

Human-in-the-Loop Processes as Enabler for Data Analytics in Digitalized Organizations

Thomas Thiele, Thorsten Sommer, Stefan Schröder, Anja Richert, Sabina Jeschke

IMA/ZLW & IfU – RWTH Aachen University, Germany

Abstract

As a key driver for innovation in science, economy and society, digitalization affects almost every aspect of our daily working and living environments. The opportunity to track data about processes, persons, and other entities in organizations allows new opportunities for digitalized working scenarios and the creation of new perspectives on matters such as inter- and intra-organizational relationships.

The aim of this paper is to elaborate on these perspectives on the basis of studies that are currently a part of our research activities. Firstly, a framework is outlined that combines topic modeling of textual data and machine learning to derive thematic synergies in the data, for example, between organizations or research projects. Secondly, classical benchmarking approaches are extended by developing a suitable text-mining process for interdisciplinary research. Lastly, a brief concept about evolution as a method for further optimizations and its implications for the human-in-the-loop process is outlined. Altogether, the approaches contribute to a comprehensive human-in-the-loop model – defined as a model that combines intelligent data technologies with human interaction – in the culture of innovation amongst modern, highly digitalized organizations.

1 Introduction

Digitalization is regarded as a main innovation driver: in vocational context examples range from new forms of human-machine interaction – such as in the digital planning of factory buildups (Büscher et al. 2016) and the interaction of product and machine tools towards cyber-physical systems – to new digitalized services and innovative economic models on a large scale (Carruthers 2016). All of these changes affect daily working environments and change the manner in which humans are embedded in working processes. Although the final outcome

Veröffentlicht durch die Gesellschaft für Informatik e.V. 2016 in
B. Weyers, A. Dittmar (Hrsg.):
Mensch und Computer 2016 – Workshopbeiträge, 4. - 7. September 2016, Aachen.
Copyright © 2016 bei den Autoren.
<http://dx.doi.org/10.18420/muc2016-ws11-0004>

of these changes is still to be determined, early forecasts show an impact on innovation capability not only in professional environments but also in the context of societal trends, such as demographic change (Jeschke et al. 2013).

This paper especially focuses on the enabling of data-driven innovation in professional environments. As these innovations especially arise if knowledge domains (or even boundaries) are crossed (Blackwell et al. 2009), the definition of relevant to-be-connected domains become a necessary prerequisite. Inter- and intra-organizational relationships are derived on the basis of data, which are created in organizations. As this work is especially connected to scientific communities and their interdependent knowledge domains, publications in all forms are considered as a data basis. By various techniques ranging from text mining, machine learning to evolutionary algorithms, data analytics are used to process the publication data towards a suitable form for human users, which are included in a human-in-the-loop process (see Figure 1).

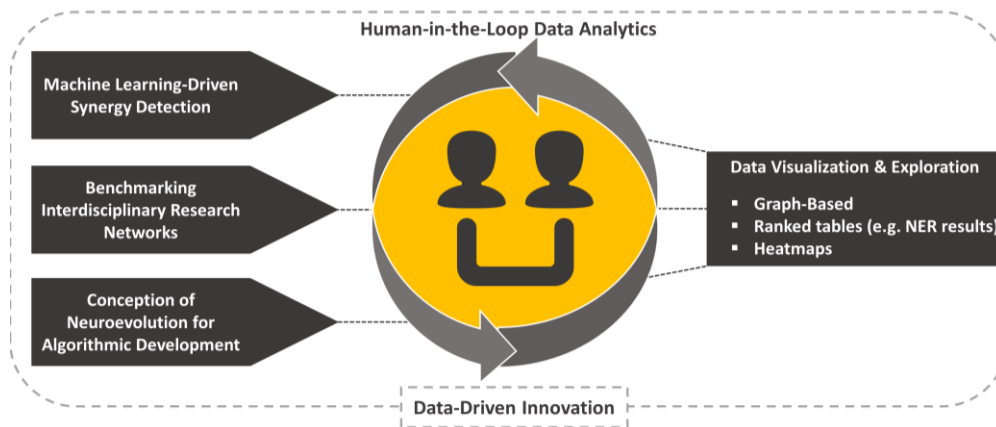


Figure 1: Data-driven innovation supported by human-in-the-loop based on three exemplary approaches

The basic concept of human-in-the-loop is connected to the idea of agent-user interactions in cognitive systems (Dautenhahn 1998). Modern interpretations combine the computational power of pattern recognition in data with the human ability to interpret and utilize the insights from these patterns (Holzinger 2016). The integration of the human user into the data analytics takes place, when relevant knowledge domains have been discovered in the data. This implies that the type of human involvement is a controlling element, which can be described as a so-called outer result control type (Piringer et al. 2014). This especially refers to the assessment of final results to initiate a discourse between the human and the machine, here based on a visualization.

