

## REVIEW

# Integrative methods for analyzing big data in precision medicine

Vladimir Gligorijević, Noël Malod-Dognin and Nataša Pržulj

Department of Computing, Imperial College London, London, UK

We provide an overview of recent developments in big data analyses in the context of precision medicine and health informatics. With the advance in technologies capturing molecular and medical data, we entered the area of “Big Data” in biology and medicine. These data offer many opportunities to advance precision medicine. We outline key challenges in precision medicine and present recent advances in data integration-based methods to uncover personalized information from big data produced by various omics studies. We survey recent integrative methods for disease subtyping, biomarkers discovery, and drug repurposing, and list the tools that are available to domain scientists. Given the ever-growing nature of these big data, we highlight key issues that big data integration methods will face.

Received: October 8, 2015

Revised: November 16, 2015

Accepted: December 9, 2015

**Keywords:**

Big data / Bioinformatics / Integration methods / Personalized medicine

## 1 Introduction

*Precision medicine*, also known as personalized, predictive, preventive, and participatory (P4) medicine [1], is an emerging approach for individualizing the practice of medicine [2]. Prevention and treatment strategies that take into account individual variability are not new; for example, blood-typing has been used to guide blood transfusion for more than a century, with a total of 35 human blood groups being recognized by the International Society of Blood Transfusion [3]. Similarly, gender, race, time of ischemia, cytomegalovirus, and sero-type are taken into account to reduce the risk of rejecting organ transplantations [4–7]. The challenge in applying the precision medicine concept to omics and clinical datasets of patient features that have become available and that cannot

be interpreted directly by medical practitioners due to their large sizes and complexities.

Big data is a broad term for datasets so large or complex that traditional data processing methods are inadequate. It is often characterized by three Vs [8]: volume, which refers to the large size of the data; velocity, which refers to the high speed at which data are generated; and variety, which refers to the heterogeneity of the data coming from different sources. All these characteristics apply to currently available biological and medical datasets. Since the beginning of the Human Genome Project [9], novel technological developments led to the era of omics sciences. Using novel high-throughput capturing technologies, we are now able to access the DNA of an individual (genetic data), the transcribed RNA over time (expression and coexpression data), proteins (protein profiles and protein interaction data), metabolism (metabolic profiles), and epigenome (DNA methylation data), among other data types [10]. The environment is also taken into account (e.g. nutrition and bacterial environment by nutriomics and metagenomics, respectively) [11, 12], and also histopathological and medical imaging data are now subject to high-throughput capturing and analysis methods [13–16].

Therefore, we are facing an increasing gap between our ability to generate big biomedical data and our ability to analyze and interpret them [17]. In this context, it is not surprising

**Correspondence:** Dr. Nataša Pržulj, Department of Computing, Imperial College London, London SW7 2AZ, UK

**E-mail:** natasha@imperial.ac.uk

**Abbreviations:** EHR, electronic health records; GNMTF, Graph-regularized non-negative matrix tri-factorization; JIVE, joint and individual variation explained; KB, kernel-based; lncRNA, long noncoding RNA; MCMC, Markov chain Monte Carlo; MSCMF, multiple similarities collaborative matrix factorization; NB, network-based; NBS, NB stratification; NGS, next-generation sequencer; PREDICT, PREDicting Drug IndiCaTions; SNF, similarity network fusion; SVM, support vector machines; TDA, topological data analysis method

\*Both authors contributed equally.

**Colour Online:** See the article online to view Figs. 1–3 in colour.

that big data and precision medicine are jointly investigated. In 2011, the “Big Data Research and Development Initiative” (<https://www.whitehouse.gov/blog/2012/03/29/big-data-big-deal>) was targeting personalized medicine through the GenI-SIS program (Genomic Information System for Integrated Science) to enhance health care for Veterans. In 2012, the US National Institutes of Health (NIH) launched the “Big Data to Knowledge” initiative, to harvest the wealth of information contained in biomedical Big Data [18]. Finally, President Obama recently announced the “Precision Medicine” initiative (<https://www.whitehouse.gov/precision-medicine>), with an ambitious goal of driving precision medicine by incorporating many different types of data, from genomes to microbiomes, with patient data collected by health care providers and patients themselves.

Out of many challenges in precision medicine, here we focus on four related problems: patient subtyping, biomarker discovery, drug repurposing, and personalized treatment prediction. We provide a review of methods capable of integrative analyses of multiple data types in addressing these problems.

### 1.1 Subtyping and biomarker discovery

Also known as patient stratification, subtyping is the task of identifying subpopulations of patients that can be used to guide treatment procedures of a given individual belonging to the subpopulation, and to predict the outcomes. Subtyping identifies *endotypes*, which refer to subtypes in which patients are related by similarities in their underlying disease mechanisms (i.e. to explain the diseases mechanisms) [19], and *verotypes*, which refer to true populations of similar patients for treatment purposes (i.e. to predict therapies for curing the patients) [20]. However, what precisely constitutes endotypes and verotypes, as well as how they should be discovered, remains open. Despite varying definitions, subtyping remains a classification task and an active and growing area of machine learning (ML) research (see Section 3.1). Diseases such as cancer, autism, autoimmune diseases, cardiovascular diseases, and Parkinson's have all been studied through the lens of subtyping [21–23].

According to FDA, a biomarker is any measurable diagnostic indicator that is used to assess the risk, or presence of a disease [24]. Biomarker discovery aims at finding features that are characteristic to particular patient subpopulations (e.g. specific gene mutations in tumor tissues, specific miRNAs, metabolites, etc.). The goal is that an individual is only tested for biomarkers to decide whether or not she/he belongs to a specific patient subtype. Biomarkers are considered key to improving healthcare and lowering medical costs [25].

### 1.2 Drug repurposing and personalized treatment

Drug repurposing refers to the identification and development of new uses for the existing or abandoned pharmacotherapies. Capitalizing on already known drugs al-

lows for reducing the cost of developing pharmacotherapies compared with de novo drug discovery and development [26]. With the availability of various omics data, computational predictions of new drug candidates for repurposing have necessitated the development of many new methods for data integration (see Section 3.2).

Drug repurposing is not only about identifying new targets for known drugs; preclinical evaluations also include predicting therapeutic regimens (i.e. dose and frequency) and safety of the treatment (i.e. side effects). Bringing together patient subtyping and precise prediction of therapeutic treatment outcomes is the key for deriving personalized treatments. For example, the American Society of Clinical Oncology estimates that testing colon cancer patients for mutations in K-RAS gene would save \$604 million in drug costs annually; since patients with these mutations do not respond well to EGF inhibitors, it is preferable to avoid giving them an inefficient and potentially toxic treatment, which is also very expensive (\$100 000 per treatment) (<http://www.asco.org/press-center/advances-treatment-gastrointestinal-cancers-0>).

In this paper, we give an overview of the available methods for analyzing large and diverse biomedical data, introduce concepts of data integration and classification, and elaborate on the successes and limitations of Big-Data approaches in precision medicine.

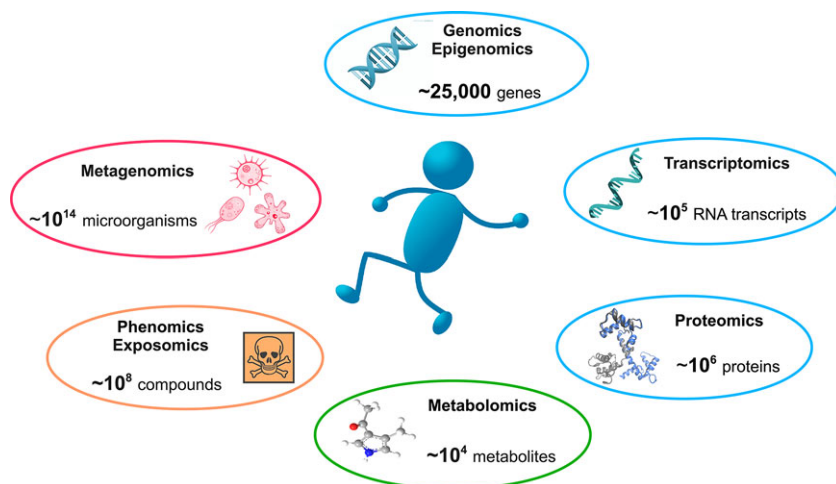
## 2 Big data

### 2.1 Avalanche of omics data

With the recent advances in biomedical data-capturing technologies, omics sciences produce ever increasing amounts of biomedical data. We briefly present key available omics data types, which are illustrated in Fig. 1.

#### 2.1.1 Genomics and exomics

Genomics is a part of genetics that focuses on capturing whole genomes. Historically, the Human Genome Project required 12 years and \$3 billion to capture the first human genome, with a final release in 2003 reporting about 20 500 genes [9]. The first commercial next-generation sequencer (NGS), the Roche GS-FLX 454 (released in 2004), allowed capturing the second human genome in 2 months [27]. In comparison, a modern NGS such as the Illumina HiSeq X is capable of producing up to 16 human genomes worth of data per 3-day run. Note that only 1–2% of a human's genetic material codes for genes, in DNA regions called exons. Exomics, which focuses on these smaller regions, leads to quicker and cheaper sequencing [28, 29]. Recently, the ability to perform sequencing of individual cells has provided novel insights into human biology and diseases [30, 31]. Heterogeneity in DNA sequence from one cell to another has unveiled the concept of *mosaicism*, i.e. the presence of two or more populations of cells



**Figure 1.** Illustration of various omics data types.

with different genotypes in one individual [32]. Cancer in particular has been studied through the lens of genomic variation to find driver mutations.

### 2.1.2 Epigenomics

Epigenomics is the study of the complete set of epigenetic modifications of the genetic material of a cell. These reversible modifications on DNA or histones affect gene expression and thus play a major role in gene regulation. High-throughput methods, such as ChipSeq and Bisulfite sequencing, allow for detection of epigenetic modifications, such as DNA methylation, histone modification, and chromatin structure [33, 34]. Epigenomics findings are cell type-specific and epigenetic reprogramming has a clear role in cancer [35, 36].

### 2.1.3 Transcriptomics

As opposed to DNA sequence, which is relatively static [37], RNA reflects the dynamic state of a cell. Transcriptomics aims at measuring the amount of transcribed genetic material over time. It includes both coding and noncoding RNAs, whose functions are sometimes unknown [38]. Coexpressed genes (i.e. with similar expression patterns over time) have been shown to be likely regulated via the same mechanisms [39] and differential expression patterns are used to identify dysregulated genes in cancer [40], predict possible drug-targets [41] and cancer outcomes [42].

### 2.1.4 Proteomics and interactomics

While transcriptomics considers all transcribed RNAs, proteomics focuses on the produced proteins, after all posttranslational sequence modifications (e.g. phosphorylation, glycosylation, and lipidation). The human proteome is several orders of magnitude larger than the human genome; because of alternative promoters, alternative splicing, and mRNA edit-

ing, the  $\approx 25\,000$  human genes lead to  $\approx 100\,000$  transcripts; with more than 300 different types of posttranslational modifications, the number of resulting proteins is estimated to be larger than 1 800 000 [43]. High-throughput capture of protein sequences is done via MS experiments [44]. Interactions amongst proteins, or between proteins and other molecules, are captured with high-throughput techniques, such as yeast-two-hybrid [45] and affinity-captured coupled with MS [46]. Interactomes and protein–protein interactions in particular, were successfully used to identify evolutionarily conserved pathways, complexes, and functional orthologs [47–49].

### 2.1.5 Metabolomics, glycomics, and fluxomics

A metabolite is any substance produced or consumed during metabolism (all chemical processes in a cell). Metabolomics studies all chemical processes involving metabolites [50]. Metabolic profiles are measured with MS and nuclear magnetic resonance spectrometry. Glycomics is the branch of metabolomics that studies glycomes, the sets of all sugars-free or in more complex molecules such as glycoproteins—in cells. Glycosylation is the most intensive and complex posttranslational modification of proteins and glycans are known to be involved in cell growth and development [51], in the immune system [52], in cell-to-cell communication [53], in cancer, and microbial diseases [54, 55]. Fluxomics refers to a range of methods in experimental and computational biology that attempt to identify, or predict the rates of metabolic reactions in biological systems [56].

