

DICE @ Rich Context Competition

Prof. Dr. Axel Ngonga
Nikit Srivastava
Richa Jalota

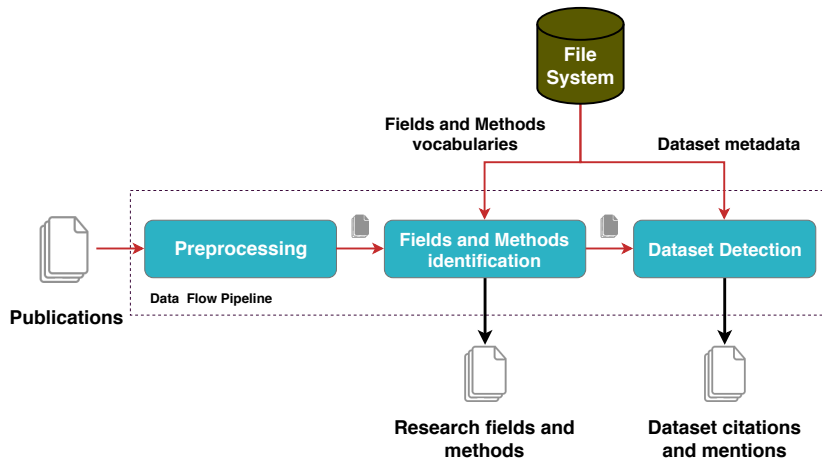


Data Science Group
Paderborn University

February 15, 2019

Section 1

Overview





- Generated text files from PDF using pdftotext^a

ISBN 3865581226 (Internetversion)

×ØØ
Ø
í× ò ×ØÓ
×Ø
×
ÓÚØØ
ØÓØ ØØØØ
× Ò Ò Ò Ò Ò Ò
Ø ×ØØØØ ØØ ØØ Ò ØØ ¹
Ø ØØ Ý Ø × ØØ Ø
× Ø Ø Ø × ØØØ¹Ø Ò Ò ØØ Ò ×Ø Ò Ø Ø ØØØØØ× Ø Ò × ÓØ ØØØ
××
Ú Ø Ø ØØ Ò
ØØ ×Ú Ø
Ø Ò Ó ØØØ×ØÝ° Ì Ý Ø ÒØ
Ø ØØ×, Ò ØØ Ò Ø ØØ ¹NØØØ
Ø Ø ØØØ n¹Ø Ò Ó Ý Ø ×, Ú Ø
ØØØ Û¹
ØØ
Ú × Ò ØØØ Ø Ø Ø Ò Ø Ø Ø Ø
Ø Ò × ÓØ Ú ØØ Ò ×Ø
Ø Ú ØÝ° Ì × × Ò
ÓØØØ ×Ø ØØ Ø Ò Ø × ØØØ¹Ø Ø ØØØ
××× Û
NØØÝ
Ø Ò Ý Ø ÒØ
Ø ØØ Ì ØØ ØØ Ø× ØØ Û

Sample text



Centrality-based Capital Allocations*

Adrian Alter[‡], Ben Craig[‡], Peter Raupach[§]

December 19, 2014

Abstract

We look at the effect of capital rules on a banking system that is connected through correlated credit exposures and interbank lending. The rules, which combine individual bank characteristics and interconnectivity measures of interbank lending, are to minimize a measure of system-wide losses. Using the detailed German Credit Register for estimation, we find capital rules based on eigenvectors to dominate any other centrality measure, followed by closeness. Compared to the baseline case, capital reallocation based on the Adjacency Eigenvector saves 14.6% in system losses as measured by expected bankruptcy costs.

Keywords: Capital Requirements, Centrality Measures, Contagion, Financial Stability

JEL classification: G21, G28, C15, C81.

