New analysis features of the CRExplorer for identifying influential publications

<u>Andreas Thor</u>¹, Lutz Bornmann² Werner Marx³, Rüdiger Mutz⁴

¹ University of Applied Sciences for Telecommunications Leipzig, thor@hft-leipzig.de
 ² Administrative Headquarters of the Max Planck Society, bornmann@gv.mpg.de
 ³ Max Planck Institute for Solid State Research, w.marx@fkf.mpg.de
 ⁴ ETH Zürich, mutz@gess.ethz.ch

23rd International Conference on Science and Technology Indicators (STI 2018), 12-14 September 2018, Leiden

CRExplorer: A Tool for Cited References Analysis

- What are the **most important** publications in a research field? On which shoulders does the research stand?
- Identify those publications in a research field or on a specific topic ...
 - ... which have been influential over many years in the past
 - ... were highly cited over a longer time period or at certain time points (shortly or several years after publication)

Cited References Explorer

- Workflow-based Data Analysis (GUI + Script language)
- Import / export data formats: Scopus, Web of Science, CrossRef, CSV
- Automatic Data Extraction and Data Cleaning
- Different types of analysis
 - Visual (Reference Publication Year Spectroscopy)
 - Indicator-based (Top-N, Sequence)
- http://www.crexplorer.net
 - Run (Java Web Start) or download (JAR), Guide + Datasets, Papers (e.g., use cases)

Cited references analysis: Workflow



Pre-Processing

• **Data Extraction** of bibliographic information of Cited References (Strings)

- Regular Expressions (Patterns), Integrity checks

CR	RPY	AL	J_L		J_N	VOL	PAG	DO	I	
Hirsch JE, 2005, P NATL ACAD SCI USA, V102, P16569, DOI 10.1073/pnas.0507655102	2005	Hirsch	ı	P NATL ACA	D SCI USA	102	16569	10.1073/PNAS.05	07655	5102
smulski M., 2006, ISSI NEWSLETTER, V2, P4 2		Kosm	ulski	ISSI NEWSL	ETTER	2 4	4			
Jin BH, 2007, CHINESE SCI BULL, V52, P855, DOI 10.1007/s11434-007-0145-9 2007 Jin CHINE				CHINESE SC	I BULL	52	355	10.1007/S11434-	007-0	145-9
CR			RPY	AU_L	J_	N		TI	VOL	PAG
Bornmann, L., Daniel, HD., What do citation counts measure? A review of studies on citing behavior (200	08) Journal (of Doc	2008	Bornmann	Journal of Documentation What do		citation counts me	64	45	
Bornmann, L., Scientific peer review (2011) Annual Review of Information Science and Technology, 45, pp. 199-245			2011	Bornmann	Annual Review of Inform Scientific		peer review	45	199	
Hirsch, J.E., An index to quantify an individual's scientific research output (2005) Proceedings of the National Academy of S 20			2005	Hirsch	Proceedings of	of the Natio	. An index	to quantify an indi	102	16569

• Data Filtering

- by Citing Publication Year (PY)
- by Reference Publication Year (RPY)
- by Number of Citations (N_CR)

• Data Cleaning (Deduplication)

- Detecting and merging duplicates is important for high-quality data analysis

Deduplication (Disambiguation): Clustering + Merge

- Different variants of the same Cited Reference
 - due to typos, missing bibliographic information, different abbreviation styles, ...
- **Clustering** based on string similarity (author, title) and year
 - Configuration: Threshold (e.g., 80%) + use of DOI, Volume and Page Number

ID	CR	RPY	N_CR
95	Hirsch JE, 2005, P NATL ACAD SCI USA, V102, P16569, DOI 10.1073/pnas.0507655102	2005	155
8664	Hirsch J. E., 2005, P NATL ACAD SCI, V102, P16569	2005	1
8898	Hirsch J. E, 2005, P NATL ACAD SCI USA, V102, P16569	2005	1
6465	Hirsch J. E., 2005, P NATL ACAD SCI USA, V102, P16569	2005	1
8896	Jin B. H, 2007, CHINESE SCI BULL, V52, P855	2007	1
13	Jin BH, 2007, CHINESE SCI BULL, V52, P855, DOI 10.1007/s11434-007-0145-9	2007	37
65	Kosmulski M., 2006, ISSI NEWSLETTER, V2, P4	2006	16
8453	Komulski M., 2006, ISSI NEWSLETTER, V2, P4	2006	1

• **Merging**: Cluster representative + Accumulation of N_CR

ID	CR	RPY	N_CR
95	Hirsch JE, 2005, P NATL ACAD SCI USA, V102, P16569, DOI 10.1073/pnas.0507655102	2005	158
13	Jin BH, 2007, CHINESE SCI BULL, V52, P855, DOI 10.1007/s11434-007-0145-9	2007	38
65	Kosmulski M., 2006, ISSI NEWSLETTER, V2, P4	2006	17

