

Disease networks and their contribution to disease understanding and drug repurposing. A survey of the state of the art

Authors

Eduardo P. García del Valle^{1*}, Gerardo Lagues García^{1,2}, Lucía Prieto Santamaría²,
Massimiliano Zanin^{2†}, Ernestina Menasalvas Ruiz^{1,2†}, Alejandro Rodríguez-González^{1,2†}.

†Equally contributing senior authors

*To whom correspondence should be addressed: ep.garcia@alumnos.upm.es

Centro de Tecnología Biomédica. Campus de Montegancedo. Pozuelo de Alarcon, 28223,
Madrid. +34 913364663

Affiliations

¹ ETS de Ingenieros Informáticos. Universidad Politécnica de Madrid. Boadilla del Monte,
Madrid, Spain.

² Centro de Tecnología Biomédica, ETS Ingenieros Informáticos. Universidad Politécnica de
Madrid. Pozuelo de Alarcón, Madrid, Spain.

About the authors

Eduardo P. García del Valle is a PhD student at the Faculty of Computer Science of the Universidad Politécnica de Madrid (UPM). His research areas are Knowledge Recovery, Artificial Intelligence and Bioinformatics.

Gerardo Lagunes García is a PhD student in the Medical Data Analytics Laboratory at the Center for Biomedical Technology (CTB) of the Technical Universidad Politécnica de Madrid (UPM). The areas of research of interest are data mining, knowledge recovery, web development and bioinformatics.

Lucía Prieto Santamaría is a biotechnology graduate student of the Universidad Politécnica de Madrid (UPM).

Ernestina Menasalvas Ruiz is a Full Professor of Universidad Politécnica de Madrid. Her research activities are on various aspects of data mining project development and in the last years her research is focused on data mining on the medical field. She leads the Data Mining and Simulation research group at UPM.

Alejandro Rodríguez-González, PhD, is an Associate Professor Universidad Politécnica de Madrid (UPM). His main research interests are the Semantic Web, Artificial Intelligence and Biomedical informatics field. He leads the Medical Data Analytics laboratory at the Center for Biomedical Technology (CTB).

Massimiliano Zanin is a Postdoctoral Researcher at the Center for Biomedical Technology (CTB) of Universidad Politécnica de Madrid (UPM). He is a member of the editorial team of Nature Scientific Reports, the European Journal of Social Behaviour, PeerJ and PeerJ Computer Science

Abstract

Over a decade ago, a new discipline called network medicine emerged as an approach to understand human diseases from a network theory point-of-view. Disease networks proved to be an intuitive and powerful way to reveal hidden connections among apparently unconnected biomedical entities such as diseases, physiological processes, signaling pathways, and genes. One of the fields that has benefited most from this improvement is the identification of new opportunities for the use of old drugs, known as drug repurposing. The importance of drug repurposing lies in the high costs and the prolonged time from target selection to regulatory approval of traditional drug development. In this document we analyze the evolution of disease network concept during the last decade and apply a data science pipeline approach to evaluate their functional units. As a result of this analysis, we obtain a list of the most commonly used functional units and the challenges that remain to be solved. This information can be very valuable for the generation of new prediction models based on disease networks.

Keywords: Disease Networks, Disease Similarity, Disease Understanding, Drug Repurposing, Data Science Pipeline

Introduction

The study of diseases as non-isolated elements and the understanding of how they resemble and relate to each other are crucial to provide novel insights into pathogenesis and etiology, as well as in the identification of new targets and applications for drugs [1]. The complete sequencing of the human genome at the beginning of the 21st century represented a revolution in the study of the relationships between diseases. In combination with the growing availability of transcriptomic, proteomic, and metabolomic data sources, it should help to improve the classification of diseases [2]. However, the use of these sources raised new problems such as their fragmentation, heterogeneity, availability and different conceptualization of their data [3, 4].

Recent developments in network theory provide a way to address this challenge by representing these complex relationships as a collection of linked nodes [5]. Complex networks theory is a statistical physics interpretation of the old graph theory, aimed at describing and understanding the structures created by the relationships between the elements of a complex system [6–9].

Those elements are represented by nodes, pairwise connected by links whenever a relationship is observed between the corresponding elements. The resulting structure can then be described by means of a plethora of topological metrics [10], or be used as a base for modelling the system.

The application of this field to biological problems has been named "network biology", while its use in biomedical problems is known as "network medicine" [11]. Following this approach, disease networks express the relationship between diseases as nodes and edges in a graph in $G = (D, W)$, where D represents the set of diseases (nodes) and W the set of their relationships (edges) based upon their similarity. The meaning of similarity varies depending on the data used

to build the network, which may be biological (genes or common proteins) or phenotypic (comorbidity, similar symptoms) [12], among other approaches.

During the past decade, numerous studies have been proposed to improve our understanding of the functioning of diseases and their relationships by creating disease networks based on different disease-disease association models and large-scale data exploitation. Of them, a significant number was oriented to exploit the new discovered relationships between diseases in the reassignment of known compounds for their treatment, the so-called "drug repurposing". In the first part of this document, we will thoroughly review this previous work, analyzing the evolution of the methodologies used in the creation of disease networks from a timeline perspective up to the state of the art.

Despite their different approaches and methodologies, in the studies dedicated to the improvement of the disease understanding and particularly to drug repositioning, the typical phases of a data science pipeline are observed, such as data extraction, data integration model, validation and presentation. In the second part of the document, these common parts are analyzed and their existing implementations are compared taking into account their use and performance. Finally, based on the previous analysis, new studies are proposed by improving or combining the phases of the pipeline.

Evolution of disease networks

Early studies proposing the use of disease networks for the analysis of their underlying relationships used data of biological origin. In 2007, Goh et al. constructed a disease-gene bipartite graph called "Diseasome" using information from OMIM database [1]. From the diseasome they derived the Human Disease Network (HDN), in which pairs of disorders are

connected if they have common genes. The study revealed that diseases tend to cluster by disease classes and that their degree of distribution follows a power law; that is, only a few diseases connect to a large number of diseases, whereas most diseases have few links to others. Aiming to reduce the bias of the HDN towards diseases transmitted in a Mendelian manner [13], subsequent studies used other sources of biological data. In 2008 year, Lee et al. constructed a metabolic disease network in which two disorders are connected if the enzymes associated with them catalyze adjacent reactions [14]. In 2009, Barrenas et al. [15] derived a complex disease-gene network (CDN) using GWAs (Genome Wide Association studies). The complex disease network showed that diseases belonging to the same disease class do not always share common disease genes. Complex disease genes are less central than the essential and monogenic disease genes in the human interactome.

The abundance of new biological data did not make researches overlook the existence of another important resource: the highest level clinical phenotypes, that is, symptoms. As one of the first and most obvious forms of diagnosis, the relationship between symptoms and diseases is widely documented in clinical records. In 2007, Rzhetsky et al. used the disease history of 1.5 million patients at the Columbia University Medical Center to infer the comorbidity links between disorders and prove that phenotypes form a highly connected network of strong pairwise correlation [16]. In 2009, Hidalgo et al. built a Phenotypic Disease Network (PDN) summarizing the connections of more than 10 thousand diseases obtained from pairwise comorbidity correlations reconstructed from over 30 million records from Medicare patients. The PDN is blind to the mechanism underlying the observed comorbidity, but it shows that patients tend to develop diseases in the network vicinity of diseases they have already had. Also disease progression was found to be different across genders and ethnicities [17]. More recently, Jiang et

al. [18] used data from the Taiwan National Health Insurance Research Database to construct the epidemiological HDN (eHDN), where two diseases are concluded as connected if their probability of co-occurring in clinics deviates from what expected under independence. However, despite their demonstrated potential in pathological analysis, the access and use of clinical records in medical research is limited by several issues, including the heterogeneity of sources [19], ethical and legal restrictions and the disparity of regulations between countries [20].

The analysis of open text sources has been used as an alternative to medical records. One of the reasons is the improvement in the techniques for Named Entity Recognition (NER) for the extraction of medical terms. Okumura et al. [21] performed an analysis of the mapping between clinical vocabularies and findings in medical literature using OMIM as a knowledge source and MetaMap as the NLP tool. Following this idea, Rodríguez et al. [22] used web scraping and a combination of NLP techniques to extract diagnostic clinical findings from MedlinePlus articles about infectious diseases using MetaMap tool. In a further study, the same team compared the performance of MetaMap and cTakes in the same task [23]. The increasing availability of retrieval engines such as PubMed or UKPMC, maintained by the US National Center for Biotechnology Information (NCBI) and the European Bioinformatics Institute (EBI), respectively, has also boosted this approach [24]. In 2014 Zhou et al. extracted symptom information from PubMed to construct the Human Symptoms Disease Network (HSDN). In the HSDN, the link weight between two diseases quantifies the similarity of their respective symptoms [25]. In 2015, Hoehndorf et al. created yet another Human Disease Network using a proposed similarity measure for text-mined phenotypes [26]. In both cases, these studies compare their results with gene-based networks, finding that symptom-based similarity of two

diseases strongly correlates with the number of shared genetic associations. They also demonstrated that not only Mendelian diseases tend to be grouped into classes, but also common ones.

Due to the intrinsic complexity of the relationships between diseases, the consideration of a single factor (shared genes or common symptoms) is a limiting point. In his review of the HDN in 2012, Goh. et al. proposed that each and every disease-contributing factor such as molecular links from interactome, co-expression and metabolism, as well as genetic interactions and phenotypic comorbidity links, will have to be integrated in a context-dependent manner. Furthermore, drug chemical information and non-biological environmental factors such as toxicity information altogether must also be incorporated [13]. The result will be a combination of general and bipartite network representations into a single, complex, k-partite heterogeneous network referred as the complete Diseasome.

In line with this idea, Sun [27] and Albornoz [28] combined multiple data sources to create tripartite networks of gene-disease-PPI and gene-disease-pathways, respectively, to predict disease-disease associations. The latter study proved that for two diseases sharing a certain number of genes, the level of inclusion can be different between both diseases due to the different pool of genes and metabolic pathways involved in each disease. In 2012, Chen et al. created an heterogeneous network from 17 public data sources relating to drugs, chemical compounds, protein targets, diseases, side effects and pathways [29]. In 2013, Žitnik et al. integrated molecular interaction and ontology data of 11 different types to create another heterogeneous network. When evaluating the predictive capacity of the network, genetic interactions proved to be the most informative feature, as they tend to be causative as opposed to

correlative and may therefore have less noise associated [4]. In both studies, the authors leveraged semantic ontology-level information to annotate the edges, as shown in **Figure 1**.