### 2.1.6 Phenomics and exposomics

Phenomics is an area of biology measuring phenomes—physical and biochemical traits of organisms—as they change in response to genetic mutation and environmental influences. Genome-wide association studies are commonly used for detecting associations between single-nucleotide polymorphisms and common diseases such as heart disease,

diabetes, autoimmune diseases, and psychiatric disorders [57]. Exposomics encompasses all human environmental (i.e. nongenetic) exposures from conception onward. It includes, amongst others, exposure to toxic molecules, drugs, and radiation. Exposomics benefits from continuous tracking that is now available for most of the key physiological metrics (blood pressure, heart rhythm, brain waves, etc.) and environmental indices, such as air pollution, pollen count, and radiation. Even medical imaging, which was traditionally manually investigated, is now a subject of high-throughput capturing [14, 15]. For example, radiomics (the high-throughput capturing and analysis of medical radio images) recently lead to connectomics, which captures and analyses brain connectivity maps.

### 2.1.7 Metagenomics

Metagenomics aims at capturing human microbiomes, usually through 16S rRNA sequencing. Our bacterial flora has been shown to play an important role in various medical conditions [12]; for example, the bacterial flora of the intestine is known to modulate the effects of drugs involved in cancer treatments [58]. However, taking into account microbiota is challenging, as human microbiome consists of circa 100 trillion microbial cells, which is about ten times the number of human cells [59].

## 2.2 Biomedical data gets more complex

The complexity of biomedical data grows in two directions: in terms of the number of samples and in terms of heterogeneity.

### 2.2.1 The growing number of samples

As capturing technologies are becoming faster and cheaper, the number of individuals for whom data are available is quickly increasing. For example, the number of available human genomes/exomes increased almost exponentially during the last decade: the first human exome was released in 2003 [9], while in 2012, 1092 human genomes were available [60]. Nowadays, the Exome Aggregation Consortium contains 60 706 unrelated human exomes (<http://exact.broadinstitute.org>). The UK government recently announced the project to map 100 000 human genomes by 2017 (<https://www.gov.uk/government/news/human-genome-uk-to-become-world-number-1-in-dna-testing>) and the precision medicine initiative in the United States plans to map 1 million human genomes. Note that this increasing number of genome samples will also come at the price of increasing variations in terms of genome quality. NGSs produce short reads that need to be assembled into genomes. The quality of the assembled genome highly depends on the ratio

between the sum of the short read lengths and of the target genomic sequence length. This ratio is called the depth of the sequencing and it is expressed in terms of X (e.g. 2X sequencing means that on average each nucleotide is covered by two short reads). While current sequencing uses  $\approx 30X$ , a recent study argues that high-quality genomes may require  $\approx 126X$  (refereed as deep sequencing) [61].

Moreover, for the same individual, an increasing number of samples is captured; data can be collected over different tissues, by using single-cell genomics [62], or on different conditions (e.g. before and after treatment). Finally, the time span of available samples is increasing. For example, gene expression can be measured over time to assess the effect of drugs. Recent developments of noninvasive capturing techniques (e.g. fetal exome sequencing from maternal blood [63] and magnetic resonance imaging, capturing brain connectivity maps from unborn babies to adults; Developing Human Connectome Project, <http://www.developingconnectome.org/project/>) will allow collecting information over the whole life span of an individual, which paves the way to personalized medicine from womb to tomb.

### 2.2.2 Increasing heterogeneity of captured data

The number of different biological entities (e.g. genes, RNAs, proteins, metabolites, drugs, diseases, etc.) for which data can be collected is increasing. The variety of available data is illustrated in Table 1, which presents some of the well-established large-scale biomedical databases. The collected data are so large that even basic data management is becoming challenging. US healthcare was already storing 150 exabytes ( $10^{18}$  Bytes) of data in 2011 and is expected to handle yottabytes of data ( $10^{24}$  Bytes) in the next few years (Institute for Health Technology Transformation, <http://ihealthtran.com/big-data-in-healthcare>). These datasets are highly heterogeneous; data from the same type can be captured with different technologies having varying coverage, bias, and noise robustness (e.g. the different technologies for capturing protein–protein interactions [64]), and the same applies across data types. Moreover, the large number of data sources poses data collection issues coming from the lack of standard format in data repositories (so-called data-extraction problem in Big Data [65]).

## 3 ML techniques

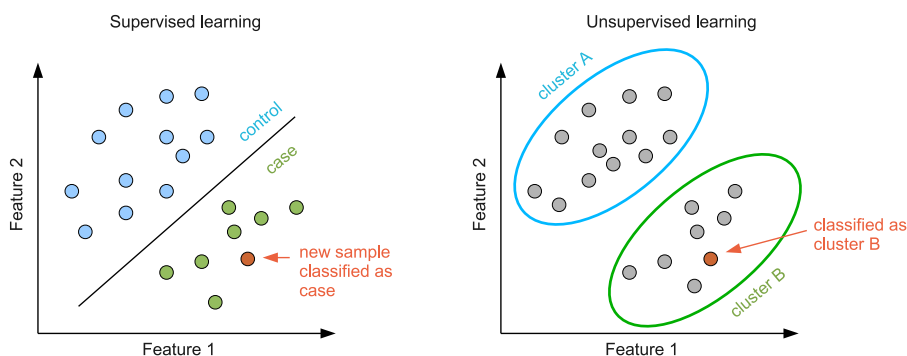
As described in the previous section, Big Data are of large-scale, diversity and complexity, and as such they require efficient algorithms for extracting knowledge hidden in them. Computational techniques that are used to analyze Big Data are either based on statistical, ML, or network-based (NB) methods [104]. These methods have already demonstrated great potential in bridging the gap between production and interpretation of big data in precision medicine, but there is still a lot of room for their improvements.

**Table 1.** Available data for human

	Database	Link	Content
Genomic	NCBI gene [66] GOA [67] ENCODE [68]	<a href="http://www.ncbi.nlm.nih.gov/gene">www.ncbi.nlm.nih.gov/gene</a> <a href="http://www.ebi.ac.uk/GOA">www.ebi.ac.uk/GOA</a> <a href="http://www.encodeproject.org">www.encodeproject.org</a>	Atlas of 59 500 human genes 487 409 gene ontology annotations for 48 569 human gene products Functional annotations of coding/noncoding DNA elements
Epigenomic	NCBI epigenomics [69] 4DGenome [70] HEA MethylomeDB [71]	<a href="http://www.ncbi.nlm.nih.gov/epigenomics">www.ncbi.nlm.nih.gov/epigenomics</a> <a href="http://4dgenome.int-med.uiowa.edu/">4dgenome.int-med.uiowa.edu/</a> <a href="http://www.genboree.org/epigenomeatlas">www.genboree.org/epigenomeatlas</a> <a href="http://www.neuroepigenomics.org/methylomedb">www.neuroepigenomics.org/methylomedb</a>	5110 epigenetic modifications 3 095 881 experimental and predicted chromatin interactions Atlas of reference epigenomes DNA methylomes of human brain cells
Transcriptomic	NCBI GEO [72] Expression atlas [73] CMAP [74] COXPRESdb [75] GeneFriends [76]	<a href="http://www.ncbi.nlm.nih.gov/geo">www.ncbi.nlm.nih.gov/geo</a> <a href="http://www.ebi.ac.uk/gxa">www.ebi.ac.uk/gxa</a> <a href="http://www.broadinstitute.org/cmap">www.broadinstitute.org/cmap</a> <a href="http://coxpresdb.jp">coxpresdb.jp</a> <a href="http://genefriends.org">genefriends.org</a>	1912 human gene expression datasets Differential and baseline gene expression data ~7000 expression profiles for 1309 perturbation compounds Coexpression of 19 803 human genes Coexpression of 159 184 human genes and transcripts
Proteomic	UniProt [77] NextProt [78] RCSB PDB [79] HPA [80] IntAct [81] BioGrid [82] I2D [83] STRING [84]	<a href="http://www.uniprot.org">www.uniprot.org</a> <a href="http://www.nextprot.org">www.nextprot.org</a> <a href="http://www.rcsb.org/pdb">www.rcsb.org/pdb</a> <a href="http://www.thehppp.org">www.thehppp.org</a> <a href="http://www.ebi.ac.uk/intact">www.ebi.ac.uk/intact</a> <a href="http://thebiogrid.org">thebiogrid.org</a> <a href="http://ophid.utoronto.ca">ophid.utoronto.ca</a> <a href="http://string-db.org">string-db.org</a>	Information about human proteome (69 693 proteins) Knowledgebase on 20 066 human proteins Portal to 113 494 biological macromolecular 3D structures Maps of human proteome on 44 normal and 20 cancer type tissues 209 852 human protein–protein interactions 215 952 human protein–protein interactions 183 524 (+55 985 predicted) protein–protein interactions 8 548 005 interactions between 20 457 proteins
Metabolomic	HMDB [85] KEGG pathway [86] SMPD [87] Reactome [88] SugarBindDB [89] UniCarbKB [90] KEGG glycan [91]	<a href="http://www.hmdb.ca">www.hmdb.ca</a> <a href="http://www.genome.jp/kegg/pathway">www.genome.jp/kegg/pathway</a> <a href="http://www.smpdb.ca">www.smpdb.ca</a> <a href="http://www.reactome.org">www.reactome.org</a> <a href="http://sugarbind.expasy.org">sugarbind.expasy.org</a> <a href="http://www.unicarbkb.org">www.unicarbkb.org</a> <a href="http://www.genome.jp/kegg/glycan/">www.genome.jp/kegg/glycan/</a>	Atlas of 41 993 human metabolites 298 human pathways ~700 human metabolic and disease pathways 8770 reactions in 1887 human pathways 1256 interactions between 200 glycans and 551 pathogenic agents 3740 glycan structure entries and 400 glycoproteins Glycan metabolic pathways
Phenomic, exposomic, metagenomic	OMIM [92] NCBI dbGaP [93] GWAS catalog [94] COSMIC [95] TCGA [96] DrugBank [97] PubChem [98] T3DB [99] FoodB [100] UMCD [101] HCP [102] HMP [103]	<a href="http://www.omim.org">www.omim.org</a> <a href="http://www.ncbi.nlm.nih.gov/gap">www.ncbi.nlm.nih.gov/gap</a> <a href="http://www.ebi.ac.uk/gwas/">www.ebi.ac.uk/gwas/</a> <a href="http://cancer.sanger.ac.uk/cosmic">cancer.sanger.ac.uk/cosmic</a> <a href="http://cancergenome.nih.gov">cancergenome.nih.gov</a> <a href="http://www.drugbank.ca">www.drugbank.ca</a> <a href="http://pubchem.ncbi.nlm.nih.gov/">pubchem.ncbi.nlm.nih.gov/</a> <a href="http://www.t3db.org">www.t3db.org</a> <a href="http://www.foodb.ca">www.foodb.ca</a> <a href="http://umcd.humanconnectomeproject.org">umcd.humanconnectomeproject.org</a> <a href="http://www.humanconnectome.org/data">www.humanconnectome.org/data</a> <a href="http://hmpdacc.org">hmpdacc.org</a>	Catalog of mendelian disorders and over 15 000 genes Database of genotypes and phenotypes Genome-wide association studies, assaying ~100 000 SNPs Somatic mutations in cancer, with 3 480 051 coding mutations Somatic mutations and expression data for ~7000 human tumors ~1600 approved/illicit/experimental drugs with known gene targets ~2 × 10 <sup>8</sup> compounds and substances, with 57 335 gene targets ~3600 common toxins and environmental pollutants ~28 000 food components/additives, with presumptive health effects 1887 brain connectivity matrices from neuroimaging data MRI captured brain connectivity maps of 500 adult individuals 11 000 samples of human microbiomes from 300 adult individuals

SNP: single-nucleotide polymorphism.