Sample text

- Generated text files from PDF using pdftotext^a
- Handled words that got split by hyphens
- Removed irrelevant data (references, acknowledgments...)
- Extracted main sections (Abstract, JEL-Classification code, keywords, methodology/data, summary, discussion)



Centrality-based Capital Allocations*

Adrian Alter[‡], Ben Craig[‡], Peter Raupach[§]

December 19, 2014

Abstract

We look at the effect of capital rules on a banking system that is connected through correlated credit exposures and interbank lending. The rules, which combine individual bank characteristics and interconnectivity measures of interbank lending, are to minimize a measure of system-wide losses. Using the detailed German Credit Register for estimation, we find capital rules based on eigenvectors to dominate any other centrality measure, followed by closeness. Compared to the baseline case, capital reallocation based on the Adjacency Eigenvector saves 14.6% in system losses as measured by expected bankruptcy costs.

Keywords: Capital Requirements, Centrality Measures, Contagion, Financial Stability

JEL classification: G21, G28, C15, C81.

Sample text

- Generated text files from PDF using pdftotext^a
- Handled words that got split by hyphens
- Removed irrelevant data (references, acknowledgments...)
- Extracted main sections (Abstract, JEL-Classification code, keywords, methodology/data, summary, discussion)
- Extracted noun phrases from these sections

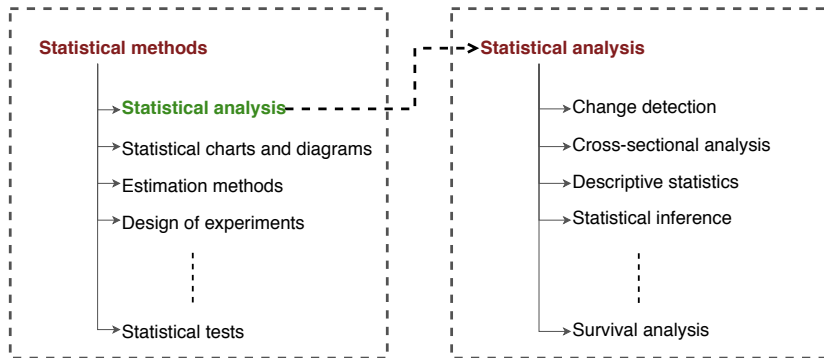
^apoppler-utils/pdfinfo.1.en.html

Section 2

Identification of Research Methods and Fields



- Research methods vocabulary from DBpedia
 - Curated statistical methods from Wikipedia¹
 - Scraped their description from corresponding DBpedia resources
 - Extracted noun phrases from description



Hierarchy of Wikipedia Categories



- Research methods vocabulary from DBpedia
 - Curated statistical methods from Wikipedia¹
 - Scraped their description from corresponding DBpedia resources
 - Extracted noun phrases from description

Statistical methods

- Bayesian model reduction
- DeMix
- Genetic assignment methods
- Kneser–Ney smoothing
- ⋮
- Synthetic control method


Statistical analysis

- Boolean analysis
- Fenwick (statistic)
- Functional data analysis
- InfoQ
- ⋮
- Weighting

Resources in Wikipedia Category



- Research methods vocabulary from DBpedia
 - Curated statistical methods from Wikipedia¹
 - Scraped their description from corresponding DBpedia resources
 - Extracted noun phrases from description
- SAGE Research fields vocabulary
 - Created blacklist of irrelevant fields
 - Extracted noun phrases from description




Meta Analysis
Mixed Methods
Narrative Analysis
Case Study and Narrative Analysis
...

Blacklisted terms

¹https://en.wikipedia.org/wiki/Category:Statistical_methods



- Research methods vocabulary from DBpedia
 - Curated statistical methods from Wikipedia¹
 - Scraped their description from corresponding DBpedia resources
 - Extracted noun phrases from description
- SAGE Research fields vocabulary
 - Created blacklist of irrelevant fields
 - Extracted noun phrases from description
- Generated vector models for both Research Fields and Methods



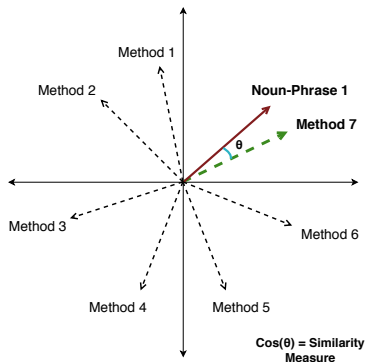
Meta Analysis
Mixed Methods
Narrative Analysis
Case Study and Narrative Analysis
...