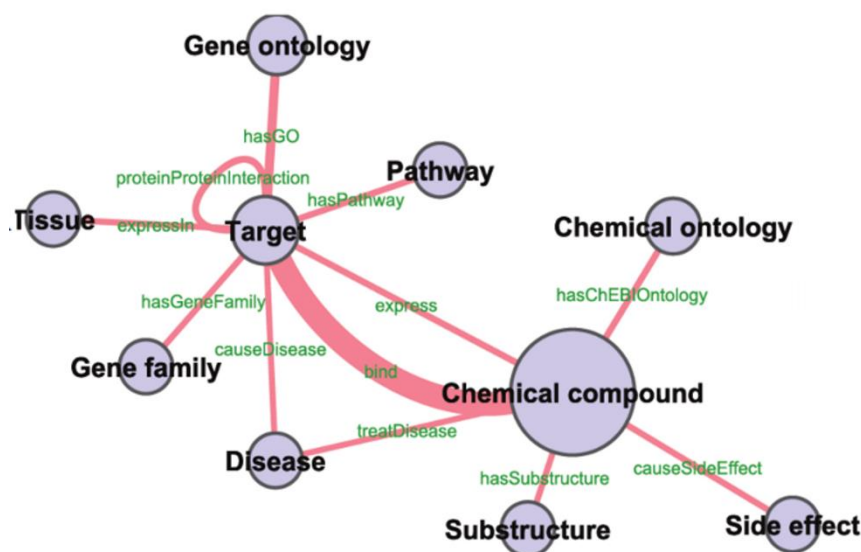


Figure 1. An example of a semantically annotated heterogeneous network. Every node and edge was semantically annotated using a systems chemical biology/chemogenomics ontology. Nodes were grouped into 10 classes which are linked by 12 types of edges. Two nodes are linked by one or more number of annotated paths. Retrieved from <https://journals.plos.org/ploscompbiol/article/figure?id=10.1371/journal.pcbi.1002574.g001>. Copyright: 2012 Chen et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License.

The evolution of these heterogeneous networks has resulted in the generation of complex tools for the study of disease associations based on multiple sources and types of relationships. A notable example is Hetionet [30], an integrative network encoding knowledge from millions of biomedical studies. Its data were integrated from 29 public resources to connect compounds, diseases, genes, anatomies, pathways, biological processes, molecular functions, cellular

components, pharmacologic classes, side effects, and symptoms. The completeness of the network is depicted in **Figure 2**.

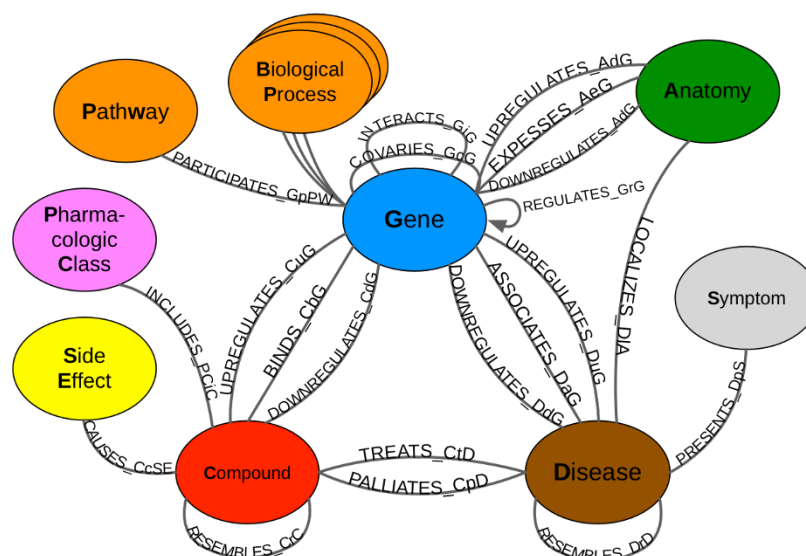


Figure 2. Representation of the metagraph in the Hetionet heterogeneous network. The schema shows the variety of data types (metanodes, depicted as circles) and connection types (metaedges, depicted as links) semantically annotated. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5640425/figure/fig1/>. Copyright: 2017, Himmelstein et al. This article is distributed under the terms of the Creative Commons Attribution License.

Application to drug repurposing

The constant improvement in disease association prediction through the use of network theory has fostered its application to drug repurposing. Drug repurposing is the utilization of known drugs and compounds to treat new indications [31]. Since the repositioned drug has already passed a significant number of toxicity and other tests, its safety is known and the risk of failure for reasons of adverse toxicology are reduced [32]. As a result, the cost and time needed to bring

a drug to market is significantly reduced compared to traditional drug development. The commercial applications of drug repositioning and the interest shown by pharmaceutical companies have led to a growing academic activity in this field. This fact is reflected in the evolution of the results for the search by “Drug Repurposing” in Google Scholar, as seen in **Figure 3**.

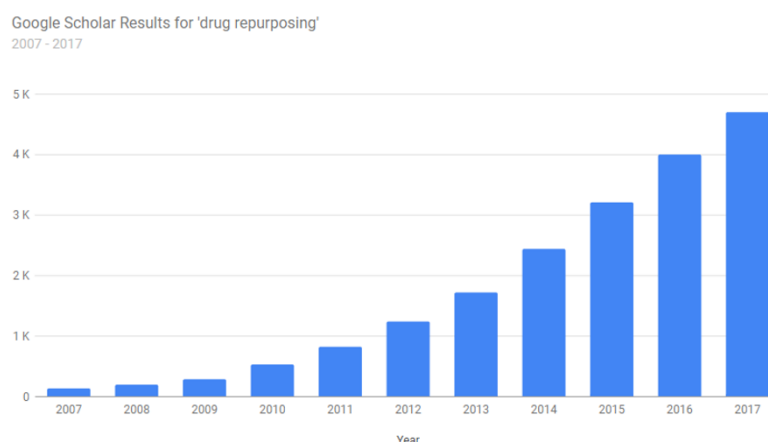


Figure 3. Evolution of the number of articles in Google Scholar containing the term “drug repurposing” within the last 10 years. Retrieved from <https://csullender.com/scholar/>. Copyright: 2017, Colin Sullender. Use authorized by the copyright owner.

First studies in drug repurposing were based on the “guilt-by-association” assumption, that is, similar drugs may share similar targets and vice-versa [33]. In 2007, Yildirim et al. created a graph composed of US Food and Drug Administration–approved drugs and proteins linked by drug–target binary associations [34]. Similar studies were carried out by Ma’ayan [35] in 2007 and Chiang [36] and Bleakley [37] in 2009. In 2008, Nacher Schwartz compiled a drug-therapy network with all US-approved drugs and associated human therapies. From this bipartite network they constructed two other networks: a drug network and a therapy network. Therapies are

closely linked to diseases, therefore the therapy network gave insights about the relations between diseases as well, making this work comparable to previous studies on human disease networks [38].

The above mentioned studies followed a drug-centric approach, that is, they discovered new indications for existing drugs based on drug-drug similarities. Other studies followed a disease-centric approach, in which effective drugs were identified based on disease-disease similarity. In 2008, Campillos et al. predicted new targets for drugs by calculating similarities between diseases based on side effect that appears from injection of drug [39]. In 2009, Guanghui Hu et al. performed a systematic, large-scale analysis of genomic expression profiles of human diseases and drugs to create a disease-drug network [40]. Suthram in 2010 [41], Mathur in 2012 [42] or Zhou in 2014 [25] also predicted new uses of existing drugs based on disease-disease associations calculated from mRNA expression similarity, biological process semantic similarity or phenotypic similarity, respectively.

As was the case in disease classification, focusing purely on drug-disease relations with no consideration of other underlying genetic or pharmacological mechanisms at play is a limiting factor in accuracy of drug repurposing prediction, due to the lack of completeness of individual information [31]. Therefore, incorporating heterogeneous data sources can potentially solve this issue. In 2011 Gottlieb made use of a broader collection of data sources to create five drug-drug similarity measures and two disease-disease similarity measures. These similarity measures were then used by PREDICT, an algorithm to infer novel drug indications [43]. Daminelli in 2012 [44] and Wang in 2014 [45] built tripartite drug-target-disease networks to predict repurposing candidate drugs.

Ultimately, advances towards more comprehensive networks have resulted in tools for the prediction of new treatments given a certain disease. This is the case of Rephetio [46], a project based on Hetionet [30] that predicted repurposing candidates by applying an algorithm originally developed for social network analysis [47]. Similarly, in the context of drug discovery, one can leverage on identifying potential associations between compounds and protein targets. To cope with the noisy, incomplete and high-dimensional nature of large-scale biological data, Luo et al. proposed DTINet [48], a Drug Target Indications (DTI) prediction system based on learning low-dimensional feature vectors that capture the context information of individual networks. DTINet showed better performance than other state-of-the-art DTI prediction methods and discovered the potential application of cyclooxygenase inhibitors in preventing inflammatory diseases.

A data science pipeline to build disease networks

Throughout the previous section, we have seen how the rise of network medicine studies has resulted in a expanding variety of innovative methods for the construction and exploitation of disease networks. However, despite using different strategies, these methods are generally based on determining the similarities and relationships between diseases and their treatments at phenotypic level (comorbidity, side-effects) or biological level (common genes, proteins, compounds). Furthermore, they clearly share common phases such as data ingestion, data processing, analysis, modeling or visualization that can be represented as functional units of a data science pipeline, as shown in **Figure 4**.

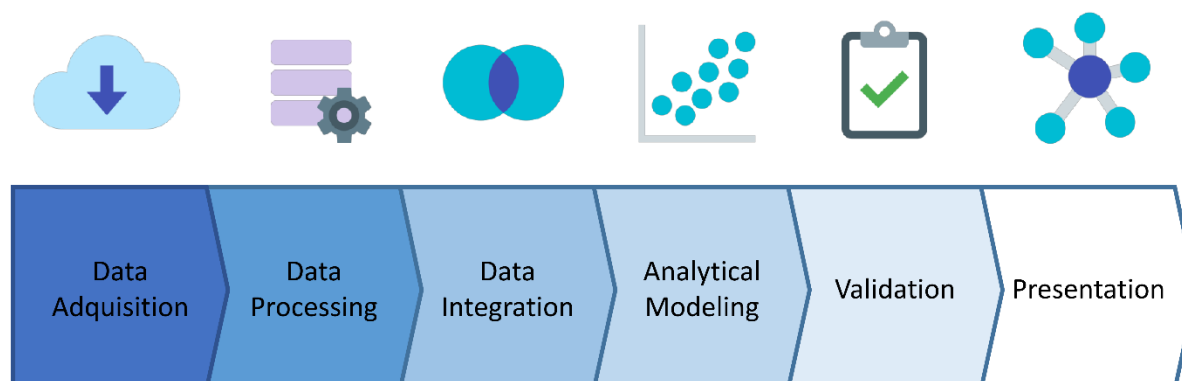


Figure 4. Sequence of functional units of a data science pipeline, including: data acquisition, data processing, data integration, analytical modeling, validation and presentation.

The data science pipeline consists in a sequence of stages or functional units that sequentially process some input data in order to solve a certain problem [49]. This concept applies to disease networks, where disease information is processed to discover how diseases relate to each other or how drugs can be repositioned. The pipeline representation also facilitates the reproducibility and the comparison among studies as a whole and also at phase level. Most importantly, it also enhances the reusability and the recombination of the functional units to build new drug-repurposing. Throughout the following sections we will describe the process of construction and exploitation of a disease network through the functional units of a data science pipeline.

