cluster structure information from  $\text{CLUSTERS}$ .

$\text{FUNC}$  denotes the method used for the selection of the number of factors from each cluster. Possible values for  $\text{FUNC}$ :

- 'f': Selection using a heuristic method from [1] (see  $\text{KS\_SELECTION}$ ).
- 'f1': Same as 'f' but use at least one factor from each cluster.
- 'equal': Use the same number of factors from each cluster.

$\text{ALPHA\_VAL}$  is a value in  $[0, 1]$  used in the number of factors selection heuristic [1]. Finally,  $\text{SVD\_METHOD}$  defines the method used for the computation of the SVD (svds or propack).

### REFERENCES:

[1] D. Zeimpekis and E. Gallopoulos. CLSI: A Flexible Approximation Scheme from Clustered Term-Document Matrices. In Proc. 5th SIAM International Conference on Data Mining, pages 631635, Newport Beach, California, 2005.

clustering\_gui

**CLUSTERING\_GUI**

CLUSTERING\_GUI is a graphical user interface for all clustering functions of the Text to Matrix Generator (TMG) Toolbox.

cm

CM - computes a rank-L approximation of the input matrix using the Centroids Method [1]

[X, Y]=CM(A, CLUSTERS) computes the rank-K approximation  $X*Y$  of the input matrix A with the Centroids Method [1], using the cluster structure information from CLUSTERS.

REFERENCES:

[1] H. Park, M. Jeon, and J. Rosen. Lower Dimensional Representation of Text Data Based on Centroids and Least Squares. BIT Numerical Mathematics, 43(2):427448, 2003.

col\_normalization

COL\_NORMALIZATION - normalizes the columns of the input matrix.



col\_rearrange

**COL\_REARRANGE** - reorders a matrix using a clustering result  
[A, N\_COLS, COL\_INDS]=COL\_REARRANGE(A, CLUSTERS) reorders  
the columns of matrix A using the clustering result represented  
by the structure CLUSTERS. N\_COLS stores the last column index  
for each column block, while COL\_INDS contains the permuted  
column indices.

column\_norms

COLUMN\_NORMS - returns the column norms of a matrix

`compute_fro_norm`

`COMPUTE_FRO_NORM` - returns the frobenius norm of a rank-1 matrix  $A = W * H$ .

`compute_scat`

**COMPUTE\_SCAT** - computes the cluster selection criterion value of PDDP

**SCAT=COMPUTE\_SCAT(A, C)** returns the square of the frobenius norm of  $A-C*\text{ones}(1, \text{size}(A, 2))$ .

`create_kmeans_response`

`CREATE_KMEANS_RESPONSE` returns an html response for k-means

`CREATE_KMEANS_RESPONSE(CLUSTERS, TITLES)` creates a summary html file containing information for the result of the k-means algorithm, defined by `CLUSTERS`, when applied to the dataset with document titles defined in the `TITLES` cell array.

`CREATE_KMEANS_RESPONSE(CLUSTERS, TITLES, VARIANT)` defines additionally the k-means variant (possible values 'k-means' and 'skmeans'). The result is stored in the "results" directory and displayed using the default web browser.

`create_pddp_response`

`CREATE_PDDP_RESPONSE` returns an html response for PDDP  
`CREATE_PDDP_RESPONSE(TREE_STRUCT, CLUSTERS, L, TITLES)`  
creates a summary html file containing information for  
the result of the PDDP algorithm, defined by `TREE_STRUCT`  
and `CLUSTERS`, when applied to the dataset with document  
titles defined in the `TITLES` cell array. `L` defines the  
maximum number of principal directions used by PDDP.  
The result is stored in the "results" directory and  
displayed using the default web browser.

`create_retrieval_response`

`CREATE_RETRIEVAL_RESPONSE` returns an html response for a query  
`CREATE_RETRIEVAL_RESPONSE(DATASET, IDS, SIMILARITY, QUERY)`  
creates an html file containing information for the text of  
documents of `DATASET` stored in MySQL defined by `IDS` and  
having `SIMILARITY` similarity coefficients against `QUERY`.  
The result is stored in the "results" directory and displayed  
using the default web browser.

`diff_vector`

**DIFF\_VECTOR**

**DIFF\_VECTOR** returns the vector of differences between consecutive elements of the input vector.



dr\_gui

**DR\_GUI**

DR\_GUI is a graphical user interface for all dimensionality reduction functions of the Text to Matrix Generator (TMG) Toolbox.

## ekmeans

### EKMEANS - Euclidean k-Means Clustering Algorithm

EKMEANS clusters a term-document matrix using the standard k-means clustering algorithm. `CLUSTERS=EKMEANS(A, C, K, TERMINATION)` returns a cluster structure with K clusters for the term-document matrix A using as initial centroids the columns of C (initialized randomly when it is empty).

TERMINATION defines the termination method used in k-means ('epsilon' stops iteration when objective function decrease falls down a user defined threshold - see OPTIONS input argument - while 'n\_iter' stops iteration when a user defined number of iterations has been reached).

`[CLUSTERS, Q]=EKMEANS(A, C, K, TERMINATION)` returns also the vector of objective function values for each iteration and `[CLUSTERS, Q, C]=EKMEANS(A, C, K, TERMINATION)` returns the final centroid vectors.

`EKMEANS(A, C, K, TERMINATION, OPTIONS)` defines optional parameters:

- `OPTIONS.iter`: Number of iterations (default 10).
- `OPTIONS.epsilon`: Value for epsilon convergence criterion (default 1).
- `OPTIONS.dsp`: Displays results (default 1) or not (0) to the command window.

## entropy

ENTROPY - computes the entropy of a clustering result  
[VENTROPY, CONFUSION\_MATRIX, MISTAKES]=ENTROPY(CLUSTERS,  
LABELS) computes the entropy value of a clustering result  
represented by the CLUSTERS structure. LABELS is a vector  
of integers containing the true labeling of the objects.  
The entropy value is stored in VENTROPY, while  
CONFUSION\_MATRIX is a  $k \times r$  matrix, where  $k$  is the number  
of clusters and  $r$  the number of true classes, and  
CONFUSION\_MATRIX( $i, j$ ) records the number of objects  
of class  $j$  assigned to cluster  $i$ . Finally, MISTAKES contains  
the number of misassigned objects, measured by  $m_1 + \dots + m_k$ ,  
where  $m_i = \sum(\text{CONFUSION\_MATRIX}(i, j)), j \neq i$ .

get\_node\_scat

GET\_NODE\_SCAT - returns the PDDP node with the maximum scatter value (see PDDP)

[MAX\_SCAT\_IND, M\_SCAT]=GET\_NODE\_SCAT(TREE\_STRUCT, SPLITTED)

returns the node index and the scatter value of the PDDP

tree defined by TREE\_STRUCT. SPLITTED is a vector that

determines the active nodes.

gui

**GUI**

GUI is a simple, top graphical user interface of the Text to Matrix Generator (TMG) Toolbox. Using GUI, the user can select any of the four GUI modules (indexing, dimensionality reduction, clustering, classification) of TMG.

