

Experiments with Randomized Algorithms in the Text to Matrix Generator Toolbox

Eugenia Maria Kontopoulou, Dimitrios Zeimpekis
and Efstratios Gallopoulos

Department of
Compute Engineering and Informatics
University of Patras

6th International Conference of the ERCIM WG on
Computational and Methodological Statistics (ERCIM 2013)

London, 15/12/2013

Outline

- 1 Introduction
- 2 Retrieval Task
- 3 Experiments
- 4 Work In Progress

The Big Data Problem

Large & Complex data sets {
Storage Problems
Interpretation Problems

The Big Data Problem

Large & Complex data sets

{ Storage Problems
Interpretation Problems

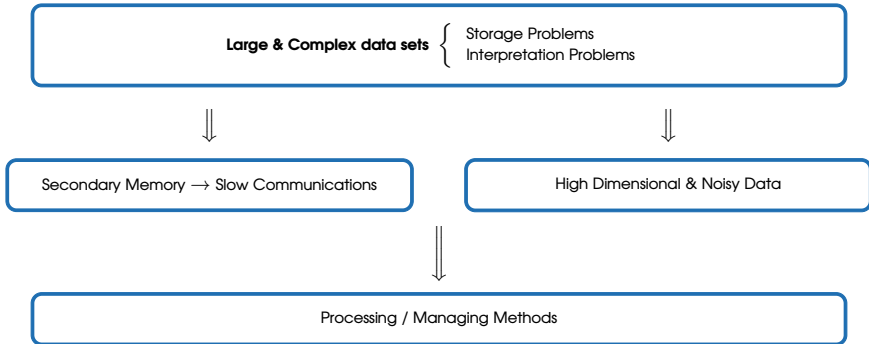


Secondary Memory → Slow Communications

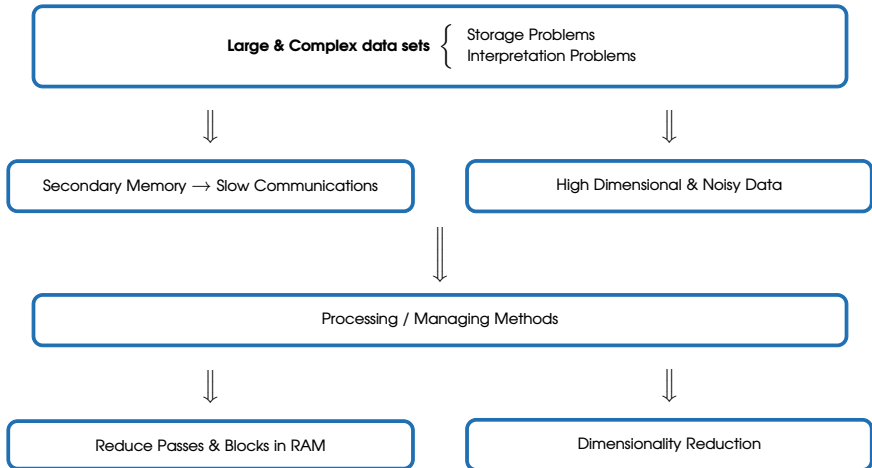


High Dimensional & Noisy Data

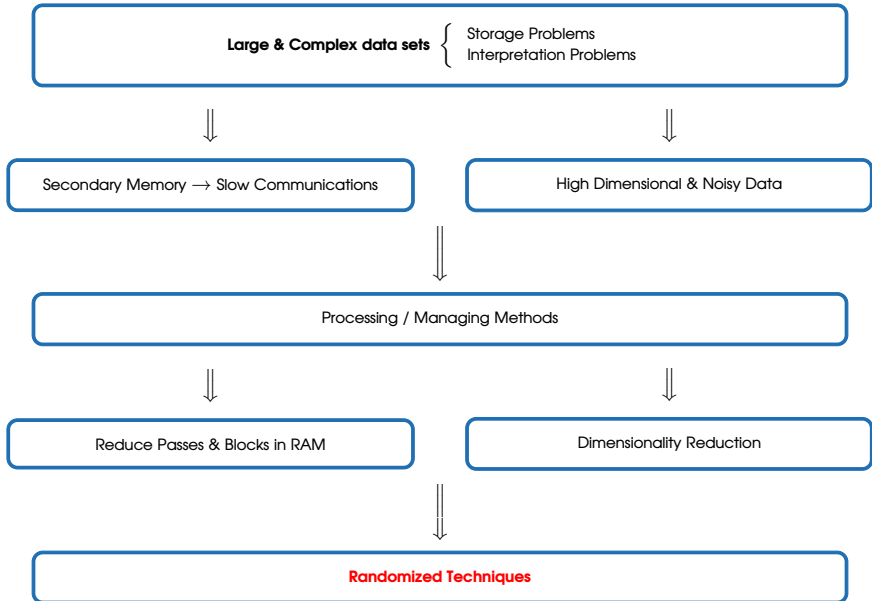
The Big Data Problem



The Big Data Problem



The Big Data Problem



Randomization inroads into Matrix Computations

BLENDENPIK: SUPERCHARGING LAPACK'S LEAST-SQUARES SOLVER*

HAIM AVRON¹, PETAR MAYMOUNKOV², AND SIVAN TOLEDO³

Abstract. Several innovative random-sampling and random-mixing techniques for solving problems in linear algebra have been proposed in the last decade, but they have not yet made a significant impact on numerical linear algebra. We show that by using a high-quality implementation of one of these techniques, we obtain a solver that performs extremely well in the traditional yardsticks of numerical linear algebra: it is significantly faster than high-performance implementations of existing

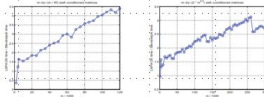


FIG. 1.1. Comparison between BLSPACK and the new solver for increasingly larger matrices. Graphs show the ratio of BLSPACK's running time to BlendenpiK's running time on random matrices with two kinds of sparsity ratios.

An efficient distributed randomized solver with application to large dense linear systems

Marc Baboulin¹, Dulceinea Becker¹, George Bosilca¹, Anthony Danalis¹, Jack Dongarra¹

Project-Team Grand-Large

Research Report n° 8043 — Août 2012 — 20 pages

