



TUTORIAL: AI research without coding: The art of fighting without fighting: Data science for qualitative researchers

Leon Ciechanowski^a, Dariusz Jemielniak^{a,*}, Peter A. Gloor^b

^a Kozminski University, Management in Networked and Digital Societies (MINDS) Department, Jagiellonska 59, 03-301 Warszawa, Poland

^b Massachusetts Institute of Technology, Center for Collective Intelligence, 245 First Street, E94, Cambridge, MA 02142, USA

ARTICLE INFO

Keywords:

Twitter
Data scraping
Sentiment analysis
Tribe finding
Wikidata

ABSTRACT

In this tutorial, we show how to scrape and collect online data, perform sentiment analysis, social network analysis, tribe finding, and Wikidata cross-checks, all without using a single line of programming code. In a step-by-step example, we use self-collected data to perform several analyses of the glass ceiling. Our tutorial can serve as a standalone introduction to data science for qualitative researchers and business researchers, who have avoided learning to program. It should also be useful for experienced data scientists who want to learn about the tools that will allow them to collect and analyze data more easily and effectively.

1. Introduction

In the 1973 classic, *Enter the Dragon*, Bruce Lee plays a formidable Shaolin kung-fu master, who describes his style as “the art of fighting without fighting,” and defeats an adversary by tricking him into entering a dinghy to go to an isolated island for a duel, only to leave him adrift behind a larger boat. We will use this philosophy in our tutorial. Even though some of us are proficient in coding, and we prepared a standalone software program for the readers of this issue, we will also show how to use free tools to collect and analyze data without needing to know any programming languages.

The social sciences in general and business studies in particular have undergone revolutionary changes because of the rapid growth of ICT platforms (Caputo & Wallezky, 2017) and big data (Lazer & Radford, 2017). Troves of data resulting from a vast expansion of social media (Kaplan, 2015) make it possible to answer questions previously beyond the reach of science and make the performance of digital analysis a necessity for most social studies. Publicly accessible open data allows the granular study of collaborative software development (Chelkowski, Gloor, & Jemielniak, 2016), the discovery of cultural differences in the perceptions of what constitutes knowledge (Jemielniak & Wilamowski, 2017), the identification of top performers through email pattern analysis (Wen, Gloor, Fronzetti Colladon, Tickoo, & Joshi, 2019), and even the prediction of crude oil prices (Elshendy, Colladon, Battistoni, & Gloor, 2018).

Tapping into the treasure chest of digital datasets is important, as

online data and social media not only reflect the social imaginary but also have a significant impact on people’s beliefs and behaviors, whether in politics (Chmielewska-Szlajfer, 2018), personal relationships (Das & Hodkinson, 2019), local activism (Stasik, 2018), smart cities management (Caputo, Wallezky, et al., 2019), organizational performance (Malone, 2018), or medical knowledge (Smith & Graham, 2019). It is difficult to imagine a sphere of life that does not have some important digital component that cannot be studied by the use of online data. Social media, digital work and digital private life have revolutionized the way we function in business and society (Jemielniak & Przegalińska, 2020).

It is not a surprise that we can observe an advancing datafication of social sciences (Millington & Millington, 2015), and of management, business, and marketing (Erevelles, Fukawa, & Swayne, 2016; Vanhala et al., 2020; Wamba et al., 2017). The data, however, never speak for themselves (Dourish & Cruz, 2018), and applying Big Data methods poses many challenges (Sivarajah, Kamal, Irani, & Weerakkody, 2017). The question about the required conditions for effective Big Data management and research remains open and timely (Caputo, Mazzoleni, et al., 2019). Interpreting big datasets requires contextualization and reflection, typical for qualitative research, but interpretive studies alone are becoming less insightful in the shadow of big datasets.

In fact, drawing from the vast resources of online data is becoming more useful than ever for qualitative researchers. Qualitative research has a long tradition in business studies, dating back to the early 20th

* Corresponding author at: Kozminski University, Management in Networked and Digital Societies (MINDS) Department, Jagiellonska 59, 03-301 Warszawa, Poland.

E-mail addresses: darekj@kozminski.edu.pl, dariuszj@mit.edu (D. Jemielniak).