**Figure 2.** A schematic illustration of the two main learning techniques in ML—supervised (left panel) and unsupervised (right panel) learning. Left: in supervised learning, a training dataset consists of samples with known class labels, e.g. cases and controls. A model is learned by maximizing the difference between cases and controls and then a label for a new sample is determined. Right: in unsupervised learning, all samples are unlabeled. A model clusters samples into different groups based on their similarity.

ML methods came into focus of Big-Data analysis due to their prominent ability to *collectively mine (integrate)* large-scale, diverse, and heterogeneous biomedical data types, a foremost challenge in precision medicine and medical informatics [105]. Thus, in this section, we mostly focus on ML methods for data integration, but we also mention some recent statistical and NB methods for data integration.

ML methods can be divided into the following classes (see Fig. 2 for an illustration):

*Supervised methods*, such as classification and regression, take as input training data samples with known labels. A model is learned through a training process that maximizes the accuracy of its performance on the training dataset. The model is then used for mapping new data samples to existing labels. For example, an input data can comprise patients classified as cases and controls. A model is learned to maximize the difference between cases and controls and then it is applied in classification of new patients. Some of the widely used supervised techniques include Support Vector Machines (SVMs) [106], kernel-based (KB) methods [107], and Logistic regression [108].

*Unsupervised methods*, such as clustering and dimensionality reduction, take as input unlabeled dataset. A model is learned by revealing hidden patterns in the data and organizing the data into meaningful subsets. These methods are often used in molecular subtyping of cancer patients, or in discovering of patterns in gene expression data. Some of the widely used unsupervised methods in precision medicine include hierarchical clustering [109], K-means [109], and its generalizations including matrix factorization methods [110].

*Semisupervised methods* take as input a mixture of labeled and unlabeled samples. A model is learned to explain the structure in the data as well as to make new predictions of unlabeled samples. For example, in predicting new drug-disease associations, semisupervised methods learn known drug-disease associations from labeled samples (i.e. prior knowledge), to predict novel drug-disease associations. This strategy is particularly suitable for data integration, as is can incorporate various data types as prior knowledge. One of the most widely used such method is network-regularized matrix factorization [111].

Based on the type of data they integrate, the integration methods can be divided into *homogeneous*, where the same

type of data, but across multiple perspectives (e.g. experimental studies) is integrated, and *heterogeneous*, where multiple data types in different formats are integrated. The latter is computationally more challenging, because it requires a framework that can deal with heterogeneous data without transforming it and losing any information through the transformation. A majority of the existing frameworks cannot cope with this issue and they require a preprocessing step prior to integration, where they transform the data into a single representation. In Section 3.2, we discuss this issue in more detail and identify methods capable of addressing this problem.

We survey recent integrative methods for disease subtyping, biomarker discovery, and drug repurposing, and provide a summary listing computational tools that can be used by domain scientists for analysis of Big Data (see Table 2 for the list of methods). The presented methods are chosen based on the following criteria: (i) the method is integrative (i.e. it considers more than one data type) and is applied on biomedical Big Data; (ii) the method is predominantly based on ML techniques, although we also consider couple of NB methods; and (iii) the method has been used to address one of the four different precision medicine challenges (see Section 1).

### 3.1 Computational methods for disease subtyping and biomarker discovery

Disease subtyping (or disease stratification) is a task of grouping patients into subgroups based on genomic, transcriptomic, epigenomic, and clinical data. The main goal of subtyping is achieving more accurate prognoses of individuals' expected outcomes that can be used to improve treatment decisions. Treatments of many diseases have benefited from subtyping, including Parkinson's, cardiovascular, autoimmune diseases, and cancer [112].

Cancer is one of the most studied diseases by subtyping. It is a disease in which genome aberrations are accumulating and eventually leading to dysregulation of the cellular system. Histologically similar cancers are composed of many molecular subtypes with significantly different clinical behaviors and molecular complexity at the genomic, epigenomic, transcriptomic, and proteomic levels. Many subtypes have been identified by utilizing techniques for data integration

**Table 2.** Summary of methods for integrative analyzes in precision medicine

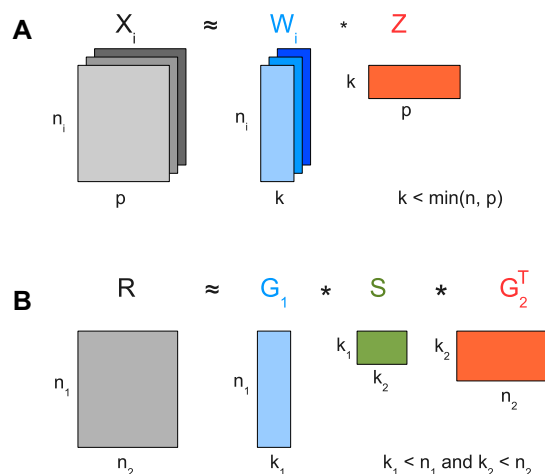
PARADIGM [121]	Inference of patient-specific pathways and patient stratification by integrating DNA copy number variations and mRNA expression data.	Matrix factorization	Unsupervised	Unsupervised
iCluster [124]	Cancer patient stratification by integrating copy number variation and mRNA expression data.	Matrix factorization	Unsupervised	Unsupervised
Joint Bayesian factor [129]	Driver genes identification by integrating mRNA expression and methylation data.	Matrix factorization	Unsupervised	Unsupervised
JIVE [130]	Cancer patient stratification by integrating mRNA expression and miRNA expression data.	NB	Unsupervised	Unsupervised
SNF [131]	Patient subtyping by integrating patient similarity networks constructed from mRNA expression, DNA methylation, and miRNA expression data.	Matrix factorization	Semisupervised	Semisupervised
NBS [133]	Cancer patient stratification by integrating somatic mutation data with molecular networks.	Matrix factorization	Semisupervised	Semisupervised
GNMTF [136]	Patient stratification, drug repurposing, and identifications of driver mutations by integrating of somatic mutations, molecular networks, drug-target interactions, and drug chemical similarity data.	Kernel-based	Supervised	Supervised
Joint kernel matrices [138]	Drug repurposing by integration of drug chemical structures, PPI network, and drug-induced gene expression data.	Kernel-based	Supervised	Supervised
PreDR [139]	Drug repurposing and prediction of novel drug-disease associations by integrating drug chemical structures, drug side-effects, and protein target structures.	Matrix factorization	Semisupervised	Semisupervised
MSCMF [140]	Drug-target interaction prediction by integrating known drug-target interactions along with multiple drug and target similarities.	Matrix factorization	Semisupervised	Semisupervised
DDR [141]	Drug-disease association prediction by integrating known drug-disease association along with multiple drug and target similarities.	Logistic regression	Supervised	Supervised
PREDICT [142]	Inference of novel drug indications by integrating multiple drug and target similarities.	NB	Unsupervised	Unsupervised
Coupled network propagation [143]	Drug-disease network inference by integrating drug, disease, and gene interaction network, as well as drug-gene and gene-disease association network.	Kolmogorov–Smirnov	Unsupervised	Unsupervised
Network completion [144]	Drug repurposing by integrating drug-target, drug-disease, and disease-target networks.	NB	Unsupervised	Unsupervised
Smirn [145]	Inference of drug-miRNA network by integrating cancer-related miRNA target gene expression and transcriptional responses to drug compounds.	Hyper geometric test	Unsupervised	Unsupervised
HGLDA [146]	Inference lncRNA-disease network by integrating miRNA-disease associations and lncRNA-miRNA interactions.	Matrix factorization	Semisupervised	Semisupervised
Regularized NMF [147]	Disease causing lncRNA prioritization by integrating lncRNA-disease associations, along with lncRNA and coding gene expression data and lncRNA-coding gene-association data.	NB	Unsupervised	Unsupervised

The first group of methods is used for subtyping and biomarker discovery; the second group is used for drug repurposing and therapy prediction. Some methods can belong to both categories (e.g. GNMTF).

for various cancer types, including colon and rectal [113], breast [114], and ovarian cancer [115].

Unsupervised clustering ML methods, such as hierarchical clustering [116], K-means [117], consensus clustering [118], and non-negative matrix factorization [119] have mostly been applied to gene expression data, by comparing expression levels of disease genes across different samples to identify meaningful subgroups. The most recent of such methods propose to divide patients into clinically relevant subtypes by comparing differentially expressed genes (based on normal and cancer tissue samples) [116]. Based on the selected set of differentially expressed genes, they calculate the distance between patients and perform hierarchical clustering [109]. Using mRNA expression data of breast and lung cancer patients, they identified four breast cancer and five lung cancer subtypes with significantly different survival rates. Moreover, instead of identifying individual driver mutations, they identify driver mutation modules for each individual subtype. Namely, by using the PPI (protein-protein interaction) network and by mapping the top 15 most frequently mutated genes of each identified subtype onto the network, they search for an optimally connected subnetwork covering these genes. The identified subnetworks are postulated as driver modules that can serve as new targets for repurposing of known drugs and their combinations [116]. Many other studies have also focused on developing methods for identifying aberrant network modules and pathways by utilizing molecular networks and other omics data. For example, Alcaraz et al. [120] developed KeyPathwayMiner, a method for extraction of aberrant network modules from PPI network by integrating gene expression and DNA methylation data. The authors demonstrated the performance of KeyPathwayMiner on the cancer genome atlas (TCGA) colorectal cancer patients. The method uses heuristic techniques based on *ant colony optimization* to extract maximally connected subnetworks with a certain number of differentially expressed genes in all patients. The resulting subnetworks were shown to be enriched in genes with overactive signaling in colorectal cancer that can be interpreted as potential therapeutic targets. Similarly, Vaske et al. [121] developed PARADIGM, a method for inferring patient-specific altered molecular pathways. The methods also allow for identification of common altered pathways among different patients and thus provide patient subtyping. The authors applied PARADIGM on TCGA gene expression and DNA copy number variations data of glioblastoma multiform patients; based on the significant pathway perturbations, the authors divide patient into four different subgroups with significantly different survival outcome.

