

# Personalized Medicine by Means of Complex Networks – a Big Data Challenge

## Inhaltsverzeichnis

<b>I. Introduction</b> .....	37
<b>II. Network Biology</b> .....	40
<b>III. Personalized Medicine by Using Network- Based Approaches</b> .....	42
<b>IV. Discussion and Outlook</b> .....	43

This conceptional note deals with the problem of performing personalized medicine by employing network-based techniques. We demonstrate this by way of example when detecting complex diseases.

Acknowledgments: We thank Florent Thouvenin for stimulating discussions concerning this work.

Key words: Personalized Medicine, Complex Networks, Biomarker, Big Data

## I. Introduction

The research on personalized medicine has been still intricate to explore medical decision models and products for individual patients<sup>1</sup>. To do so, several decision models by using statistical techniques have been used and optimized<sup>2</sup>. A major

---

<sup>\*</sup> Ao. Professor, Division for Bioinformatics and Translational Research UMIT – The Health and Life Sciences University, Tirol.

<sup>\*\*</sup> Assistenzprofessor, IBM Watson Think Group, Research Unit HCI, Institute for Medical Informatics, Medical University Graz.

<sup>\*\*\*</sup> Ao. Professor, Computational Biology and Machine Learning Laboratory, Center for Cancer Research and Cell Biology, School of Medicine, Dentistry and Biomedical Sciences, Faculty of Medicine, Health and Life Sciences, Queen's University Belfast.

<sup>1</sup> SHASTRY B. S., Pharmacogenetics and the concept of individualized medicine. *Pharmacogenomics J.*, 6:16–21, 2006.

<sup>2</sup> EMMERT-STREIB F., The chronic fatigue syndrome: A comparative pathway analysis. *Journal of Computational Biology*, 14(7), 2007; SIEBERT U., ROCHAU U. / CLAXTON K., When is

application of personalized medicine has been cancer management and its stratification, i.e., to estimate the cancer risk for the clinical management. In this paper, we sketch a conceptual approach to perform personalized medicine by using complex network analysis<sup>3</sup> where the pathways (networks) are inferred from DNA microarray cancer data<sup>4</sup>.

In order to start, we briefly outline the most recent developments in network biology<sup>5</sup>. Afterwards, we state the crucial idea how to represent biomarker as complex networks. Also, we briefly mention mathematical techniques to analyze pseudo pathways quantitatively. The paper ends with a discussion and outlook.

Due to the increasing trend towards personalized, molecular and precision medicine (P4 medicine: Predictive, Preventive, Participatory, Personalized<sup>6</sup>), biomedical data today results from various sources in different structural dimensions, ranging from the microscopic world, and in particular from the omics world (e.g., from genomics, proteomics, metabolomics, lipidomics, transcriptomics, epigenetics, microbiomics, fluxomics, phenomics, etc.) to the macroscopic world (e.g., disease spreading data of populations in public health informatics<sup>7</sup>).

This «Big Data» is not a buzz word, actually we are at a paradigmatic shift of scientific work and medicine is turning into a data intensive science. Whilst, the traditional method of turning data into knowledge relied on manual analysis and interpretation by a medical doctor in order to find useful patterns in data for decision support, now The systematic exploration of this big data can be seen as a

---

enough evidence enough? using systematic decision analysis and value-of-information analysis to determine the need for further evidence. *Zeitschrift für Evidenz, Fortbildung und Qualität im Gesundheitswesen*, 107: 575–584, 2013.

<sup>3</sup> EMMERT-STREIB F./DEHMER M., (eds), *Analysis of Microarray Data: A Network-based Approach*. Wiley VCH Publishing, 2010; EMMERT-STREIB F./DEHMER M., *Networks for systems biology: Conceptual connection of data and function*. *IET Systems Biology*, 5: 185–207, 2011; JUNKER B. H./SCHREIBER F., *Analysis of Biological Networks*. Wiley Series in Bioinformatics. Wiley-Interscience, 2008.