<https://doi.org/10.1016/j.jbusres.2020.06.012>

Received 17 January 2020; Received in revised form 3 June 2020; Accepted 5 June 2020

0148-2963/ © 2020 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

century and the famous Elton Mayo experiments (Mahoney & Baker, 2002). Qualitative studies have the unique advantage of getting insight into the internal logic of a group, organization or culture (Ciesielska & Jemielniak, 2018a, 2018b). In short, while quantitative studies can plausibly show what is happening and seek correlations and measurable effects, qualitative studies are better at dealing with social backgrounds and explaining their reasons and rationales (Ravitch & Carl, 2019).

Even though qualitative studies in business research were treated with suspicion in the past, they have found an equal place in the methodological repertoire of management scholars (Eriksson & Kovalainen, 2015; Ghauri, Grønhaug, & Strange, 2020; Myers, 2019). However, the weaknesses of qualitative studies are also widely recognized; these include questions of reliability, validity, generalizability, and objectivity (Sinkovics, Penz, & Ghauri, 2008; Symon & Cassell, 2012). They also rely on the researcher's own interpretation and processing of fieldwork material (Queirós, Faria, & Almeida, 2017). Grounding and framing the scope of the study then is inherently subjective. For this reason, taking advantage of quantitative approaches to inform qualitative studies has been advocated as particularly useful (Creswell & Plano Clark, 2017; Ivankova & Wingo, 2018)

2. Aim and objectives: Big data meets thick data

Although the advantages of using mixed methods for business research have already been recognized (Harrison, 2013), and qualitative-quantitative crossovers are known to bring many benefits (Polsa, 2013), allowing qualitative researchers to perform pilot studies and reconnaissance by the use of data science allows the quick mapping of areas that merit deeper, long-term analysis, to frame the qualitative part more accurately. This is why combining big data with thick data, or a *Thick Big Data* approach, make so much sense (Charles & Gherman, 2019; Jemielniak, 2020), and why even diehard non-coding researchers could benefit from the use of data science and machine learning.

Luckily, these fields are rapidly changing (Haenlein & Kaplan, 2019). Many tasks that once required advanced programming skills now do not. However, knowledge about the novel and evolving tools is not readily available.

In this paper we introduce several powerful tools for data analysis to researchers, especially qualitative researchers, who have no experience in collecting large digital datasets. In this tutorial, we show how someone with no knowledge of coding can conduct research with the use of data science and AI, and then, by acquiring primary data, processing and analyzing it. In the spirit of open science, aimed at eliminating financial and knowledge-access barriers in academia (Vicente-Saez & Martinez-Fuentes, 2018), we are going to demonstrate five powerful and free tools. Together, they give business research scholars a strong arsenal. We believe that just by relying on these tools, academics who so far have not been doing digital quantitative studies, should be able to embark on their data science adventure. Creating such tutorials and beginner's guides is crucial for advancing business research (Haenlein & Kaplan, 2004; Watson, 2014). It is also an important step in the direction of closing the gap between qualitative and quantitative researchers.

3. Theoretical framework

In the tutorial—as an example of our methodology—we will collect data on the glass ceiling, a topic of vital importance for research in social media and business. It is still not fully understood in the context of rapid digital transformations, and has a multiplier effect on many other social processes online (Humprecht & Esser, 2017; Jemielniak, 2016; Yarchi & Samuel-Azran, 2018).

In the following section, we present an accessible way to scrape (download) tweets containing the phrase “glass ceiling.” Then, we

perform sentiment analysis (assessment of emotional content) of the tweets, analyze their “tribes” and study them by collecting information about the most influential authors with Wikidata tools.

Sentiment analysis or opinion mining is a method of automatic assessment of the emotional charge of text (e.g., tweets) based on algorithms using Artificial Intelligence (Liu & Zhang, 2012). Its task is to classify emotionally charged texts, both those that indicate the author's emotional state and those that may indicate the emotional effect of the text on the recipient. It can be used to analyze people's opinions, feelings or attitudes about products, services, organizations, events or individuals; this makes it potentially useful in the social and economic sciences (Agarwal, Xie, Vovsha, Rambow, & Passonneau, 2011; Liu, 2012; Medhat, Hassan, & Korashy, 2014) in addition to business (Aldahawi & Allen, 2013; Saura, Palos-Sanchez, & Grilo, 2019; Wang & Zhou, 2009). Sentiment analysis tools are built using neural networks, which is a part of AI, especially through various convolutional and recurrent neural networks (Ain et al., 2017).