`init_tmg`

**INIT\_TMg - Installation script of TMG**

INIT\_TMg is the installation script of the Text to Matrix Generator (TMG) Toolbox. INIT\_TMg creates the MySQL database and adds all TMG directories to the path.

knn.multi

**KNN\_MULTI** - k-Nearest Neighbors classifier for multi-label collections

**LABELS\_AS=KNN\_MULTI(A, Q, K, LABELS, NORMALIZED\_DOCS, THRESHOLDS)** classifies the columns of Q with the K-Nearest Neighbors classifier using the pre-classified columns of matrix A with labels LABELS (cell array of vectors of integers). THRESHOLDS is a vector of class threshold values. NORMALIZED\_DOCS defines if cosine (1) or euclidean distance (0) similarity measure is to be used. LABELS\_AS contains the assigned labels for the columns of Q.

knn\_single

**KNN\_SINGLE** - k-Nearest Neighbors classifier for single-label collections

**LABELS\_AS=KNN\_SINGLE(A, Q, K, LABELS, NORMALIZED\_DOCS)**  
classifies the columns of Q with the K-Nearest Neighbors classifier using the pre-classified columns of matrix A with labels LABELS (vector of integers). NORMALIZED\_DOCS defines if cosine (1) or euclidean distance (0) similarity measure is to be used. LABELS\_AS contains the assigned labels for the columns of Q.



ks\_selection

KS\_SELECTION - implements the heuristic method from [2] for the selection of the number of factors from each cluster used in the Clustered Latent Semantic Indexing method [1].

N\_ST=KS\_SELECTION(A, N\_COLS, ALPHA\_VAL, L) returns in N\_ST a vector of integers denoting the number of factors (sum equals L) selected from each cluster of the tdm A. N\_COLS is a vector containing the last column index for each column block, while ALPHA\_VAL is a value in [0, 1].

ks\_selection1

**KS\_SELECTION1** - implements the heuristic method from [2] for the selection of the number of factors from each cluster used in the Clustered Latent Semantic Indexing method [1]. The number of factors from each cluster is at least 1.

**N\_ST=KS\_SELECTION1(A, N\_COLS, ALPHA\_VAL, L)** returns in **N\_ST** a vector of integers denoting the number of factors

(sum equals L) selected from each cluster of the tdm A.

**N\_COLS** is a vector containing the last column index for each column block, while **ALPHA\_VAL** is a value in [0, 1].

llsf\_multi

LLSF\_MULTI - Linear Least Squares Fit for multi-label collections [2]

LABELS\_AS=LLSF\_MULTI(A, Q, CLUSTERS, LABELS, L, METHOD, THRESHOLDS, SVD\_METHOD, CLSI\_METHOD) classifies the columns of Q with the Linear Least Squares Fit classifier [2] using the pre-classified columns of matrix A with labels LABELS (cell array of vectors of integers).

THRESHOLDS is a vector of class threshold values, while CLUSTERS is a structure defining the classes. METHOD is the method used for the approximation of the rank-l truncated SVD, with possible values:

- 'clsi': Clustered Latent Semantic Indexing [3].
- 'cm': Centroids Method [1].
- 'svd': Singular Value Decomposition.

SVD\_METHOD defines the method used for the computation of the SVD, while CLSI\_METHOD defines the method used for the determination of the number of factors from each class used in Clustered Latent Semantic Indexing in case METHOD equals 'clsi'.

## llsf\_single

LLSF\_SINGLE - Linear Least Squares Fit for single-label collections [2]

LABELS\_AS=LLSF\_SINGLE(A, Q, CLUSTERS, LABELS, L, METHOD, SVD.METHOD, CLSI.METHOD) classifies the columns of Q with the Linear Least Squares Fit classifier [2] using the pre-classified columns of matrix A with labels LABELS (cell array of vectors of integers). CLUSTERS is a structure defining the classes. METHOD is the method used for the approximation of the rank-1 truncated SVD, with possible values:

- 'clsi': Clustered Latent Semantic Indexing [3].
- 'cm': Centroids Method [1].
- 'svd': Singular Value Decomposition.

SVD.METHOD defines the method used for the computation of the SVD, while CLSI.METHOD defines the method used for the determination of the number of factors from each class used in Clustered Latent Semantic Indexing in case METHOD equals 'clsi'.

lsa

LSA - Applies the Latent Semantic Analysis Model to a document collection

[SC, DOCS\_INDS] = LSA(D, P, Q, NORMALIZE\_DOCS) applies

LSA to the text collection represented by the latent semantic factors D, P of the collection's term - document matrix, for the query defined by the vector Q [1].

NORMALIZE\_DOCS defines if the document vectors are to be normalized (1) or not (0). SC contains the sorted similarity coefficients, while DOC\_INDS contains the corresponding document indices.

`make_clusters_multi`

**MAKE\_CLUSTERS\_MULTI** - auxiliary function for the classification algorithms

**CLUSTERS=MAKE\_CLUSTERS\_MULTI(LABELS)** forms the cluster structure of a multi-label collection with document classes defined by LABELS (cell array of vectors of integers).