Data acquisition and processing

The first step in the pipeline is to acquire data from a variety of sources, a process known as data acquisition or data ingestion. As seen in the section about the Evolution of Disease Networks, the

growing availability of information sources has allowed developing different approaches to improve our understanding of diseases and to predict new drug applications.

A significant number of studies use biological data, such as KEGG (genes and pathways) [14], BioGRID (protein interactions) [4] or OMIM (genes and phenotypes) [1, 42], among many others. **Supplementary Table 1** contains some of the most important sources of biological information, including their type and description. Studies on disease networks focusing on drug repositioning exploit drug databases and their relation to genes, phenotypes and compounds, such as those offered by the FDA [34–37] or DrugBank [25, 39–42], for instance.

Supplementary Table 2 collects the most common drug data sources. Finally, an increasingly significant number of studies use data obtained by mining medical literature sources (e.g. articles, clinical trials) such as PubMed [25, 26, 50] or the GWAS Catalog [27]. **Supplementary Table 3** contains some of the most relevant sources of medical literature.

A second step in the pipeline consists in transforming and mapping data into a format with the intent of making it more appropriate to work (usually referred as data processing, data wrangling or data munging). Recent studies combine multiple databases to provide more accurate prediction models [4, 29, 30]. However, this poses a challenge when relating identifiers or terms obtained from different sources. To address this problem, researchers use thesauri of terms such as MeSH, SNOMED CT or UMLS; code listings such as ICD or HGNC; and ontologies such as DO, PO, GO or Uberon [26, 51]. Being a valuable source of semantic and hierarchical information themselves, these resources allow mapping data such as disease codes or medical terms. In the case of medical literature sources, the use of metadata (such as MeSH headers in the case of Pubmed, for example) is often combined with terms extraction tools such as

MetaMap or cTakes [23]. **Supplementary Table 4** lists some of the sources used for data mapping.

The way to exploit the information in these databases varies greatly from one source to another. Largest databases offer online advanced search and provide developers with application programming interfaces (APIs) to facilitate intensive access to data. For example, the NCBI provides the E-utilities, a public API to access all the Entrez databases including PubMed, PMC, Gene, Nucleotide and Protein. The Japanese KEGG also provides REST APIs for data consumption. DisGeNET provides an SPARQL endpoint that allows exploration of the DisGeNET-RDF data set and query federation to expand gene-disease association information with data on gene expression, drug activity and biological pathways, among other. In some cases, data can also be downloaded for their consumption through on-premise applications, as in the case of the Disease Ontology or the Gene Ontology, for example. This disparity complicates the use of different sources in research projects. To alleviate this problem, initiatives such as Biopython¹ offer common libraries to access multiple sources reducing code duplication in computational biology. Finally, it is very important to know the limitations imposed by each source regarding the volume and use of the data. **Supplementary Tables 1-4** also include information in this regard.

Data integration and modeling

In the next steps of the data science pipeline, data previously acquired and processed are integrated and analyzed in order to answer the matter of our study. In other words, a disease network is built by combining the output of the previous stage and a model is constructed from

¹ <https://biopython.org>

it. Disease networks consist of a set of nodes (mainly, but not only, representing diseases) and a set of edges (connecting diseases directly or through other related node types). Depending on the type of node they connect, network edges can be directed or undirected, weighted or unweighted. As described in previous sections, over the past decade successive studies based on disease networks have proposed different models of data integration.

Homogeneous networks

Homogeneous disease networks (i.e. those where nodes represent diseases and edges represent direct connections among them) are the simplest type of disease networks. In many studies these networks are built as a projection of a heterogeneous disease network (i.e. a network in which diseases are connected to other types of nodes) [1, 28]. For example, in **Figure 5**, the gene-disease bipartite network is projected onto the disease similarity network (DSN) by relating two diseases that have a gene in common. The disease–disease network can then be analysed by using standard network based methods [1, 52]. In a simplistic approach, the link weights in the resulting disease–disease network represent the link multiplicity resulting from the projection. More complex methods, such as hyperbolic weighting or resource allocation weighting, have been proposed as an alternative [53, 54].

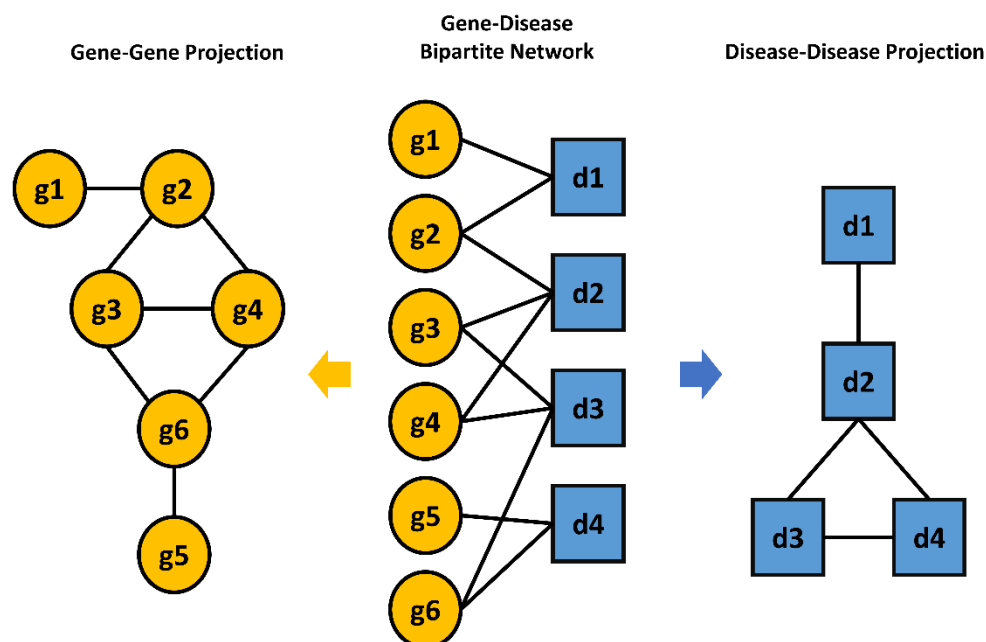


Figure 5. A representation of a heterogeneous network composed of gene-disease interactions. Homogeneous gene-gene or disease-disease networks are obtained via gene or disease projections, respectively.

In other studies, homogeneous disease networks are built as similarity networks. In these networks, if the similarity score between disease i and j is more than zero, the corresponding vertices are linked by an edge in the network. The weight of this edge is the corresponding disease similarity score. Several computation methods for the disease similarity score have been proposed, being Vector Space Model (VSM) [55] among the most popular ones. For instance, in 2006 Van Driel et al. represented diseases as vectors of features (viz. disease associated MeSH terms extracted from OMIM records) weighted by their inverse document frequency [56]. The similarity between diseases was then computed as the cosine of the disease vector angles (i.e. cosine similarity). A similar approach was followed by Zhou to build the HSDN [25] and by Sun to build the Integrated Disease Network [57]. Hoehndorf et al. proposed Normalized Pointwise

Mutual Information (NPMI) for disease phenotypic term weighting and later used the PhenomeNET system to compute similarity between diseases using a Jaccard index based measured [26]. Similarity measures based on the term hierarchy in the Disease Ontology and the Gene Ontology have been proposed by Resnik, Lin, Wang, Mathur and Cheng [42, 58–61], and have been integrated in online tools like DisSim or DisSetSim [62, 63]. Okumura et al. described alternative similarity measures based on standardized disease classification, probabilistic calculation, and machine learning [64].

Heterogeneous networks

The projection of heterogeneous networks into homogeneous disease-disease networks allows applying simpler network analysis techniques on the resulting network. However, it often results in information loss. For instance, in **Figure 5** by projecting the gene-gene network onto the disease-disease network, the information about gene interactions and their structure is lost. In contrast, heterogeneous networks make it easy to predict relationship between entities of different types, such as diseases, genes or drugs, following a guilt-by-association paradigm [33]. For example, a drug that regulates a gene associated to a disease could be repurposed for diseases associated to the same gene. Data fusion by matrix factorization and network topology based techniques, such as diffusion and meta-path, are the most common methods for edge prediction in heterogeneous networks.

Matrix Factorization methods are closely related to clustering (unsupervised) algorithms. Non-Negative Matrix Factorization (NNMF) decompose matrices of heterogeneous data and data relationships to obtain low-dimensional matrix factors. These factors are then used to reconstruct the data matrices, adding new unobserved data obtained from the latent structure captured by the low-dimensional matrix factors. Hence, NNMF provides a mechanism to integrate heterogeneous

data of any number, type and size. In 2013 Žitnik et al. applied a variant of NNMF called non-negative matrix tri-factorization to discover new disease-disease association by fusing 11 data sources on four type of objects including drugs, genes, DO terms and GO terms [4]. In 2015 Dai et al. integrated drug-disease associations, drug-gene interactions, and disease-gene interactions with a matrix factorization model to predict novel drug indications [65]. More recently, Zhang et al. proposed a similarity constrained matrix factorization method for the drug-disease association prediction using data of known drug-disease associations, drug features and disease semantic information [66].

Methods based on diffusion (i.e. information spreading across network links) have also been extensively proposed to estimate the strength of the connection between nodes of heterogeneous networks. An advantage of such approaches, also called network propagation methods, over matrix factorization is that they preserve the network structure. Chen et al developed the method of Network-based Random Walk with Restart on the Heterogeneous network (NRWRH), a variation of a ranking algorithm, to predict potential drug-target interactions on heterogeneous networks [67]. Further variations of random walk algorithms, such Bi-Random Walk (BiRW) have been applied to predict novel disease-gene [68], disease-MiRNA [69] or disease-lncRNA associations [70], among others.

Metapath-based approaches also preserve the network structure, and additionally provide an intuitive framework and interpretable models and results. A meta-path P is a path defined over the general schema of the heterogeneous network $G = (A, R)$, where A represents the set of nodes and R the set of their relationships. The metapath is denoted by $A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_l} A_{l+1}$,

where l is an index indicating the corresponding metapath [47]. **Figure 6** shows the metapaths extracted from an annotated heterogeneous network.