In sum, our theoretical framework relies on a methodology of exploring “glass ceiling” tweets and encompasses the following steps:

1. Retrieval of tweets with Twitter Archiver
2. Sentiment extraction with MeaningCloud
3. Twitter and sentiment analysis using our original and free software FWF: Fighting Without Fighting:
 - a. Sentiment analysis
 - b. Tags Cloud generation
4. Wikidata scraping with Wikipedia and Wikidata Tools
5. Tribe analysis with TribeFinder

4. Methodology: A complex quantitative exploration of a research problem without complex tools

1. Tweet retrieval with Twitter Archiver

We collected tweets containing the phrase “glass ceiling” using *Twitter Archiver*, a Google Sheets add-on, which even in its free version is fully functional. The setup is user-friendly: after registering and authorizing our Twitter account (we recommend using a prop one for data harvesting) it is enough to type in the requested parameters in the form (Fig. 1).

After approving the choice of topics (“update search rule”) the script starts collecting tweets directly into Google Sheets and keeps adding newly found tweets every hour until turned off.

From 26 October through 16 December 2019, we collected 6651 tweets containing the exact phrase “glass ceiling.” The database that was generated contains well-structured information on the tweets with that phrase, including date, screen name, full name, tweet text, tweet ID (linking to the original tweet), link(s) from the tweet, media, location, the number of retweets, the number of likes, app used, the number of followers, or the number of accounts followed. There are also several other useful cells, such as informing whether the account is verified, the user location and time zone, the date of registering the account, or a link to the user's website and a profile image.

The Google Sheet we were working on is available at bit.ly/GlassCeilingTweets. We added the collected and processed data in separate sheets.

Of the tweets, 2757 received a retweet or a like; 1431 were retweeted, 1252 were both retweeted and liked at least once; One hundred twenty six were retweeted at least 10 times, and 389 were liked at least 10 times. The most-often liked tweet (1042 likes) was the following, supporting the breaking of the glass ceiling, authored by @jLLOL (Fig. 2).

The second most-often liked tweet (731 likes) ridiculed Hillary Clinton, authored by @DFBHarvard (Fig. 3).

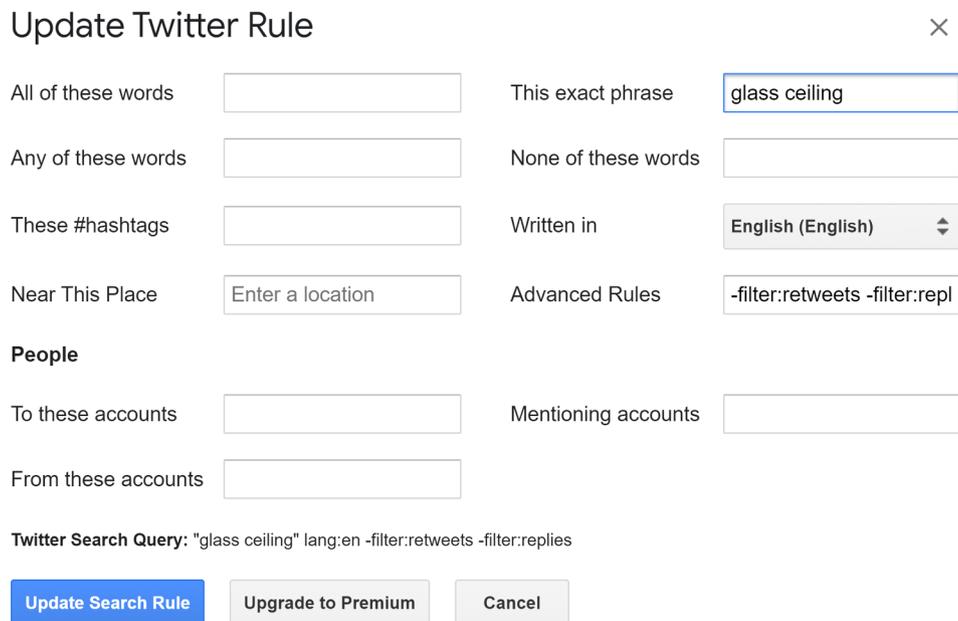


Fig. 1. Twitter Archiver main window.