`make_clusters_single`

**MAKE\_CLUSTERS\_SINGLE** - auxiliary function for the classification algorithms

**CLUSTERS=MAKE\_CLUSTERS\_SINGLE(LABELS)** forms the cluster structure of a single-label collection with document classes defined by LABELS (vector of integers).

make\_labels

**MAKE\_LABELS** - creates a label vector of integers for the input cell array of string  
[LABELS, UNIQUE\_LABELS]=MAKE\_LABELS(INPUT\_LABELS) creates a vector of integer labels (LABELS) for the input cell array of strings INPUT\_LABELS. UNIQUE\_LABELS contains the strings of unique labels of the input cell array.



make\_val\_inds

MAKE\_VAL\_INDS - auxiliary function for the classification algorithms

INDS=MAKE\_VAL\_INDS(LABELS) constructs an index vector used during the thresholding phase of any classifier for the multi-label collection with document classes defined by LABELS (cell array of vectors of integers).

merge\_dictionary

**MERGE\_DICTIONARY** - merges two cell arrays of chars and returns only the distinct elements of their union (used by `tmg.m`, `tmg_query.m`, `tdm_update.m`)

[**ALL\_WORDS**, **ALL\_DOC\_IDS**]=

**MERGE\_DICTIONARY**(**ALL\_WORDS**, **NEW\_WORDS**, **ALL\_DOC\_IDS**, **NEW\_DOC\_IDS**) returns in **ALL\_WORDS** all distinct elements of the union of the cell arrays of chars **ALL\_WORDS**, **NEW\_WORDS** corresponding to two document collections. **ALL\_DOC\_IDS** and **NEW\_DOC\_IDS** contain the inverted indices of the two collections. Output argument **ALL\_DOC\_IDS** contains the inverted index of the whole collection.

merge\_tdms

MERGE\_TDMS - Merges two document collections

[A, DICTIONARY]=MERGE\_TDMS(A1, DICTIONARY1, A2, DICTIONARY2)  
merges the tdms A1 and A2 with corresponding dictionaries  
DICTIONARY1 and DICTIONARY2.

MERGE\_TDS(A1, DICTIONARY1, A2, DICTIONARY2, OPTIONS) defines  
optional parameters:

- OPTIONS.min\_local\_freq: The minimum local frequency for a term (default 1)
- OPTIONS.max\_local\_freq: The maximum local frequency for a term (default inf)
- OPTIONS.min\_global\_freq: The minimum global frequency for a term (default 1)
- OPTIONS.max\_global\_freq: The maximum global frequency for a term (default inf)
- OPTIONS.local\_weight: The local term weighting function (default 't'). Possible values (see [1, 2]):
  - 't': Term Frequency
  - 'b': Binary
  - 'l': Logarithmic
  - 'a': Alternate Log
  - 'n': Augmented Normalized Term Frequency
- OPTIONS.global\_weight: The global term weighting function (default 'x'). Possible values (see [1, 2]):
  - 'x': None
  - 'e': Entropy
  - 'f': Inverse Document Frequency (IDF)
  - 'g': GfIdf
  - 'n': Normal
  - 'p': Probabilistic Inverse
- OPTIONS.normalization: Indicates if we normalize the document vectors (default 'x'). Possible values:
  - 'x': None
  - 'c': Cosine

myperms
---------

<p>MYPERMS - computes all possible combinations of the input V=MYPERMS[P, L] returns all possible combinations of the input vector of integers with L numbers.</p>
--

nmmf\_gui

**NNMF\_GUI**

nmmf\_gui is a graphical user interface for all non-negative dimensionality reduction techniques implemented in the Text to Matrix Generator (TMG) Toolbox.

`nmmf_mul_update`

**NNMF\_MUL\_UPDATE** - Applies the multiplicative update algorithm of Lee and Seung.

**NNMF\_MUL\_UPDATE** applies the multiplicative update algorithm of Lee and Seung for non-negative factorizations.  $[W, H, S] = \text{nmmf\_mul\_update}(A, W, H, \text{NIT}, \text{DSP})$  produces a non-negative factorization of  $A$ ,  $W^*H$ , using as initial factors  $W$  and  $H$ , applying  $\text{NIT}$  iterations.

**REFERENCES:**

[1] D. Lee, S. Seung, Algorithms for Non-negative Matrix Factorization, NIPS (2000), 556-562.

`open_file`

**OPEN\_FILE**

OPEN\_FILE is a graphical user interface for selecting a file, directory or variable from the workspace. The function returns the name of the selected file, directory or variable.

opt\_2means

OPT\_2MEANS - a special case of k-means for k=2

OPT\_2MEANS(A, X) returns the clustering that optimizes the objective function of the k-means algorithm based on the ordering of vector X.

[CLUSTERS, S]=OPT\_2MEANS(A, X) returns the cluster structure as well as the value of the objective function.



pca

PCA - Principal Component Analysis

[U, S, V]=PCA(A, C, K, METHOD) computes the K-factor Principal Component Analysis of A, i.e. SVD of  $A-C \cdot \text{ones}(\text{size}(A, 2), 1)$ , using either the svds function of MATLAB or the PROPACK package [1].

REFERENCES:

[1] R.M.Larsen, PROPACK: A Software Package for the Symmetric Eigenvalue Problem and Singular Value Problems on Lanczos and Lanczos Bidiagonalization with Partial Reorthogonalization, Stanford University, <http://sun.stanford.edu/~rmunk/PROPACK>.

pca.mat

PCA\_MAT - Principal Component Analysis with MATLAB (svds)  
[U, S, V]=PCA\_MAT(A, C, K) computes the K-factor Principal  
Component Analysis of A, i.e. SVD of A-C\*ones(size(A, 2),  
1), using the svds function of MATLAB.

pca\_mat\_afun

PCA\_MAT\_AFUN - Auxiliary function used in PCA\_MAT.

pca\_propack

PCA\_PROPACK - Principal Component Analysis with PROPACK

[U, S, V]=PCA\_PROPACK(A, C, K) computes the K-factor  
Principal Component Analysis of A, i.e. SVD of  
A-C\*ones(size(A, 2), 1), using the PROPACK package [1].

REFERENCES:

[1] R.M.Larsen, PROPACK: A Software Package for the Symmetric Eigenvalue  
Problem and Singular Value Problems on Lanczos and Lanczos Bidiagonalization  
with Partial Reorthogonalization, Stanford University,  
<http://sun.stanford.edu/~rmunk/PROPACK>.

pca\_propack.Atransfunc

PCA\_PROPACK\_ATRANSFUNC - Auxiliary function used in PCA\_PROPACK.

pca\_propack\_afun

PCA\_PROPACK\_AFUN - Auxiliary function used in TMG\_PCA\_PROPACK.

pca\_update

PCA\_UPDATE - Principal Component Analysis of a rank-1 updated matrix with MATLAB (eigs)

[U, S, V]=PCA\_UPDATE(A, W, H, C, K) computes the K-factor Principal Component Analysis of  $A - W * H$ , i.e. SVD of  $(A - W * H) - C * \text{ones}(\text{size}(A, 2), 1)$ , using the svds function of MATLAB.

pca\_update\_afun

PCA\_UPDATE\_AFUN - Auxiliary function used in PCA\_UPDATE.



