

Searching for Research Fraud in OpenAlex with Graph Data Science

Adam Day - SAGE Publishing Ebru Cucen - OpenCredo



OC - Sage R&D Project



Fahran Wallace



Sebastian Margineanu



Helen King







Key Takeaways

- 1. Graph analytics is a good approach to identify organised misconduct
- 2. To do this efficiently, we need reliable data structures, such as a knowledge graph built on well-curated data
 - OpenAlex is a good data set for this
- 3. We can detect unusual co-authorships. Either with machine-learning, or a simple rule:
 - It is unusual for an author to co-author with a different institution on their first paper

The value of graphs

Imagine what happens when we link these accounts to:

- Authors who recommended them
- Coauthors
- Past papers....



The different kinds of paper mills



Candal-Pedreira et al BMJ 2022; 379 DOI: https://doi.org/10.1136/bmj-2022-071517 (Published 28 November 2022)

AI isn't really helping!

- Already a number of AI-based tools to detect paper mills.
 - But these don't solve the problem on their own
- Already examples of paper mills using generated text from easy-to-detect models like GPT-2
 - Interestingly, NO examples of detections of GPT-3 or larger models
- Generating fake scientific images with Stable Diffusion is trivial
- Meta recently released a model specifically for generating scientific text.
 - Swiftly pulled!



Update: Meta's Galactica Al Criticized as 'Dangerous' for Science

Model quickly pulled after renowned experts said results were 'statistical nonsense' and 'wrong'



Retraction Watch

Tracking retractions as a window into the scientific process





A company in Russia hawks its wares

← Tweet



Anna Abalkina @AbalkinaAnna

This is my new preprint on @arxiv "Publication and collaboration anomalies in academic papers originating from a paper mill: evidence from Russia". At least 303 papers were identified. and a set of predictors of a Russian paper mill were proposed.

...

arxiv.org/abs/2112.13322

Example of the identification of the offer Date of issue August September 2019 Тема доступна только клиентам. которые оплатили Август Сентябрь 2019 выход #187 курнала five co-authors 26-04-2019 Набор в Google Translate of the topic: журнал до Superhero Cinema: Refraction 5 чел. (авторов) в этой статье 126000 py of Topical Issues in the Modern Epic В статье рассмотрены такие вопросы. COABT 2-й 3-й 5-**й** Супергеройское кино: преломление OD злободневных проблем в современном свобод свобод свобод свобод свобод эпосе HO HO HO HO HO Рубли 32200 28700 25200 21700 18200 00 1 место - свободно (продается) No Nº No No No 2 место - свободно (продается) 187 1 18 2 187.3 187.4 187.5 3 место - своболно (продается) Country of the journal 4 место - свободно (продается) Название журнала доступно только 5 место - свободно (продается) клиентам, которые оплатили Venezuela Номер заказа № 187 O3 Scopus SJR=0.199 Самый быстрый и простой способ заказать или задать вопрос, написать Специализация журнала: Arts and Humanities (miscellaneous) на WhatsApp 😕 или Viber 🏁 или article@ Social Sciences (miscellaneous) 123mi ru или по телефону +7 (968) 655-29-49 или Нажмите Перейти в контакты и свяжитесь по почте, телефону или другим способом с любым менеджером. Scopus Superhero movie: Breaking the challenges of topics in the Akim, K., Kara-2019 Opcion modern epos l [Película de superhéroes: Rompiendo los Murza, G., Saenko, 15(Special Issue 22), c. 1408-1428 desafios de los temas en los epos modernos] N., Suharyanto, A., Kalimullin, D. Scimago Opcion Universidad del Zulia

Questions:

- Can we take the publicly-available data for these papers and model the co-authorships?
- Can we work out when a co-authorship is unusual and worth looking into more closely?



Data



- Free, Hundreds of millions of entities
- Inspired by the ancient Library of Alexandria
- Launched at Sep '22
- Updated regularly
- 4 ways to access:
 - Web
 - API for ad-hoc queries
 - AWS Snapshot : 1.7 TB
 - Snowflake Marketplace

- of Authors : people who create works
- Works : papers, books, datasets, etc; they cite other works
- 💡 Concepts : keywords associated with works
- 🚛 Venues : journals and repositories that host works
- institutions : universities and other orgs that are affiliated with works (via authors)