Abstract: Randomized algorithms are gaining ground in high-performance computing as they have the potential to outperform deterministic methods, while still providing good results. We propose a randomized algorithm for distributed multicore architectures solving large dense symmetric indefinite linear systems that are encountered, for instance, in transfer estimation problems or electromagnetic simulations. This solver combines

RANDOMIZED ALGORITHMS FOR ESTIMATING THE TRACE OF AN IMPLICIT SYMMETRIC POSITIVE SEMI-DEFINITE MATRIX

HAIM AVRON AND SIVAN TOLEDO

Accelerating linear system solutions using randomization techniques

Marc Baboulin¹, Jack Dongarra^{2,3,4}, Julien Hermant¹, and Stanimir Tomov²

Randomized Extended Kaczmarz for Solving Least Squares

Anastasiou Zouzias¹ Nikolaos M. Papanikolaou²

Site Map

Proceedings of the National Academy of Sciences of the United States of America

PNAS

CUR matrix decompositions for improved data analysis

original data set. We present an algorithm that preferentially chooses columns and rows that exhibit high "statistical leverage" and, thus, in a very precise statistical sense, exert a disproportionately large "influence" on the best low-rank "fit" of the data matrix. By selecting columns and rows in this manner, we obtain a fast relative-error and constant-factor approximation guaranteed by deterministic



Subscribe to our technology newsletter, www.newscientist.com, or us.journals@pnas.org

random matrix theory. Unimagined developments more than 20 years ago to describe the energy levels of atomic nuclei, the theory is turning up in everything from inflation rates to the behaviour of solids. So much so that many researchers believe that it points to some kind of deep pattern in nature that we don't yet understand. "It really does feel like the ideas of random matrix theory are somehow buried deep in the heart of nature," says electrical engineer Raj Nadakuditi of the University of Michigan. Ann Arbor

All of this, oddly enough, emerged from an effort to tame physicists' ignored into an advantage. In 1955, when ...