**PDDP - Principal Direction Divisive Partitioning Clustering****Algorithm**

PDDP clusters a term-document matrix (tdm) using the Principal Direction Divisive Partitioning clustering algorithm [1, 2].

CLUSTERS=PDDP(A, K, L) returns a cluster structure with K clusters for the tdm A formed using information from the first L principal components of the tdm.

[CLUSTERS, TREE\_STRUCTURE]=PDDP(A, K, L) returns also the full PDDP tree, while [CLUSTERS, TREE\_STRUCTURE, S]=PDDP(A, K, L) returns the objective function of PDDP.

PDDP(A, K, L, SVD\_METHOD) defines the method used for the computation of the PCA (svds - default - or propack), while PDDP(A, K, L, SVD\_METHOD, DSP) defines if results are to be displayed to the command window (default 1) or not (0).

**REFERENCES:**

- [1] D.Boley, Principal Direction Divisive Partitioning, Data Mining and Knowledge Discovery 2 (1998), no. 4, 325-344.
- [2] D.Zeimpekis, E.Gallopoulos, PDDP(1): Towards a Flexible Principal Direction Divisive Partitioning Clustering Algorithm, Proc. IEEE ICDM'03 Workshop on Clustering Large Data Sets (Melbourne, Florida), 2003.

pddp\_2means

**PDDP\_2MEANS - Hybrid Principal Direction Divisive Partitioning**

**Clustering Algorithm and k-means**

PDDP\_2MEANS clusters a term-document matrix (tdm) using a combination of the Principal Direction Divisive Partitioning clustering algorithm [1] and k-means [2].

CLUSTERS=PDDP\_2MEANS(A, K) returns a cluster structure with K clusters for the tdm A.

[CLUSTERS, TREE\_STRUCTURE]=PDDP\_2MEANS(A, K) returns also the full PDDP tree, while [CLUSTERS, TREE\_STRUCTURE, S]=PDDP\_2MEANS(A, K) returns the objective function of PDDP.

PDDP\_2MEANS(A, K, SVD\_METHOD) defines the method used for the computation of the PCA (svds - default - or propack).

PDDP\_2MEANS(A, K, SVD\_METHOD, DSP) defines if results are to be displayed to the command window (default 1) or not (0).

Finally, PDDP\_2MEANS(A, K, SVD\_METHOD, DSP, EPSILON) defines the termination criterion value for the k-means algorithm.

**REFERENCES:**

[1] D.Boley, Principal Direction Divisive Partitioning, Data Mining and Knowledge Discovery 2 (1998), no. 4, 325-344.

[2] D.Zeimpekis, E.Gallopoulos, k-means Steering of Spectral Divisive Clustering Algorithms, Proc. of Text Mining Workshop, Minneapolis, 2007.

`pddp_extract_centroids`

`PDDP_EXTRACT_CENTROIDS` - returns the cluster centroids of a PDDP clustering result.

pddp\_optcut

**PDDP\_OPTCUT - Hybrid Principal Direction Divisive**

**Partitioning Clustering Algorithm and k-means**

PDDP\_OPTCUT clusters a term-document matrix (tdm) using a combination of the Principal Direction Divisive Partitioning clustering algorithm [1] and k-means [2].

CLUSTERS=PDDP\_OPTCUT(A, K) returns a cluster structure with K clusters for the tdm A.

[CLUSTERS, TREE\_STRUCTURE]=PDDP\_OPTCUT(A, K) returns also the full PDDP tree, while [CLUSTERS, TREE\_STRUCTURE, S]=PDDP\_OPTCUT(A, K) returns the objective function of PDDP.

PDDP\_OPTCUT(A, K, SVD\_METHOD) defines the method used for the computation of the PCA (svds - default - or propack).

PDDP\_OPTCUT(A, K, SVD\_METHOD, DSP) defines if results are to be displayed to the command window (default 1) or not (0). Finally,

PDDP\_OPTCUT(A, K, SVD\_METHOD, DSP, EPSILON) defines the termination criterion value for the k-means algorithm.

**REFERENCES:**

[1] D.Boley, Principal Direction Divisive Partitioning, Data Mining and Knowledge Discovery 2 (1998), no. 4, 325-344.

[2] D.Zeimpekis, E.Gallopoulos, k-means Steering of Spectral Divisive Clustering Algorithms, Proc. of Text Mining Workshop, Minneapolis, 2007.

pddp-optcut\_2means

**PDDP\_OPTCUT\_2MEANS - Hybrid Principal Direction Divisive**

**Partitioning Clustering Algorithm and k-means**

PDDP\_OPTCUT\_2MEANS clusters a term-document matrix (tdm) using a combination of the Principal Direction Divisive Partitioning clustering algorithm [1] and k-means [2].

CLUSTERS=PDDP\_OPTCUT\_2MEANS(A, K) returns a cluster structure with K clusters for the tdm A.

[CLUSTERS, TREE\_STRUCTURE]=PDDP\_OPTCUT\_2MEANS(A, K) returns also the full PDDP tree, while [CLUSTERS, TREE\_STRUCTURE, S]=PDDP\_OPTCUT\_2MEANS(A, K) returns the objective function of PDDP.

PDDP\_OPTCUT\_2MEANS(A, K, SVD\_METHOD) defines the method used for the computation of the PCA (svds - default - or propack).

PDDP\_OPTCUT\_2MEANS(A, K, SVD\_METHOD, DSP) defines if results are to be displayed to the command window (default 1) or not (0). Finally, PDDP\_OPTCUT\_2MEANS(A, K, SVD\_METHOD, DSP, EPSILON) defines the termination criterion value for the k-means algorithm.

**REFERENCES:**

[1] D.Boley, Principal Direction Divisive Partitioning, Data Mining and Knowledge Discovery 2 (1998), no. 4, 325-344.

[2] D.Zeimpekis, E.Gallopoulos, k-means Steering of Spectral Divisive Clustering Algorithms, Proc. of Text Mining Workshop, Minneapolis, 2007.

**PDDP\_OPTCUTPD - Hybrid Principal Direction Divisive****Partitioning Clustering Algorithm and k-means**

PDDP\_OPTCUTPD clusters a term-document matrix (tdm) using a combination of the Principal Direction Divisive

Partitioning clustering algorithm [1, 2] and k-means [3].

CLUSTERS=PDDP\_OPTCUT\_PD(A, K, L) returns a cluster structure with K clusters for the tdm A formed using information from the first L principal components of the tdm.