However, a majority of recent methods use integrative approaches to combine multiple types of molecular data, such as DNA copy number alteration, DNA methylation, mRNA and protein expression, and molecular interaction data, accounting for different levels of variations among affected individuals and thereby providing more accurate subtyping [122, 123]. For example, Shen et al. [124] developed iCluster, an unsupervised learning framework that can simulta-



**Figure 3.** Illustration of MF-based methods. (A) Matrix factorization of multiple data matrices,  $X_i$ , representing different data types (e.g. mRNA expression, DNA methylation, copy number variation, etc.) over the same number of samples  $p$ . The matrices are decomposed into a common feature space, represented by matrix  $Z$ , which is also a cluster indicator matrix; it is used for assigning  $p$  samples into  $k$  clusters. Matrices  $W_i$  called coefficient matrices are specific to each dataset  $i$ . (B) Tri-factorization of the data matrix  $R$  representing relations between two datasets of sizes  $n_1$  and  $n_2$  (e.g. drug-target interactions) into three low-dimensional matrices. Matrices  $G_1$  and  $G_2$  are cluster indicator matrices for the first and second datasets, respectively; matrix  $G_1$  ( $G_2$ ) is used for assigning  $n_1$  ( $n_2$ ) data points to  $k_1$  ( $k_2$ ) clusters. Matrix  $S$  is the low-dimensional representation of  $R$ .

neously perform clustering, data integration, feature selection, and dimension reduction of multiple data types. It uses a probabilistic matrix factorization approach to simultaneously decompose data matrices, representing different data types (e.g. DNA methylation, DNA copy number variations, mRNA expression data) over the same number of samples (patients), into a common feature space represented by two low-dimensional matrices (Fig. 3A). Specifically, they decompose the data matrices by simultaneously factorizing each data matrix into a product of two low-dimensional matrices. The dimensionality of the low-dimensional matrices represents the number of cancer subtypes and it is a predefined parameter. The first matrix, also called *the coefficient matrix*, is specific to each data type, while the second matrix, also called *the cluster indicator matrix*, is shared across the decomposition. The second matrix captures the dependencies across the data types, and based on its entries it is used for a single, integrated assignment of tumor samples to clusters (subtypes). The authors applied iCluster on DNA copy number variation and gene expression data to stratify breast and lung cancer patients. After obtaining the probabilistic representation of the low-dimensional, cluster indicator matrix, they assign tumor samples to different subgroups. In both the breast and lung cancer data examples, they identify novel subgroups with statistically different clinical outcomes as a result of combined information from both the data types [124].



iCluster is a widely used tool and it has been applied for subtyping of various cancers. For example, Curtis et al. [125], applied it to breast cancer patients from METABRIC (Molecular Taxonomy of Breast Cancer International Consortium) cohort and identified ten subgroups with significantly different outcomes. Moreover, they reported significant correlation between genome variations and gene expression data and based on that, they discovered novel putative genes for breast cancer [125]. iCluster was also applied on TCGA glioblastoma multiforme (the most common and most aggressive malignant brain tumor) dataset by simultaneous clustering of DNA copy number variation, methylation, and gene expression data [126]. The authors reveal three distinct tumor subtypes of glioblastoma multiforme, as opposed to the four distinct subtypes reported by previous studies that used solely gene expression data [22]. This demonstrates the power of integrative analysis over analyses of single data types in characterizing, classifying, and predicting clinical outcomes of cancer patients.

The first method that deals with detection of contradictory signals across different data types is proposed by Yuan et al. [127]. They propose a Patient-Specific Data Fusion method based on nonparametric Bayesian approach to integrate gene expression and copy number variation data of prostate and breast cancer patients [127]. A Bayesian approach is a statistical ML approach that builds a model of data by constructing conditional dependencies between data variables represented by conditional probabilities. One of the widely used methods for learning conditional probabilities is Markov chain Monte Carlo (MCMC) technique [128]. Unlike other methods, this method successfully detects contradictory signals between different data types arising from different measurement errors. Specifically, a latent variable is assigned to each patient; it measures whether or not the patient's data are concordant (i.e. in agreement) across different data types. This approach allows for contradictory data information to be suppressed in the patient clustering assignment. The biggest drawback of this approach is that it does not scale well with the number of data types and thus, the authors restrict their analysis only on two data types. Namely, the MCMC step is computationally the most intensive and requires around 48 h for a single MCMC chain to complete. Despite this drawback, the authors report a novel subtype of prostate cancer patients with extremely poor survival outcome [127].

To further take into account data inconsistency across data types, iCluster was further generalized by Ray et al. [129] by introducing Bayesian joint factor model built upon iCluster framework. Namely, instead of having a single cluster indicator matrix common for all data types, they further decompose it into shared and data-specific matrix components. Specifically, the cluster indicator matrix is represented as a sum of data type-specific and common low-dimensional feature matrices. The common and specific low-dimensional matrices are learned jointly via simultaneous decomposition of all data matrices. This generalization was shown to be particularly useful for joint analysis of multiplatform genomic

data, as it allows more flexibility in the decomposition of distinct data types. Moreover, the authors reported better performance of their model compared to iCluster, because unlike iCluster, that enforces all tumor samples to be included into the clustering procedure, the proposed model can selectively choose between more and less correlated samples across data types when performing clustering assignment. The authors demonstrated their method on TCGA gene expression, copy number variation, and methylation data of ovarian cancer patients, particularly for uncovering key driver genes in ovarian cancer [129]. Similarly, Lock et al. [130] introduced JIVE (Joint and Individual Variation Explained), a method which instead of having the same coefficient matrices for shared and data-specific components proposed a model with different coefficient matrices corresponding to joint and data-specific components capturing low-dimensional joint variations across data types, as well as variations specific to each data type. With this extension, JIVE performed a better characterization of tumor subtypes, as well as a better understanding of the biological interactions between different data types [130].

To overcome scalability drawbacks of the previous ML clustering methods that operate with high-dimensional gene  $\times$  patient matrices, Wang et al. [131] proposed an NB method that integrates data represented by patient  $\times$  patient matrices. This method, called Similarity Network Fusion (SNF), combines mRNA expression, DNA methylation, and microRNA expression data for the same set of cancer patients. First, for each data type, it constructs a weighted network of patients, with nodes being patients and weighted links being similarities between patients. The similarities are computed based on their gene profiles for a particular data type. Second, it normalizes weights of each network by taking into account the networks from all data types. Finally, it fuses all the networks into a single network by performing a diffusion of information within each network and across different networks. After the convergence of the diffusion process, the authors use a spectral clustering method [132] on the final fused network to group patients into clusters. Unlike the previous methods, SNF is more scalable. Namely, instead of processing large-scale matrices constructed over a large number of genes, SNF method fuses much smaller matrices representing networks constructed over patients (i.e. samples), which makes the convergence faster. SNF is shown to be robust to noise and when applied on five different cancer types from TCGA database, it was shown to be effective in prediction of patient survival outcomes [131].

A majority of studies are based on analyzing mRNA expression data from RNA sequencing and microarrays, and DNA copy number alteration data. Because of noisiness of these data, the patient stratification studies for cancer types often do not produce patient subgroups that agree well with any clinical, or survival data [113]. To overcome these shortcomings, Hofree et al. [133] recently proposed the use of somatic mutation data as a new source of information for cancer patient stratification. However, highly heterogeneous somatic mutation profiles between different patients make the use

of somatic mutations for patient stratification into subtypes much harder [115, 133, 134]. Namely, two clinically identical tumors rarely have a large set of common mutated genes. Moreover, very few genes are frequently mutated across tumor samples. However, despite this genetic diversity between tumor samples, the perturbed pathways are often similar [134]. Therefore, Hofree et al. [133] proposed to address this problem by integrating somatic mutations with molecular networks that contain pathways. Their method, called NB Stratification (NBS), is based on network-regularized non-negative matrix factorization [135]. Namely, they factorize patient-gene binary matrix, encoding patients' somatic mutation profiles, into a product of two low-dimensional, non-negative matrices; the second of which being the cluster indicator matrix. The non-negativity constraint provides an easier interpretation of clustering assignment of tumor samples. They further incorporate molecular networks into the clustering procedure by constraining the construction of the cluster indicator matrix to respect the local network connectivity. This semisupervised approach uses molecular networks as prior knowledge about clusters, ensuring that the patients are grouped not only based on the similarity of their somatic mutation profiles, but also on the proximity of their mutated genes in the molecular network. Using the consensus clustering method [118] applied on the final cluster indicator matrix, the authors stratify patients into different subgroups. The method was applied on ovarian, uterine, and lung cancer patients from TCGA database, and it yielded cancer subtypes with different clinical outcomes, response to therapies, and tumor histologies.

MF-based methods are promising for mining heterogeneous datasets. These methods have a potential to incorporate any number and type of heterogeneous data and to perform comprehensive analyses. We recently made a step toward this goal and extended the NBS method to incorporate drug data into the framework [136]. Unlike the previous, our method is more comprehensive because it can simultaneously perform three tasks: cancer patient subtyping, drug repurposing, and biomarker discovery (driver gene identification). We used Graph-regularized Non-negative Matrix Tri-Factorization (GNMTF) [111] (see Fig. 3B for an illustration) approach to integrate somatic mutation profiles of ovarian cancer patients, molecular networks, drug-target interactions, and drug chemical similarity data. We simultaneously tri-factorize patient-gene and drug-target matrix by sharing common low-dimensional matrix factors representing cluster indicator matrices. We compute three different cluster indicator matrices used for clustering assignment of genes, patients, and drugs, respectively. The computation of the gene cluster indicator matrix is constrained by connectivity of integrated molecular network, whereas the computation of the drug cluster indicator matrix is constrained by drug chemical similarities. The integrated network is composed of three different molecular networks, namely PPI, genetic, and metabolic interaction networks. Given that GNMTF is both a coclustering and dimensionality-

reduction approach, we use GNMTF to perform the following three tasks: (i) we use the patient cluster indicator matrix to stratify ovarian cancer patients into different subgroups with different clinical outcomes; (ii) we use the gene cluster indicator matrix to uncover gene modules enriched in driver mutations and postulate new genes as drivers of tumor progression; and (iii) we use the matrix completion property of the drug-target matrix to predict novel drug-target interactions and discover new drug candidates that can be repurposed to treat ovarian cancer patients.

### 3.1.1 Challenges and open questions

Identification of disease subtypes has been shown to be both data and method dependent. Moreover, there is no consensus in the literature about the number of subtypes of a particular cancer type. Depending on the methods and data types they use, different studies report different numbers of subtypes of a particular cancer type (e.g. breast cancer). Also, unsupervised methods require the number of subtypes to be predetermined. Determining the number of subtypes is not a straightforward task and different approaches can be used to discover the correct number of clusters in the data. For example, iCluster uses a cross-validation technique [124], while NBS determines the number of subtypes based on the stability of the consensus clustering [133]. Furthermore, there is an urgent need for a reference dataset that should be used in future studies for systematic evaluation and comparisons of methods.

Moreover, many of the above-mentioned integrative methods for subtyping are incapable of simultaneously considering different data types. For example, SNF method can only integrate data types given by continuous variables (e.g. mRNA expression levels), as they can be easily used for construction of similarity networks. However, SNF cannot incorporate somatic mutation profiles, as it cannot construct a similarity network from highly heterogeneous somatic mutation profiles. Namely, due to the small overlap between somatic mutation profiles across different patients, it is difficult to define a proper similarity measure between patients. Approaches such as NBS and GNMTF are more convenient for integration of somatic mutation profiles. Very few studies integrate somatic mutation data with mRNA and methylation data, due to the difficulty in integrating binary with continuous data types [137].

A proper normalization of different data types is another issue in integrative data analyses. If not properly accounted for, it often results in cases where the largest dataset wins. Unlike iCluster, JIVE properly takes into account the data normalization problem [130].

### 3.2 Computational methods for drug repurposing and personalized treatments

Various computational methods for drug repurposing have been proposed and they can be classified under different

criteria. For example, from the data viewpoint, Dudley et al. [148] suggested classification into *drug-based* and *disease-based* methods. The first group of methods uses some notion of similarity between drugs (e.g. chemical similarity [149], similarity between gene expressions induced by drug actions [74], or drug side-effect similarity [150]) to group drugs and infer a novel drug candidate for repurposing from the group that can perform the same action as other drugs in the group. The second group of methods uses similarities between diseases (e.g. phenotype similarity [151], or similarity between disease symptoms [152]) to group diseases and to infer a novel drug for repurposing by expanding known associations between the drug and some members of the group to the rest of the group. Other approaches use *target-based* similarities [153], i.e. protein sequence similarity [154], or 3D structural similarity [155], to infer novel drugs. On the other hand, all three approaches can be classified as *similarity-based approaches* [153]. They often use either machine-learning, or NB methods in the drug inference process. Other computational approaches include *molecular docking simulation* approaches that deal with prediction of a binding place of a drug within protein 3D structure [156]. However, the biggest limitations of these methods are the lack of knowledge of 3D structures for many protein targets and extensive computational costs for testing a single drug-target interaction.

A full review of similarity-based and molecular docking approaches for single data type analyses is beyond the scope of this article and we refer the reader to recent review articles by Li et al. [157] and Ding et al. [153]. Here, we focus on integrative methods capable of integrating various similarities from different data types containing complementary information, such as pharmacological, chemical, genetic, and clinical data. Namely, due to heterogeneity and complexity of many diseases characterized with different subtypes, drugs are not always equally efficient in treatment of the same disease. Thus, the overarching goal of precision medicine is to take into account molecular diversity between individuals when diagnosing patients and prescribing drugs specific to each individual [158]. With the Big-Data initiative (see Section 2), integrative computational approaches have started attracting more attention due to their ability to address this goal.