Reference Publication Year Spectroscopy (RPYS)

- Method to analyze historical roots based on cited references within single research fields
- Analysis of the frequency with which references are cited in the publications in terms of the publication years of these CRs.
- Spectrogram
- Origins show up in the form of more or less pronounced peaks mostly caused by individual publications that are cited particularly frequently

Number of Cited References

Reference Publication Year (RPY)

RPYS: Example

• Scientometrics papers (2007-2015) \rightarrow 9,375 Cited References



Number of Publication Years (PYs)

- N_PYEARS = No. of PYs in which the CR has been cited
- N_PY_{RPY} = No. of PYs in which a CR from RPY has been cited
- **PERC_PYEAR** = N_PYEARS / N_PY_{RPY}



Top-N

- **Relative comparison** of Cited References w.r.t. the Reference Publication Year (RPY) and the Publication Year (PY) of citing publications
- N_TOP10 = Number of publication years in which the CR has been in the Top-10% of all CRs of the same RPY



Sequence Types

- **SEQUENCE** and **TYPE**: Distribution of citations over time and identification of common time-series patterns
- Example: Number of citations per publication year ...

Is the number of citations per publication year increasing or decreasing over time?

	2001	2002	2003	2004	2005	2006	Σ
Α	2	6	6	4	2	2	22
В	2	3	3	4	4	5	21
С	0	0	1	4	7	8	20
Σ	4	9	10	12	13	15	63

• ... induce time sequence patterns for cited references



Sequence Computation

• **Observed** values (number of citation per publication year)

	2001	2002	2003	2004	2005	2006	Σ
Α	2	6	6	4	2	2	22
В	2	3	3	4	4	5	21
С	0		1	4	7	8	20
Σ	4	9	10	12	13	15	63

• Expected values

	2001	2002	2003	2004	2005	2006	Σ
Α	1.4	3.1	3.5	4.2	4.5	5.2	22
В	1.3	3.0	3.3	4.0	4.3	5.0	21
С	1.3	2.9	3.2	3.8	4.1	4.8	20
Σ	4	9	10	12	13	15	63

• z-value: Standard Normal Distribution (mean=0, std. dev.=1)

	2001	2002	2003	2004	2005	2006
Α	0.5	1.6	1.3	-0.1	-1.2	-1.4
В	0.6	0.0	-0.2	0.0	-0.2	0.0
С	-1.1	-1.7	-1.2	0.1	1.4	1.5

$$22 \cdot \frac{9}{63} \approx 3.1$$

$$\frac{6-3.1}{\sqrt{3.1}} \approx 1.6$$

2

3

1

-3

-2

-1

0

Sequence Types

• Classification of Cited References based on z-value patterns

	2001	2002	2003	2004	2005	2006
Α	0.5	1.6	1.3	-0.1	-1.2	-1.4
В	0.6	0.0	-0.2	0.0	-0.2	0.0
С	-1.1	-1.7	-1.2	0.1	1.4	1.5



• Scientometrics dataset (1978-2016)

CR	RPY	N_CR	SEQUENCE	ТҮРЕ
Lotka, A.J., The frequency distribution of scientific productivi	1926	155	000000000000000000000000000000000000000	Constant performer
Garfield, E., (1979) Citation Indexing: Its Theory and Applicati	1979	151	00-00000000-000+0++00000-0+00+0+00-0	Constant performer + Life cycle
Small, H., Co-citation in the scientific literature: A new meas	1973	162	-00-0000-00-0000++00+0+++0++	Sleeping beauty
Katz, J.S., Martin, B.R., What is research collaboration? (1997)	1997	171	00-00000++0+++	Sleeping beauty

CRExplorer: User Interface



Mon Aug 20 14:34:15 CEST 2018 Clustering done

#CRs: 33812 (33812 shown), #Clusters: 33812, RPY: 1900-2005 (1900-2005 shown), PY: 1978-2016