[CLUSTERS, TREE\_STRUCTURE]=PDDP\_OPTCUTPD(A, K, L) returns also the full PDDP tree, while [CLUSTERS, TREE\_STRUCTURE, S]=

PDDP\_OPTCUTPD(A, K, L) returns the objective function of PDDP.

PDDP\_OPTCUTPD(A, K, L, SVD\_METHOD) defines the method used for the computation of the PCA (svds - default - or

propack). Finally, PDDP\_OPTCUTPD(A, K, L, SVD\_METHOD, DSP)

defines if results are to be displayed to the command window

(default 1) or not (0).

**REFERENCES:**

[1] D.Boley, Principal Direction Divisive Partitioning, Data Mining and Knowledge Discovery 2 (1998), no. 4, 325-344.

[2] D.Zeimpekis, E.Gallopoulos, PDDP(1): Towards a Flexible Principal Direction Divisive Partitioning Clustering Algorithm, Proc. IEEE ICDM'03 Workshop on Clustering Large Data Sets (Melbourne, Florida), 2003.

[3] D.Zeimpekis, E.Gallopoulos, k-means Steering of Spectral Divisive Clustering Algorithms, Proc. of Text Mining Workshop, Minneapolis, 2007.

ps\_pdf2ascii

PS\_PDF2ASCII - converts the input ps or pdf file to ASCII  
RESULT = PS\_PDF2ASCII(FILENAME) converts the input ps or pdf files to ASCII, using ghostscript's utility 'ps2ascii'.  
RESULT returns a success indicator, e.g. -2 if the input file does not exist or has a wrong format, -1 if gs is not installed or the path isn't set, 0 if 'ps2ascii' didn't work properly, and 1 if the conversion was successful.

retrieval\_gui

**RETRIEVAL\_GUI**

RETRIEVAL\_GUI is a graphical user interface for all retrieval functions of the Text to Matrix Generator (TMG) Toolbox.



rocchio\_multi

ROCCHIO\_MULTI - Rocchio classifier for multi-label collections

LABELS\_AS=KNN\_MULTI(A, CLUSTERS, BETA, GAMMA, Q, LABELS, NORMALIZED\_DOCS, THRESHOLDS) classifies the columns of Q with the Rocchio classifier using the pre-classified columns of matrix A with labels LABELS (vector of integers).

THRESHOLDS is a vector of class threshold values. BETA and GAMMA define the weight of positive and negative examples in the formation of each class centroid. NORMALIZED\_DOCS defines if cosine (1) or euclidean distance (0) similarity measure is to be used. LABELS\_AS contains the assigned labels for the columns of Q.

rocchio\_single

ROCCHIO\_SINGLE - Rocchio classifier for single-label collections

LABELS\_AS=KNN\_SINGLE(A, CLUSTERS, BETA, GAMMA, Q, LABELS, NORMALIZED\_DOCS) classifies the columns of Q with the Rocchio classifier using the pre-classified columns of matrix A with labels LABELS (vector of integers).

BETA and GAMMA define the weight of positive and negative examples in the formation of each class centroid.

NORMALIZED\_DOCS defines if cosine (1) or euclidean distance (0) similarity measure is to be used. LABELS\_AS contains the assigned labels for the columns of Q.

scut\_knn

SCUT\_KNN - implements the Scut thresholding technique from [1]

for the k-Nearest Neighbors classifier

THRESHOLD=SCUT\_KNN(A, Q, K, LABELS\_TR, LABELS\_TE, MINF1, NORMALIZE, STEPS) returns the vector of thresholds for the k-Nearest Neighbors classifier for the collection [A Q]. A and Q define the training and test parts of the validation set with labels LABELS\_TR and LABELS\_TE respectively. MINF1 defines the minimum F1 value and NORMALIZE defines if cosine (1) or euclidean distance (0) measure of similarity is to be used. Finally, STEPS defines the number of steps used during thresholding. [THRESHOLD, F, THRESHOLDS]=SCUT\_KNN(A, Q, K, LABELS\_TR, LABELS\_TE, MINF1, NORMALIZE, STEPS) returns also the best F1 value as well as the matrix of thresholds for each step (row i corresponds to step i).

REFERENCES:

[1] Y. Yang. A Study of Thresholding Strategies for Text Categorization. In Proc. 24th ACM SIGIR, pages 137145, New York, NY, USA, 2001. ACM Press.

SCUT\_LLSF - implements the Scut thresholding technique from [2] for the Linear Least Squares Fit classifier [3]

THRESHOLD=SCUT\_LLSF(A, Q, CLUSTERS, K, LABELS\_TR, LABELS\_TE, MINF1, L, METHOD, STEPS, SVD\_METHOD, CLSI\_METHOD) returns the vector of thresholds for the Linear Least Squares Fit classifier for the collection [A Q]. A and Q define the training and test parts of the validation set with labels LABELS\_TR and LABELS\_TE respectively. CLUSTERS is a structure defining the classes, while MINF1 defines the minimum F1 value and STEPS defines the number of steps used during thresholding.

METHOD is the method used for the approximation of the rank-1 truncated SVD, with possible values:

- 'clsi': Clustered Latent Semantic Indexing [4].
- 'cm': Centroids Method [1].
- 'svd': Singular Value Decomposition.

SVD\_METHOD defines the method used for the computation of the SVD, while CLSI\_METHOD defines the method used for the determination of the number of factors from each class used in Clustered Latent Semantic Indexing in case METHOD equals 'clsi'.

[THRESHOLD, F, THRESHOLDS]=SCUT\_LLSF(A, Q, CLUSTERS, K, LABELS\_TR, LABELS\_TE, MINF1, L, METHOD, STEPS, SVD\_METHOD, CLSI\_METHOD) returns also the best F1 value as well as the matrix of thresholds for each step (row i corresponds to step i).