For example, Napolitano et al. [138] used a KB method [106] to integrate drug chemical similarity, PPI network, and drug-induced gene expression data after a patient treatment. Each data are represented by a kernel matrix in a drug-centered feature space. Particularly, the three kernel matrices represent drug–drug similarities based on: (i) drug chemical structures from DrugBank; (ii) proximity of their targets in the PPI network; and (iii) correlations between gene profiles under the drug's influence retrieved from CMap database. After combining these kernel matrices into a single kernel matrix, the authors applied an SVM, a supervised ML method for classification. They trained the SVM on the existing drug classification achieving 78% of classification accuracy and they used the top scoring misclassified drugs as new candidates for repurposing [138]. A similar approach was used

by Wang et al. [139], who developed a PreDR (Predict Drug Repurposing) method where drug-centered kernel matrices represent: (i) drug chemical similarities obtained from PubChem database; (ii) target (protein) sequence similarities retrieved from KEGG BRITE and DrugBank; and (iii) drug side-effect similarities for SIDER database. The disease-centered kernel matrix represents disease similarities measured by their semantic similarity of disease phenotypes retrieved from OMIM database. The authors trained the SVM classifier on the combined kernel matrix and reported accuracy in identifying novel drug-disease interactions.

Zheng et al. [140] developed an integrative framework called Multiple Similarities Collaborative Matrix Factorization (MSCMF) for drug-target prediction. It takes as an input a matrix representing drug-target interactions, as well as multiple matrices representing different types of similarities between drugs and targets constructed from various databases. MSCMF projects drugs and targets into a common low-dimensional feature space by factorizing the drug-target matrix into a product of two low-dimensional matrices representing drug and target low-dimensional feature vectors, respectively. The computation of low-dimensional matrices of drugs and targets is done in a semisupervised manner by constraining their values to be consistent with drug–drug and target–target similarity matrices, respectively. Namely, the similarity between two drugs is approximated by the inner product of their corresponding feature vectors. The same is applied on target feature vectors. The authors mathematically formulated the factorization condition and constraints within the same objective function, which they minimize by applying the Alternating Least Squares algorithm [159]. After convergence, they reconstructed the drug-target matrix from the obtained low-dimensional matrices (i.e. from matrix completion) and extracted new, previously unobserved entries representing predicted drug-target interactions. MSCMF is shown to perform better than the previous state-of-the-art methods for drug-target prediction. Moreover, the big advantage of MSCMF over the previous methods is the fact that it can integrate similarities from multiple data sources over the same set of drugs or targets and estimate their influence onto the quality of the drug-target prediction.

Similar to MSCMF, Zhang et al. [141] proposed drug-disease repositioning, a semisupervised, matrix tri-factorization-based framework for novel drug-disease association prediction. It takes as input known drug-disease associations, as well as multiple drug and multiple disease similarity networks and generates new drug-disease associations. In particular, it constructs three drug similarity matrices based on their chemical structures, side effects, and target proteins and three disease similarity matrices based on their phenotypes, Disease Ontology, and disease genes. The predicted associations are validated in clinical trial databases. Unlike MSCMF, drug-disease repositioning factorizes drug-disease associations into a product of three low-dimensional matrices, where the first and the last matrices can be interpreted as cluster assignment matrices of drugs and diseases,

respectively. These matrices can be used to identify subgroups of highly correlated drugs and diseases, thus providing additional insights into drug repurposing by identifying a group of similar drug candidates that can be used in clinical trials.

Gottlieb et al. [142] developed a supervised method, called PREDICT (PREdicting Drug IndiCaTions). First, it computes drug–drug and disease–disease similarity measures from five and six different drug and disease data sources, respectively. Second, based on these similarities, it constructs an overall similarity for each drug–disease pair. Finally, based on the drug–disease similarity, it trains a logistic regression classifier on correctly classifying known drug–disease associations. The authors demonstrated a great accuracy of PREDICT in identifying novel indications of drugs with area under the ROC curve [160] of 0.92. Moreover, they propose PREDICT as a general framework that can be used in future personalized drug treatments by incorporating gene expression data of disease patients into the framework.

All previous methods use either supervised, or semisupervised strategy in predicting drug–target, or drug–disease associations and thus, they require a gold standard (i.e. a set of known associations) to train their models. For many specific diseases, that dataset is unknown, or incomplete, which makes the use of the methods more difficult. To overcome this, Huang et al. [143] proposed a completely unsupervised integrative method that can infer drug–disease associations without any prior associations. They used coupled network propagation [161] on drug–drug chemical similarity, disease–disease phenotype similarity, and gene–gene coexpression similarity homogeneous networks, connected by drug–gene and gene–disease heterogeneous networks. They applied their method on data for prostate and colorectal cancer patients. They identified top scoring drugs predicted to be used in treatment of these groups of patients. Another unsupervised, NB method for heterogeneous network integration and drug repurposing was introduced by Daminelli et al. [144]. They predicted novel drug–target associations by completing incomplete bi-cliques in the integrated drug–target–disease network. They demonstrate the power of their method by predicting novel drugs for cardiovascular and parasitic diseases, as well as by predicting novel drugs for cancer-related kinases. For other NB methods for drug repurposing, we refer a reader to a recent review paper by Wu et al. [162].

Noncoding RNAs, in particular microRNAs (or miRNAs) and long noncoding RNAs (lncRNAs), have recently started attracting attention due to their involvement in various diseases, including cancer and autoimmune disorders [163] and thus, have been proposed as potential biomarkers [146, 164] and drug targets [165, 166]. Due to large collections of transcriptional and drug data being available, new computational methods for identification of miRNAs as potential drug targets have recently been proposed. For example, Jiang et al. [145] proposed a framework for construction of a network, SMirN, of interactions between small drug molecules (compounds) and miRNAs using data from different human can-

cers. Specifically, they used transcriptional responses to compounds and differentially expressed miRNA target genes in 23 different human cancers. For each miRNA, they partitioned their target genes into GO modules, and for each GO module they evaluated the association between its differentially expressed target genes and the transcriptional response to the compound by using Kolmogorov–Smirnov test. If these associations are confirmed for a significant number of GO modules of a particular miRNA, then the authors hypothesized a link between the miRNA and the corresponding drug compound. They analyzed the SMirN network and separately grouped miRNAs and compounds into modules, based on which they infer novel potential miRNA targets, as well as novel drug compounds that can be used in drug repurposing for cancer therapy. Chen [167] developed a novel model of HyperGeometric distribution for lncRNA–Disease Associations inference. The model integrates known miRNA–disease associations and lncRNA–miRNA interactions and without a gold standard dataset, it infers a network of lncRNA–disease associations with AUC of 0.76 in the leave-one-out cross-validation. Based on the top 19 predicted associations, they reported novel lncRNAs involved in breast, lung, and colorectal cancer that can be used as novel biomarkers for diagnosis of these cancers. A more sophisticated integrative method, based on non-negative matrix factorization, was recently proposed by Biswas et al. [147]. They factorize lncRNA–disease association matrix into a product of two non-negative, low-dimensional matrices specific to lncRNAs and diseases, respectively. The non-negativity of the obtained, low-dimensional matrices allows for easier extraction of lncRNA and disease subgroups in the data. They can also be interpreted as cluster assignment matrices for lncRNAs and diseases, respectively. The factorization of the lncRNA–disease association matrix was done in a semisupervised way, by constraining the construction of the low-dimensional matrices with additional data, including coding gene and lncRNA expression data, as well as lncRNA–coding gene–association network. The authors identified several biologically relevant lncRNA and disease groups. Based on the membership scores in the lncRNA low-dimensional matrix, they ranked disease causing lncRNAs for each particular disease. They identified a prominent group of lncRNAs associated with heart diseases, as well as a group of lncRNAs strongly associated with neurological disorders that can be used in future experimental testing as biomarkers of these disorders.

### 3.2.1 Challenges and open questions

Many of the methods presented in this section require different data types to be represented in common feature space. For example, KB methods (e.g. PreDR) require the matrices of all data types to be constructed over the same set of entities (e.g. drugs or diseases). This often requires transforming the data that may lead to information loss. On the other hand, MF-based methods (e.g. MSCMF) can handle these



heterogeneous data without any data transformation and thus, without any information loss. Also, many methods require choosing an appropriate similarity measure to integrate various data types. This is not always a straightforward task and different measures may result in different final conclusions. Similar to the methods described in Section 3.1, the methods for drug-target (and drug-disease) prediction and drug repurposing are lacking a reference corpus of data for comparing their performances.

#### 4 Challenges and perspectives

As presented in Section 2, biomedical data are increasingly becoming available and dealing with their “three V” components will impose many challenges and open questions. For example, in addressing Big Data’s volume (i.e. high dimensionality), many dimensionality-reduction techniques have been devised, reviewed in Sections 3.1 and 3.2. However, they are all computationally intensive on large-scale datasets and devising techniques that are both efficient and accurate in revealing hidden substructures in them is still an open question. One of the possible solutions to addressing this question might be Topological Data Analysis methods (TDAs) [168, 169]. TDAs use mathematical concepts developed in algebraic topology. TDAs analyze Big Data by converting them into low-dimensional geometric representations from which they extract shapes (patterns) and obtain insight into them. These methods have been shown to be more efficient in finding substructures in large-scale datasets than standard methods, such as clustering, or principal component analysis methods. Moreover, they succeed in finding hidden structures in the data that standard methods failed to discover [169].

Dealing with Big Data’s velocity (i.e. coping with its growth over time) is particularly challenging and poorly addressed in the literature on precision medicine. One of the possible future directions in addressing this challenge is the utilization of so-called “anytime algorithms” [170] that can learn from streaming data (e.g. time-dependent Bayesian classifiers) [171] and that still return a valuable result if their execution is interrupted at any time. Moreover, in the future, we will have access to more and more time series data. At the moment, such time series are either preprocessed to find patterns, e.g. time series of expression data are either used to find genes with time-correlated expression (coexpression network), or used to study the effect of drugs on short time scales by differential expression analysis. With the increasing number of measured features and the increasing time span of the measurements, a key challenge will be to find a data integration model that will directly mine time series measurements for which the time spans and frequencies of measurements vary greatly.

The Big Data’s variety (i.e. heterogeneity) has been addressed by many methods as presented in Section 3.2. MF-

based methods are promising for mining heterogeneous datasets. Although GNMTF is a versatile data integration framework [136], its computational complexity increases with the number of data types to be integrated. Thus, integrating large numbers of heterogeneous data types within the MF-based framework necessitates novel algorithmic improvements. Extracting the complementary information conveyed in data of different formats and types is another challenge that is partially addressed by the presented integrative methods. For example, proteomics data have been shown to be a good complement to other omics data. Namely, many studies have confirmed that proteins having physical interactions in a PPI network are more likely to have correlated coexpression profiles of their corresponding genes [172]. On the contrary, protein physical interactions are less likely to coincide to genetic interactions of their corresponding genes [173]. Thus, integrating genetic interaction network with PPI network and other molecular networks has been shown to be beneficial in many biological problems [133, 136, 174].

Moreover, many data types including exposomic and metagenomic data are yet to be analyzed and their integration with other data will be a focus of future studies. For example, much of an individual’s health data, such as demographic data, personal and family medical history, vaccination records, laboratory tests, and imaging results are systematically being collected and stored in Electronic Health Records (EHR). EHR data are increasingly becoming available for academic research purposes and they present numerous computational challenges that are yet to be addressed. Two major computational challenges include developing algorithms for: (i) individual *phenotyping* (i.e. annotating patient records with disease conditions) [175] and (ii) integration of EHR data with omics data for better understanding of disease mechanisms and treatments [176]. The biggest obstacles of the first challenge are nosiness and incompleteness of the EHR data that need to be properly taken into account. On the other hand, the biggest obstacles of the second challenge are heterogeneity and different format types of EHR and genomic data. Some steps toward addressing these challenges have been made [175, 176], but developing methods that can overcome these obstacles are yet to come.