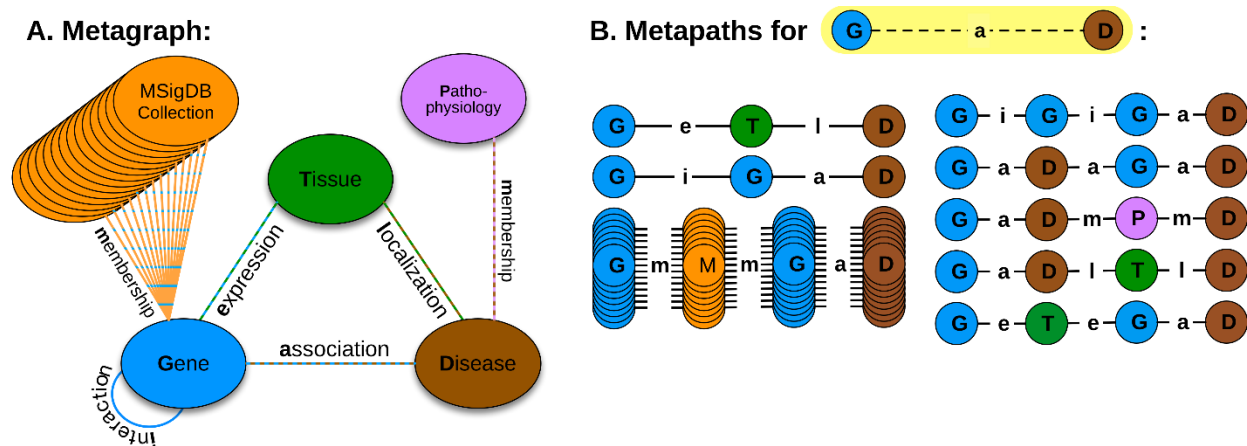


Figure 6. A) The Hetionet annotated heterogeneous network is constructed according to the metagraph schema, which is composed of metanodes (node types) and metaedges (edge types). B) The network topology connecting a gene and disease node is measured along metapaths (types of paths). Starting on Gene and ending on Disease, all metapaths length three or less are computed by traversing the metagraph. Retrieved from <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004259>. Copyright: 2015 Himmelstein, Baranzini. This is an open access article distributed under the terms of the Creative Commons Attribution License.

In their 2012 study, Chen et al. developed a meta-path based statistical model called Semantic Link Association Prediction (SLAP) to assess the association of drug target pairs and to predict missing links [29]. In 2016 Gang Fu et al. proposed an alternative DTI approach to the SLAP algorithm taking advantage of machine learning methods such as Random Forest and Support Vector Machine [63]. To quantify the prevalence of the meta-paths, Himmelstein adapted an existing method developed for social network analysis (PathPredict) and developed a new metric

called degree-weighted path count (DWPC). The DWPC downweights paths through high-degree nodes when computing meta-path prevalence [30].

Despite maintaining and exploiting the structure of heterogeneous networks, methods based on diffusion or meta-paths present some scalability limitations, such as the bias introduced by the noise and high-dimensionality of biological data or the effort in feature engineering. Recently, Luo et al. designed DTINet, a novel network integration pipeline for DTI prediction. DTINet integrates information from heterogeneous sources (e.g., drugs, proteins, diseases and side-effects) and copes with the noisy, incomplete and high-dimensional nature of large-scale biological data by learning low-dimensional but informative vector representations of features for both drugs and proteins [48].

Model validation

In this analysis of the reconstruction of disease networks, we wanted to give a special relevance to the validation process. Ensuring that the computational pipeline is producing correct and valid results is critical, particularly in a clinical setting [71]. As previously explained, disease networks are used in studies as diverse as the discovery of new disease-disease relationships, the prediction of gene-disease relationships (GDA) or the repositioning of drugs. The validation of the network depends, therefore, on the type of study in question. In general, the validation can be done experimentally or by computational techniques.

Approaches and sources

Experimental validation includes the verification of the predictions in a controlled environment outside of a living organism (*in vitro*) or using a living organism (*in vivo*). Animal studies and clinical trials are two forms of *in vivo* research. For example, in their drug repositioning study

based on heterogeneous networks, Luo et al. validated the bioactivities of the COX inhibitors predicted by DTINet experimentally. They tested their inhibitory potencies on the mouse kidney lysates using the COX fluorescent activity assays [48]. Jodeleit et al. validated their disease network of inflammatory processes in humanized NOD/SCID/IL2R γ (NSG) mice [72]. While experimental validation studies have the potential to offer more conclusive results about the performance of disease networks, they have several limitations. First, animal studies and clinical trials require expensive lab work and are long and costly. In addition, their conclusions can be misleading. For example, a therapy can offer a short-term benefit, but a long-term harm. Also, it is debatable that genomic responses in mouse models mimic human inflammatory disease [73].

In silico is an expression used to mean “performed on computer or via computer simulation.” In silico tests have the potential to speed the validation process while reducing the need for expensive lab work. In silico validation requires a point of reference for evaluating the model performance, also known as Criterion Standard or Gold Standard. It is noteworthy that in the field of biomedicine usually the Criterion Standard is actually the best performing test available under reasonable conditions [74]. For example, in this sense, a MRI is the gold standard for brain tumour diagnosis, though it is not as good as a biopsy [75]. Hence, the most recurrent benchmarks used in the validation in silico of disease networks include consolidated data biomedical sources and medical literature.

Sources of biological, phenotypic or chemical data as well as several available ontologies and code standards (see Data extraction section) are used for validation in many studies focusing on disease networks. For instance, their performance to discover disease-disease relationships has been validated with the disease classifications in the Disease Ontology [4, 26] or in the ICD codes [28], as well as with comorbidity associations downloaded from the Human Disease

Network (HuDiNe) [27]. DisGeNET has been used to validate de novo gene-disease associations [76], as it integrates data from expert curated repositories with information gathered through text-mining of the scientific literature, GWAS catalogues and animal models [77]. For the validation of drug repositioning predictions, sources such as PharmacotherapyDB and DrugCentral were exploited [46].

The aforementioned sources are inevitably biased towards consolidated knowledge, and therefore they might suffer some limitations in corroborating new discoveries. As an alternative (or usually, as a complement) to these sources, medical literature (i.e. studies, medical trials, clinical histories) are used to validate disease network based studies. For instance, Mathur and Paik used previous studies to validate disease-disease and drug-target associations [42, 78]. In some cases, the validation process also combined human (i.e. medical experts) action to corroborate the discoveries [25].

Methods

Leaving aside the particularities of biomedical research and its sources, the validation of classification or prediction methods based on disease networks does not differ from other validation cases. Therefore, in the analyzed studies we found validation methods widely used. For example, k-fold cross-validation is often used to check whether the model is an overfit or not [79, 80]. Overfitting is one of the typical problems of validation, especially when limited data sets are available.

To quantify the predictive power of their network-based model, many studies use the Area Under the Curve of the Receiver Operating Characteristic (AUC-ROC), another frequently used method in validation problems [26, 81, 82]. The AUC-ROC is the plot between sensitivity and (1-

specificity). (1- specificity) is also known as false positive rate and sensitivity is also known as true positive rate. The p -value Is the probability that the observed sample AUC-ROC could actually correspond to a model of no predictive power (null hypothesis), i.e. to a model whose population AUC-ROC is 0.5. If p -value is small, then it can be concluded that the AUC-ROC is significantly different from 0.5 and that therefore there is evidence that the model actually discriminates between significant and non-significant results [83]. Typically, a threshold value (called significance level) of p -value < 0.05 is used. However, biomedicine studies often use more restrictive values like 0.005 [42] or even 0.001 [4]. As an alternative to the AUC-ROC, the p -value can be obtained for other tests such as chi-squared or Fisher's exact, depending on the case of study [84]. Finally, to control the familywise error rate associated with multiple testing, a correction algorithm like Benjamini–Hochberg or Bonferroni is applied.

Presentation

Last but not least, at the end of the pipeline the results obtained should come out in a format that can be consumed by the audience (e.g. the scientific community, the media or even ourselves to inform the next iteration). One of the major advantages of disease networks is the intuitive access to the underlying complex interactions between diseases and other diseases, genes or drugs. Thus, publishing not only the data but also means to explore and exploit the network is key to ensure reproducibility and extensibility of the study [85]. Early studies lacked this option, although access to their data allowed the construction of visualization tools a posteriori. For example, Ramiro Gómez created an interactive view of the Human Disease Network proposed by Goh in 2007 using the graph visualization software Gephi² and the original dataset from the

² <https://gephi.org>

Figure 7. A representation of gene-disease associations with the DisGeNET application for Cytoscape. Retrieved from <http://apps.cytoscape.org/apps/disgenetapp>. Copyright: 2016, Anna Bauer-Mehren et al. The DisGeNET plugin is distributed under the GNU GPL 3.0 license.

On their side, Himmelstein et al. accompanied their study based on heterogeneous disease networks with a powerful visualization tool built with Neo4j⁶ [30] that provides browsing and querying on Hetionet (see **Figure 8**). Being a remarkable example of data accessibility, not only the data but also the code of this tool is publicly available. Different studies of the University of Rome, such as SIGNOR⁷ and DISNOR⁸, also provides a disease network visualization tool that includes intuitive representations of the interactions between biological entities at different complexity levels (see **Figure 9**). This visualization tool was developed ad-hoc for these projects [90–92].

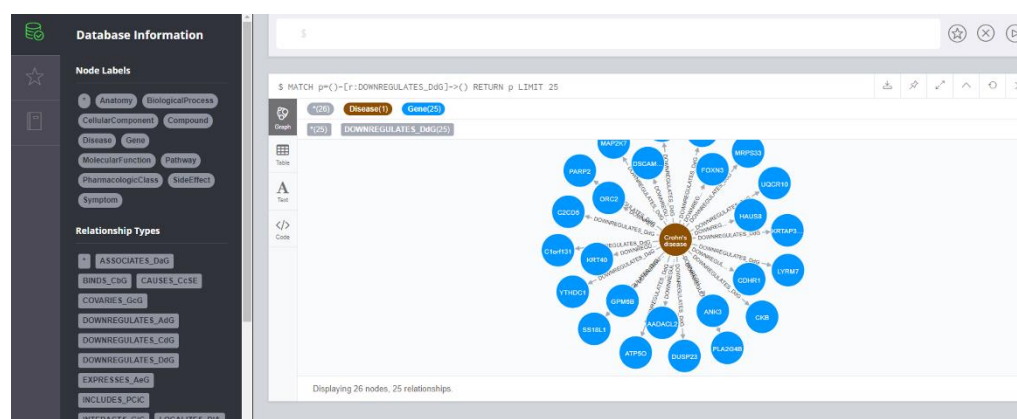


Figure 8. Hetionet Neo4j browser displaying disease-gene associations for Crohn's disease in an interactive way. Retrieved from <https://neo4j.het.io/browser>. Copyright: 2015

⁶ <https://neo4j.com>

⁷ <https://signor.uniroma2.it>

⁸ <https://disnor.uniroma2.it>

Himmelstein, Baranzini. This is an open access project distributed under the terms of the Creative Commons Universal License.

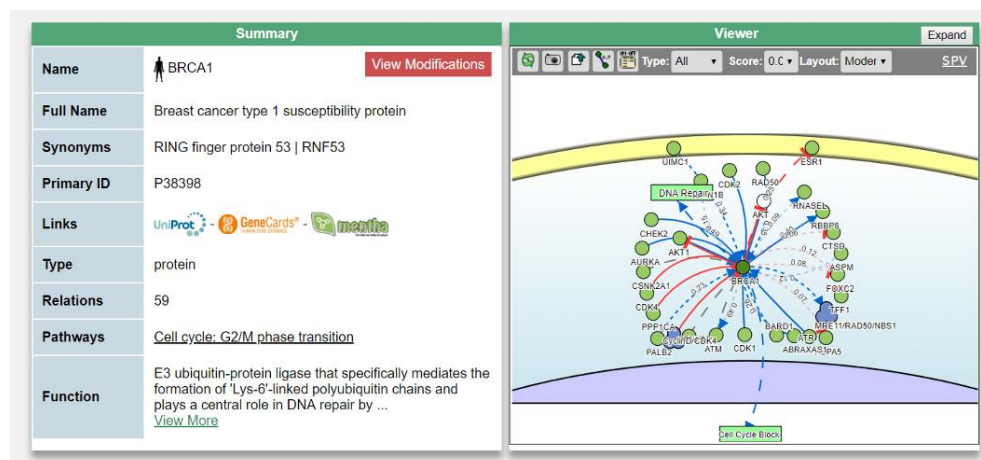


Figure 9. DISNOR signaling network browser showing an interactive graph for a protein related with Breast Cancer. Copyright: 2018, SIGNOR. This is an open access project distributed under the terms of the Creative Commons Attribution-ShareAlike 4.0 International License.