<sup>4</sup> EMMERT-STREIB F./DEHMER M., *Analysis of Microarray Data* (Fn. 3).

<sup>5</sup> BARABÁSI A. L./OLTVAI Z. N., *Network biology: Understanding the cell's functional organization*. *Nature Reviews. Genetics*, 5(2):101–113, 2004; EMMERT-STREIB F./DEHMER M., *Networks for systems biology* (Fn. 3).

<sup>6</sup> HOOD L./FRIEND S. H., *Predictive, personalized, preventive, participatory (P4) cancer medicine*. *Nature Reviews Clinical Oncology*, 8, (3), 2011, 184–187.

<sup>7</sup> HOLZINGER, A., *Biomedical Informatics: Discovering Knowledge in Big Data*, New York 2014, Springer; HOLZINGER A./STOCKER C./DEHMER M., *Big complex biomedical data: Towards a taxonomy of data*, 2014. Springer; *Communications in Computer and Information Science*. Berlin, Heidelberg, New York 2014, Springer, pp. in print; HOLZINGER A./DEHMER M./JURISICA I., 2014, *Knowledge Discovery and Interactive Data Mining in Bioinformatics – State-of-the-Art, Future challenges and Research Directions*. *BMC Bioinformatics*, to appear.

completely new paradigm in the investigation of nature<sup>8</sup>: Machine learning approaches provide a mechanism for data driven hypotheses generation, optimized experiment planning, precision medicine and evidence-based medicine. The challenge is not only to extract meaningful information from this data, but to gain knowledge, to discover previously unknown insights, to look for patterns, and to make sense of the data<sup>9</sup>. Diverse approaches, including statistical and graph theoretical methods, data mining, and computational pattern recognition, have been applied to this task in the past<sup>10</sup> and although machine learning approaches are indispensable current trends are increasingly putting the human-into-the-loop<sup>11</sup>, leaving the grand goals of the 1970ies artificial intelligence to make everything automatic. However, an increasingly important issue is the limited time which medical doctors have in their daily clinical routine: in average a medical doctor in a public hospital has only five minutes to make a decision<sup>12</sup>, hence interactive tools for decision support are a necessity, which calls for powerful grid-based computing approaches. Inferring that the future of personalized medicine will be based on mobile technologies and cloud-based technologies, one of the most important future aspects are in privacy, data protection, data security and data safety

- 
- <sup>8</sup> BELL G./HEY T./SZALAY A., Beyond the data deluge. *Science*, 323, (5919), 2009, 1297–1298.
- <sup>9</sup> HOLZINGER A./STOCKER C./BRUSCHI M./AUINGER A./SILVA H./FRED A., On Applying Approximate Entropy to ECG Signals for Knowledge Discovery on the Example of Big Sensor Data. In: R. HUANG, E. A. E. (ed.), *Active Media Technologies AMT 2012*, LNCS 7669, Macau 2012, Springer, pp. 646–657; HOLZINGER A./SCHERER R./SEEBER M./WAGNER J./MÜLLER-PUTZ G., Computational Sensemaking on Examples of Knowledge Discovery from Neuroscience Data: Towards Enhancing Stroke Rehabilitation. In: BÖHM C./KHURI S./LHOTSKÁ L./RENDA M. (eds.), *Information Technology in Bio- and Medical Informatics*, Lecture Notes in Computer Science, LNCS 7451. Heidelberg, New York 2012, Springer, pp. 166–168.
- <sup>10</sup> RAYMER M. L./DOOM T. E./KUHN L. A./PUNCH W. F., Knowledge discovery in medical and biological datasets using a hybrid Bayes classifier/evolutionary algorithm. *IEEE Transactions on Systems Man and Cybernetics Part B-Cybernetics*, 33, (5), 2003, 802–813; JURISICA I./MYLOPOULOS J./GLASGOW J./SHAPIRO H./CASPER R. F., Case-based reasoning in IVF: prediction and knowledge mining. *Artificial intelligence in medicine*, 12, (1), 1998, 1–24.
- <sup>11</sup> HOLZINGER A., Human–Computer Interaction & Knowledge Discovery (HCI-KDD): What is the benefit of bringing those two fields to work together? In: CUZZOCREA A./DIMITRIS C. K./SIMOS E./WEIPPL E./XU L. (ed.) *Multidisciplinary Research and Practice for Information Systems*, Springer Lecture Notes in Computer Science LNCS 8127. Heidelberg, Berlin, New York 2013, Springer, pp. 319–328.
- <sup>12</sup> GIGERENZER G., *GUT FEELINGS: Short Cuts to Better Decision Making* London, Penguin 2008; GIGERENZER G./GAISSMAIER W., Heuristic Decision Making. In: Fiske S. T./Schacter, D. L./Taylor S. E. (eds.), *Annual Review of Psychology*, Vol. 62., 2011, Palo Alto: Annual Reviews, pp. 451–482.