FAST MONTE CARLO ALGORITHMS FOR MATRICES I: APPROXIMATING MATRIX MULTIPLICATION*

PETROS DRINEAS¹, RAVI KANNAN², AND MICHAEL W. MAHONEY¹

ST MONTE CARLO ALGORITHMS FOR MATRICES III: COMPUTING A COMPRESSED APPROXIMATE MATRIX DECOMPOSITION*

PETROS DRINEAS¹, RAVI KANNAN², AND MICHAEL W. MAHONEY¹

SAMPLING FROM LARGE MATRICES: AN APPROACH THROUGH GEOMETRIC FUNCTIONAL ANALYSIS

MARK DEDELLON AND ROMAN VERSHIYAN

Effective Resistances, Statistical Leverage, and Applications to Linear Equation Solving

Petros Drineas¹ Michael W. Mahoney¹

FINDING STRUCTURE WITH RANDOMNESS: PROBABILISTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE MATRIX DECOMPOSITIONS

N. BALOGH, P. G. MARTINSSON¹, AND J. A. TRÖPP

RELATIVE-ERROR CUR MATRIX DECOMPOSITIONS*

PETROS DRINEAS¹, MICHAEL W. MAHONEY¹, AND S. MUTHUKRISHNAN²

Randomization inroads into Matrix Computations

" One great and underused technique at this scale is sampling. **For most things that you want to do with data, 100.000 randomly selected rows is as good as 10.000.000 rows** and working at R or SPSS scale allows for a much faster analysis cycle "

*Lukas Biewald
CEO of CrowdFlower*

Text Data

From document collections . . .

Documents



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

Text Data

... to Term-Document structures ...

Term-Document Matrix (TDM)

33 × 5

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

✓ tf-idf

✓ Stemming

Text Data

... for text mining tasks

Retrieval



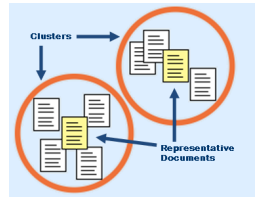
Text Data

... for text mining tasks



Retrieval

Clustering



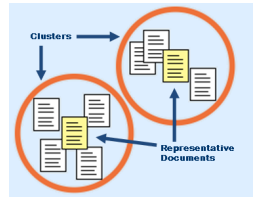
Text Data

... for text mining tasks

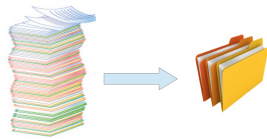


Retrieval

Clustering



Classification



Text to Matrix Generator

What is TMG:

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

TMG: A MATLAB Toolbox for Generating
Term-Document Matrices from Text Collections
(ZG06)



Grouping Multidimensional Data

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

Grouping Multidimensional Data
2006, pp 187-210

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

D. Zaimpekis, E. Gallopoulos

Text to Matrix Generator

What is TMG:

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

Implementation:

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

TMG: A MATLAB Toolbox for Generating Term-Document Matrices from Text Collections (ZG06)



Grouping Multidimensional Data

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

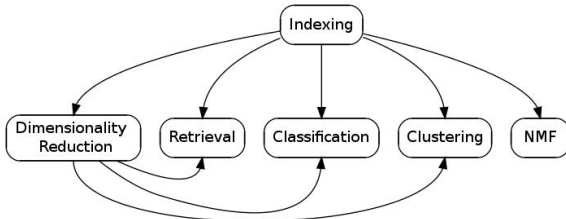
Grouping Multidimensional Data
2006, pp 187-210

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

D. Zeimpekis, E. Gallopoulos

Six basic modules:

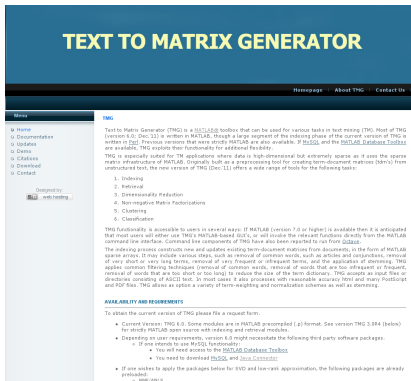
- 1 Indexing
- 2 Dimensionality Reduction
- 3 Non-Negative Matrix Factorizations
- 4 Retrieval
- 5 Clustering
- 6 Classification



How can I find TMG?

Free under request from:

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



TEXT TO MATRIX GENERATOR

Homepage About TMG Contact Us

Home

TMG

Text to Matrix Generator (TMG) is a [SUSLISS](#) toolbox that can be used for various tasks in text mining (TM). Most of TMG (version 8.0, Dec '11) is written in MATLAB, though a large segment of the indexing phase of the current version of TMG is written in C++. Previous versions that were strictly MATLAB are also available. If [EUSQL](#) and the [MATLAB Database Toolbox](#) are available, TMG exploits their functionality for additional flexibility.

TMG is especially suited for TM applications where data is high-dimensional but extremely sparse as it uses the sparse matrix infrastructure of MATLAB. Originally built as a preprocessing tool for creating term-document matrices (tdm's) from unstructured text, the new version of TMG (Dec '11) offers a wide range of tools for the following tasks:

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

TMG functionality is accessible to users in several ways: If MATLAB (version 7.0 or higher) is available then it is anticipated that most users will either use TMG's MATLAB-based GUI's, or will invoke the relevant functions directly from the MATLAB command-line interface. Command line components of TMG have also been repaired to run from [Shell](#).

The indexing process constructs new and updates existing term-document matrices from documents, in the form of MATLAB sparse arrays. Errors include various things, such as removal of common words, such as articles and conjunctions, removal of very short or very long terms, removal of very frequent or infrequent terms, and the application of stemming. TMG applies common filtering techniques (removal of common words, removal of words that are too infrequent or frequent, removal of words that are too short or too long) to reduce the size of the term dictionary. TMG accepts as input files or directories consisting of ASCII text. In most cases it also processes with reasonable accuracy HTML and many Postscript and PDF files. TMG allows as option a variety of term-weighting and normalization schemes as well as stemming.

AVAILABILITY AND REQUIREMENTS

To obtain the current version of TMG please file a request form.

- Current version of TMG 8.0. Some modules are in MATLAB precompiled (.p) format. See version TMG 3.084 (below) for study MATLAB open source with indexing and retrieval modules.
- Depending on user requirements, version 8.0 might necessitate the following third party software packages.
 - If one intends to use MATLAB functionality:
 - You will need access to the [MATLAB Database Toolbox](#)
 - You need to download [EUSQL](#) and [Java Connector](#)
- If one wishes to apply the packages below for EUC and low-rank approximation, the following packages are already preloaded:
 - [SVD/SVD++](#)

More than 4000 requests worldwide . . .

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

Outline

- 1 Introduction
- 2 Retrieval Task**
- 3 Experiments
- 4 Work In Progress

Retrieval Latent Semantic Analysis

Data Explosion



Big Data Collections \Rightarrow

Parsing & Processing



Large & Sparse
Term-Document Matrices



Dimensionality Reduction



Difficult Management



Low Rank
Approximation

Low Rank Approximation

Given an $m \times n$ matrix A and a rank parameter $k \ll \min\{m, n\}$, the Low-Rank Approximation problem is to find a matrix Z of rank k such that $\|A - Z\|_{2,F}$ is sufficient small.