A recent study by Pavlopoulos et al. performs an empirical comparison of visualization tools for large-scale network analysis [93].

Discussion

The analysis of the evolution of the disease networks carried out in the first part of the document shows how these models have become increasingly complex and allow to address arduous problems such as the improvement of our disease understanding or the repositioning of drugs with promising results. However, as a side effect of this growing complexity, new challenges have emerged that need to be addressed.

The growing availability of biological sources, key in the improvement of disease networks, is ballasted by their fragmentation, heterogeneity, availability and different conceptualization of their data [3]. Furthermore, these sources are intrinsically biased towards consolidated knowledge, which complicates the discovery of novel findings. The exploitation of textual sources such as clinical histories or scientific articles - more abundant and faster growing - allows researchers to compensate for these limitations. As an example of the abundance and potential of these alternative sources, in a recent study Westergaard extracted and analyzed 15 million English scientific full-text articles published during the period 1823–2016 [94].

Despite this demonstrated potential, the exploitation of medical literature is hindered by factors such as its limited access and heterogeneity. In the aforementioned study by Westergaard, the team could only access a subset of the Medline articles in full-text mode, while for the rest only the abstracts were available. In addition, depending on the source, they had to process documents with different structures and format. As an alternative, a recent study proposed the use of Wikipedia as a source of structured and free-access text data, evaluating its usefulness in the detection of relations between diseases based on its symptoms/diagnosis elements, and comparing its performance with that of PubMed. The obtained results showed that Wikipedia can be as relevant a source as PubMed for this type of analysis [95].

Another limiting factor when integrating new sources to enhance the predictive capacity of disease networks is noise [96]. Adding new sources does not necessarily imply an improvement, since some databases are more informative than others. For example, Žitnik et al. evaluated the impact of removing sources in the performance of the proposed model to validate their informativeness. They observed that while the absence of some sources significantly affected the performance, in other case the impact was minimum [4]. It is therefore necessary to counteract

this effect by choosing algorithms that eliminate irrelevant sources or features before constructing the model [48].

Validation is yet another challenge in the studies based on disease networks. In some cases, the absence of a Standard Criterion leads to the use of previous studies for the validation of the new models [42, 78]. This might ultimately result in the propagation of errors from one study to another. The use of curated sources and of sufficiently contrasted studies, combined when possible with in-vitro and in-vivo validations, helps to alleviate solution to this problem [48, 72]. Related with the challenge of validation, the difficulty in accessing data from some studies prevents their reproducibility and verification by other teams, which makes them less reliable as references for future studies or as benchmarks. However, the effort of some researchers in making available the results of their work is worth to mention. Study cases such as Hetionet, Rephetio, SIGNOR and DisNOR [30, 46, 90, 91], which offer advanced search and visualization tools, undoubtedly represent the path to follow.

The review of the process of creating a disease network from the point of view of a data science pipeline carried out in the second part of the document allows to compare how each study has faced these challenges. **Supplementary Table 5** lists some of the most notable studies related to disease networks of the last decade, breaking down each of its phases. It also contains information on the type of problem addressed and the characteristics of the obtained network. This table could be considered an extension / update of the one compiled by Sun K. et al. [27].

Conclusion and future work

Research studies on based disease networks have significantly advanced over the last decade. From the initial simple undirected networks that associated diseases with symptoms or genes in a

way, we have moved to complex networks that relate the disease to dozens of features from different sources in a semantic, directional and weighted way. The growing availability of biological and textual sources, the improvement in techniques and processing capacity and the use of new models have contributed fundamentally to this progress. As can be concluded from the analysis in the first part of the document, the contribution of disease networks to fields of disease understanding and drug repositioning is increasingly notable.

Nevertheless, an exhaustive analysis of the phases in the process of creating disease networks carried out in the second half of the document reveals important challenges. First, biological sources suffer from fragmentation, heterogeneity, lack of availability and different conceptualization, that can only be alleviated in part with the aggregation of textual sources. Second, the combination of sources involves the introduction of noise that can affect the performance of the model, which makes it necessary to take preventive measures in this regard. Finally, the scarcity of reference data and verifiable studies hinders the validation of the new models.

In addition to detecting these challenges, the analysis of disease networks from the point of view of their functional units allows for a more precise comparison of studies, highlighting their differences and common points. This study and the presented analyses, reflected in the summary tables, can serve to inspire future work. For example, a performance comparison of the prediction models in the different studies might lead to deduce which functional units offer better results. In a next phase, based on the obtained results, alternative combinations of these functional units could be proposed to build new pipelines and obtain more precise models based on disease networks.

Funding

Horizon 2020 research and innovation programme under grant agreement No. 727658, project IASIS (Integration and analysis of heterogeneous big data for precisionmedicine and suggested treatments for different types of patients).

Conflicts of interest

Authors declare no conflict of interest.

Keypoints

- Disease networks have proved to be an intuitive and powerful way to address arduous problems such as the improvement of our disease understanding or the repositioning of drugs.
- Over the last decade, disease networks have evolved from initial simple and undirected homogeneous networks, to complex, semantic, directional and weighted heterogeneous networks.
- Depite their increasing complexity, studies on disease networks share common phases that can be represented as functional units of a data science pipeline for a better analysis and comparison.
- The heterogeneity and fragmentation of biological and textual sources, the noise introduced by their combination and the scarcity of validation datasets are some of the challenges discovered through this analysis.

References

1. Goh K-I, Cusick ME, Valle D, Childs B, Vidal M, Barabási A-L. The human disease network. *PNAS*. 2007;104:8685–90.
2. Yang J, Wu S-J, Dai W-T, Li Y-X, Li Y-Y. The human disease network in terms of dysfunctional regulatory mechanisms. *Biol Direct*. 2015;10. doi:10.1186/s13062-015-0088-z.
3. Loscalzo J, Kohane I, Barabasi A-L. Human disease classification in the postgenomic era: A complex systems approach to human pathobiology. *Mol Syst Biol*. 2007;3:124.
4. Žitnik M, Janjić V, Larminie C, Zupan B, Pržulj N. Discovering disease-disease associations by fusing systems-level molecular data. *Scientific Reports*. 2013;3:3202.
5. Chan SY, Loscalzo J. The emerging paradigm of network medicine in the study of human disease. *Circ Res*. 2012;111:359–74.
6. Strogatz SH. Exploring complex networks. *Nature*. 2001;410:268–76.
7. Albert R, Barabási A-L. Statistical mechanics of complex networks. *Rev Mod Phys*. 2002;74:47–97.
8. Newman M. The Structure and Function of Complex Networks. *SIAM Rev*. 2003;45:167–256.
9. Boccaletti S, Latora V, Moreno Y, Chavez M, Hwang D-U. Complex networks: Structure and dynamics. *Physics Reports*. 2006;424:175–308.
10. Costa L da F, Rodrigues FA, Travieso G, Boas PRV. Characterization of complex networks: A survey of measurements. *Advances in Physics*. 2007;56:167–242.
11. Barabási A-L. Network Medicine — From Obesity to the “Diseasome.” *New England Journal of Medicine*. 2007;357:404–7.
12. Park S, Lee D, Shin H. Network mirroring for drug repositioning. *BMC Med Inform Decis Mak*. 2017;17 Suppl 1. doi:10.1186/s12911-017-0449-x.
13. Goh K-I, Choi I-G. Exploring the human diseasome: the human disease network. *Brief Funct Genomics*. 2012;11:533–42.
14. Lee D-S, Park J, Kay KA, Christakis NA, Oltvai ZN, Barabási A-L. The implications of human metabolic network topology for disease comorbidity. *PNAS*. 2008;105:9880–5.
15. Barrenas F, Chavali S, Holme P, Mobini R, Benson M. Network Properties of Complex Human Disease Genes Identified through Genome-Wide Association Studies. *PLoS One*. 2009;4. doi:10.1371/journal.pone.0008090.
16. Rzhetsky A, Wajngurt D, Park N, Zheng T. Probing genetic overlap among complex human phenotypes. *PNAS*. 2007;104:11694–9.
17. Hidalgo CA, Blumm N, Barabási A-L, Christakis NA. A Dynamic Network Approach for the Study of Human Phenotypes. *PLOS Computational Biology*. 2009;5:e1000353.
18. Jiang Y, Ma S, Shia B-C, Lee T-S. An Epidemiological Human Disease Network Derived from Disease Co-occurrence in Taiwan. *Scientific Reports*. 2018;8:4557.

19. Conway M, Berg RL, Carrell D, Denny JC, Kho AN, Kullo JJ, et al. Analyzing the heterogeneity and complexity of Electronic Health Record oriented phenotyping algorithms. *AMIA Annu Symp Proc.* 2011;2011:274–83.
20. Yip C, Han N-LR, Sng BL. Legal and ethical issues in research. *Indian J Anaesth.* 2016;60:684–8.
21. Okumura T, Aramaki E, Tateisi Y. Clinical Vocabulary and Clinical Finding Concepts in Medical Literature. In: *The First Workshop on Natural Language Processing for Medical and Healthcare Fields.* Nagoya: Asian Federation of Natural Language Processing; 2013. p. 7–13. <http://www.aclweb.org/anthology/W13-4602>. Accessed 3 Sep 2018.
22. Rodríguez-González A, Martínez-Romero M, Costumero R, Wilkinson MD, Menasalvas-Ruiz E. Diagnostic Knowledge Extraction from MedlinePlus: An Application for Infectious Diseases. In: Overbeek R, Rocha MP, Fdez-Riverola F, De Paz JF, editors. *9th International Conference on Practical Applications of Computational Biology and Bioinformatics.* Springer International Publishing; 2015. p. 79–87.
23. Rodríguez González A, Costumero Moreno R, Martínez Romero M, Wilkinson MD, Menasalvas Ruiz E. Extracting diagnostic knowledge from MedLine Plus: a comparison between MetaMap and cTAKES Approaches. *Current Bioinformatics.* 2015;375:1–7.
24. Rebholz-Schuhmann D, Oellrich A, Hoehndorf R. Text-mining solutions for biomedical research: enabling integrative biology. *Nat Rev Genet.* 2012;13:829–39.
25. Zhou X, Menche J, Barabási A-L, Sharma A. Human symptoms–disease network. *Nature Communications.* 2014;5:4212.
26. Hoehndorf R, Schofield PN, Gkoutos GV. Analysis of the human diseasome using phenotype similarity between common, genetic, and infectious diseases. *Sci Rep.* 2015;5:10888.
27. Sun K, Gonçalves JP, Larminie C, Pržulj N. Predicting disease associations via biological network analysis. *BMC Bioinformatics.* 2014;15:304.
28. Garcia-Albornoz M, Nielsen J. Finding directionality and gene-disease predictions in disease associations. *BMC Syst Biol.* 2015;9. doi:10.1186/s12918-015-0184-9.
29. Chen B, Ding Y, Wild DJ. Assessing Drug Target Association Using Semantic Linked Data. *PLoS Comput Biol.* 2012;8. doi:10.1371/journal.pcbi.1002574.
30. Himmelstein DS, Lizee A, Hessler C, Brueggeman L, Chen SL, Hadley D, et al. Systematic integration of biomedical knowledge prioritizes drugs for repurposing. *eLife.* 6. doi:10.7554/eLife.26726.
31. Mullen J, Cockell SJ, Woollard P, Wipat A. An Integrated Data Driven Approach to Drug Repositioning Using Gene-Disease Associations. *PLoS One.* 2016;11. doi:10.1371/journal.pone.0155811.
32. Hernandez JJ, Prysizlak M, Smith L, Yanchus C, Kurji N, Shahani VM, et al. Giving Drugs a Second Chance: Overcoming Regulatory and Financial Hurdles in Repurposing Approved Drugs As Cancer Therapeutics. *Front Oncol.* 2017;7:273.
33. Altshuler D, Daly M, Kruglyak L. Guilt by association. *Nat Genet.* 2000;26:135–7.
34. Yildirim MA, Goh K-I, Cusick ME, Barabási A-L, Vidal M. Drug-target network. *Nat Biotechnol.* 2007;25:1119–26.
35. Ma'ayan A, Jenkins SL, Goldfarb J, Iyengar R. Network Analysis of FDA Approved Drugs and their Targets. *Mt Sinai J Med.* 2007;74:27–32.