By data visualization the human user is able to explore and interpret the dataset. A deeper integration of the human into, for example, the data mining process is possible by targeting the inner result control, where humans assess intermediate results. Possible kinds of integration are client-driven and algorithm-driven (Piringer et al. 2014). A client-driven integration is an approach, where e.g. the interactive visualization (i.e. client) controls also the computation or mining process. Likewise, in case of an algorithm-driven integration, the client serves only as

observer and the algorithm takes the computation control. Within this paper, three approaches are presented that especially address human-in-the-loop as an algorithm-driven integration of an outer-result control.

All presented approaches tackle the augmentation of the human perspective in innovation processes, by, for example, revealing future cooperation partners or organizational developments due to new topics. On the one hand, they might be caused by internal developments in organizations. On the other hand, they may occur due to external constraints such as trends in scientific communities. As the former is addressed in chapter 2 via machine learning-driven connection of modeled topics, chapter 0 refers to the benchmarking of scientific communities. Chapter 4 conceptualizes an approach for the evolution of algorithms.

2 Machine Learning-Driven Topic Connection

Machine learning is weighted as one of the fastest growing areas of computer science over the past years (Jordan and Mitchell 2015). Google's autonomous cars or Amazon shopping basket analyses are just some prominent examples that heavily rely on machine learning. In the context of our work this raises the question of whether machine learning could be applied to connect different knowledge domains in organizations. Figure 2 displays the process used in this chapter. Process iterations (at the end of the process) are conducted in order to include new data, retrain the neural network and especially involve the Human-in-the-Loop for the assessment of results (data exploration) and the enhancement of machine learning training.

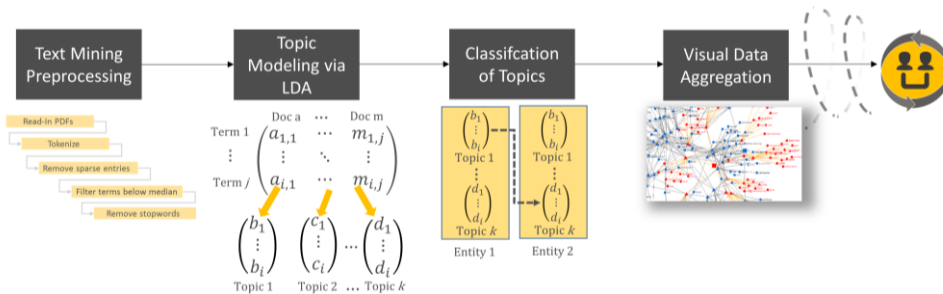


Figure 2: Process of our approach towards knowledge domain connection

2.1 Modeling the Digital Organization

A necessary prerequisite for the application of machine learning can be seen in the data modeling of the to-be-connected domains. The data, which are used to describe in particular scientific organizations, are publications and their meta data, e.g. dates, authors and organizational affiliations. In order to create a suitable format for these data, text mining offers a variety of methods to transform unstructured to structured data which are processed by further analysis (see e.g. Aggarwal and Zhai 2012). The composition of the text corpora already reflects the aim of the analysis: as we target the connection of knowledge domains, the publications of a

group and the publications related to the group's work are labeled with each of the domains. Based on an exemplary use case within a research network, these knowledge domains are represented by projects (Thiele et al. 2015). Hence, each corpus consists of project-specific publications that define the knowledge domains.

This leads to the question, how should these knowledge domains be revealed? Topic modeling offers the possibility to discover hidden layers in textual data, and Latent Dirichlet Allocation is a representative algorithm for this field of action (Blei 2012). Based on probabilistic models the algorithm determines the affiliation of each word towards a semantic space within the text data. As a result, the text data are processed from a bag of words towards a topic model that groups together in a generative process depending on their common usage in semantic topics within the data.