Blacklisted terms

¹https://en.wikipedia.org/wiki/Category:Statistical_methods

Identification of Research Methods and Fields

Methods Identification - Word2Vec

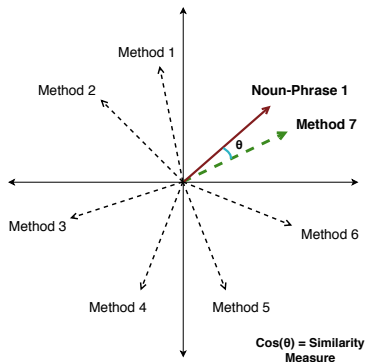


- Closest research method vector found by measuring cosine similarity between noun phrase vectors and method vectors

Noun-phrases from text and methods in Embedding Space

Identification of Research Methods and Fields

Methods Identification - Word2Vec

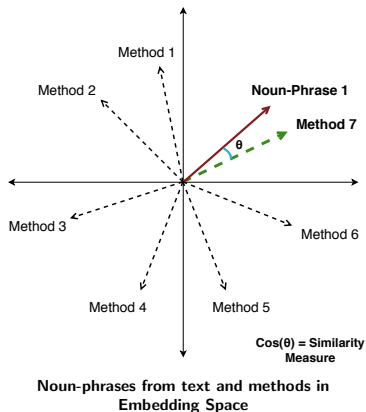


Noun-phrases from text and methods in
Embedding Space

- Closest research method vector found by measuring cosine similarity between noun phrase vectors and method vectors
- Computed the significance of recurring methods using IDF
- Each research method assigned a weightage

Identification of Research Methods and Fields

Methods Identification - Word2Vec



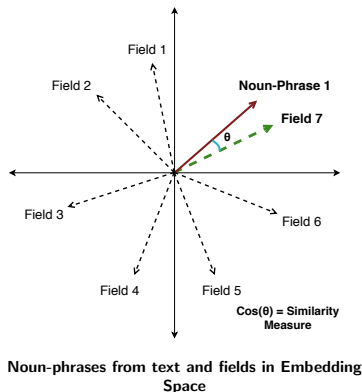
- Closest research method vector found by measuring cosine similarity between noun phrase vectors and method vectors
- Computed the significance of recurring methods using IDF
- Each research method assigned a weightage
- Re-ran the algorithm to find closest method vector and then sorted the pairs based on weighted cosine similarity
- Pair with highest cosine similarity chosen

Identification of Research Methods and Fields

Fields Identification - Word2Vec



- Top 10 closest research field vectors found using cosine similarity between noun phrase vectors and field vectors
- Pairs with similarity score < 0.9 filtered out
- If not blacklisted, top-ranked term marked as Research Field



Section 3

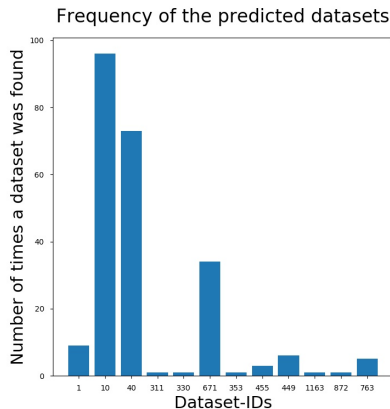
Dataset Detection

Dataset Detection

Simple dataset mention search



- Unique datasets selected from Dataset metadata
- Datasets searched in publication text
- Frequency of dataset mentions calculated
- Removed mentions $>$ threshold-value
* median of dataset-frequency





- Trained Entity Extraction model using **Rasa NLU**²
 - Training data : 2500 labeled publications (phase-1 training data)
 - Training data : 7500 labeled publications (entire phase-1 corpus)
- **Only those entities that**
 - Had confidence-score > threshold-value
 - Belonged to the Research Field of the article**considered as datasets.**

²<https://rasa.com/docs/nlu/>

Dataset Detection

Combining the two approaches



- Removed irrelevant mentions from Rasa-identified datasets
- Took a **union of results** from two approaches