Fig. 2. The most-often liked tweet about the glass ceiling.

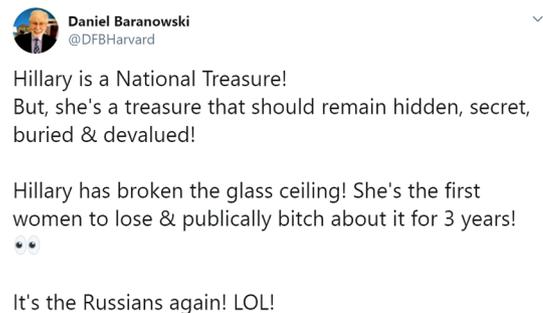


Fig. 3. The second most-often liked tweet.

Among the 10 most-liked tweets, two criticized Hillary Clinton. The remaining eight supported gender equality and breaking the glass ceiling.

2. Sentiment extraction with MeaningCloud

For the purpose of sentiment extraction from the collected tweets, we used MeaningCloud (meaningcloud.com), which also offers a fully functional free Google Sheets add-on (Fig. 4). Processing more than 6000 tweets was not possible in one go: the add-on got stuck a couple of times and we had to restart it on the remaining cell range and combine the results manually. Nevertheless, the process was easy, and batches of 400–500 tweets were analyzed.

MeaningCloud attempts to determine whether the analyzed content was very positive, positive, neutral, negative, very negative, or undetermined. It also estimates whether the content expresses agreement or disagreement, whether it is objective or subjective, ironic or non-ironic, and gives algorithmic confidence in the estimation (Fig. 5).

From our database, 2038 tweets had content classified as negative and 222 as very negative. In addition, 2176 were positive and 382 very positive. Surprisingly, positive comments noticeably outnumbered negative ones, and very positive ones significantly outnumbered very negative ones. All results are added as a separate sheet in our demo link.

3. Twitter and Sentiment analysis using our software

After using Twitter Archiver and MeaningCloud, we ended up with a database that can be easily analyzed with the usage of any statistical

program like IBM SPSS, Statistica, or free and open-source ones like Jamovi or Jasp. One can also use a Python-based program FWF: Fighting Without Fighting¹ that we have developed for this paper, for readers who are not well versed in programming or inferential statistics. The program functionalities allow for automatic statistical analyses and visualization of the data generated with Twitter Archiver and MeaningCloud. The program is available for free unrestricted use, as a

¹ <http://bit.ly/JBRsoftware>.

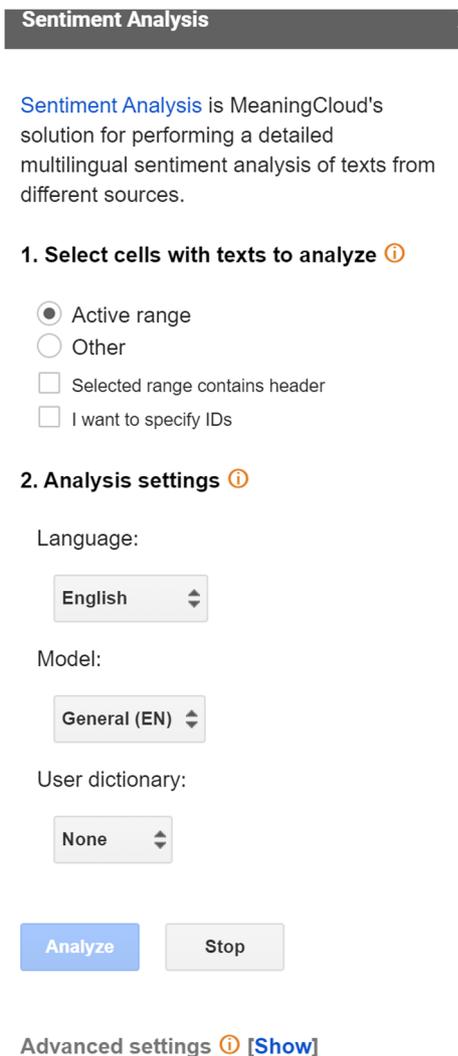


Fig. 4. MeaningCloud, a free Google Sheets add-on for automatic sentiment extraction.

license requirement imposing only the obligation to cite this source article. Below we present an exemplary use of FWF.