Finally, while we focus on the four challenges of precision medicine, big data integration also opens novel opportunities in bioinformatics and in other data sciences. For example, it can be used to reprocess raw data in a more coherent way, or with novel research questions in mind [177].

*This work was supported by the European Research Council (ERC) Starting Independent Researcher Grant 278212, the National Science Foundation (NSF) Cyber-Enabled Discovery and Innovation (CDI) OIA-1028394, the ARRS project J1-5454, and the Serbian Ministry of Education and Science Project III144006.*

*The authors have declared no conflict of interest.*



## 5 References

- [1] Hood, L., Friend, S. H., Predictive, personalized, preventive, participatory (p4) cancer medicine. *Nat. Rev. Clin. Oncol.* 2011, *8*, 184–187.
- [2] Mirnezami, R., Nicholson, J., Darzi, A., Preparing for precision medicine. *N. Engl. J. Med.* 2012, *366*, 489–491.
- [3] Table of blood group systems v4.0. International Society of Blood Transfusion, Amsterdam 2014.
- [4] Smits, J., De Meester, J., Persijn, G., Claas, F., Vanrenterghem, Y., Long-term results of solid organ transplantation. Report from the eurotransplant international foundation. *Clin. Transplant.* 1995, 109–127.
- [5] Takemoto, S., Terasaki, P. I., Cecka, J. M., Cho, Y. W., Gjertson, D. W., Survival of nationally shared, hla-matched kidney transplants from cadaveric donors. *N. Engl. J. Med.* 1992, *327*, 834–839.
- [6] Thorogood, J., Persijn, G. G., Schreuder, G. M., Zantvoort, F. A. et al., The effect of hla matching on kidney graft survival in separate post-transplantation intervals. *Transplantation* 1990, *50*, 146–149.
- [7] Mitsuishi, Y., Terasaki, P., HLA matching effect on five-year graft survival and half-life in the cyclosporine era. *Kidney Int.* 1992, Suppl *38*, S172–S175.
- [8] Beyer, M. A., Laney, D., *The Importance of 'Big Data': A Definition*, Gartner, Stamford, CT 2012.
- [9] International Human Genome Sequencing Consortium, I. H. G. S. et al., Finishing the euchromatic sequence of the human genome. *Nature* 2004, *431*, 931–945.
- [10] McDermott, J. E., Wang, J., Mitchell, H., Webb-Robertson, B. J. et al., Challenges in biomarker discovery: combining expert insights with statistical analysis of complex omics data. *Expert Opin. Med. Diagn.* 2013, *7*, 37–51.
- [11] Kato, H., Takahashi, S., Saito, K., Omics and integrated omics for the promotion of food and nutrition science. *J. Tradit. Complement. Med.* 2011, *1*, 25–30.
- [12] Cho, I., Blaser, M. J., The human microbiome: at the interface of health and disease. *Nat. Rev. Genet.* 2012, *13*, 260–270.
- [13] Yuan, Y., Failmezger, H., Rueda, O. M., Ali, H. R. et al., Quantitative image analysis of cellular heterogeneity in breast tumors complements genomic profiling. *Sci. Transl. Med.* 2012, *4*, 143–157.
- [14] Kumar, V., Gu, Y., Basu, S., Berglund, A. et al., Radiomics: the process and the challenges. *Magn. Reson. Imaging* 2012, *30*, 1234–1248.
- [15] Lambin, P., Rios-Velazquez, E., Leijenaar, R., Carvalho, S. et al., Radiomics: extracting more information from medical images using advanced feature analysis. *Eur. J. Cancer* 2012, *48*, 441–446.
- [16] Ahrens, M. B., Orger, M. B., Robson, D. N., Li, J. M., Keller, P. J., Whole-brain functional imaging at cellular resolution using light-sheet microscopy. *Nat. Methods* 2013, *10*, 413–420.
- [17] Mardis, E. R., The \$1,000 genome, the \$100,000 analysis. *Genome Med.* 2010, *2*, 84–86.
- [18] Margolis, R., Derr, L., Dunn, M., Huerta, M. et al., The national institutes of health's big data to knowledge (bd2k) initiative: capitalizing on biomedical big data. *J. Am. Med. Inform. Assoc.* 2014, *21*, 957–958.
- [19] Lötvall, J., Akdis, C. A., Bacharier, L. B., Bjermer, L. et al., Asthma endotypes: a new approach to classification of disease entities within the asthma syndrome. *J. Allergy Clin. Immunol.* 2011, *127*, 355–360.
- [20] Boland, M. R., Hripcsak, G., Shen, Y., Chung, W. K., Weng, C., Defining a comprehensive verotype using electronic health records for personalized medicine. *J. Am. Med. Inform. Assoc.* 2013, *20*, e232–e238.
- [21] Schulam, P., Wigley, F., Saria, S., Clustering longitudinal clinical marker trajectories from electronic health data: applications to phenotyping and endotype discovery. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, Austin, Texas, USA 2015, pp. 2956–2964.
- [22] Verhaak, R. G., Hoadley, K. A., Purdom, E., Wang, V. et al., Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in *pdgfra*, *idh1*, *egfr*, and *nf1*. *Cancer Cell* 2010, *17*, 98–110.
- [23] Lewis, S., Foltynie, T., Blackwell, A. D., Robbins, T. W. et al., Heterogeneity of Parkinson's disease in the early clinical stages using a data driven approach. *J. Neurol. Neurosurg. Psychiatry* 2005, *76*, 343–348.
- [24] Gutman, S., Kessler, L. G., The us food and drug administration perspective on cancer biomarker development. *Nat. Rev. Cancer* 2006, *6*, 565–571.
- [25] Davis, J. C., Furstenthal, L., Desai, A. A., Norris, T. et al., The microeconomics of personalized medicine: today's challenge and tomorrow's promise. *Nat. Rev. Drug Discov.* 2009, *8*, 279–286.
- [26] Ashburn, T. T., Thor, K. B., Drug repositioning: identifying and developing new uses for existing drugs. *Nat. Rev. Drug Discov.* 2004, *3*, 673–683.
- [27] Wheeler, D. A., Srinivasan, M., Egholm, M., Shen, Y. et al., The complete genome of an individual by massively parallel dna sequencing. *Nature* 2008, *452*, 872–876.
- [28] Ng, S. B., Turner, E. H., Robertson, P. D., Flygare, S. D. et al., Targeted capture and massively parallel sequencing of 12 human exomes. *Nature* 2009, *461*, 272–276.
- [29] Choi, M., Scholl, U. I., Ji, W., Liu, T. et al., Genetic diagnosis by whole exome capture and massively parallel dna sequencing. *Proc. Natl. Acad. Sci. USA* 2009, *106*, 19096–19101.
- [30] Owens, B., The single life. *Nature* 2012, *491*, 27–29.
- [31] Shapiro, E., Biezuner, T., Linnarsson, S., Single-cell sequencing-based technologies will revolutionize whole-organism science. *Nat. Rev. Genet.* 2013, *14*, 618–630.
- [32] Lupski, J. R., Genome mosaicism-one human, multiple genomes. *Science* 2013, *341*, 358–359.
- [33] Ziller, M. J., Gu, H., Müller, F., Donaghey, J. et al., Charting a dynamic dna methylation landscape of the human genome. *Nature* 2013, *500*, 477–481.
- [34] Rivera, C. M., Ren, B., Mapping human epigenomes. *Cell* 2013, *155*, 39–55.

- [35] Dawson, M. A., Kouzarides, T., Cancer epigenetics: from mechanism to therapy. *Cell* 2012, *150*, 12–27.
- [36] Suvà, M. L., Riggi, N., Bernstein, B. E., Epigenetic reprogramming in cancer. *Science* 2013, *339*, 1567–1570.
- [37] Vishwanathan, N., Le, H., Le, T., Hu, W.-S., Advancing biopharmaceutical process science through transcriptome analysis. *Curr. Opin. Biotechnol.* 2014, *30*, 113–119.
- [38] Marian, A., Sequencing your genome: what does it mean? *Methodist DeBakey Cardiovasc. J.* 2014, *10*, 3–6.
- [39] Allocco, D. J., Kohane, I. S., Butte, A. J., Quantifying the relationship between co-expression, co-regulation and gene function. *BMC Bioinformatics* 2004, *5*, 18.
- [40] DeRisi, J., Penland, L., Brown, P. O., Bittner, M. L. et al., Use of a cDNA microarray to analyse gene expression patterns in human cancer. *Nat. Genet.* 1996, *14*, 457–460.
- [41] Volinia, S., Calin, G. A., Liu, C. G., Ambs, S. et al., A microRNA expression signature of human solid tumors defines cancer gene targets. *Proc. Natl. Acad. Sci. USA* 2006, *103*, 2257–2261.
- [42] Van't Veer, L.J., Dai, H., Van De Vijver, M. J., He, Y. D. et al., Gene expression profiling predicts clinical outcome of breast cancer. *Nature* 2002, *415*, 530–536.
- [43] Jensen, O. N., Modification-specific proteomics: characterization of post-translational modifications by mass spectrometry. *Curr. Opin. Chem. Biol.* 2004, *8*, 33–41.
- [44] Ong, S.-E., Mann, M., Mass spectrometry-based proteomics turns quantitative. *Nat. Chem. Biol.* 2005, *1*, 252–262.
- [45] Fields, S., Song, O. K., A novel genetic system to detect protein-protein interactions. *Nature* 1989, *340*, 245–246.
- [46] Ho, Y., Gruhler, A., Heilbut, A., Bader, G. D. et al., Systematic identification of protein complexes in *Saccharomyces cerevisiae* by mass spectrometry. *Nature* 2002, *415*, 180–183.
- [47] Kelley, B. P., Sharan, R., Karp, R. M., Sittler, T. et al., Conserved pathways within bacteria and yeast as revealed by global protein network alignment. *Proc. Natl. Acad. Sci. USA* 2003, *100*, 11394–11399.
- [48] Kuchaiev, O., Milenković, T., Memišević, V., Hayes, W., Pržulj, N., Topological network alignment uncovers biological function and phylogeny. *J. R. Soc. Interface* 2010, *7*, 1341–1354.
- [49] Bandyopadhyay, S., Sharan, R., Ideker, T., Systematic identification of functional orthologs based on protein network comparison. *Genome Res.* 2006, *16*, 428–435.
- [50] Daviss, B., Growing pains for metabolomics. *The Scientist* 2005, *19*, 25–28.
- [51] Lowe, J. B., Marth, J. D., A genetic approach to mammalian glycan function. *Annu. Rev. Biochem.* 2003, *72*, 643–691.
- [52] Kinjo, Y., Wu, D., Kim, G., Xing, G. W. et al., Recognition of bacterial glycosphingolipids by natural killer t cells. *Nature* 2005, *434*, 520–525.
- [53] Crocker, P. R., Siglecs: sialic-acid-binding immunoglobulin-like lectins in cell-cell interactions and signalling. *Curr. Opin. Struct. Biol.* 2002, *12*, 609–615.
- [54] Sasisekharan, R., Shriver, Z., Venkataraman, G., Narayanasami, U., Roles of heparan-sulphate glycosaminoglycans in cancer. *Nat. Rev. Cancer* 2002, *2*, 521–528.
- [55] Fry, E. E., Lea, S. M., Jackson, T., Newman, J. W. et al., The structure and function of a foot-and-mouth disease virus-oligosaccharide receptor complex. *EMBO J.* 1999, *18*, 543–554.
- [56] Winter, G., Krömer, J. O., Fluxomics-connecting 'omics analysis and phenotypes. *Environ. Microbiol.* 2013, *15*, 1901–1916.
- [57] Visscher, P. M., Brown, M. A., McCarthy, M. I., Yang, J., Five years of GWAS discovery. *Am. J. Hum. Genet.* 2012, *90*, 7–24.
- [58] Viaud, S., Saccheri, F., Mignot, G., Yamazaki, T. et al., The intestinal microbiota modulates the anticancer immune effects of cyclophosphamide. *Science* 2013, *342*, 971–976.
- [59] Rajendhran, J., Gunasekaran, P., Human microbiomics. *Indian J. Microbiol.* 2010, *50*, 109–112.
- [60] 1000 Genomes Project Consortium, An integrated map of genetic variation from 1,092 human genomes. *Nature* 2012, *491*, 56–65.
- [61] Ajay, S. S., Parker, S. C., Abaan, H. O., Fajardo, K. V. F., Margulies, E. H., Accurate and comprehensive sequencing of personal genomes. *Genome Res.* 2011, *21*, 1498–1505.
- [62] Kalisky, T., Quake, S. R., Single-cell genomics. *Nat. Methods* 2011, *8*, 311–314.
- [63] Fan, H. C., Gu, W., Wang, J., Blumenfeld, Y. J. et al., Non-invasive prenatal measurement of the fetal genome. *Nature* 2012, *487*, 320–324.
- [64] Chen, Y.-C., Rajagopala, S. V., Stellberger, T., Uetz, P., Exhaustive benchmarking of the yeast two-hybrid system. *Nat. Methods* 2010, *7*, 667–668.
- [65] Jagadish, H. V., Gehrke, J., Labrinidis, A., Papakonstantinou, Y. et al., Big data and its technical challenges. *Commun. ACM* 2014, *57*, 86–94.
- [66] Brown, G. R., Hem, V., Katz, K. S., Ovetsky, M. et al., Gene: a gene-centered information resource at NCBI. *Nucleic Acids Res.* 2015, *43*, D36–D42.
- [67] Dimmer, E. C., Huntley, R. P., Alam-Faruque, Y., Sawford, T. et al., The uniprot-go annotation database in 2011. *Nucleic Acids Res.* 2012, *40*, D565–D570.
- [68] ENCODE Project Consortium, The encode (encyclopedia of DNA elements) project. *Science* 2004, *306*, 636–640.
- [69] Bernstein, B. E., Stamatoyannopoulos, J. A., Costello, J. F., Ren, B. et al., The NIH roadmap epigenomics mapping consortium. *Nat. Biotechnol.* 2010, *28*, 1045–1048.
- [70] Teng, L., He, B., Wang, J., Tan, K., 4Dgenome: a comprehensive database of chromatin interactions. *Bioinformatics* 2015, *31*, 2560–2564.
- [71] Xin, Y., Chanrion, B., O'Donnell, A. H., Milekic, M. et al., MethylomeDB: a database of DNA methylation profiles of the brain. *Nucleic Acids Res.* 2012, *40*, D1245–D1249.
- [72] Edgar, R., Domrachev, M., Lash, A. E., Gene expression omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res.* 2002, *30*, 207–210.