Source: Fundamentals of Data Engineering by Joe Reis, Matt Housley



Fundamentals of Data Engineering by Joe Reis, Matt Housley



Source: https://www.kidslovegreece.com/greece_online/the-12-labors-of-hercules/



Source: https://www.kidslovegreece.com/greece_online/the-12-labors-of-hercules/



Fundamentals of Data Engineering by Joe Reis, Matt Housley

Learn/Optimise Aggregate / Label Explore/Transform Move/Store Collect

Data LifeCycle



Concepts

1	{'ancestors': [],							
2	'cited_by_count': 391428928,							
3 >	'counts_by_year': [{'cited_by_count': 26308154,-							
36	'created_date': '2016-06-24',							
37 >	'description': 'theoretical study of the formal foundation enabling the '							
48	'display_name': 'Computer science',							
41	'id': 'https://openalex.org/C41008148',							
42	'ids': {'mag': '41008148',							
43	'openalex': 'https://openalex.org/C41008148',							
44	'umls_cui': ['C0599726'],							
45	'wikidata': 'https://www.wikidata.org/wiki/021198'.							
46	'wikipedia': 'https://en.wikipedia.org/wiki/Computer%20science'},							
47	'image thumbnail url': 'https://upload.wikimedia.org/wikipedia/commons/thumb/6/6a/Sorting guicksort anim.gif/100px-Sorting guicksort anim.gif',							
48	'image_url': 'https://upload.wikimedia.org/wikipedia/commons/6/6a/Sorting_quicksort_anim.gif',							
49	' دراسة العصليات التي تتغايل سع بيانات' : {'ar': (description': {'ar')							
50	و التي يعكن تعقيلها كبيانات في شكل ا							
51	, 'بر امج'							
52	'az': 'elm sahəsi'.							
53 >	'bn': 'কপ্রিউটার নামক যথের সাহায্যে '							
56 >	'bs': 'naučna disciplina koja se bavi svim '							
58 >	'ca': 'estudi teòric de la base formal que '							
61 >	'cs': 'věda studující principy strojového '-							
63 >	'da': 'disciplin med rødder i matematisk '							
65 >	'de': 'Wissenschaft von der systematischen '							
67 >	'en': 'theoretical study of the formal '-							
73 >	'en-ob': 'study of the theoretical '-							
76 >	ing of stars and the charter of the stars and the im-							
78 >	'est': 'ciencia dedicada a la comutación y '-							
80 >	fair 's stants Attack a stants at the							
82 >	fil' tisteenala joka tutkii '							
85)	ifri tartivité crientificue technique et '-							
80 \	1) - L'interità dedicatà é computación e an lu							
91 >								
93 >	The the state of t							
95 \	id) - Madinate no systematical							
07	iti i rejun daan teretai i fandamanti							
100	1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.							
101								
102 >								
105 >	while 'fanfalt and rather i matematick '							
109 >	10. Identical participanti and 1							
111	The scale value of a non-net en -							
112 >	isti larta da fugamente totica da l							
115 /	pt. estado dos fundamentos teorizos da							
113 /	prebi letito ese ante en ante							
110	in , grining care se occupation and a predictarea -							
120 >	la , declamina do engonacia ;							
122 >	sti i stuta ti stuta t augusta t augusta i stuta i stuta t e cauda aziani							
125	sor accenterize on protected opprodult "							
127 >								
120 >	su , statini i texilotogi se se "							
123 /	ti : Detrit evien kurataria dayat,							
134 >	'un · naynuba ghuqunulana ;							
137)	Vec - Stenso che ca stauta - l'olimantati							
140	vi i ngaming nguan cu ve cu so ky "							
141	Han a synthe Al studeye tes constructes ;							
142 >	All + W751862/Mall ##9528882892 7, Idionalay name / 1/16/1: Toformatikal							
331	Tovol's A							
332 >	related concents': [['disnlay name': 'Mathematics'							
797	Indated date: 1202-11-2017-48-02.45438-							
798	'wikidata': 'https://www.wikidata.org/wiki/021198'.							
799	'works api url': 'https://api.openalex.org/works?filter=concepts.id:C41008148'.							
800	'works count': 76576650}							
	works_count : /05/0050/							