Concerning our primal goal – to connect different knowledge domains – topic modeling allows a more detailed connection since each knowledge domain is decomposed towards several sub-domains. This allows a more precise classification between the domains since each of the topics redescribes individual thematic areas of the domain. The synthesis of matchings between the created topics is realized via another area of machine learning – namely classification – along with neural networks at an algorithmic level.

2.2 Matching Knowledge Domains

The fundamental idea of artificial neural networks reassembles concepts of the human brain (Russell and Norvig 2010). Thus, a neural network is composed out of connected neurons. These are arranged in different layers depending on the complexity of the task for the neural network. The minimum requirement for this topology are an input and an output layer. Deep neural networks are neural networks with more than one (so-called hidden) layer of neurons between input and output. The most interesting part regarding the modeling of neural network can be seen in the mathematical modeling of the connection and neurons and the topology as function of the given task. In order to adjust these parameters, a learning process is necessary. The classic learning process for deep neural networks is the backpropagation approach (Schmidhuber 2015). An advantage of this approach is the ability to get results after a short time of training. However, the approach also faces issues, such as the tendency towards overfitting, deadlocks at a local minima and the need for learning examples. This means that in order to execute backpropagation, an adequate number of well-known examples is mandatory.

Regarding our aim to derive relationships between knowledge domains, these examples can be seen in topics that are possibly eligible for connection. The deep neural network fulfills the task of processing the huge amount of topics (e.g. a topic consists of 3000 features and the estimated topic probabilities) and give a prediction towards thematic overlap to other topics.

The last step within the derivation of connections between the knowledge domains includes the human as an explorer of the data. Hence, the data has to be aggregated in a form that allows the exploration of a complex system. One approach to this challenge is represented by graph-based visualizations (Kolaczyk and Csárdi 2014). On the one hand this form allows the depiction of the hierarchical dependency of a knowledge domain (e.g. a project) and a topic. On the

other hand the connection between the different knowledge domains can be illustrated as networked data – a type of data that is especially easy to present via graphs.

3 Text Mining Based Benchmarking

Over the last decade benchmarking approaches have been developed and implemented primarily for different kinds of enterprises in order to keep in pace with competitors. Predominant in this case are databases such as surveys, documents and patents etc., which are mainly investigated by experts. Due to the increasing amount of data, the development of new ways of handling these data need to be developed. Especially in the framework of research networks, where publications are one of the major outcomes, a text mining based benchmarking approach is desired. For this reason, the benchmarking approach based on text mining, where the human in the loop plays an important role, is applied within a specific research network.

3.1 Benchmarking Scientific Research with Text Mining

The objective of benchmarking in the context of research is to help the management to improve performance and productivity within a research network. Over the last years, many organizations have adopted a range of approaches to research benchmarking (Levy and Valcik 2012). While they deal primarily with quantitative measureable facts (e.g. indicators such as grades of graduates, third party funding, quality assurance etc.), none of them integrates text mining based content analysis. However, since the beginning of the Excellence Initiative in Germany, the German Research Foundation and the German Science and Humanities Council claim a need for adopting business-driven benchmarking approaches.

Following the approach of benchmarking scientific output has been defined as a fact-based consulting (Banerjee et al. 2013). Modern commercial benchmarking usually refers to the process of identifying the best methods, practices and processes. This is justifiable, in that it improve one's own business (Banerjee et al. 2013; Levy and Ronco 2012). Thus, the identification of research priorities and trends is focused. In order to benchmark scientific output, there is a demand for measuring, comparing and judging references (Levy and Ronco 2012). If the quality of a research network's performance is to be judged, scientific output, rather than sales revenues should be consulted. Thus, scientific publication abstracts, because of their accessibility and condensed information content, were selected for analysis. Thereby all publications of the research network itself (data corpus I) were compared to the content related abstracts of the scientific community (data corpus II).

Because of the large amount of scientific publications, a partly automated, text mining based analysis was applied. Consequently, text mining is expected to provide a base for comparison (Kayser et al. 2014) in order to identify apparently unknown domain knowledge. Most text mining solutions are aimed at discovering patterns (Feldman and Sanger 2006). This offers the possibility of key words and elements extraction as well as the identification of relationships. In order to derive information, which is useful for researchers, a demand analysis was upstreamed in an exemplary use case within a research network. Required information includes

research topics and alignment, method expertise and treated substances and products within and outside the network. To fulfil this informational need, the following process was designed (see Figure 3).