Section 4

Results



Quantitative Evaluation of Datasets against Validation data

	Rasa-based Approach (2500)	Rasa-based Approach (7500)	Combined Approach (2500)	Combined Approach (7500)
Precision	0.382	0.388	0.456	0.456
Recall	0.26	0.26	0.31	0.31
F1	0.309	0.311	0.369	0.369

Numbers in brackets indicate training samples

Improvement - 19.42%



Evaluation against Phase-1 holdout

pub_id	Keywords	Phase-1	Phase-2
10328	Cycling for transport, leisure and sport cyclists	Health evaluation	Public health and health promotion
7270	Older adult drug users, harm reduction	Health Education	Correctional health care
6053	Economic conditions - crime relationship, homicide	Homicide	Gangs and crime



pub_id	Keywords	Phase-1	Phase-2
10328	Thematic content analysis	Thematic analysis	Sidak correction
7270	interviews conducted face to face, finding systematic patterns or relationships among categories identified by reading the interview transcript	Qualitative interviewing	Sampling design
6053	Autoregressive integrated moving average (ARIMA) time-series model	Methodological pluralism	Multivariate statistics

Section 5

Challenges Encountered and Future Agenda



- Appropriate extraction of text from PDFs
 - Extraction of specific sections from text



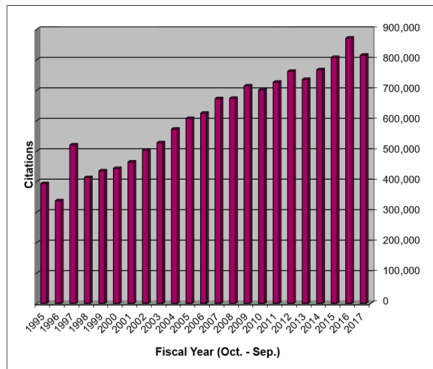
- Appropriate extraction of text from PDFs
 - Extraction of specific sections from text
- Not enough uniformity in labeled data to detect datasets



- Appropriate extraction of text from PDFs
 - Extraction of specific sections from text
- Not enough uniformity in labeled data to detect datasets
- Some polysemous and high-level terms in the given SAGE vocabularies



- Appropriate extraction of text from PDFs
 - Extraction of specific sections from text
- Not enough uniformity in labeled data to detect datasets
- Some polysemous and high-level terms in the given SAGE vocabularies
- Unavailability of Open Ontologies for Social Science Fields and Methods



Number of indexed citations added to MEDLINE during each fiscal year since 1995³

³https://www.nlm.nih.gov/bsd/stats/cit_added.html



Vision: Giant Global Scientific Knowledge Graph



Vision: Giant Global Scientific Knowledge Graph

- Never ending



Vision: Giant Global Scientific Knowledge Graph

- Never ending
- Distributed



Vision: Giant Global Scientific Knowledge Graph

- Never ending
- Distributed
- Self-feeding (focused crawling)



Vision: Giant Global Scientific Knowledge Graph

- Never ending
- Distributed
- Self-feeding (focused crawling)
- Self-repairing (introspection)



Vision: Giant Global Scientific Knowledge Graph

- Never ending
- Distributed
- Self-feeding (focused crawling)
- Self-repairing (introspection)
- Humans in the loop (weak supervision)



Vision: Giant Global Scientific Knowledge Graph

- Never ending
- Distributed
- Self-feeding (focused crawling)
- Self-repairing (introspection)
- Humans in the loop (weak supervision)
- Standardized access (SPARQL, Linked Data Fragments, question answering, etc.)



ScholarBrew

Distilling scholarly data from the Web

Search for person, organization or keyword..

Example search: [Axel Ngonga](#)

[Faceted search](#)

 5382618  8303571  7210983 [More stats](#)

Check out: s2.exynize.com



Prof. Dr.
Axel Ngonga



Dr. Ricardo
Usbeck



Michael
Röder



Daniel
Vollmers



Nikit
Srivastava



Richa Jalota



Rene Speck



Thank You! Questions?

Follow us

 [@DiceResearch](https://twitter.com/DiceResearch)

 github.com/dice-group

 cs.uni-paderborn.de/ds/