We intended our software to be intuitive and self-explanatory (Fig. 6).

After loading the user's file, the dataset can undergo the following analyses.

4.1. Sentiment analysis

Users can pick the desired columns containing data produced by MeaningCloud, preprocessing the data and generating visualizations (Fig. 7) and statistical analyses (Fig. 8). It is worth mentioning that FWF software can analyze and visualize any data relationship chosen by the user, provided that he or she picks the columns containing categorical and continuous values.

We can see that the most positive tweets were most often retweeted. The most negative ones were most rarely retweeted. The program

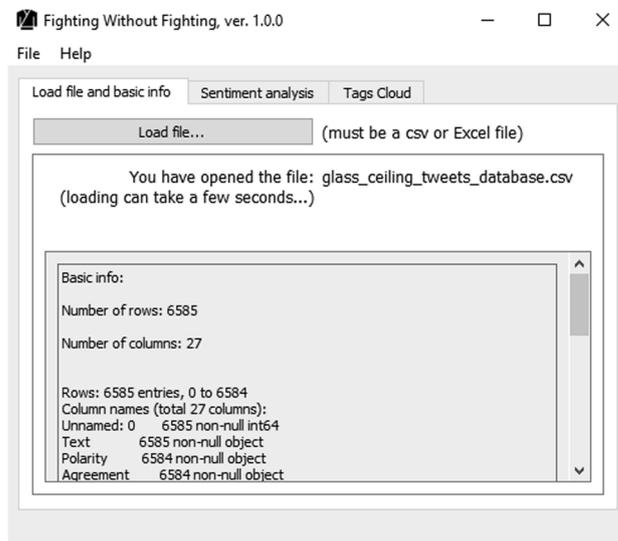


Fig. 6. Main screen from FWF: Fighting Without Fighting after loading an exemplary dataset. Source: <http://bit.ly/JBRsoftware>.

performs non-parametric analysis, since the distributions of each group of tweets with various sentiments are not normal, and the number of tweets in each group is different. A similar situation will likely take place in most Twitter projects; it is unlikely that researchers will find the same number of tweets in five different sentiments. The Kruskal–Wallis test for independent groups was statistically significant ($H = 26.006, p < 0.001$). The post-hoc analysis, with the use of Mann-Whitney test with Bonferroni adjustment for multiple comparisons, revealed that positive (P) and negative (N) tweets are significantly different in the number of retweets ($p < 0.001$), and there was also a statistically significant difference between neutral (NEU), and negative (N) tweets ($p = 0.001$).

4.2. Tags cloud

Users may pick here any column from the data that contains text, and this tab will produce Word Clouds or Tag Clouds (Fig. 9). Word Clouds present the words that appear most frequently in a set of analyzed sentences, tweets and data. The more often the word or phrase was used, the bigger it is displayed in the picture.

We can observe that the most visible difference (except for obvious differences in positive and negative word choice) is the repetition of terms like “Hillary” or “Hillary Clinton.” It is, therefore, possible that commenters' negative sentiments were directed more at the person than at the topic of glass ceiling.

4. WikiData check

We also used the [Wikipedia and Wikidata Tools](#) Google Sheets add-on to check which of the people whose tweets led to stronger reactions have biographical articles on Wikipedia. It is worth remembering that we were able to only the declared combinations of full first and last names given on Twitter, and not verify the actual identity of a person.