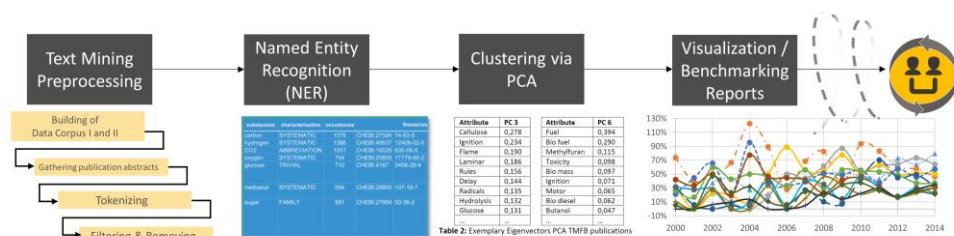


Figure 3: Exemplary process: text mining based benchmarking

3.2 Bridging the Gap: Applying Text Mining

In order to apply text mining, two data corpora (one containing publications of the research network, another gathered from relevant community data) are needed and their corresponding texts must also be converted (Schröder et al. 2016). Therefore, the following process is used to reveal information from a series of documents. Firstly, the estimation of corpora is performed, which distinguishes a set of documents as an answer to a logical query (set of key words augmented by a thesaurus). Afterwards information extraction is used to extract specific information that is analyzed for trends (Miner et al. 2011). Therefore, relevant key words are determined.

The next process deals with Named Entity Recognition (NER). NER is part of the information extraction process and seeks to identify specific keywords (in this case substances). By means of NER, keywords were able to be linked to the specific concept that is being referred to in the document. Thus, concepts can be defined as a biological entity that can be referred to by multiple keywords (Fleuren and Alkema 2015).

After this preprocessing of the texts has been completed, the individual word tokens must be transformed into a vector representation suitable for input into text mining algorithms in a fourth step. Storing a text as weighted vectors requires choosing a weighting scheme first. The most popular scheme is the TF-IDF weighting approach. To apply clustering algorithms or association rules as a next step, TF-IDF transforms the textual data into a suitable metric.

In the end, it is up to human analyst to decide whether or not the results, derived by the above described process, are useful and lead to recommendations of action for the management of a research network. Text Mining based benchmarking gives unbiased information, without judgements and considerations of valuations. This feedback has to be evaluated in a human-in-the-loop manner in order to derive useful insights from this information towards useful feedback.

4 Future Research: Towards Neuroevolution

With a rising amount of information, the human in the loop needs additional support as derived results grow in complexity. Continuing with this idea, more heterogeneous sources can be considered, e.g. social media, industrial websites, open source repositories, etc. The principal challenge is to find patterns of interest in high dimensional vector spaces and assess these by humans. On an algorithmic layer, "neuroevolution" offers the methods to enhance current classifiers (e.g. neural networks) towards a smarter result presentation. Neuroevolution is defined as machine learning technique which uses simulated nature-like evolution to construct artificial neural networks (Lehman and Miikkulainen 2013).

To adjust the composition of artificial neurons and synapses in deep neural networks, a learning process is necessary. The set of all parameters and interconnections decide whether the deep neural networks become smart. Usually, the mentioned backpropagation gets used for the learning process. Backpropagation needs known examples and a static problem space. The presented approaches in chapter 2 and 3 rely on text mining processes to create the known examples and the problem space (e.g. for topic modeling). Hence, no human can provide the required number of examples for the high dimensional vector space. Furthermore, no human even knows on which rules the vector space is based on, if neural networks are used.

Although deep neural networks were already able to show how to outperform the human capabilities (Rutkin 2016), the discovery of meaningful patterns at a high dimensional vector space is obviously a task which overstrained the human's performance. New learning processes for neural networks that are able to evolve represent one possibility towards future applications. In a holistic approach, neuroevolution e.g. through the NEAT method (Stanley and Miikkulainen 2002) provides reinforcement learning and the ability to evolve the entire network.