#### REFERENCES:

- [1] H. Park, M. Jeon, and J. Rosen. Lower Dimensional Representation of Text Data Based on Centroids and Least Squares. *BIT Numerical Mathematics*, 43(2):427448, 2003.
- [2] Y. Yang. A Study of Thresholding Strategies for Text Categorization. In *Proc. 24th ACM SIGIR*, pages 137145, New York, NY, USA, 2001. ACM Press.
- [3] Y. Yang and C. Chute. A Linear Least Squares Fit Mapping Method for Information Retrieval from Natural Language Texts. In *Proc. 14th Conference on Computational Linguistics*, pages 447453, Morristown, NJ, USA, 1992.
- [4] D. Zeimpekis and E. Gallopoulos, "Non-Linear Dimensional Reduction via Class Representatives for Text Classification". In *Proc. 2006 IEEE International Conference on Data Mining (ICDM'06)*, Hong Kong, Dec. 2006.

scut\_rocchio

SCUT\_ROCCHIO - implements the Scut thresholding technique

from [1] for the Rocchio classifier

THRESHOLD=SCUT\_ROCCHIO(A, CLUSTERS, BETA, GAMMA, Q,  
LABELS\_TR, LABELS\_TE, MINF1, NORMALIZE, STEPS) returns

the vector of thresholds for the Rocchio classifier

for the collection [A Q]. A and Q define the training

and test parts of the validation set with labels

LABELS\_TR and LABELS\_TE respectively. MINF1 defines

the minimum F1 value, while NORMALIZE defines if cosine

(1) or euclidean distance (0) measure of similarity is

to be used, CLUSTERS is a structure defining the classes

and STEPS defines the number of steps used during

thresholding. BETA and GAMMA define the weight of positive

and negative examples in the formation of each class

centroid.

[THRESHOLD, F, THRESHOLDS]=SCUT\_ROCCHIO(A, CLUSTERS, BETA,  
GAMMA, Q, LABELS\_TR, LABELS\_TE, MINF1, NORMALIZE, STEPS)

returns also the best F1 value as well as the matrix of

thresholds for each step (row i corresponds to step i).

REFERENCES:

[1] Y. Yang. A Study of Thresholding Strategies for Text  
Categorization. In Proc. 24th ACM SIGIR, pages 137145,  
New York, NY, USA, 2001. ACM Press.

sdd\_tmg

SDD\_TMGM - interface for SDDPACK

[X, D, Y]=SDD\_TMGM(A, K) computes a rank-K Semidiscrete  
Decomposition of A using the SDDPACK [1].

REFERENCES:

Tamara G. Kolda and Dianne P. O'Leary, Computation and Uses of the  
Semidiscrete Matrix Decomposition, Computer Science Department Report  
CS-TR-4012 Institute for Advanced Computer Studies Report UMIACS-TR-99-22,  
University of Maryland, April 1999.

## skmeans

### SKMEANS - Spherical k-Means Clustering Algorithm

SKMEANS clusters a term-document matrix using the Spherical k-means clustering algorithm [1]. `CLUSTERS=SKMEANS(A, C, K, TERMINATION)` returns a cluster structure with K clusters for the term-document matrix A using as initial centroids the columns of C (initialized randomly when it is empty). `TERMINATION` defines the termination method used in spherical k-means ('epsilon' stops iteration when objective function increase falls down a user defined threshold - see `OPTIONS` input argument - while 'n\_iter' stops iteration when a user defined number of iterations has been reached).  
`[CLUSTERS, Q]=SKMEANS(A, C, K, TERMINATION)` returns also the vector of objective function values for each iteration and `[CLUSTERS, Q, C]=SKMEANS(A, C, K, TERMINATION)` returns the final centroid vectors.  
`SKMEANS(A, C, K, TERMINATION, OPTIONS)` defines optional parameters:

- `OPTIONS.iter`: Number of iterations (default 10).
- `OPTIONS.epsilon`: Value for epsilon convergence criterion (default 1).
- `OPTIONS.dsp`: Displays results (default 1) or not (0) to the command window.

#### REFERENCES:

[1] I. S. Dhillon and D. M. Modha, "Concept Decompositions for Large Sparse Text Data using Clustering", *Machine Learning*, 42:1, pages 143-175, Jan, 2001.

stemmer

STEMMER - applies the Porter's Stemming algorithm [1]

S = STEMMER(TOKEN, DSP) returns in S the stemmed word of  
TOKEN. DSP indicates if the function displays the result  
of each stem (1).

REFERENCES:

[1] M.F.Porter, An algorithm for suffix stripping, Program, 14(3): 130-137,  
1980.



`strip_html`

**STRIP\_HTML** - removes html entities from an html file  
**S = STRIP\_HTML(FILENAME)** parses file **FILENAME** and removes the html entities, while the result is stored in **S** as a cell array and written in file "**FILENAME.TXT**".

svd\_tmg

**SVD\_TMG - Singular Value Decomposition**

[U, S, V]=SVD\_TMG(A, K, METHOD) computes the K-factor truncated Singular Value Decomposition of A using either the svds function of MATLAB or the PROPACK package [1].

**REFERENCES:**

[1] R.M.Larsen, PROPACK: A Software Package for the Symmetric Eigenvalue Problem and Singular Value Problems on Lanczos and Lanczos Bidiagonalization with Partial Reorthogonalization, Stanford University, <http://sun.stanford.edu/~rmunk/PROPACK>.

svd\_update

SVD\_UPDATE - Singular Value Decomposition of a rank-1 update matrix with MATLAB (eigs)  
[U, S, V]=SVD\_UPDATE(A, X, Y, K) computes the K-factor SVD of  $A-X*Y$ , using the eigs function of MATLAB.

svd\_update\_afun

SVD\_UPDATE\_AFUN - Auxiliary function used in SVD\_UPDATE.

tdm\_downdate

TDM\_DOWNDATE - renews a text collection by downdating the corresponding term-document matrix

A = TDM\_DOWNDATE(UPDATE\_STRUCT, REMOVED\_DOCS) returns the new term - document matrix of the downdated collection.

UPDATE\_STRUCT defines the update structure returned by TMG, while REMOVED\_DOCS defines the indices of the documents that is to be removed.

[A, DICTIONARY] = TDM\_DOWNDATE(UPDATE\_STRUCT, REMOVED\_DOCS) returns also the dictionary for the updated collection, while

[A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZED\_FACTORS] = TDM\_DOWNDATE(UPDATE\_STRUCT, REMOVED\_DOCS) returns the vectors of global weights for the dictionary and the

normalization factor for each document in case such a factor is used. If normalization is not used TDM\_DOWNDATE returns a

vector of all ones. [A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZATION\_FACTORS, WORDS\_PER\_DOC] =

TDM\_DOWNDATE(UPDATE\_STRUCT, REMOVED\_DOCS) returns statistics for each document, i.e. the number of terms for each document.

[A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZATION\_FACTORS, WORDS\_PER\_DOC, TITLES, FILES] = TDM\_DOWNDATE(UPDATE\_STRUCT, REMOVED\_DOCS) returns in FILES the filenames containing the collection's documents and a cell array (TITLES) that contains a declaratory title for each document, as well as the document's first line.