36. Chiang AP, Butte AJ. SYSTEMATIC EVALUATION OF DRUG-DISEASE RELATIONSHIPS TO IDENTIFY LEADS FOR NOVEL DRUG USES. *Clin Pharmacol Ther.* 2009;86:507–10.
37. Bleakley K, Yamanishi Y. Supervised prediction of drug–target interactions using bipartite local models. *Bioinformatics.* 2009;25:2397–403.
38. Nacher JC, Schwartz J-M. A global view of drug-therapy interactions. *BMC Pharmacol.* 2008;8:5.
39. Campillos M, Kuhn M, Gavin A-C, Jensen LJ, Bork P. Drug target identification using side-effect similarity. *Science.* 2008;321:263–6.
40. Hu G, Agarwal P. Human Disease-Drug Network Based on Genomic Expression Profiles. *PLoS One.* 2009;4. doi:10.1371/journal.pone.0006536.
41. Suthram S, Dudley JT, Chiang AP, Chen R, Hastie TJ, Butte AJ. Network-based elucidation of human disease similarities reveals common functional modules enriched for pluripotent drug targets. *PLoS Comput Biol.* 2010;6:e1000662.
42. Mathur S, Dinakarbandian D. Finding disease similarity based on implicit semantic similarity. *J Biomed Inform.* 2012;45:363–71.
43. Gottlieb A, Stein GY, Ruppel E, Sharan R. PREDICT: a method for inferring novel drug indications with application to personalized medicine. *Mol Syst Biol.* 2011;7:496.
44. Daminelli S, Haupt VJ, Reimann M, Schroeder M. Drug repositioning through incomplete bi-cliques in an integrated drug-target-disease network. *Integr Biol (Camb).* 2012;4:778–88.
45. Wang W, Yang S, Zhang X, Li J. Drug repositioning by integrating target information through a heterogeneous network model. *Bioinformatics.* 2014;30:2923–30.
46. Himmelstein D, Lizee A, Hessler C, Brueggeman L, Chen S, Hadley D, et al. Rephetio: Repurposing drugs on a hetnet [report]. Thinklab. 2016. doi:10.15363/thinklab.a7.
47. Sun Y, Barber R, Gupta M, Aggarwal CC, Han J. Co-author Relationship Prediction in Heterogeneous Bibliographic Networks. In: 2011 International Conference on Advances in Social Networks Analysis and Mining. 2011. p. 121–8.
48. Luo Y, Zhao X, Zhou J, Yang J, Zhang Y, Kuang W, et al. A Network Integration Approach for Drug-Target Interaction Prediction and Computational Drug Repositioning from Heterogeneous Information. *bioRxiv.* 2017;:100305.
49. Ojeda T, Murphy SP, Bengfort B, Dasgupta A. *Practical Data Science Cookbook.* Packt Publishing; 2014.
50. Lu M, Zhang Q, Deng M, Miao J, Guo Y, Gao W, et al. An Analysis of Human MicroRNA and Disease Associations. *PLOS ONE.* 2008;3:e3420.
51. Zhou X, Lei L, Liu J, Halu A, Zhang Y, Li B, et al. A Systems Approach to Refine Disease Taxonomy by Integrating Phenotypic and Molecular Networks. *EBioMedicine.* 2018;31:79–91.
52. Gligorićević V, Pržulj N. Methods for biological data integration: perspectives and challenges. *Journal of The Royal Society Interface.* 2015;12:20150571.
53. Fan Y, Li M, Zhang P, Wu J, Di Z. The effect of weight on community structure of networks. *Physica A: Statistical Mechanics and its Applications.* 2007;378:583–90.

54. Zhou T, Ren J, Medo M, Zhang Y-C. Bipartite network projection and personal recommendation. *Phys Rev E Stat Nonlin Soft Matter Phys.* 2007;76 4 Pt 2:046115.
55. Salton G, Lesk ME. Computer Evaluation of Indexing and Text Processing. *J ACM.* 1968;15:8–36.
56. van Driel MA, Bruggeman J, Vriend G, Brunner HG, Leunissen JAM. A text-mining analysis of the human phenome. *Eur J Hum Genet.* 2006;14:535–42.
57. Sun K, Buchan N, Larminie C, Pržulj N. The integrated disease network. *Integr Biol (Camb).* 2014;6:1069–79.
58. Resnik P. Using Information Content to Evaluate Semantic Similarity in a Taxonomy. In: *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 1*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.; 1995. p. 448–453. <http://dl.acm.org/citation.cfm?id=1625855.1625914>. Accessed 3 Sep 2018.
59. Lin D. An Information-Theoretic Definition of Similarity. In: *Proceedings of the 15th International Conference on Machine Learning*. Morgan Kaufmann; 1998. p. 296–304.
60. Wang JZ, Du Z, Payattakool R, Yu PS, Chen C-F. A new method to measure the semantic similarity of GO terms. *Bioinformatics.* 2007;23:1274–81.
61. Cheng L, Li J, Ju P, Peng J, Wang Y. SemFunSim: A New Method for Measuring Disease Similarity by Integrating Semantic and Gene Functional Association. *PLOS ONE.* 2014;9:e99415.
62. Cheng L, Jiang Y, Wang Z, Shi H, Sun J, Yang H, et al. DisSim: an online system for exploring significant similar diseases and exhibiting potential therapeutic drugs. *Scientific Reports.* 2016;6:30024.
63. Hu Y, Zhao L, Liu Z, Ju H, Shi H, Xu P, et al. DisSetSim: an online system for calculating similarity between disease sets. *J Biomed Semantics.* 2017;8 Suppl 1. doi:10.1186/s13326-017-0140-2.
64. Omura M, Tateishi Y, Okumura T. Disease Similarity Calculation on Simplified Disease Knowledge Base for Clinical Decision Support Systems. :6.
65. Dai W, Liu X, Gao Y, Chen L, Song J, Chen D, et al. Matrix Factorization-Based Prediction of Novel Drug Indications by Integrating Genomic Space. *Computational and Mathematical Methods in Medicine.* 2015. doi:10.1155/2015/275045.
66. Zhang W, Yue X, Lin W, Wu W, Liu R, Huang F, et al. Predicting drug-disease associations by using similarity constrained matrix factorization. *BMC Bioinformatics.* 2018;19:233.
67. Chen X, Liu M-X, Yan G-Y. Drug-target interaction prediction by random walk on the heterogeneous network. *Mol Biosyst.* 2012;8:1970–8.
68. Xie M, Hwang T, Kuang R. Prioritizing Disease Genes by Bi-Random Walk. In: Tan P-N, Chawla S, Ho CK, Bailey J, editors. *Advances in Knowledge Discovery and Data Mining*. Springer Berlin Heidelberg; 2012. p. 292–303.
69. Liu Y, Zeng X, He Z, Zou Q. Inferring MicroRNA-Disease Associations by Random Walk on a Heterogeneous Network with Multiple Data Sources. *IEEE/ACM Trans Comput Biol Bioinformatics.* 2017;14:905–915.
70. Yu G, Fu G, Lu C, Ren Y, Wang J. BRWLDA: bi-random walks for predicting lncRNA-disease associations. *Oncotarget.* 2017;8:60429–46.
71. Yang A, Troup M, Ho JWK. Scalability and Validation of Big Data Bioinformatics Software. *Comput Struct Biotechnol J.* 2017;15:379–86.

72. Jodeleit H, Palamides P, Beigel F, Mueller T, Wolf E, Siebeck M, et al. Design and validation of a disease network of inflammatory processes in the NSG-UC mouse model. *Journal of Translational Medicine*. 2017;15:265.
73. Seok J, Warren HS, Cuenca AG, Mindrinos MN, Baker HV, Xu W, et al. Genomic responses in mouse models poorly mimic human inflammatory diseases. *PNAS*. 2013;110:3507–12.
74. Versi E. “Gold standard” is an appropriate term. *BMJ*. 1992;305:187–187.
75. D S, S S, A M. Measurement error correction for logistic regression models with an “alloyed gold standard”. *American Journal of Epidemiology*. 1997;145:184–96.
76. Suratanee A, Plaimas K. Network-based association analysis to infer new disease-gene relationships using large-scale protein interactions. *PLOS ONE*. 2018;13:e0199435.
77. Piñero J, Bravo À, Queralt-Rosinach N, Gutiérrez-Sacristán A, Deu-Pons J, Centeno E, et al. DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants. *Nucleic Acids Res*. 2017;45:D833–9.
78. Paik H, Chen B, Sirota M, Hadley D, Butte AJ. Integrating Clinical Phenotype and Gene Expression Data to Prioritize Novel Drug Uses. *CPT Pharmacometrics Syst Pharmacol*. 2016;5:599–607.
79. Zhang X, Yuan Z, Ji J, Li H, Xue F. Network or regression-based methods for disease discrimination: a comparison study. *BMC Medical Research Methodology*. 2016;16:100.
80. Liu H, Song Y, Guan J, Luo L, Zhuang Z. Inferring new indications for approved drugs via random walk on drug-disease heterogeneous networks. *BMC Bioinformatics*. 2016;17:539.
81. Carson MB, Lu H. Network-based prediction and knowledge mining of disease genes. *BMC Medical Genomics*. 2015;8:S9.
82. Gu C, Liao B, Li X, Li K. Network Consistency Projection for Human miRNA-Disease Associations Inference. *Sci Rep*. 2016;6. doi:10.1038/srep36054.
83. Detector Performance Analysis Using ROC Curves | Receiver Operating Characteristic | Signal To Noise Ratio. Scribd. <https://es.scribd.com/document/339719122/Detector-Performance-Analysis-Using-ROC-Curves>. Accessed 3 Sep 2018.
84. du Prel J-B, Röhrig B, Hommel G, Blettner M. Choosing statistical tests: part 12 of a series on evaluation of scientific publications. *Dtsch Arztebl Int*. 2010;107:343–8.
85. Bustin SA. The reproducibility of biomedical research: Sleepers awake! *Biomolecular Detection and Quantification*. 2014;2:35–42.
86. Kobourov SG. Spring Embedders and Force Directed Graph Drawing Algorithms. arXiv:12013011 [cs]. 2012. <http://arxiv.org/abs/1201.3011>. Accessed 3 Sep 2018.
87. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks. *Genome Res*. 2003;13:2498–504.
88. Le D-H, Pham V-H. HGPEC: a Cytoscape app for prediction of novel disease-gene and disease-disease associations and evidence collection based on a random walk on heterogeneous network. *BMC Systems Biology*. 2017;11:61.
89. Piñero J, Queralt-Rosinach N, Bravo À, Deu-Pons J, Bauer-Mehren A, Baron M, et al. DisGeNET: a discovery platform for the dynamical exploration of human diseases and their genes. *Database (Oxford)*. 2015;2015:bav028.