and the fair use of data<sup>13</sup>. This calls for data accessibility, legal restrictions, confidentiality and data provenance. Consequently, the field of knowledge discovery and data mining is challenged by a number of non-trivial future aspects to properly address these complex restrictions in dealing with big data.

## II. Network Biology

Network-based techniques have been investigated in various disciplines such as chemistry<sup>14</sup>, biology<sup>15</sup>, ecology<sup>16</sup> and finance<sup>17</sup>. In particular, network biology<sup>18</sup> deals with analyzing biological data sets in the context of systems biology<sup>19</sup>. In systems biology, the underlying biological systems are usually understood as dynamical systems which are investigated holistically<sup>20</sup>.

Towards network analysis, two main categories of network analysis may be distinguished: One may use methods from so-called descriptive graph theory<sup>21</sup> or approaches from quantitative graph theory<sup>22</sup>. Methods which belong to the first category describe structural features of a network deterministically. The second category is concerned with developing techniques to quantify structural information by using a measurement approach<sup>23</sup>. In addition, so-called machine learning

- 
- <sup>13</sup> WEIPPL E./HOLZINGER A./TJOA A. M., Security aspects of ubiquitous computing in health care. *Springer Elektrotechnik & Informationstechnik, e&i*, 123, (4), 2006, 156–162.
- <sup>14</sup> D. BONCHEV, *Information Theoretic Indices for Characterization of Chemical Structures*. Research Studies Press, Chichester, 1983.
- <sup>15</sup> DIUDEA M. V./GUTMAN I./JÄNTSCHI L., *Molecular Topology*. Nova Publishing, 2001. New York, NY, USA; JUNKER B. H./SCHREIBER F., *Analysis of Biological Networks*. Wiley Series in Bioinformatics. Wiley-Interscience, 2008.
- <sup>16</sup> ULANOWICZ R. E., Quantitative methods for ecological network analysis. *Computational Biology and Chemistry*, 28: 321–339, 2004.
- <sup>17</sup> BOGINSKI V./BUTENKO S./PARDALOS P., Statistical analysis of financial networks. *Computational Statistics and Data Analysis*, 48(2): 431–443, 2005.
- <sup>18</sup> BARAB'ASI A. L./OLTVAI Z. N., Network biology: Understanding the cell's functional organization. *Nature Reviews. Genetics*, 5(2): 101–113, 2004.
- <sup>19</sup> PALSSON B. O., *Systems Biology: Properties of Reconstructed Networks*. Cambridge University Press, 2006.
- <sup>20</sup> EMMERT-STREIB F./DEHMER M., *Networks for systems biology* (Fn. 3).
- <sup>21</sup> HARARY F., *Graph Theory*. Addison Wesley Publishing Company, 1969, Reading, MA, USA.
- <sup>22</sup> DEHMER M./EMMERT-STREIB F., *Quantitative Graph Theory. Theory and Applications*. CRC Press, 2014. to appear.
- <sup>23</sup> DEHMER M./EMMERT-STREIB F./MEHLER M. (eds.), *Towards an Information Theory of Complex Networks: Statistical Methods and Applications*. Birkhäuser Publishing, 2011; DEHMER M./EMMERT-STREIB F. (Fn. 22); EMMERT-STREIB F./DEHMER M., *Analysis of Microarray Data* (Fn. 3).