Finally [A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZATION\_FACTORS, WORDS\_PER\_DOC, TITLES, FILES, UPDATE\_STRUCT] =

TDM\_DOWNDATE(UPDATE\_STRUCT, REMOVED\_DOCS) returns the update structure that keeps the essential information for the collection's update (or downdate).

TDM\_DOWNDATE(UPDATE\_STRUCT, REMOVED\_DOCS, OPTIONS) defines optional parameters:

- OPTIONS.dsp: Displays results (default 1) or not (0) to the command window.

## tdm\_update

TDM.UPDATE renews a text collection by updating the corresponding term-document matrix.

A = TDM.UPDATE(FILENAME, UPDATE\_STRUCT) returns the new term - document matrix of the updated collection. FILENAME defines the file (or files in case a directory is supplied) containing the new documents, while UPDATE\_STRUCT defines the update structure returned by TMG. In case FILENAME variable is empty, the collection is simply updated using the options defined by UPDATE\_STRUCT (for example, use another term-weighting scheme).

[A, DICTIONARY] = TDM.UPDATE(FILENAME, UPDATE\_STRUCT) returns also the dictionary for the updated collection, while [A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZED\_FACTORS] = TDM.UPDATE(FILENAME, UPDATE\_STRUCT) returns the vectors of global weights for the dictionary and the normalization factor for each document in case such a factor is used.

If normalization is not used TDM.UPDATE returns a vector of all ones.

[A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZATION\_FACTORS, WORDS\_PER\_DOC] = TDM.UPDATE(FILENAME, UPDATE\_STRUCT) returns statistics for each document, i.e. the number of terms for each document.

[A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZATION\_FACTORS, WORDS\_PER\_DOC, TITLES, FILES] = TDM.UPDATE(FILENAME, UPDATE\_STRUCT) returns in FILES the filenames contained in directory (or file) FILENAME and a cell array (TITLES) that contains a declaratory title for each document, as well as the document's first line.

Finally [A, DICTIONARY, GLOBAL\_WEIGHTS, NORMALIZATION\_FACTORS, WORDS\_PER\_DOC, TITLES, FILES, UPDATE\_STRUCT] = TDM.UPDATE(FILENAME, UPDATE\_STRUCT) returns the update structure that keeps the essential information for the collection's update (or downdate).

TDM.UPDATE(FILENAME, UPDATE\_STRUCT, OPTIONS) defines optional parameters:

- OPTIONS.delimiter: The delimiter between documents within the same file. Possible values are 'emptyline' (default), 'none\_delimiter' (treats each file as a single document) or any other string.
- OPTIONS.line\_delimiter: Defines if the delimiter takes a whole line of text (default, 1) or not.

- `OPTIONS.update_step`: The step used for the incremental built of the inverted index (default 10,000).
- `OPTIONS.dsp`: Displays results (default 1) or not (0) to the command window.

**TMG - Text to Matrix Generator**

TMG parses a text collection and generates the term - document matrix.

$A = \text{TMG}(\text{FILENAME})$  returns the term - document matrix, that corresponds to the text collection contained in files of directory (or file) FILENAME.

Each document must be separated by a blank line (or another delimiter that is defined by OPTIONS argument) in each file.

$[A, \text{DICTIONARY}] = \text{TMG}(\text{FILENAME})$  returns also the dictionary for the collection, while  $[A, \text{DICTIONARY}, \text{GLOBAL\_WEIGHTS}, \text{NORMALIZED\_FACTORS}] = \text{TMG}(\text{FILENAME})$  returns the vectors of global weights for the dictionary and the normalization factor for each document in case such a factor is used. If normalization is not used TMG returns a vector of all ones.

$[A, \text{DICTIONARY}, \text{GLOBAL\_WEIGHTS}, \text{NORMALIZATION\_FACTORS}, \text{WORDS\_PER\_DOC}] = \text{TMG}(\text{FILENAME})$  returns statistics for each document, i.e. the number of terms for each document.

$[A, \text{DICTIONARY}, \text{GLOBAL\_WEIGHTS}, \text{NORMALIZATION\_FACTORS}, \text{WORDS\_PER\_DOC}, \text{TITLES}, \text{FILES}] = \text{TMG}(\text{FILENAME})$  returns in FILES the filenames contained in directory (or file)

FILENAME and a cell array (TITLES) that contains a declaratory title for each document, as well as the document's first line. Finally

$[A, \text{DICTIONARY}, \text{GLOBAL\_WEIGHTS}, \text{NORMALIZATION\_FACTORS}, \text{WORDS\_PER\_DOC}, \text{TITLES}, \text{FILES}, \text{UPDATE\_STRUCT}] = \text{TMG}(\text{FILENAME})$  returns a structure that keeps the essential information for the collection's update (or downdate).

$\text{TMG}(\text{FILENAME}, \text{OPTIONS})$  defines optional parameters:

- OPTIONS.use\_mysql: Indicates if results are to be stored in MySQL.
- OPTIONS.db\_name: The name of the directory where the results are to be saved.
- OPTIONS.delimiter: The delimiter between documents within the same file. Possible values are 'emptyline' (default), 'none\_delimiter' (treats each file as a single document) or any other string.
- OPTIONS.line\_delimiter: Defines if the delimiter takes a whole line of text (default, 1) or not.
- OPTIONS.stoplist: The filename for the stoplist,



i.e. a list of common words that we don't use for the indexing (default no stoplist used).

- `OPTIONS.stemming`: Indicates if the stemming algorithm is used (1) or not (0 - default).
- `OPTIONS.update_step`: The step used for the incremental build of the inverted index (default 10,000).
- `OPTIONS.min_length`: The minimum length for a term (default 3).
- `OPTIONS.max_length`: The maximum length for a term (default 30).
- `OPTIONS.min_local_freq`: The minimum local frequency for a term (default 1).
- `OPTIONS.max_local_freq`: The maximum local frequency for a term (default inf).
- `OPTIONS.min_global_freq`: The minimum global frequency for a term (default 1).
- `OPTIONS.max_global_freq`: The maximum global frequency for a term (default inf).
- `OPTIONS.local_weight`: The local term weighting function (default 't'). Possible values (see [1, 2]):
  - 't': Term Frequency
  - 'b': Binary
  - 'l': Logarithmic
  - 'a': Alternate Log
  - 'n': Augmented Normalized Term Frequency
- `OPTIONS.global_weight`: The global term weighting function (default 'x'). Possible values (see [1, 2]):
  - 'x': None
  - 'e': Entropy
  - 'f': Inverse Document Frequency (IDF)
  - 'g': GfIdf
  - 'n': Normal
  - 'p': Probabilistic Inverse
- `OPTIONS.normalization`: Indicates if we normalize the document vectors (default 'x'). Possible values:
  - 'x': None
  - 'c': Cosine
- `OPTIONS.dsp`: Displays results (default 1) or not (0) to the command window.