90. Perfetto L, Briganti L, Calderone A, Cerquone Perpetuini A, Iannuccelli M, Langone F, et al. SIGNOR: a database of causal relationships between biological entities. *Nucleic Acids Res.* 2016;44:D548–54.
91. Lo Surdo P, Calderone A, Iannuccelli M, Licata L, Peluso D, Castagnoli L, et al. DISNOR: a disease network open resource. *Nucleic Acids Res.* 2018;46:D527–34.
92. Calderone A, Cesareni G, Stegle O. SPV: a JavaScript Signaling Pathway Visualizer. *Bioinformatics.* 2018;34:2684–6.
93. Pavlopoulos GA, Paez-Espino D, Kyrpides NC, Iliopoulos I. Empirical Comparison of Visualization Tools for Larger-Scale Network Analysis. *Advances in Bioinformatics.* 2017. doi:10.1155/2017/1278932.
94. Westergaard D, Stærfeldt H-H, Tønsberg C, Jensen LJ, Brunak S. A comprehensive and quantitative comparison of text-mining in 15 million full-text articles versus their corresponding abstracts. *PLOS Computational Biology.* 2018;14:e1005962.
95. Valle EPG del, García GL, Santamaría LP, Zanin M, Ruiz EM, González AR. Evaluating Wikipedia as a Source of Information for Disease Understanding. In: 2018 IEEE 31st International Symposium on Computer-Based Medical Systems (CBMS). 2018. p. 399–404.
96. Grewal N, Singh S, Chand T. Effect of Aggregation Operators on Network-Based Disease Gene Prioritization: A Case Study on Blood Disorders. *IEEE/ACM Transactions on Computational Biology and Bioinformatics.* 2017;14:1276–87.

Supplementary Tables

Supplementary Table 1. Biological data sources

Data Source	Type	URL	Description	Access ⁹	API ¹⁰
NCBI Databases	Various	https://www.ncbi.nlm.nih.gov	The National Center for Biotechnology Information (NCBI) advances science and health by providing access to biomedical and genomic information. Major biological databases include GenBank for DNA sequences, RefSeq for reference sequences and PheGenI for Phenotype-Genotype integration, among others.	Free	Yes ¹¹
KEGG Pathway Database	Various	http://www.genome.jp/kegg	The Kyoto Encyclopedia of Genes and Genome (KEGG) is a database resource for understanding high-level functions and utilities of the biological system, such as the cell, the organism and the ecosystem, from genomic and molecular-level information.	Free	Yes
ArrayExpress	Genomic	https://www.ebi.ac.uk/arrayexpress/browse.html	Functional Genomics Data from high-throughput functional genomics experiments. It is part of the European Bioinformatics Institute ¹² (EMBL-EBI) and the ELIXIR infrastructure ¹³ .	Free	Yes
MetaCyc	Biological Pathways	https://metacyc.org	Curated database of experimentally elucidated metabolic pathways from all domains of life. MetaCyc contains pathways involved in both primary and secondary metabolism. MetaCyc is part of the BioCyc database collection ¹⁴ .	Free ¹⁵	Yes
WikiPathways	Biological Pathways	http://www.wikipathways.org	A database maintained by and for the scientific community dedicated to the curation of biological pathways.	Free (CC)	Yes
Reactome	Biological Pathways	http://reactome.org	An open-source, open access, manually curated and peer-reviewed pathway database.	Free (CC)	Yes
BioGRID	PPI	https://thebiogrid.org	Biological General Repository for Interaction Dataset with data compiled through comprehensive curation efforts.	Free	Yes
STRING	PPI	https://string-db.org	A database of known and predicted protein-protein interactions. STRING is part of the ELIXIR infrastructure.	Free (CC)	Yes

⁹ To resources, online search, downloads and/or API, least for academic purposes. Databases with a Creative Commons License type are marked with CC. For commercial use, please refer to licensing details on the Database URL.

¹⁰ The database provides a consumable API for extensive use, via tools and/or web services. See details on database URL.

¹¹ Entrez: API key required (as of May 1, 2018) to reach a request rate up to 10 per second. See <https://www.ncbi.nlm.nih.gov/books/NBK25500>.

¹² <https://www.ebi.ac.uk/>

¹³ <https://www.elixir-europe.org/>

¹⁴ <https://biocyc.org>

¹⁵ Access via BioCyc web services require paid subscription, but MetaCyc search online and downloads are free access.

UniProt	Protein sequence	http://www.uniprot.org	A database of protein sequence and functional information, many entries being derived from genome sequencing projects. It contains a large amount of information about the biological function of proteins derived from the research literature. It is part of the European Bioinformatics Institute (EMBL-EBI) and the ELIXIR infrastructure.	Free (CC)	Yes
Human Protein Atlas	Protein / Anatomy / Phenotype	http://www.proteinatlas.org	Contains information for a large majority of all human protein-coding genes regarding the expression and localization of the corresponding proteins based on both RNA and protein data. It is part of the ELIXIR infrastructure.	Free (CC)	No ¹⁶
CMAP	Gene expression	https://portals.broadinstitute.org/cmap	The Connectivity Map is a collection of genome-wide transcriptional expression data from cultured human cells treated with bioactive small molecules and simple pattern-matching algorithms that together enable the discovery of functional connections between drugs, genes and diseases through the transitory feature of common gene-expression change	Free	Yes ¹⁷
JASPAR	Gene expression	http://jaspar.genereg.net	The high-quality transcription factor binding profile database (JASPAR) (regulatory).	Free (CC)	Yes
Expression Atlas	Gene expression	https://www.ebi.ac.uk/gxa/home	An open science resource with information about gene and protein expression across species. It is part of the European Bioinformatics Institute (EMBL-EBI) and the ELIXIR infrastructure.	Free (CC)	No ¹⁸
DisGeNET	Gene / Phenotype	http://www.disgenet.org	One of the largest publicly available collections of genes and variants associated to human diseases. It is part of the ELIXIR infrastructure.	Free (CC)	Yes ¹⁹
OMIM	Gene / Phenotype	http://www.omim.org	A comprehensive, authoritative compendium of human genes and genetic phenotypes. OMIM is currently biocurated at the McKusick-Nathans Institute of Genetic Medicine, The Johns Hopkins University School of Medicine.	Free	Yes ²⁰

¹⁶ Data is accessible via downloads. A programmatic access to filter downloadable data is provided. <https://www.proteinatlas.org/about/help/dataaccess>

¹⁷ Clue API Key required for unlimited use. <https://clue.io/api>

¹⁸ Tools for R available, but data must be downloaded first.

¹⁹ API-like access through SPARQL endpoint <http://rdf.disgenet.org/sparql/>

²⁰ Registration is required. <https://omim.org/api>

Supplementary Table 2. Drug data sources

Data Source	Type	URL	Description	Access	API
FDA Databases	Various	https://www.fda.gov/Drugs/InformationOnDrugs	Information about FDA-approved brand name and generic prescription and over-the-counter human drugs and biological therapeutic products.	Free	Yes ²¹
PubChem	Compound	https://pubchem.ncbi.nlm.nih.gov	A database of chemical molecules and their activities against biological assays. The system is maintained by the National Center for Biotechnology Information (NCBI), a component of the National Library of Medicine, which is part of the United States National Institutes of Health (NIH).	Free	Yes ²²
DrugBank	Phenotype	http://www.drugbank.ca	This database combines detailed drug (i.e. chemical, pharmacological and pharmaceutical) data with comprehensive drug target (i.e. sequence, structure, and pathway) information.	Free (CC)	Yes ²³
Orphanet	Phenotype	http://www.orphadata.org	A unique resource, gathering and improving knowledge on rare diseases so as to improve the diagnosis, care and treatment of patients with rare diseases. The Orphanet Rare Disease ontology (ORDO) provides a structured vocabulary for rare diseases capturing relationships between diseases, genes and other relevant features.	Free (CC)	No ²⁴
SIDER	Phenotype	http://sideeffects.embl.de	Contains information on marketed medicines and their recorded adverse drug reactions. The information is extracted from public documents and package inserts. The available information includes side effect frequency, drug and side effect classifications as well as links to further information, for example drug–target relations.	Free (CC)	No ²⁵
PharmaGKB	Gene / Phenotype	http://www.pharmgkb.org	Pharmacogenomics knowledge resource that encompasses clinical information including dosing guidelines and drug labels, potentially clinically actionable gene-drug associations and genotype-phenotype relationships	Free (CC)	Yes ²⁶
ChEMBL	Compound	https://www.ebi.ac.uk/chembl	Database of bioactive drug-like small molecules, it contains 2-D structures, calculated properties and abstracted bioactivities. It is part of the European Bioinformatics Institute (EMBL-EBI) and the ELIXIR infrastructure.	Free	Yes

²¹ Open FDA Access Key is needed to get the maximum rate of requests per minute. <https://open.fda.gov/apis/authentication/>

²² Limited to 5 requests per second. <https://pubchemdocs.ncbi.nlm.nih.gov/pug-rest>

²³ A paid subscription is required to access the API. Data is available through online search and downloads.

²⁴ Downloads available. SPARQL endpoint available to consume the ontology.