Eckart-Young Theorem

The minimization problem:

$$\min_{\text{rank}(Z)=k} \|A - Z\|_{2,F}$$

has a solution given by the truncated SVD:

$$Z = A_k = U_k S_k V_k^T$$

Truncated Versions of SVD

- + reveals latent semantic structure
- + construct orthogonal bases for the terms (rows) and documents (columns)

Low Rank Approximation

Given an $m \times n$ matrix A and a rank parameter $k \ll \min\{m, n\}$, the Low-Rank Approximation problem is to find a matrix Z of rank k such that $\|A - Z\|_{2,F}$ is sufficient small.

Eckart-Young Theorem

The minimization problem:

$$\min_{\text{rank}(Z)=k} \|A - Z\|_{2,F}$$

has a solution given by the truncated SVD:

$$Z = A_k = U_k S_k V_k^T$$

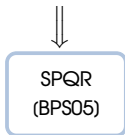
Truncated Versions of SVD

- + reveals latent semantic structure
- + construct orthogonal bases for the terms (rows) and documents (columns)
- requires the entire matrix in RAM
- lacks interpretability
- results in dense factors
- is time inefficient

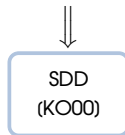
Replacing SVD

Basic Concerns

① Time Efficient Algorithms



② Storage Efficient Algorithms



③ Interpretable Algorithms

"Matrices are about their columns and rows" (G. Strang)

Concerning

NMF

Deterministic

Pseudo-Skeleton
(GTZ97)

SCRA
(BPS05)

Randomized

CUR
(DMM08)

CX
(DMM08)

Basic Citation:

SIAM J. MATRIX ANAL. APPL.
Vol. 30, No. 2, pp. 844-881

© 2008 Society for Industrial and Applied Mathematics

RELATIVE-ERROR CUR MATRIX DECOMPOSITIONS*

PETROS DRINEAS¹, MICHAEL W. MAHONEY², AND S. MUTHUKRISHNAN³

Basic Idea:

Randomly select columns (and) rows of A based on probability vectors constructed by dominant right/left singular vectors (2nd Generation).

CUR

$$\begin{array}{c} \boxed{A} \\ m \times n \end{array} = \begin{array}{c} \boxed{C} \\ m \times c \end{array} \begin{array}{c} \boxed{U} \\ c \times r \end{array} \begin{array}{c} \boxed{R} \\ r \times n \end{array}$$

CX

$$\begin{array}{c} \boxed{A} \\ m \times n \end{array} = \begin{array}{c} \boxed{C} \\ m \times c \end{array} \begin{array}{c} \boxed{X} \\ c \times n \end{array}$$

CUR/CX Algorithms in TMG Graphical Interface

Text to Matrix Generator - Dimensionality Reduction

Window Help

Text to Term-Document Matrix (tdm) Generator

Select Dataset:

Method

- Singular Value Decomposition (SVD)
- Principal Component Analysis (PCA)
- Clustered Latent Semantic Indexing (CLSI)
- Centroid Method (CM)
- Semidiscrete Decomposition (SDD)
- SPQR
- SCRA
- CUR/CX
- CGR

Compute SVD with

- MATLAB (svds)
- Propack

Compute Submatrix with

- MATLAB
- Proved SP

CUR/CX Algorithms

Number of Columns:

Number of Rows:

Number of Runs:

Select Columns/Rows: Exact Expected

Clustering Algorithm

- Euclidean K-means
- Spherical K-means
- Fuzzy

Initialize Centroids: At random

Termination Criterion: Epsilon 1

Principal Directions: 1 Variant: Basic

Maximum num. of PCs:

Automatic Determination of Num. of factors for each cluster:

Number of Clusters: Display Results

Select at least one factor from each cluster (Recommended for classification)