Interactive Workflow using GUI

File	Edit Vie	w Disambiguation Help		▼ Matching							
Open	. Ctrl+O	Cluster equivalent Cited References		50 60	70 80	90 100	Volume Page DOI	Same	Different	Extract	Unde
Import	: •	Merge clustered Cited References		ID	CR	RPY	AU	N_CR	▼ N_PYEARS	SEQUENCE	ТҮР
~	C 1 C			27396	Hirsch, J.E., An ind	2005	Hirsch, J.E.	403	12	-0++++0-0	Sleeping b
Save	Ctrl+S			10221	De Celle Déce D	1063	De Celle Déce D	512	25	0.00.0000.0000	· · · · · · ·
Save A	S										
Export	•	Web of Science		Volume		Samo	Diffe	ront	Extract		ada
			100	Page		Same	Dille	rent	EXHACL		100
Setting	JS	Scopus									
			4	3041	menton, r.k., rne	1300	merton, n.n.	110	50	000-00-000000	Life cycle
Exit		CSV (Graph)		4354	Glänzel, W., Natio	2001	Glänzel, W.	114	15	00-0-00+0+-+0+0-	Life cycle
			V	1416	Garfield, E., Citati	1972	Garfield, E.	109	31	00-0000-0000	Life cycle
	4500	CSV (Cited References)		17173	Schubert, A., Glän	1989	Schubert, A.	99	25	000++0+000+00	Life cycle
				29836	Wasserman, S., Fa	1994	Wasserman, S.	94	15	00-00+0	Sleeping
		CSV (Citing Publications)	1	1070	Newman, M.E.J., T	2001	Newman, M.E.J.	91	13	00000+00++0	
			7.	50435	White, H.D., McCa	1998	White, H.D.	90	16	0-0000+-+00	Sleeping
	4000	CSV (Cited References + Citing Publications)		4087	Seglen, P.O., Why	1997	Segien, P.O.	89	10	00+++000	Steeping
	4000	1	V	770	Kercler M.M. Pibl	1063	Kerrler M.M.	97	24	00-0-00-000-00	Constan
00		1986 Number of Cited References: 1 498		3075	Moed H.F. De Br	1995	Moed H.F.	87	21	0-000-000000000	Constan
		Deviation from the 5-Year-Median: 197		1331	King, D.A., The sci	2004	King, D.A.	83	12	00-000++00-00	Constan
		Ň Ň		1176	van Raan, A.F.J., F	2005	van Raan, A.F.J.	82	12	0+00-000++00	Constan
00			1	200	Bradford, S.C., So	1934	Bradford, S.C.	81	28	0000000000000	
		Λf		32	Etzkowitz, H., Ley	2000	Etzkowitz, H.	81	15	00-0-+00+0+0-	Sleeping
00		1		991	Melin, G., Persson,	1996	Melin, G.	79	19	0000-0-++00+	Sleeping
				15452	Narin, F., Stevens,	1991	Narin, F.	78	24	0000-000-0000+	Constan
		A A A A A A A A A A A A A A A A A A A	. A t	14489	Small, H., Sweene	1985	Small, H.	75	27	+00-00+0+000	Constan
0	******************	and vertices see accesses that had by have have have by have	A.M. M.	22224	Cole, J.R., Cole, S.,	1973	Cole, J.R.	74	32	0+00+00+000000	Constan
		Ý	V V	37478	Newman, M.E.J., C	2004	Newman, M.E.J.	73	11	00000+0+	Sleeping
00				5690	Gibbons, M., Limo	1994	Gibbons, M.	72	18	00-0000+000+	Sleeping
				5778	Price, D.D., Netwo	1965	Price, D.D.	72	19	0-0-000000-0	
0.0				195	Barabási, A.L., Jeo	2002	Barabási, A.L.	72	12	000++0-0++	Sleeping
1900	1910 192	0 1930 1940 1950 1960 1970 1980 199	0 2000	7094	Narin, F., Hamilto	1997	Narin, F.	71	18	00-000+00+000+	Life cycle
		Reference Publication Year		3255	Schubert, A., Brau	1990	Schubert, A.	71	23	0-++000+00000	Constant
				3073	Martin, B.R., Irvine	1983	Martin, B.R.	70	26	0000000+000++	Life cycle
	Nu	mber of Cited References Deviation from the 5-Year-Median		31950	Borgatti, S.P., Ever	2002	Borgatti, S.P.	70	14	000-000000-00++	Constant

CRExplorer's Script Language

• Script language for workflow automation

- **Reproducibility** of results
- Same analysis procedure for different publication sets
- Processing large files

```
importFile (
   dir: "E:/data/input/",
   type: "WOS",
   sampling: "RANDOM",
   maxCR: 1000
cluster(
   threshold: 0.8,
   volume: true,
   page: true,
   DOI: false
merge ()
removeCR (RPY: [0, 1995])
exportFile (
   file: "E:/data/output.csv",
   type: "CSV CR"
```

Summary + Future Work

- CRExplorer: A Tool for Cited References Analysis
- Data Extraction + Data Cleaning (Deduplication)
- Reference Publication Year Spectroscopy (RPYS)
- Indicators (TOP-N, Sequence)
- Script-based Automation

- New import / export formats
- User-defined indicators
- Help wanted 🙂

Website: http://www.crexplorer.net Source: https://github.com/andreas-thor/cre