methods have been used to analyze network data efficiently<sup>24</sup>. On the one hand, network-based techniques have been successfully applied to tackle various questions in systems biology<sup>25</sup>. On the other hand, network-based methods have been proven amazingly efficient and flexible to examine problems related to the function and structure of biological systems in the postgenomics era<sup>26</sup>.

In the following, we list some still outstanding and intricate problems in network biology and related areas in brief:

- Network Biology for Medicine
  - Inferring complex networks from High-Throughput data sets, e.g. cancer<sup>27</sup>
  - Structural Analysis of biological networks<sup>28</sup>
  - Developing information-theoretic and statistical techniques for complex network analysis<sup>29</sup>
- Structure-oriented Drug Design
  - Representing molecule structures as graphs<sup>30</sup>
  - Quantifying the structure of molecules using so-called topological descriptors<sup>31</sup>
  - Predicting chemical or biological activities of the molecules using graph-theoretical and statistical techniques<sup>32</sup>

---

<sup>24</sup> GÄRTNER T./FLACH P. A./WROBEL S., On graph kernels: Hardness results and efficient alternatives. In COLT, pages 129–143, 2003; MÜLLER L. A. J./KUGLER K. G./GRABER A./DEHMER M., A networkbased approach to classify the three domains of life. *Biology Direct*, 6:140–141, 2011.

<sup>25</sup> EMMERT-STREIB F./DEHMER M., *Networks for systems biology* (Fn. 3).

<sup>26</sup> EMMERT-STREIB F./DEHMER M., *Networks for systems biology* (Fn. 3).

<sup>27</sup> EMMERT-STREIB F./DEHMER M., *Analysis of Microarray Data* (Fn. 3); EMMERT-STREIB F. (Fn. 2); EMMERT-STREIB F./DEHMER M., *Networks for systems biology* (Fn. 3).

<sup>28</sup> EMMERT-STREIB F./DEHMER M., *Analysis of Microarray Data*, (Fn. 3); EMMERT-STREIB F. (Fn. 2).

<sup>29</sup> DEHMER M., Information processing in complex networks: Graph entropy and information functionals. *Appl. Math. Comput.*, 201: 82–94, 2008; DEHMER M./EMMERT-STREIB F./MEHLER M. (Fn. 23).

<sup>30</sup> DIUDEA M. V./GUTMAN I./JÄNTSCHI L. (Fn. 15).

<sup>31</sup> DEHMER M./VARMUZA K./BORGERT S./EMMERT-STREIB F., On entropy-based molecular descriptors: Statistical analysis of real and synthetic chemical structures. *J. Chem. Inf. Model.*, 49: 1655–1663, 2009; DIUDEA M. V./GUTMAN I./JÄNTSCHI L. (Fn. 15); TODESCHINI R., CONSONNI V., MANNHOLD R., *Handbook of Molecular Descriptors*. Wiley-VCH, Weinheim 2002.

<sup>32</sup> BASAK S. C./MAGNUSON V. R., Molecular topology and narcosis. *Arzeim.-Forsch./Drug Design*, 33(I): 501–503, 1983; DEHMER M./VARMUZA K./BONCHEV D. (eds.), *Statistical Modelling of Molecular Descriptors in QSAR/QSPR. Quantitative and Network Biology*. Wiley-Blackwell, 2012; TODESCHINI R./CONSONNI V./MANNHOLD R. (Fn. 31).