Number of Factors: Store Results

GUI usage

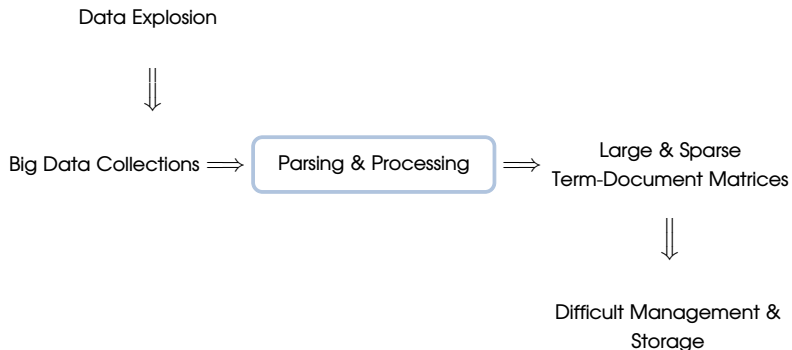
- 1 Select CUR/CX radio button
- 2 Block CUR/CX Algorithm is activated
- 3 Number of Columns & Number of Rows fields define the algorithm

Algorithm	Number of Columns	Number of Rows
CX	✓	×
CUR	✓	✓
Error	×	✓
Error	×	×

- 4 Number of Runs field defines the repetitions of the algorithm
- 5 Random selection policies:

Selection Method	Description
Exactly	select the exact desired number of rows/columns
Expected	select in expectation the desired number of rows/columns

Retrieval Vector Space Model



Retrieval Vector Space Model

Consider a document-term matrix A and a query vector q . We seek documents similar to the query.

Similarity Measures

① Euclidean Distance

$$Aq$$

② Cosine

$$\frac{Aq}{\|A\| \|q\|}$$



$$Aq$$

Data Intensive



Reduction of the size of A and q

BMM/PIP Algorithms

Basic Matrix Multiplication (DKM04)

SIAM J. COMPUT.
Vol. 30, No. 1, pp. 132–157

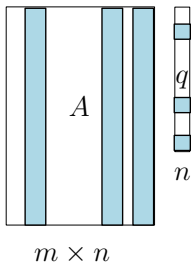
© 2006 Society for Industrial and Applied Mathematics

FAST MONTE CARLO ALGORITHMS FOR MATRICES I: APPROXIMATING MATRIX MULTIPLICATION*

PETROS DRINEAS¹, RAVI KANNAN¹, AND MICHAEL W. MAHONEY²

Basic Idea:

Randomly select columns of A and elements of q based on one probability vector constructed by the euclidean norms of the columns of A and q .



Probabilistic Inner Product (EB+11)

SIAM J. SCI. COMPUT.
Vol. 33, No. 4, pp. 1889–1706

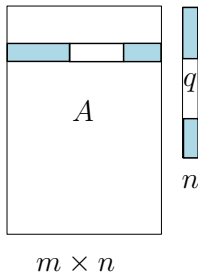
© 2011 Society for Industrial and Applied Mathematics

IMPORTANCE SAMPLING FOR A MONTE CARLO MATRIX MULTIPLICATION ALGORITHM, WITH APPLICATION TO INFORMATION RETRIEVAL*

SYLVESTER ERIKSSON-BIQUE¹, MARY SOLBRIG², MICHAEL STEFANELLI³,
SARAH WARKENTIN⁴, RALPH ABBEY¹, AND ILSE C. F. IPSEN⁵

Basic Idea:

Randomly select columns of A and elements of q based on n probability vectors constructed by the euclidean norms of the columns of A and q .



BMM/PIP Algorithms in TMG Graphical Interface

Text to Matrix Generator - Retrieval

Window Help

Text to Term-Document Matrix (tdm) Generator

Select Dataset:

Insert query:

Alternative Global Weights: Use Stored Global Weights

Stoplist:

Local Term Weighting:

Vector Space Model

Latent Semantic Analysis by:

Number of Factors:

Similarity Measure:

Retrieve Documents:

- Number of most relevant:
- Similarity measure exceeds:

Probabilistic:

- Probabilistic Inner Product
Num. Iterations:
- Basic Matrix Multiplication
Num. Samples:

BMM/PIP Algorithms in TMG Graphical Interface

GUI usage

- 1 Select Vector Space Model radio button
- 2 Block Probabilistic is activated
- 3 Algorithm Selection