- [73] Petryszak, R., Burdett, T., Fiorelli, B., Fonseca, N. A. et al., Expression atlas update—a database of gene and transcript expression from microarray-and sequencing-based functional genomics experiments. *Nucleic Acids Res.* 2014, *42*, D926–D932.
- [74] Lamb, J., Crawford, E. D., Peck, D., Modell, J. W. et al., The connectivity map: using gene-expression signatures to connect small molecules, genes, and disease. *Science* 2006, *313*, 1929–1935.
- [75] Okamura, Y., Aoki, Y., Obayashi, T., Tadaka, S. et al., CoXpresDB in 2015: coexpression database for animal species by DNA-microarray and RNAseq-based expression data with multiple quality assessment systems. *Nucleic Acids Res.* 2014, *43*, D82–D86.
- [76] van Dam, S., Craig, T., de Magalhães, J. P., Genefriends: a human RNA-seq-based gene and transcript co-expression database. *Nucleic Acids Res.* 2015, *43*, D1124–D1132.
- [77] UniProt Consortium, Uniprot: a hub for protein information. *Nucleic Acids Res.* 2014, *43*, D204–D212.
- [78] Gaudet, P., Michel, P. A., Zahn-Zabal, M., Cusin, I. et al., The nextprot knowledgebase on human proteins: current status. *Nucleic Acids Res.* 2015, *43*, D764–D770.
- [79] Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G. et al., The protein data bank. *Nucleic Acids Res.* 2000, *28*, 235–242.
- [80] Uhlén, M., Fagerberg, L., Hallström, B. M., Lindskog, C. et al., Tissue-based map of the human proteome. *Science* 2015, *347*, 1260419.
- [81] Kerrien, S., Aranda, B., Breuza, L., Bridge, A. et al., The intact molecular interaction database in 2012. *Nucleic Acids Res.* 2011, *40*, D841–D846.
- [82] Chatr-aryamontri, A., Breitkreutz, B. J., Oughtred, R., Boucher, L. et al., The BioGRID interaction database: 2015 update. *Nucleic Acids Res.* 2014, *43*, D470–D478.
- [83] Brown, K. R., Jurisica, I., Online predicted human interaction database. *Bioinformatics* 2005, *21*, 2076–2082.
- [84] Szklarczyk, D., Franceschini, A., Wyder, S., Forslund, K. et al., String v10: protein-protein interaction networks, integrated over the tree of life. *Nucleic Acids Res.* 2014, *43*, D447–D452.
- [85] Wishart, D. S., Jewison, T., Guo, A. C., Wilson, M. et al., HMDB 3.0—the human metabolome database in 2013. *Nucleic Acids Res.* 2012, *41*, D801–D807.
- [86] Ogata, H., Goto, S., Sato, K., Fujibuchi, W. et al., KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* 1999, *27*, 29–34.
- [87] Jewison, T., Su, Y., Disfany, F. M., Liang, Y. et al., Smpdb 2.0: big improvements to the small molecule pathway database. *Nucleic Acids Res.* 2014, *42*, D478–D484.
- [88] Croft, D., Mundo, A. F., Haw, R., Milacic, M. et al., The REACTOME pathway knowledgebase. *Nucleic Acids Res.* 2014, *42*, D472–D477.
- [89] Shakhsher, B., Anderson, M., Khatib, K., Tadoori, L. et al., Sugarbind database (sugarbinddb): a resource of pathogen lectins and corresponding glycan targets. *J. Mol. Recognit.* 2013, *26*, 426–431.
- [90] Campbell, M. P., Peterson, R., Mariethoz, J., Gasteiger, E. et al., UnicarbKB: building a knowledge platform for glycoproteomics. *Nucleic Acids Res.* 2013, *42*, D215–D221.
- [91] Hashimoto, K., Goto, S., Kawano, S., Aoki-Kinoshita, K. F. et al., KEGG as a glycome informatics resource. *Glycobiology* 2006, *16*, 63R–70R.
- [92] Hamosh, A., Scott, A. F., Amberger, J. S., Bocchini, C. A., McKusick, V. A., Online mendelian inheritance in man (OMIM), a knowledgebase of human genes and genetic disorders. *Nucleic Acids Res.* 2005, *33*, D514–D517.
- [93] Mailman, M. D., Feolo, M., Jin, Y., Kimura, M. et al., The NCBI dbgap database of genotypes and phenotypes. *Nat. Genet.* 2007, *39*, 1181–1186.
- [94] Welter, D., MacArthur, J., Morales, J., Burdett, T. et al., The nhgri GWAS catalog, a curated resource of snp-trait associations. *Nucleic Acids Res.* 2014, *42*, D1001–D1006.
- [95] Forbes, S. A., Beare, D., Gunasekaran, P., Leung, K. et al., Cosmic: exploring the world's knowledge of somatic mutations in human cancer. *Nucleic Acids Res.* 2015, *43*, D805–D811.
- [96] Weinstein, J. N., Collisson, E. A., Mills, G. B., Shaw, K. R. M. et al., The cancer genome atlas pan-cancer analysis project. *Nat. Genet.* 2013, *45*, 1113–1120.
- [97] Law, V., Knox, C., Djoumbou, Y., Jewison, T. et al., Drugbank 4.0: shedding new light on drug metabolism. *Nucleic Acids Res.* 2014, *42*, D1091–D1097.
- [98] Wang, Y., Xiao, J., Suzek, T. O., Zhang, J. et al., PubChem: a public information system for analyzing bioactivities of small molecules. *Nucleic Acids Res.* 2009, *37*, W623–W633.
- [99] Lim, E., Pon, A., Djoumbou, Y., Knox, C. et al., T3DB: a comprehensively annotated database of common toxins and their targets. *Nucleic Acids Res.* 2010, *38*, D781–D786.
- [100] Scalbert, A., Andres-Lacueva, C., Arita, M., Kroon, P. et al., Databases on food phytochemicals and their health-promoting effects. *J. Agric. Food Chem.* 2011, *59*, 4331–4348.
- [101] Brown, J. A., Rudie, J. D., Bandrowski, A., Van Horn, J. D., Bookheimer, S. Y., The UCLA multimodal connectivity database: a web-based platform for brain connectivity matrix sharing and analysis. *Front. Neuroinform.* 2012, *6*.
- [102] Van Essen, D. C., Smith, S. M., Barch, D. M., Behrens, T. E. et al., The WU-Minn human connectome project: an overview. *Neuroimage* 2013, *80*, 62–79.
- [103] Human Microbiome Project Consortium, Structure, function and diversity of the healthy human microbiome. *Nature* 2012, *486*, 207–214.
- [104] Greene, C. S., Tan, J., Ung, M., Moore, J. H., Cheng, C., Big data bioinformatics. *J. Cell. Physiol.* 2014, *229*, 1896–1900.
- [105] Gligorijević, V., Pržulj, N., Methods for biological data integration: perspectives and challenges. *J. R. Soc. Interface* 2015, *12*, 20150571.
- [106] Vapnik, V. N., Vapnik, V., *Statistical Learning Theory*, Vol. 1, Wiley, New York 1998.
- [107] Scholkopf, B., Smola, A. J., *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, MIT Press, Cambridge, MA, USA 2001.
- [108] Freedman, D. A., *Statistical Models: Theory and Practice*, Cambridge University Press, New York, USA 2009.