### III. Personalized Medicine by Using Network-Based Approaches

In this section, we describe the problem of performing personalized medicine by using cancer data and complex networks by way of example. More precisely, we consider so-called complex diseases which have been defined due to Strohman<sup>33</sup> for inferring cancer networks.

**Definition 3.1** A complex disease can not be defined by a single biomarker (gene) but a group of interacting genes (pathways).

For instance, the Coronary heart disease, Hypertension, Diabetes, Obesity, Various cancers, Alzheimer's disease, and Parkinson's disease are known to be complex diseases. Non-complex diseases, that is Mendelian diseases, are for example Huntington (D), BRCA (breastcancer, D), Cysticfibrosis, and Thrombophilia.

Now we briefly recall the classical view to describe Mendelian diseases on a gene level. If the expression of a gene changes, then the function of the gene changes. Also if the function of a gene changes, the phenotype of the organism changes accordingly. Consequently, we obtain a genotype-phenotype mapping in the Mendelian sense and the phenotype corresponds to the disease. Finally, we obtain that differentially expressed genes are markers of diseases which corresponds to the classical view of biological data analysis.

The new view is based on the hypothesis that biological processes in a cell form dynamical systems that maintain biological functions. That means the biological function is in general a collective effort. This implies that groups of genes participating in a biological process are important and groups of genes interact with each other. Otherwise we were unable to infer any dynamical system.

Ideally, this aims to infer pathways from an individual patient by using high-throughput data sets. Therefore this leads to a procedure to perform personalized medicine. In the following, we give a concrete example to do so when detecting complex diseases<sup>34</sup> and consider Figure 1. We see that two groups of patients exist, here sick versus non-sick. Then pseudo pathways are inferred from the underlying DNA microarray data, e.g., cancer. Finally Figure 1 says that one can detect complex diseases by comparing the inferred pseudo pathways quantitatively. For instance, we could pairwise compute the structural similarity between the pathways and then calculate the distance between the obtained similarity distributions.

---

<sup>33</sup> STROHMAN R., Maneuvering in the complex path from genotype to phenotype. *Science*, 296(5568): 701–703, 2002.

<sup>34</sup> EMMERT-STREIB F./DEHMER M., *Analysis of Microarray Data* (Fn. 3).

Figure 1: Detecting complex diseases by means of complex networks and comparative network analysis<sup>35</sup>.

Key Challenges along with big data include Data in the Cloud, mobile solutions, the trend towards software-as-a-service, and the massive increase in the amount of data, consequently lot of future effort must be spent in Privacy, Data Protection, Security and Safety. The challenges of data integration, data fusion and the increased use of data for secondary use put these issues from a «nice-to-have» into the key interest; just an example: In January 2013, the US Department of Health released the so-called Omnibus Final Rule, which significantly modified the privacy and security standards under the Health Insurance Portability and Accountability Act (HIPAA). These new regulations were driven by a need to ensure the confidentiality, integrity, and security of patients' protected health information (PHI) in electronic health records (EHRs) and addresses these concerns by expanding the scope of regulations and increasing penalties for health information violations<sup>36</sup>.

A further big issue – particularly in the context of P4 medicine – is the secondary use of data, providing patient data for clinical and/or medical research. For most secondary data use, it is possible to use de-identified data, but for the remaining data protection issues are very important<sup>37</sup>. The secondary use of data involves the linkage of data sets to bring different modalities of data together, which raises more concerns over the privacy of the data. A classic example is the publication of the Human Genome, which raised new ways of finding relationships between clinical disease and human genetics. The increasing use and storage of P4 relevant data also impacts the use of familial records, since the information about the patient also provides information on the patient's relatives, which can indeed be very useful. Another big issue is the production of anonymized open data sets to support international joint research efforts.