5 Conclusions

We presented two approaches that allow a derivation of inter- and intra-organizational relationships. Chapter 4 concluded these approaches by conceptualizing neuroevolution as a principle to advance these approaches in future research. All concepts show different approaches where the human-in-the-loop is an important factor for the digitalized organizations. In terms of data analysis, the human fulfills the role of an initiator of further measures as a result of the data exploration. On an algorithmic layer the human-in-the-loop can cover the role of an enhanced trainer for the algorithms associated with the iterations of our concepts. Although algorithms continue to increase in complexity, human-in-the-loop is still able to provide the chance to capture the interpretational capacity of the human brain and include this into the results.

The current approaches only allow human feedback during analysis steps (e.g. of visualizations). In order to address this challenge, human feedback could also be included in the processing steps of the data, for example, when topics are estimated. In addition, the human-in-the-loop could assess the long-term effect. It would be easier to let humans assess impulsively

and directly after they receive a recommendation. Currently, when the individual gets recommendations to cross its knowledge domains only based on visualizations, this is obviously not included. Instead, the human should try the recommendation and assess it in later follow-ups, such as meetings.

Although the evolution of (classifying) algorithms sounds rational in terms of the technical outlook, this leads to another question: Does the system derive reasonable results? A deep neural network will produce an output, regardless of the input. But is a pattern, detected by a neural network and evolved by neuroevolution, semantic useful and is it causal to humans? For a human, the process is not comprehensible due to the scale of data and the amount of execution steps. Thus, our current approaches neglect individual human feedback. If a system's task is to analyze larger data corpora (e.g. a whole community), a crowd-in-the-loop can be helpful to assess interpretability and usefulness of results. Thus, a single human must trust the machine, but only a crowd is able to enhance the process.

Acknowledgements

This research has been funded by the German Research Foundation (DFG) as part of the Clusters of Excellences “Tailor-Made Fuels from Biomass” and “Integrative Production Technology for High-Wage Countries” at RWTH Aachen University, as well as the Federal Ministry of Education and Research as part of the project “ELLI - Excellent Teaching and Learning in Engineering Sciences”.

References

- Aggarwal, C. C. and Zhai, C. (2012). *Mining text data*: Springer Science & Business Media.
- Banerjee, A., Bandyopadhyay, T. & Acharya, P. (2013). Data Analytics Hyped Up Aspirations or True Potential?, *VIKALPA*, Vol. 38, no. 4. pp. 1–11.
- Blackwell, A. F., Wilson, L., Street, A., Boulton, C. & Knell, J. (2009). Radical innovation: crossing knowledge boundaries with interdisciplinary teams, *University of Cambridge*.
- Blei, D. M. (2012). Probabilistic Topic Models, *Commun. ACM*, vol. 55, no. 4. pp. 77–84.
- Büscher, C., Voet, H., Krunke, M., Burggräf, P., Meisen, T. & Jeschke, S. (2016). Semantic Information Modelling for Factory Planning Projects, *Procedia CIRP*, vol. 41. pp. 478–483.
- Carruthers, K. (2016). *Internet of Things and Beyond: Cyber-Physical Systems*, IEEE Internet of Things (Accessed 30 May 2016).
- Dautenhahn, K. (1998). The Art of Designing Socially Intelligent Agents – Science, Fiction, and the Human in the Loop, *Applied Artificial Intelligence*, vol. 12, 7-8. pp. 573–617.
- Feldman, R. and Sanger, J. (2006). *The Text Mining Handbook. Advanced Approaches in analyzing Unstructured Data*. Cambridge: Cambridge University Press.
- Fleuren, W. & Alkema, W. (2015). Application of text mining in the biomedical domain, *Methods (San Diego, Calif.)*, vol. 74. pp. 97–106.
- Holzinger, A. (2016). Interactive machine learning for health informatics – When do we need the human-in-the-loop?, *Brain Informatics*, vol. 3, no. 2. pp. 119–131.