²⁵ Data are available via online search and downloads

²⁶ Beta version. Limited to 2 requests per second. <https://api.pharmgkb.org/>

Drug Repurposing Hub	Compound	https://clue.io/repurposing	Contains extensive curated annotations for each drug, including details about commercial sources of all compounds	Free	Yes ²⁷
----------------------	----------	---	---	------	-------------------

²⁷ Clue API Key required for unlimited use. <https://clue.io/api>

Supplementary Table 3. Medical textual sources

Data Source	Type	URL	Description	Access	API
PubMed/PMC	Journals	https://www.nlm.nih.gov/bsd/pmresources.html	PubMed contains more than 28 million medical journal citations from MEDLINE indexed journals and NCI Bookshelf. These citations may have links to full-text articles or manuscripts in PMC (PubMed Central).	Free	Yes ²⁸
MedlinePlus	Medical Encyclopedia	https://medlineplus.gov	Provides encyclopedic information on health and drug issues, and provides a directory of medical services. The service provides curated consumer health information in English and Spanish.	Free	Yes
ClinicalTrials	Clinical Trials	http://www.clinicaltrials.gov	A database of privately and publicly funded clinical studies conducted around the world. It is a resource provided by the U.S. National Library of Medicine.	Free	No ²⁹
Europe PMC	Journals / Clinical Trials	http://europepmc.org	Provides access to worldwide life sciences articles, books, patents and clinical guidelines. Europe PMC provides links to relevant records in databases such as Uniprot, European Nucleotide Archive (ENA), Protein Data Bank Europe (PDBE) and BioStodie. It is part of the ELIXIR infrastructure.	Free	Yes
GWAS Catalog	GWAS Studies	https://www.ebi.ac.uk/gwas/	The NHGRI-EBI Catalog of published genome-wide association studies. Eligible studies are curated within 1-2 months of publication, dependent on the availability of literature, and the data is released on a weekly cycle	Free	Yes

²⁸ Entrez: API key required (as of May 1, 2018) to reach a request rate up to 10 per second. See <https://www.ncbi.nlm.nih.gov/books/NBK25500>.

²⁹ Not strictly an API. The combination of Advanced Search query building and results download provides intensive data access.

Table 4. Mapping sources

Data Source	Type	URL	Description	Access	API
MeSH	Medical Thesaurus	https://www.ncbi.nlm.nih.gov/mesh	Medical Subject Headings (MeSH) is the NLM controlled vocabulary thesaurus used for indexing articles for PubMed. It is also used by ClinicalTrials.gov registry to classify which diseases are studied by trials.	Free	Yes ³⁰
SNOMED CT	Medical Thesaurus	http://www.snomed.org/snomed-ct	Systematized Nomenclature of Medicine Clinical Terms (SNOMED CT) is a comprehensive and precise clinical health terminology product.	Licensed ³¹	Yes
UMLS	Medical Thesaurus	https://www.nlm.nih.gov/research/umls/	Unified Medical Language System (UMLS) integrates and distributes key terminology, classification and coding standards, and associated resources to promote creation of more effective and interoperable biomedical information systems and services, including electronic health records	Licensed ³²	Yes
ICD-CM	Disease Codes	http://www.who.int/classifications/icd	International Classification of Diseases (ICD) is the classification used to code and classify mortality data from death certificates. ICD-Clinical Modification (CM) is used to code and classify morbidity data from the inpatient and outpatient records, physician offices, and most National Center for Health Statistics (NCHS) surveys.	Free	Yes ³³
HGNC	Gene Codes	https://www.genenames.org/	HGNC is responsible for approving unique symbols and names for human loci, including protein coding genes, ncRNA genes and pseudogenes, to allow unambiguous scientific communication.	Free	Yes
DO	Disease Ontology	https://bioportal.bioontology.org/ontologies/DOID	The Disease Ontology (DO) provides the biomedical community with consistent, reusable and sustainable descriptions of human disease terms, phenotype characteristics and related medical vocabulary disease concepts. The DO semantically integrates disease and medical vocabularies through extensive cross mapping of DO terms to MeSH, ICD, NCI's thesaurus, SNOMED and OMIM.	Free (CC)	Yes
PO	Phenotype Ontology	https://hpo.jax.org/app/	The Human Phenotype Ontology (HPO) provides a standardized vocabulary of phenotypic abnormalities encountered in human disease. The HPO is currently being developed using the medical literature, Orphanet, DECIPHER, and OMIM.	Free	Yes ³⁴
GO	Gene Ontology	http://www.geneontology.org	The Gene Ontology (GO) provides a computational representation of our evolving	Free (CC)	Yes

³⁰ API-like access through the SPARQL endpoint. <https://hhs.github.io/meshrdf/sparql-and-uri-requests>

³¹ Free license for member countries. See other fee exemptions: <https://www.snomed.org/snomed-ct/get-snomed-ct>

³² Free license available under certain conditions. <https://uts.nlm.nih.gov/license.html>

³³ Provided by NLM Clinical Tables Search Service: <https://clinicaltables.nlm.nih.gov/apidoc/icd10cm/v3/doc.html>

³⁴ SPARQL endpoint available at http://www.orphadata.org/cgi-bin/inc/sparql_hoom.inc.php

		logy.org	knowledge of how genes encode biological functions at the molecular, cellular and tissue system levels. GO terms are organised in three domains: cellular component, molecular function and biological process.		
Uberon	Anatomy Ontology	http://uberon.github.io/	An integrated cross-species ontology covering anatomical structures in animals. See the about page for more info, or read the Uberon paper in Genome Biology	Free	Yes ³⁵

³⁵ SPARQL endpoint available: <http://sparql.hegroup.org/sparql>

Supplementary Table 5. Studies on disease networks

Author. Study (year)	Addressed problem	Sources	Methods	Validation	Network facts	Access/ Visualization
Mathur et al. Finding disease similarity based on implicit semantic similarity (2012)	Disease similarity	OMIM, SwissProt, GeneRif	Similarity of disease gene and processes (Jaccard based)	Comparison with previous studies using KEGG data	1,477 disease-disease linked by GO process	Validation data set and results available as supplementary material ³⁶ .
Zhou et al. Human symptoms–disease network (2014)	Disease similarity	PubMed, OMIM, PharmaGKB, BioGrid	Similarity of disease symptoms (Cosine based)	Manual validation. Overlapping of predicted similarities with HPO	HSDN with 147,978 connections between 4,219 diseases and 322 symptoms.	Results available as supplementary material in the article ³⁷ . Network analysis and visualization with Gephi.
Sun et al. Predicting disease associations via biological network analysis (2014)	Disease similarity	BioGrid, OMIM, CTD, HuGeNet	Similarity of disease genes and PPIs (Jaccard based). Disease gene topology.	Correlation of similarity scores with ICD-9 groups, comorbidity from HuDiNet and GWAS associations. Manual confirmation with medical literature.	Disease-disease associations predicted for 543 diseases.	
Hoehndorf et al. Analysis of the human diseaseome using phenotype similarity between common, genetic, and infectious diseases (2015)	Disease similarity and drug repurposing.	PubMed	Similarity of disease phenotypes (Jaccard based).	Correlation of predicted disease clusters and top-level DO groups. Comparison against models of the Mouse Genome Informatics database and gene-disease associations from OMIM. Drug-disease associations validated with data from SIDER.	5,030 disease nodes and 65,795 disease-disease association edges.	Results and visualization available online ³⁸ .
Garcia-Albornoz et al. Finding directionality and gene-disease predictions in disease associations (2015)	Disease similarity	OMIM, KEGG	Mapping of disease codes to OMIM and KEGG data to obtain directional disease-gene-pathway associations using the	Correlation of predicted associations with ICD-10 categories.	880 diseases with 3,430 disease-disease gene-based associations and 112,956 disease-disease pathway-based associations.	Top-rated gene disease pairs available as additional info in the article.

³⁶ <https://ars.els-cdn.com/content/image/1-s2.0-S1532046411002073-mmcl.doc>

³⁷ <https://www.nature.com/articles/ncomms5212#s1>

³⁸ <http://aber-owl.net/aber-owl/diseasephenotypes>

			level-of-inclusion.			Cytoscape used for network analysis ³⁹ .
Paik et al. Integrating Clinical Phenotype and Gene Expression Data to Prioritize Novel Drug Uses (2016)	Drug repurposing	SIDER2,, DrugBank, CMap	Similarity of drug-phenotype indications and drug-gene expressions and (Cosine based)	Comparison of results with previous studies.	Two bipartite networks: Drug-phenotype network with 1,631 drug nodes, 1,587 phenotype nodes and 72,848 edges. Drug-gene expression network with 756 drugs, 8,101 gene nodes and 17,000 edges.	Results available as supporting information of the article ⁴⁰ .
Himmelstein et al. Systematic integration of biomedical knowledge prioritizes drugs for repurposing (2016)	Drug repurposing	DrugBank, SIDER, DrugCentral, Gene, WikiPathways, Reactome, UniChem, ChEMBL, PharmacotherapyD B, DisGeNET, PubMed, GWAS Catalog, among others	Metapath based approach with Degree-Weighted Path Count for metapath weighting. Machine learning techniques to translate the network connectivity between a compound and a disease into a probability of treatment.	Comparison of drug predictions with new indications extracted from DrugCentral and ClinicalTrials.	Metagraph with 47,031 metanodes of 11 types and 2,250,197 metaedges of 24 types. A derived drug-disease bipartite network with 136 disease nodes, 1,538 drugs and 209,168 edges.	Network accessible online ⁴¹ . Neo4j used for network analysis and visualization.
Guney et al. Network-based in silico drug efficacy screening (2016)	Drug repurposing	DrugBank, OMIM, GWAS, UniProtKB, PheneGenI, Orphanet	Mapping of drugs and diseases using associated genes. Drug-disease proximity based on shortest-path.	Comparison with drug-disease associations on DailyMed.	Bipartite disease-drug network with 402 drug-disease associations between 238 drugs and 78 diseases.	Results available as supplementary information ⁴² .
Luo et al. A network integration approach for drug-target interaction prediction and computational drug	Drug repurposing	DrugBank, HPRD, Comparative Toxicogenomics, SIDER	Low-dimensional vector representation of features and vector space projection to predict drug-disease associations.	Comparison with known drug-disease associations and in-vivo validation of some of the obtained predictions.	Heterogeneous network with 12,015 nodes of 4 types and 1,895,445 edges of six types.	Data and source code available online ⁴³ .

³⁹ https://static-content.springer.com/esm/art%3A10.1186%2Fs12918-015-0184-9/MediaObjects/12918_2015_184_MOESM1_ESM.xlsx

⁴⁰ <https://ascpt.onlinelibrary.wiley.com/doi/full/10.1002/psp4.12108>

⁴¹ <http://het.io>

⁴² <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4740350/#S1>

⁴³ <https://github.com/luoyunan/DTINet>

repositioning from heterogeneous information (2017)						
Luo et al. Computational Drug Repositioning using Low-Rank Matrix Approximation and Randomized Algorithms (2018)	Drug repurposing	DrugBank, OMIM	Similarity of drug-chemical structure and disease-phenotypes. Singular Value Thresholding (SVT) algorithm to complete the drug-disease adjacency matrix with predicted scores for unknown drug-disease pairs.	Comparison of drug predictions with previous studies, mainly Gottlieb et al. (2011).	Bipartite disease-drug network with 593 drugs, 313 diseases and 1,933 interactions.	Drug repositioning recommendation system and data available online ⁴⁴ .

⁴⁴ <http://bioinformatics.csu.edu.cn/resources/softs/DrugRepositioning/DRRS/index.html>