## IV. Discussion and Outlook

Inferring and analyzing networks by using DNA microarray data has been proven useful for at least a decade. Consequently, genome-wide high-throughput tech-

---

<sup>35</sup> EMMERT-STREIB F./DEHMER M., Analysis of Microarray Data (Fn. 3).

<sup>36</sup> WANG C./HUANG D. J., The HIPAA conundrum in the era of mobile health and communications. JAMA, 310, (11), 2013, 1121–1122.

<sup>37</sup> SAFRAN, C./BLOOMROSEN, M./HAMMOND W. E./LABKOFF S./MARKEL-FOX S./TANG P. C./DETMER D. E., Toward a national framework for the secondary use of health data: an American Medical Informatics Association white paper. Journal of the American Medical Informatics Association, 14, (1), 2007, 1–9.

nologies have been proven useful for performing personalized medicine. In fact, applying the underlying mathematical methods bears enormous potential to perform personalized medicine in the future. As huge data sets for inferring and analyzing pathways exist, the issue of data protection is crucial too. This encompasses the whole range of issues from privacy, data protection, safety and security which always must be taken into account, see e.g. [hci4all.at](#)<sup>38</sup>.

Also, it has been shown that not only the analysis of next-generation sequencing technologies to generate, e.g., DNA sequence data has been useful<sup>39</sup>. In this context, Emmert-Streib and Dehmer<sup>40</sup> recently explored the hypothesis that dynOmics data can be important for enabling personalized medicine by complementing genotype data<sup>41</sup>.

A further grand challenge and big issue in research is the interaction with large graphs for graph-based data mining and knowledge discovery, where Holzinger and Dehmer do some pioneering future work<sup>42</sup>. During the last decades, many biological data sets have been growing exponentially in their size. Also, the complexity of those data sets (e.g., high-throughput data) bears potential problems when analyzing the data. For example, the networks might not be deterministically inferrable and, hence, one needs statistical techniques to cope structural errors. This example gives an idea about the complexity of the problem of selecting the right method.

---

<sup>38</sup> KIESEBERG, P./HOBEL, H./SCHRITTWIESER, S./WEIPPL, E./HOLZINGER, A. 2014. Protecting Anonymity in the Data-Driven Medical Sciences. In: Holzinger, A./Jurisica, I. (eds.) *Interactive Knowledge Discovery and Data Mining: State-of-the-Art and Future Challenges in Biomedical Informatics*, Springer Lecture Notes in Computer Science LNCS 8401. Berlin, Heidelberg: Springer, pp. in print.

<sup>39</sup> EMMERT-STREIB F./DEHMER M., Enhancing systems medicine beyond genotype data by dynamic patient signatures: Having information and using it too. *Frontiers in Genetics*, 4(241), 2013.

<sup>40</sup> EMMERT-STREIB F./DEHMER M. (Fn. 38).

<sup>41</sup> EMMERT-STREIB F./DEHMER M. (Fn. 38).

<sup>42</sup> HOLZINGER, A./OFNER, B./STOCKER, C./VALDEZ, A. C./SCHAAR, A. K./ZIEFLE, M./DEHMER, M. 2013. On Graph Entropy Measures for Knowledge Discovery from Publication Network Data. In: Cuzzocrea, A./Kittl, C./Simos, D. E./Weippl, E./Xu, L. (eds.) *Multidisciplinary Research and Practice for Information Systems*, Springer Lecture Notes in Computer Science LNCS 8127. Heidelberg, Berlin: Springer, pp. 354–362. HOLZINGER, A./OFNER, B./DEHMER, M. 2014. Multi-touch Graph-Based Interaction for Knowledge Discovery on Mobile Devices: State-of-the-Art and Future Challenges. In: Andreas Holzinger, I. J. (ed.) *Interactive Knowledge Discovery and Data Mining: State-of-the-Art and Future Challenges in Biomedical Informatics*, Springer Lecture Notes in Computer Science LNCS 8401. Berlin, Heidelberg: Springer, pp. in print.