Concepts

1	{'ancestors': [],							
2	'cited_by_count': 391428928,							
3 >	'counts_by_year': [{'cited_by_count': 26308154,							
36	'created_date': '2016-06-24',							
37 >	'description': 'theoretical study of the formal foundation enabling the '							
40	'display_name': 'Computer science',							
41	'id': 'https://openalex.org/C41008148',							
42	'ids': {'mag': '41008148',							
43	'openalex': 'https://openalex.org/C41008148',							
44	'umls_cui': ['C0599726'],							
45	'wikidata': 'https://www.wikidata.org/wiki/Q21198',							
46	'wikipedia': 'https://en.wikipedia.org/wiki/Computer%20science'},							
47	'image_thumbnail_url': 'https://upload.wikimedia.org/wikipedia/commons/thumb/6/6a/Sorting_quicksort_anim.gif/100px-Sorting_quicksort_anim.gif',							
48	'image_url': 'https://upload.wikimedia.org/wikipedia/commons/6/6a/Sorting_quicksort_anim.gif',							
49	' دراسة الععليات التي تتفاعل سع بيانات' :'international': {'ar':							
50	' و التي يعكن تعتينها كبيانات في شكل'							
51	ابراسج'							
52	'az': 'elm sahəsi',							
53 >	'bn': 'কন্সিউটার নামক যন্ধের সাহায্যে '							
56 >	'bs': 'naučna disciplina koja se bavi svim '							
58 >	'ca': 'estudi teòric de la base formal que '							
61 >	'cs': 'věda studující principy strojového '							
63 >	'da': 'disciplin med rødder i matematisk '							
65 >	'de': 'Wissenschaft von der systematischen '-							
67 >	'en': 'theoretical study of the formal '							
73 >	'en-gb': 'study of the theoretical '							
76 >	'eo': 'scienco pri teoria fondado de '							
78 >	'es': 'ciencia dedicada a la computación y '							
80 >	··· علم اطلاعات و يردازش اطلاعات و "··							
82 >	'fi': 'tieteenala, joka tutkii '							
85 >	'fr': 'activité scientifique, technique, et '-							
89 >	'gl': 'ciencia dedicada á computación e ao '							
91 >	"he': "חחום מחקר של תיאוריות בסיסיות"							
93 >	'hsb': 'wédomosć wo systematiskim '							
95 >	'id': 'kajian dasar teoretis dari informasi '-							
97 >	'it': 'scienza che studia i fondamenti '							
100	'ja': '情報と計算の理論的基礎、及びそのコンピュータ上への実装と応用に関する研究分野',							
101 >	'kn': 'ಮಾಹಿತಿ ಮತ್ತು ಗಣನೆಯ ಸೈದ್ಯಾಂತಿಕ '							
103 >	'ko': '전산 이론, 소프트웨어, 하드웨어의 중점을 둔\xa0정보과학의 한 '~							
105 >	'nb': 'fagfelt med røtter i matematisk '							
108 >	'nl': 'de studie van informatie en '							
111 >	'pl': 'nauka zajmująca się przetwarzaniem '							
113 >	'pt': 'estudo dos fundamentos teóricos da '							
115 >	"pt-br": "estudo dos fundamentos teoricos							
117 >	'ro': 'stlința care se ocupa cu prelucrarea							
119	ги:: дисциплина оо информации,							
120 >	sch : "scienza chi studia t'eladourazzoni "-							
123 2	sto: strencerte an practice approach ""							
125 /	st: veda o racunakniski oddelavi							
120 >	sq , studini i tenintugjise se							
129 /	ci : betrit evien nulatarina dayati, "							
134)	lugri 'Sianga cha za studia i fondaminti '							
137 >	'vi' ' 'nahnh hog nghiện cứu về cơ sở lý '-							
148	"wit 'typice ki studeve les conjutrees'.							
141	「カト」、「研究信頼TFF留約理論基礎」)。							
142 >	'display name': {'af': 'Informatika'							
331	'level': 0.							
332 >	'related concepts': [{'display name': 'Mathematics'							
797	'updated date': '2022-11-29107:48:02.454438'.							
798	'wikidata': 'https://www.wikidata.org/wiki/021198',							
799	'works_api_url': 'https://api.openalex.org/works?filter=concepts.id:C41008148',							
888	'works_count': 76576650}							

"Property"	"Value"					
"cited_by_count"	"231013792"					
"description"	"theoretical study of the formal foundation enabling the automated pro cessing or computation of information, for example on a computer or ov er a data transmission network"					
"display_name"	"Computer science"					
"id"	"https://openalex.org/C41008148"					
"image_thumbnail_url"	"https://upload.wikimedia.org/wikipedia/commons/thumb/6/6a/Sorting_qui cksort_anim.gif/100px-Sorting_quicksort_anim.gif"					
"image_url"	"https://upload.wikimedia.org/wikipedia/commons/6/6a/Sorting_quicksort _anim.gif"					
"level"	"0"					
"updated_date"	"2022-04-12"					
"wikidata"	"https://www.wikidata.org/wiki/Q21198"					
"works_api_url"	"https://api.openalex.org/works?filter=concepts.id:C41008148"					
"works_count"	"27478212"					

Authors



Categoric	al Features	Chart to show					
count	missing	unique	top	freq top	avg len	custom	
last_known 653k	_institution 46.56%	-	-	-	-	data type: struct	0.2
most_cited_ 1.20M	_work 1.5%	1.07M Sea	arch f	2,915	95.3	data type: string	SHOW RAW DATA
orcid 48.1k	96.06%	46.0k htt	ps://or	2	37	data type: string	10 30 50 70 90 SHOW RAW DATA
updated_da 1.22M	te 0%	1.15M 20	22-10	1	26	data type: string	10 30 50 70 90 SHOW RAW DATA
works_api_u 1.22M	url 0%	1.21M htt	ps://a	1	58.98	data type: string	10 30 50 70 90 SHOW RAW DATA
x_concepts 1.22M	0%	-	-	-	-	data type: array max size: 126 min size: 0 avg size: 20.53	

https://t... https://t... https://t..