Algorithm	Fields	Characteristics
Probabilistic Inner Product	Num .	repetitions
	Iterations	sampling 1% of elements
Basic Matrix Multiplication	Num . Samples	samples

Outline

- 1 Introduction
- 2 Retrieval Task
- 3 Experiments**
- 4 Work In Progress

CRANFIELD Collection (CRAN) (2568 × 1398)

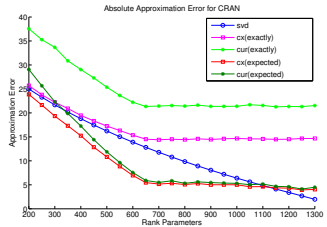
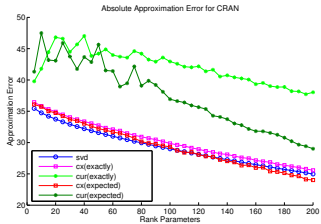
Cranfield collection of 1398 documents, 225 queries

Parsing Options			
Terms		Min Global Frequency	Weights (NGL)
numerics removal	stemming	2	cae
DR Parameters			
Rank	$k = [5 : 5 : 200]$		
	$k = [200 : 50 : 1300]$		
CUR/CX	Num. Runs	Num. Rows	Num. Columns
	5	4k (2500 if $2k > 2568$)	2k (1300 if $2k > 1398$)
VSM Parameters			
PIP	Num. Iterations		
	10		
BMM	Num. Samples		
	[10 : 50 : 2500]		

System Specifications

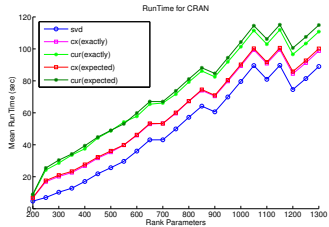
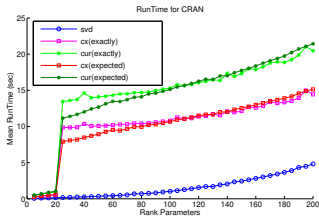
System Specifications			
Processor	RAM	MATLAB	OS
Intel Core i5 2500 (4Cores) @3.3GHz	16 GB	R2012b	Debian 3.2.12-1

DR Approximation Error

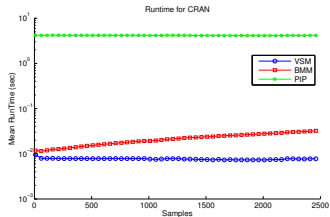


Running Time

Latent Semantic Indexing

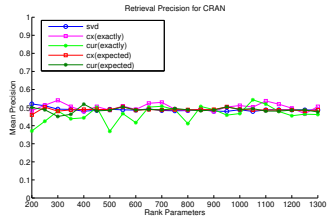
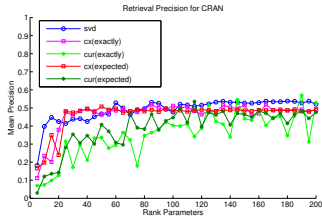


Vector Space Model

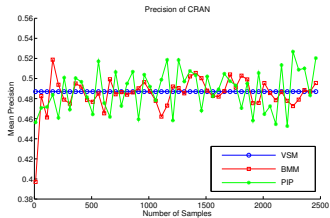


Retrieval Accuracy

Latent Semantic Indexing

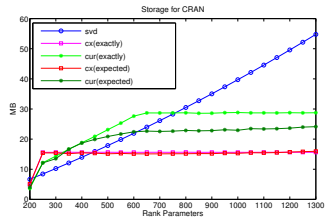
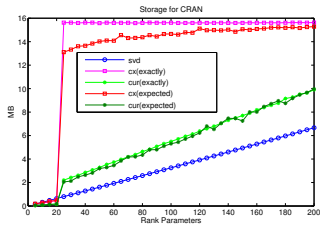