- [109] Hartigan, J. A., *Clustering Algorithms* (99th Edn.), John Wiley & Sons, Inc., New York, NY, USA 1975.
- [110] Lee, D. D., Seung, H. S., Learning the parts of objects by non-negative matrix factorization. *Nature* 1999, *401*, 788–791.
- [111] Wang, F., Li, T., Zhang, C., Semi-supervised clustering via matrix factorization, in: SDM, SIAM, Atlanta, Georgia, USA 2008, 1–12.
- [112] Saria, S., Goldenberg, A., Subtyping: what it is and its role in precision medicine. *Intell. Syst. IEEE* 2015, *30*, 70–75.
- [113] Network, C. G. A., Comprehensive molecular characterization of human colon and rectal cancer. *Nature* 2012, *487*, 330–337.
- [114] Cancer Genome Atlas Network, Comprehensive molecular portraits of human breast tumours. *Nature* 2012, *490*, 61–70.
- [115] Cancer Genome Atlas Research Network, Integrated genomic analyses of ovarian carcinoma. *Nature* 2011, *474*, 609–615.
- [116] Wang, L., Li, F., Sheng, J., Wong, S. T., A computational method for clinically relevant cancer stratification and driver mutation module discovery using personal genomics profiles. *BMC Genomics* 2015, *16*, S6.
- [117] de Souto, M. C., Costa, I. G., de Araujo, D. S., Ludermir, T. B., Schliep, A., Clustering cancer gene expression data: a comparative study. *BMC Bioinformatics* 2008, *9*, 497.
- [118] Monti, S., Tamayo, P., Mesirov, J., Golub, T., Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data. *Mach. Learn.* 2003, *52*, 91–118.
- [119] Brunet, J.-P., Tamayo, P., Golub, T. R., Mesirov, J. P. Meta-genes and molecular pattern discovery using matrix factorization. *Proc. Natl. Acad. Sci. USA* 2004, *101*, 4164–4169.
- [120] Alcaraz, N., Pauling, J., Batra, R., Barbosa, E. et al., Key-pathwayminer 4.0: condition-specific pathway analysis by combining multiple omics studies and networks with cytoscape. *BMC Syst. Biol.* 2014, *8*, 99.
- [121] Vaske, C. J., Benz, S. C., Sanborn, J. Z., Earl, D. et al., Inference of patient-specific pathway activities from multi-dimensional cancer genomics data using paradigm. *Bioinformatics* 2010, *26*, i237–i245.
- [122] List, M., Hauschild, A. C., Tan, Q., Kruse, T. A. et al., Classification of breast cancer subtypes by combining gene expression and dna methylation data. *J. Integr. Bioinform.* 2014, *11*, 236.
- [123] Kristensen, V. N., Lingjærde, O. C., Russnes, H. G., Vøllan, H. K. M. et al., Principles and methods of integrative genomic analyses in cancer. *Nat. Rev. Cancer* 2014, *14*, 299–313.
- [124] Shen, R., Olshen, A. B., Ladanyi, M., Integrative clustering of multiple genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioinformatics* 2009, *25*, 2906–2912.
- [125] Curtis, C., Shah, S. P., Chin, S. F., Turashvili, G. et al., The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* 2012, *486*, 346–352.
- [126] Shen, R., Mo, Q., Schultz, N., Seshan, V. E. et al., Integrative subtype discovery in glioblastoma using icluster. *PLoS One* 2012, *7*, e35236.
- [127] Yuan, Y., Savage, R. S., Markowetz, F., Patient-specific data fusion defines prognostic cancer subtypes. *PLoS Comput. Biol.* 2011, *7*, e1002227.
- [128] Ben-Gal, I., *Bayesian Networks, Chapter 1*, John Wiley & Sons, Ltd, Hoboken, New Jersey, USA 2008.
- [129] Ray, P., Zheng, L., Lucas, J., Carin, L., Bayesian joint analysis of heterogeneous genomics data. *Bioinformatics* 2014, *30*, 1370–1376.
- [130] Lock, E. F., Hoadley, K. A., Marron, J., Nobel, A. B., Joint and individual variation explained (JIVE) for integrated analysis of multiple data types. *Ann. Appl. Stat.* 2013, *7*, 523–542.
- [131] Wang, B., Mezlini, A. M., Demir, F., Fiume, M. et al., Similarity network fusion for aggregating data types on a genomic scale. *Nat. Methods* 2014, *11*, 333–337.
- [132] Ng, A. Y., Jordan, M. I., Weiss, Y., On spectral clustering: analysis and an algorithm. *Adv. Neural Inf. Process. Syst.* 2002, *2*, 849–856.
- [133] Hofree, M., Shen, J. P., Carter, H., Gross, A., Ideker, T., Network-based stratification of tumor mutations. *Nat. Methods* 2013, *10*, 1108–1115.
- [134] Vogelstein, B., Papadopoulos, N., Velculescu, V. E., Zhou, S. et al., Cancer genome landscapes. *Science* 2013, *339*, 1546–1558.
- [135] Cai, D., He, X., Wu, X., Han, J., Non-negative matrix factorization on manifold. In *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on*. IEEE, Pisa, Italy 2008, pp. 63–72.
- [136] Gligorijević, V., Malod-Dognin, N., Pržulj, N., Patient-specific data fusion for cancer stratification and personalised treatment. *Pac. Symp. Biocomput.* 2015, 321–332.
- [137] Mo, Q., Wang, S., Seshan, V. E., Olshen, A. B. et al., Pattern discovery and cancer gene identification in integrated cancer genomic data. *Proc. Natl. Acad. Sci. USA* 2013, *110*, 4245–4250.
- [138] Napolitano, F., Zhao, Y., Moreira, V. M., Tagliaferri, R. et al., Drug repositioning: a machine-learning approach through data integration. *J. Cheminformatics* 2013, *5*, 30–39.
- [139] Wang, Y., Chen, S., Deng, N., Wang, Y., Drug repositioning by kernel-based integration of molecular structure, molecular activity, and phenotype data. *PLoS One* 2013, *8*, e78518.
- [140] Zheng, X., Ding, H., Mamitsuka, H., Zhu, S., Collaborative matrix factorization with multiple similarities for predicting drug-target interactions. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining* ACM 2013, pp. 1025–1033.
- [141] Zhang, P., Wang, F., Hu, J., Towards drug repositioning: a unified computational framework for integrating multiple aspects of drug similarity and disease similarity. In *AMIA Annual Symposium Proceedings, Vol. 2014* American Medical Informatics Association 2014, p. 1258.
- [142] Gottlieb, A., Stein, G. Y., Ruppin, E., Sharan, R., Predict: a method for inferring novel drug indications with application to personalized medicine. *Mol. Syst. Biol.* 2011, *7*, 496.



- [143] Huang, Y.-F., Yeh, H.-Y., Soo, V.-W., Inferring drug-disease associations from integration of chemical, genomic and phenotype data using network propagation. *BMC Med. Genomics* 2013, 6, 1–14.
- [144] Daminelli, S., Haupt, V. J., Reimann, M., Schroeder, M., Drug repositioning through incomplete bi-cliques in an integrated drug–target–disease network. *Integr. Biol.* 2012, 4, 778–788.
- [145] Jiang, W., Chen, X., Liao, M., Li, W. et al., Identification of links between small molecules and miRNAs in human cancers based on transcriptional responses. *Sci. Rep.* 2012, 2, 1–9.
- [146] Chung, S., Nakagawa, H., Uemura, M., Piao, L. et al., Association of a novel long non-coding rna in 8q24 with prostate cancer susceptibility. *Cancer Sci.* 2011, 102, 245–252.
- [147] Biswas, A. K., Kang, M., Kim, D. C., Ding, C. H. et al. Inferring disease associations of the long non-coding RNAs through non-negative matrix factorization. *Netw. Model. Anal. Health Inform. Bioinform.* 2015, 4, 1–8.
- [148] Dudley, J. T., Deshpande, T., Butte, A. J., Exploiting drug-disease relationships for computational drug repositioning. *Brief. Bioinform.* 2011, 4, 303–311.
- [149] Keiser, M. J., Setola, V., Irwin, J. J., Laggner, C. et al., Predicting new molecular targets for known drugs. *Nature* 2009, 462, 175–181.
- [150] Campillos, M., Kuhn, M., Gavin, A.-C., Jensen, L. J., Bork, P., Drug target identification using side-effect similarity. *Science* 2008, 321, 263–266.
- [151] Van Driel, M. A., Bruggeman, J., Vriend, G., Brunner, H. G., Leunissen, J. A., A text-mining analysis of the human phenome. *Eur. J. Hum. Genet.* 2006, 14, 535–542.
- [152] Zhou, X., Menche, J., Barabási, A.-L., Sharma, A. Human symptoms-disease network. *Nat. Commun.* 2014, 5, 1–10.
- [153] Ding, H., Takigawa, I., Mamitsuka, H., Zhu, S., Similarity-based machine learning methods for predicting drug–target interactions: a brief review. *Brief. Bioinform.* 2013, 15, 734–747.
- [154] Yamanishi, Y., Araki, M., Gutteridge, A., Honda, W., Kanehisa, M., Prediction of drug-target interaction networks from the integration of chemical and genomic spaces. *Bioinformatics* 2008, 24, i232–i240.
- [155] Minai, R., Matsuo, Y., Onuki, H., Hirota, H., Method for comparing the structures of protein ligand-binding sites and application for predicting protein-drug interactions. *Proteins* 2008, 72, 367–381.
- [156] Chen, Y., Zhi, D., Ligand-protein inverse docking and its potential use in the computer search of protein targets of a small molecule. *Proteins* 2001, 43, 217–226.
- [157] Li, J., Zheng, S., Chen, B., Butte, A. J. et al., A survey of current trends in computational drug repositioning. *Brief. Bioinform.* 2015, 1, 11.
- [158] Li, Y., Jones, S., Drug repositioning for personalized medicine. *Genome Med.* 2012, 4, 27.
- [159] Berry, M. W., Browne, M., Langville, A. N., Pauca, V. P., Plemmons, R. J., Algorithms and applications for approximate nonnegative matrix factorization. *Comput. Stat. Data Anal.* 2007, 52, 155–173.
- [160] Fawcett, T., An introduction to ROC analysis. *Pattern Recogn. Lett.* 2006, 27, 861–874.
- [161] Vanunu, O., Magger, O., Ruppin, E., Shlomi, T., Sharan, R., Associating genes and protein complexes with disease via network propagation. *PLoS Comput. Biol.* 2010, 6, e1000641.
- [162] Wu, Z., Wang, Y., Chen, L., Network-based drug repositioning. *Mol. BioSyst.* 2013, 9, 1268–1281.
- [163] Del Vescovo, V., Grasso, M., Barbareschi, M., Denti, M. A., Micrnas as lung cancer biomarkers. *World J. Clin. Oncol.* 2014, 5, 604–620.
- [164] Hayes, J., Peruzzi, P. P., Lawler, S., Micrnas in cancer: biomarkers, functions and therapy. *Trends Mol. Med.* 2014, 20, 460–469.
- [165] Wen, X., Deng, F.-M., Wang, J., Micrnas as predictive biomarkers and therapeutic targets in prostate cancer. *Am. J. Clin. Exp. Urol.* 2014, 2, 219–230.
- [166] Mercer, T. R., Dinger, M. E., Mattick, J. S., Long non-coding RNAs: insights into functions. *Nat. Rev. Genet.* 2009, 10, 155–159.
- [167] Chen, X., Predicting lncRNA-disease associations and constructing lncrna functional similarity network based on the information of miRNA. *Sci. Rep.* 2009, 5, 13186.
- [168] Carlsson, G., Topology and data. *B. Am. Math. Soc.* 2009, 46, 255–308.
- [169] Lum, P. Y., Singh, G., Lehman, A., Ishkanov, T. et al., Extracting insights from the shape of complex data using topology. *Sci. Rep.* 2013, 3.
- [170] Hoens, T. R., Polikar, R., Chawla, N. V., Learning from streaming data with concept drift and imbalance: an overview. *Prog. Artificial Intelligence* 2012, 1, 89–101.
- [171] Stella, F., Amer, Y., Continuous time bayesian network classifiers. *J. Biomed. Inform.* 2012, 45, 1108–1119.
- [172] Ge, H., Walkout, A. J., Vidal, M., Integrating “omic” information: a bridge between genomics and systems biology. *Trends Genet.* 2003, 19, 551–560.
- [173] Mani, R., St.Onge, R. P., Hartman, J. L., Giaever, G., Roth, F. P., Defining genetic interaction. *Proc. Natl. Acad. Sci. USA* 2008, 105, 3461–3466.
- [174] Gligorijević, V., Janjić, V., Pržulj, N., Integration of molecular network data reconstruct gene ontology. *Bioinformatics* 2014, 30, i594–i600.
- [175] Ho, J. C., Ghosh, J., Sun, J., Marble: high-throughput phenotyping from electronic health records via sparse nonnegative tensor factorization. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining ACM*, New York, USA 2014, pp. 115–124.
- [176] Gottesman, O., Kuivaniemi, H., Tromp, G., Faucett, W. A. et al., The electronic medical records and genomics (emerge) network: past, present, and future. *Genet. Med.* 2013, 15, 761–771.
- [177] Vaudel, M., Verheggen, K., Csordas, A., Raeder, H. et al., Exploring the potential of public proteomics data. *Proteomics* 2016, 16, 214–225.