We analyzed a subset of 712 tweets that were retweeted at least twice. We copied a list of full names of authors of these tweets,

Don't be afraid of the glass ceiling. It's just an illusion. We'll guide you through. Click here for more: #bt	P	AGREEMENT OBJECTIVE	92	NONIRONIC
Shout to @trenchmag for including @ChildrenOfZeus' 'Ghost' (from the 'Excess Baggage' EP) in their r	P+	AGREEMENT SUBJECTIVE	100	NONIRONIC
How Over-40 Female Writers Are Trying to Shatter Hollywood's Last Glass Ceiling via @thr	NONE	AGREEMENT OBJECTIVE	100	NONIRONIC
"I can't say how many times I've been up for jobs for books that I should adapt. Just very clearly, 'I'm thr	N+	AGREEMENT OBJECTIVE	92	NONIRONIC
I fall asleep from the sound of rain and the glass ceiling at my job doesn't help	N	AGREEMENT OBJECTIVE	92	NONIRONIC

Fig. 5. Part of a Google Spreadsheet with a sentiment of each tweet.

- guide to social research. SAGE.
- Ghauri, P., Grønhaug, K., & Strange, R. (2020). *Research methods in business studies*. Cambridge University Press.
- Gloor, P. A., Fronzetti Colladon, A., de Oliveira, J. M., & Rovelli, P. (2019). Put your money where your mouth is: Using deep learning to identify consumer tribes from word usage. *International Journal of Information Management*. <https://doi.org/10.1016/j.ijinfomgt.2019.03.011>.
- Haenlein, M., & Kaplan, A. M. (2004). A beginner's guide to partial least squares analysis. *Understanding Statistics*, 3(4), 283–297.
- Haenlein, M., & Kaplan, A. (2019). A brief history of artificial intelligence: On the past, present, and future of artificial intelligence. *California Management Review*, 61(4), 5–14.
- Harrison, R. L. (2013). Using mixed methods designs in the Journal of Business Research, 1990–2010. *Journal of Business Research*, 66(11), 2153–2162.
- Humprecht, E., & Esser, F. (2017). A glass ceiling in the online age? Explaining the underrepresentation of women in online political news. *European Journal of Disorders of Communication: The Journal of the College of Speech and Language Therapists, London*, 32(5), 439–456.
- Ivankova, N., & Wingo, N. (2018). Applying mixed methods in action research: Methodological potentials and advantages. *The American Behavioral Scientist*, 62(7), 978–997.
- Jemielniak, D. (2016). breaking the glass ceiling on Wikipedia. *Feminist Review*, 113(1), 103–108.
- Jemielniak, D. (2020). *Thick big data: Doing digital social sciences*. Oxford University Press.
- Jemielniak, D., & Przegalińska, A. (2020). *Collaborative society*. MIT Press.
- Jemielniak, D., & Wilamowski, M. (2017). Cultural diversity of quality of information on Wikipedias. *Journal of the Association for Information Science and Technology*. <http://onlinelibrary.wiley.com/doi/10.1002/asi.23901/full>.
- Kaplan, A. M. (2015). Social media, the digital revolution, and the business of media. *International Journal on Media Management*, 17(4), 197–199.
- Lazer, D. M. J., & Radford, J. (2017). Data ex Machina: Introduction to big data. *Annual Review of Sociology*, 43(1), 19–39. <https://doi.org/10.1146/annurev-soc-060116-053457>.
- Liu, B. (2012). *Sentiment analysis and opinion mining*. Morgan & Claypool Publishers.
- Liu, B., & Zhang, L. (2012). A survey of opinion mining and sentiment analysis. In *Mining text data* (pp. 415–463). https://doi.org/10.1007/978-1-4614-3223-4_13.
- Mahoney, K. T., & Baker, D. B. (2002). Elton Mayo and Carl Rogers: A tale of two techniques. *Journal of Vocational Behavior*, 60(3), 437–450.
- Malone, T. W. (2018). *Superminds: The surprising power of people and computers thinking together*. Brown: Little.
- Medhat, W., Hassan, A., & Korashy, H. (2014). Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal*, 5(4), 1093–1113.
- Millington, B., & Millington, R. (2015). “The datafication of everything”: Toward a sociology of sport and Big Data. *Sociology of Sport Journal*, 32(2), 140–160.
- Myers, M. D. (2019). *Qualitative research in business and management*. Sage Publications Limited.
- Polsa, P. (2013). The crossover-dialog approach: The importance of multiple methods for international business. *Journal of Business Research*, 66(3), 288–297.
- Przegalińska, A., Ciechanowski, L., Stroz, A., Gloor, P., & Mazurek, G. (2019). In bot we trust: A new methodology of chatbot performance measures. *Business Horizons*, 62(6), 785–797.