Modelling



ETL / ELT

ELTTLTLTL ...

- Challenges
 - Database provisioning is easy
 - Ingestion is still expecting you transformation*
 - Storage is cheap
 - Not the database licensing model
 - Not the compute power you will need
 - Pipeline as code helps next iteration
 - Data is oil vs data should be curated



ETL/ELT ... ETLTLTL



https://datacreatio n.substack.com/p/ why-data-contract s-are-obviously





Similarity

Node level similarity :

Regular

Edge level similarity

Automorphic





Co-authorship

MATCH p=(a:Author)←[:HAS_AUTHOR]-(w:Work)-[:HAS_AUTHOR]→(a2:Author)
WHERE a ⇔a2
RETURN p



Co-authorship

 $\begin{array}{l} \mbox{MATCH } p=(a:Author) \leftarrow [:HAS_AUTHOR] - (w:Work) - [:HAS_AUTHOR] \rightarrow (a2:Author) \\ \mbox{WHERE } a \diamondsuit a2 \\ \mbox{RETURN } p \end{array}$

How many times have they co-authored before?

MATCH (a1:Author)←[:HAS_AUTHOR]-(work:Work)-[:HAS_AUTHOR]→(a2:Author)
WHERE id(a1) < id(a2)
WITH a1, a2, count(distinct work) as worksCount
CREATE (a1)-[:IS_COAUTHOR_WITH {works_count: worksCount}]→(a2)</pre>



Subgraph

```
CALL gds.graph.project(
    'coauthors',
    'Author',
    {IS_COAUTHOR_WITH: {
        type: 'IS_COAUTHOR_WITH',
        orientation: 'UNDIRECTED'
    }})
```



New Labels

- 1 MATCH (c:Concept)
- 2 CALL apoc.create.addLabels (id(c), ['L'+c.level])
- 3 YIELD node
- 4 RETURN node





Concept Similarity

Explored the concepts of the work co-authors published

The similarity between the concepts gave us the idea how much concepts they share:

Set(AuthorA)∩Set(AuthorB)

Set(AuthorA)USet(AuthorB)

The example pic has low similarity (0!), so possibly fake authors: 0/6



Fake Authors: Igor & Simon



Fake Authors: Ben & Alex



Fake Authors: Charles & Tom



Lessons learned...



- Work on small dataset first to prove data flow
- Work as Data product team
- Communicate data constraints as early as possible
- Change model rapidly when there is no hope
- Scale infrastructure when needed is cost-effective

Data Analysis



- Compare large set of suspect co-authorships with large set of related articles.
- Limit to coauthors at different institutions



suspect

• 1

Related = 0, Suspect = 1

Findings



Machine learning

- Does work well...
- However,
 - a simple rule is almost as good
 - more data and testing would be required to make something usable in the real world.

Conclusions

- 1. Graph analytics are a good approach to identifying organised misconduct
- 2. To do this efficiently, we need graph data structures, such as a knowledge graph built on well-curated data.
 - OpenAlex is a good data set for this
- 3. We can detect unusual co-authorships. Either with machine-learning, or a simple rule:
 - It is unusual for an author to co-author with a different institution on their first paper.

References

- Graph Machine Learning ebook by Ebru Cucen Publication and collaboration anomalies in academic papers originating from a paper mill: evidence from a Russia-based paper mill Anna Abalkina | @AnnaAbalkina | ArXiv: 2202.03310,
- Exploratory analysis of text duplication in peer-review reveals peer-review fraud and paper mills Adam Day | @AdamSci12 | ArXiv: 2202.03310
- For more on Papermills, see Adam's blog: <u>https://publisherad.medium.com/</u> _
- Retraction Watch story: https://retractionwatch.com/2019/07/18/exclusive-russian-site-savs-it-has-brokered-authorships-for-more-than-10000-res earchers/
- Candal-Pedreira et al BMJ 2022; 379 doi: https://doi.org/10.1136/bmj-2022-071517 (Published 28 November 2022)