Vector Space Model

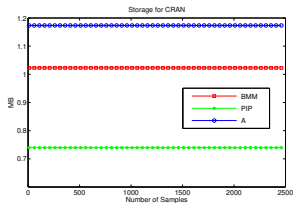


Storage

Latent Semantic Indexing



Vector Space Model





Outline

- 1 Introduction
- 2 Retrieval Task
- 3 Experiments
- 4 Work in Progress**

Summary

Goals:

-  rapid familiarization of randomized techniques
-  prototyping and testing new algorithms

Summary

Goals:

- 😊 rapid familiarization of randomized techniques
- 😊 prototyping and testing new algorithms

Landscape:

	Tasks	Categories
Term Document Matrix	Dimensionality Reduction	General DRM
		NMF
	Retrieval	Vector Space Model
		Latent Semantic Analysis
	Clustering	
	Classification	

■ Present Work

■ Future Work/
Work in Progress

Work in Progress:

- Parallel implementations
- Construction of a complete & extensible Randomization Module
 - ✓ Randomized NMF
 - ✓ Randomized Clustering
 - ✓ Randomized Classification
- Incorporating methods for handling efficiently the coupling

Text Data + Algorithms + Friendly MATLAB Interface

Work in Progress:

- Parallel implementations
- Construction of a complete & extensible Randomization Module
 - ✓ Randomized NMF
 - ✓ Randomized Clustering
 - ✓ Randomized Classification
- Incorporating methods for handling efficiently the coupling

Text Data + Algorithms + Friendly MATLAB Interface

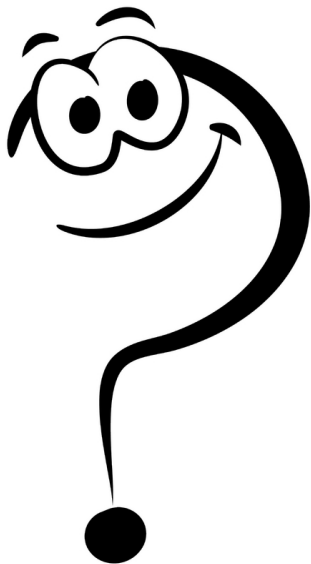
Your Algorithms are welcome!



Thank You!

Thank you!

Questions ?



Bibliography I

- (BPS05) M. W. Berry, S. A. Pulatova, and G. W. Stewart. "Algorithm 844: Computing sparse reduced-rank approximations to sparse matrices". In: *ACM TOMS* 31.2 (June 2005), 252–269.
- (DKM04) P. Drineas, R. Kannan, and M. W. Mahoney. "Fast Monte Carlo algorithms for matrices I: Approximating Matrix Multiplication". In: *SISC* 36.1 (2004), 132–157.
- (DMM08) P. Drineas, M. Mahoney, and S. Muthukrishnan. "Relative-Error CUR Matrix Decompositions". In: *SIMAX* 30.2 (Sept. 2008), 844–881.
- (EB+11) S. Eriksson-Bique et al. "Importance Sampling for a Monte Carlo Matrix Multiplication Algorithm, with Application to Information Retrieval". In: *SISC* 33.4 (2011), 1689–1706.
- (GTZ97) S. Goreinov, E. Tyrtyshnikov, and N. Zamarashkin. "A theory of pseudoskeleton approximations". In: *Linear Algebra and its Applications* 261.1-3 (1997), 1–21.
- (KO00) T. Kolda and D. O'Leary. "Algorithm 805: Computation and Uses of the Semidiscrete Matrix Decomposition". In: *ACM TOMS* 26.3 (Sept. 2000), 415–435.

Bibliography II

- (ZG06) D. Zeimpekis and E. Gallopoulos. ``TMG: A MATLAB toolbox for generating term document matrices from text collections''. In: *Grouping Multidimensional Data: Recent Advances in Clustering*. Ed. by J. Kogan, C. Nicholas, and M. Teboulle. Berlin: Springer, 2006, 187–210.