- Queirós, A., Faria, D., & Almeida, F. (2017). Strengths and limitations of qualitative and quantitative research methods. *European Journal of Education Studies*. <https://oapub.org/edu/index.php/ejes/article/view/1017>.
- Ravitch, S. M., & Carl, N. M. (2019). *Qualitative research: Bridging the conceptual, theoretical, and methodological*. SAGE Publications.
- Saura, J. R., Palos-Sanchez, P., & Grilo, A. (2019). Detecting indicators for startup business success: Sentiment analysis using text data mining. *Sustainability: Science Practice and Policy*, 11(3), 917.
- Sinkovics, R. R., Penz, E., & Ghauri, P. N. (2008). Enhancing the trustworthiness of qualitative research in international business. *Management International Review*, 48(6), 689–714.
- Sivarajah, U., Kamal, M. M., Irani, Z., & Weerakkody, V. (2017). Critical analysis of Big Data challenges and analytical methods. *Journal of Business Research*, 70, 263–286.
- Smith, N., & Graham, T. (2019). Mapping the anti-vaccination movement on Facebook. *Information, Communication and Society*, 22(9), 1310–1327.
- Stasik, A. (2018). Global controversies in local settings: Anti-fracking activism in the era of Web 2.0. *Journal of Risk Research*, 21(12), 1562–1578.
- Symon, G., & Cassell, C. (2012). *Qualitative organizational research: Core methods and current challenges*. SAGE.
- Vanhala, M., Lu, C., Peltonen, J., Sundqvist, S., Nummenmaa, J., & Järvelin, K. (2020). The usage of large data sets in online consumer behaviour: A bibliometric and computational text-mining-driven analysis of previous research. *Journal of Business Research*, 106, 46–59.
- Vicente-Saez, R., & Martinez-Fuentes, C. (2018). Open Science now: A systematic literature review for an integrated definition. *Journal of Business Research*, 88, 428–436.
- Wamba, S. F., Gunasekaran, A., Akter, S., Ren, S. J.-F., Dubey, R., & Childe, S. J. (2017). Big data analytics and firm performance: Effects of dynamic capabilities. *Journal of Business Research*, 70, 356–365.
- Wang, W., & Zhou, Y. (2009). E-business websites evaluation based on opinion mining. *International Conference on Electronic Commerce and Business Intelligence*, 2009, 87–90.
- Watson, H. J. (2014). Tutorial: Big data analytics: Concepts, technologies, and applications. *Communications of the Association for Information Systems*, 34(1), 65.
- Wen, Q., Gloor, P. A., Fronzetti Colladon, A., Tickoo, P., & Joshi, T. (2019). Finding top performers through email patterns analysis. *Journal of Information Science and Engineering* 0165551519849519.
- Yarchi, M., & Samuel-Azran, T. (2018). Women politicians are more engaging: Male versus female politicians' ability to generate users' engagement on social media during an election campaign. *Information, Communication and Society*, 21(7), 978–995.

Leon Ciechanowski is a research assistant at Management in Networked and Digital Environments (MINDS) department, Kozminski University.

Dariusz Jemielniak is Full Professor and head of Management in Networked and Digital Environments (MINDS) department, Kozminski University, visiting scholar at he Center for Collective Intelligence at MIT's Sloan School of Management, and associate faculty at Berkman-Klein Center for Internet and Society, Harvard University. He is a corresponding member of the Polish Academy of Sciences. His recent books include Collaborative Society (2020, MIT Press, with A. Przegalińska), Thick Big Data (2020, Oxford University Press), Common Knowledge? (2014, Stanford University Press).

Peter A. Gloor is a Research Scientist at the Center for Collective Intelligence at MIT's Sloan School of Management where he leads a project exploring Collaborative Innovation Networks (COIN). He is also Founder and Chief Creative Officer of software company galaxyadvisors where he puts his academic insights to practical use, helping clients to coolhunt by analyzing social networking patterns on the Internet – spot the next big thing by finding the trendsetters, and to coolfarm – increase organizational happiness, creativity and performance through workforce analytics. In addition Peter is a Honorary Professor at University of Cologne and a Honorary Professor at Jilin University, Changchun China.