- Jeschke, S., Vossen, R., Leisten, I., Welter, F., Fleischer, S. & Thiele, T. (2013). Industrie 4.0 als Treiber der demografischen Chancen. In Jeschke, S. (ed.): *Innovationsfähigkeit im demografischen Wandel*: Campus Verlag. pp. 9–19.
- Jordan, M. I. & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects, *Science (New York, N.Y.)*, vol. 349, no. 6245. pp. 255–260.
- Kayser, V., Goluchowicz, K. & Bierwisch, A. (2014). Text Mining For Technology Roadmapping — The Strategic Value of Information, *International Journal of Innovation Management*, no. 18. pp. 1–23.
- Kolaczyk, E. D. and Csárdi, G. (2014). *Statistical analysis of network data with R*: Springer.
- Lehman, J. & Miikkulainen, R. (2013). *Neuroevolution* [Online]. Available at <http://www.scholarpedia.org/article/Neuroevolution> (Accessed 25 March 2016).
- Levy, G. D. & Ronco, S. L. (2012). How Benchmarking and Higher Education came together, *New Directions for Institutional Research*, vol. 2012, no. 156. pp. 5–13.
- Levy, G. D. and Valcik, N. A., eds. (2012). *Benchmarking in Institutional Research*.
- Miner, G., Delen, D., Elder, J., Fast, A., Hill, T. and Nisbet, R. (2011). *Practical Text Mining and Statistical Analysis for Non-structured Text Data Applications*: Elsevier.
- Piringer, H., Streit, M., Sedlmair, M., Gratzl, S. & Mühlbacher, T. (2014). Opening the Black Box: Strategies for Increased User Involvement in Existing Algorithm Implementations, *IEEE Transactions on Visualization and Computer Graphics* [Online]. Available at <http://eprints.cs.univie.ac.at/4163/>.
- Russell, S. J. and Norvig, P. (2010). *Artificial intelligence – A modern approach*, 3. Auflage. Upper Saddle River, NJ: Prentice-Hall.
- Rutkin, A. (2016). Anything you can do..., *New Scientist*, vol. 229, no. 3065. pp. 20–21.
- Schmidhuber, J. (2015). Deep learning in neural networks: an overview, *Neural networks : the official journal of the International Neural Network Society*, vol. 61. pp. 85–117.
- Schröder, S., Tummel, C., Isenhardt, I., Jeschke, S. & Richert, A. (2016). Benchmarking of Scientific Research Clusters by Use of Text Mining Algorithms on Textual Artefacts, *Proceedings of the ICISE (IEEE Computer Society)*. pp. 22–28.
- Stanley, K. O. & Miikkulainen, R. (2002). Evolving Neural Networks through Augmenting Topologies, *Evolutionary Computation*, vol. 10, no. 2. pp. 99–127 [Online]. DOI: 10.1162/106365602320169811.
- Thiele, T., Jooß, C., Richert, A. & Jeschke, S. (2015). Terminology Based Visualization of Interfaces in Interdisciplinary Research Networks. In *19th Triennial Congress of the IEA*.

Authors



Thiele, Thomas

Thomas Thiele is a scientific researcher at IMA/ZLW, RWTH Aachen University. Within his doctoral thesis, he is focused on data science based methods for the detection of synergies in scientific publications, in particular text mining, topic modeling and machine learning.



Sommer, Thorsten

Thorsten Sommer is a scientific researcher at the IMA/ZLW of RWTH Aachen University. His PhD thesis is about artificial intelligence by using enhanced neuroevolution. Further research topics are machine learning, data mining and e-learning.



Schröder, Stefan

Stefan Schröder is a research group leader at the IMA/ZLW of RWTH Aachen University. His doctoral thesis is about adopting benchmarking approaches to university research networks e.g. by establishing a text mining based process.



Richert, Anja

Anja Richert is Junior Professor for Agile Management in the Faculty of Mechanical Engineering and Director of the Center for Learning and Knowledge Management at RWTH Aachen University. Her research focuses on agile management and learning and knowledge processes, developing mixed reality learning concepts and developing and testing data-science-based socio-technological research designs.



Jeschke, Sabina

Sabina Jeschke is head of the Cybernetic-Cluster IMA/ZLW & IfU and Vice Dean of the Faculty of Mechanical Engineering at the RWTH Aachen University. Her main research areas are: Complex IT-systems, robotics and automation, traffic and mobility and virtual worlds for research alliances and education. For more information, visit: <http://www.ima-zlw-ifu.rwth-aachen.de/en/Sabina.Jeschke>

Image Copyright: Marcus Gerards RWTH Aachen University