#### REFERENCES:

- [1] M.Berry and M.Browne, Understanding Search Engines, Mathematical Modeling and Text Retrieval, Philadelphia, PA: Society for Industrial and Applied Mathematics, 1999.
- [2] T.Kolda, Limited-Memory Matrix Methods with Applications, Tech.Report CS-TR-3806, 1997.

tmg\_gui

**TMG\_GUI**

TMG\_GUI is a graphical user interface for all indexing routines of the Text to Matrix Generator (TMG) Toolbox. For a full documentation type 'help tmg', 'help tmg\_query', 'help tdm\_update' or 'help tdm\_downdate'. For a full documentation of the GUI's usage, select the help tab to the GUI.

## tmg\_query

TMG\_QUERY - Text to Matrix Generator, query vector constructor

TMG\_QUERY parses a query text collection and generates the query vectors corresponding to the supplied dictionary.

Q = TMG\_QUERY(FILENAME, DICTIONARY) returns the query vectors, that corresponds to the text collection contained in files of directory FILENAME. DICTIONARY is the array of terms corresponding to a text collection.

Each query must be separated by a blank line (or another delimiter that is defined by OPTIONS argument) in each file.

[Q, WORDS\_PER\_QUERY] = TMG\_QUERY(FILENAME, DICTIONARY) returns statistics for each query, i.e. the number of terms for each query.

Finally, [Q, WORDS\_PER\_QUERY, TITLES, FILES] = TMG\_QUERY(FILENAME) returns in FILES the filenames contained in directory (or file) FILENAME and a cell array (TITLES) that contains a declaratory title for each query, as well as the query's first line.

TMG\_QUERY(FILENAME, DICTIONARY, OPTIONS) defines optional parameters:

- OPTIONS.delimiter: The delimiter between queries within the same file. Possible values are 'emptyline' (default), 'none\_delimiter' (treats each file as a single query) or any other string.
- OPTIONS.line\_delimiter: Defines if the delimiter takes a whole line of text (default, 1) or not.
- OPTIONS.stoplist: The filename for the stoplist, i.e. a list of common words that we don't use for the indexing (default no stoplist used).
- OPTIONS.stemming: Indicates if the stemming algorithm is used (1) or not (0 - default).
- OPTIONS.update\_step: The step used for the incremental built of the inverted index (default 10,000).
- OPTIONS.local\_weight: The local term weighting function (default 't'). Possible values (see [1, 2]):
  - 't': Term Frequency
  - 'b': Binary
  - 'l': Logarithmic
  - 'a': Alternate Log
  - 'n': Augmented Normalized Term Frequency
- OPTIONS.global\_weights: The vector of term global weights (returned by tmg).
- OPTIONS.dsp: Displays results (default 1) or not (0).

**REFERENCES:**

- [1] M.Berry and M.Browne, Understanding Search Engines, Mathematical Modeling and Text Retrieval, Philadelphia, PA: Society for Industrial and Applied Mathematics, 1999.
- [2] T.Kolda, Limited-Memory Matrix Methods with Applications, Tech.Report CS-TR-3806, 1997.

`tmg_save_results`

**TMG\_SAVE\_RESULTS**

TMG\_SAVE\_RESULTS is a graphical user interface used from TMG\_GUI, for saving the results to a (or multiple) .mat file(s).

tmg\_template

TDM\_TEMPLATE - demo script

This is a template script demonstrating the use of TMG, as well as the application of the resulting TDM'S in two IR tasks, quering and clustering. The quering models used is the Vector Space Model (see vsm.m) and LSI (see lsi.m), while two versions of the k-means algorithm (euclidean and spherical, see ekmeans.m and skmeans.m) cluster the resulting matrix (see pddp.m). The user can edit this code in order to change the default OPTIONS of TMG, as well as to apply other IR tasks or use his own implementations regarding these tasks.

two\_means\_1d

TWO\_MEANS\_1D - returns the clustering that optimizes the objective function of the k-means algorithm for the input vector.

[CUTOFF, CLUSTERS, DISTANCE, OF, MEAN1, MEAN2]=  
TWO\_MEANS\_1D(A) returns the cutoff value of the clustering, the cluster structure, the separation distance, the value of the objective function and the two mean values.

unique\_elements

UNIQUE\_ELEMENTS - detects all distinct elements of a vector  
[ELEMENTS, N] = UNIQUE\_ELEMENTS(X) returns in ELEMENTS all  
distinct elements of vector X, and in N the number of times  
each element appears in X. A value is repeated if it appears  
in non-consecutive elements. For no repetitive elements sort  
the input vector.



unique\_words

UNIQUE\_WORDS - detects all distinct elements of a cell array of chars (used by tmg.m, tmg\_query.m, tdm\_update.m)  
[NEW\_WORDS, NEW\_DOC\_IDS]=UNIQUE\_WORDS(WORDS, DOC\_IDS, N\_DOCS)  
returns in NEW\_WORDS all distinct elements of the cell array of chars WORDS. DOC\_IDS is the vector of the document identifiers containing the corresponding words, while N\_DOCS is the total number of documents contained to the collection. NEW\_DOC\_IDS contains the inverted index of the collection as a cell array of 2 x N\_DOCS arrays.

vsm

VSM - Applies the Vector Space Model to a document collection  
[SC, DOCS\_INDS] = VSM(D, Q, NORMALIZE\_DOCS) applies the  
Vector Space Model to the text collection represented by  
the term - document matrix D for the query defined by the  
vector Q [1]. NORMALIZE\_DOCS defines if the document  
vectors are to be normalized (1) or not (0). SC contains  
the sorted similarity coefficients, while DOC\_INDS contains  
the corresponding document indices.

REFERENCES:

[1] M.Berry and M.Browne, Understanding Search Engines,  
Mathematical Modeling and Text Retrieval, Philadelphia,  
PA: Society for Industrial and Applied Mathematics, 1999.