Briefings in Bioinformatics, 2022, pp. 1-24

doi: DOI HERE

Heterogeneous data integration methods for patient similarity networks Heterogeneous data integration methods for patient similarity networks

HETEROGENEOUS DATA INTEGRATION METHODS FOR PATIENT SIMILARITY NETWORKS

Heterogeneous data integration methods for patient similarity networks

Jessica Gliozzo,^{1,2,5} Marco Mesiti,^{1,5} Marco Notaro,^{1,5} Alessandro Petrini,^{1,5} Alex Patak,² Antonio Puertas-Gallardo,² Alberto Paccanaro,^{3,4} Giorgio Valentini^{1,5,6,7} and Elena Casiraghi^{1,5 *}

FOR PUBLISHER ONLY Received on Date Month Year; revised on Date Month Year; accepted on Date Month Year

Abstract

Patient similarity networks (PSNs), where patients are represented as nodes and their similarities as weighted edges, are being increasingly used in clinical research. These networks provide an insightful summary of the relationships among patients and can be exploited by inductive or transductive learning algorithms for the prediction of patient outcome, phenotype, and disease risk. PSNs can also be easily visualized, thus offering a natural way to inspect complex heterogeneous patient data, and providing some level of explainability of the predictions obtained by machine learning algorithms. The advent of high-throughput technologies, enabling us to acquire high-dimensional views of the same patients (e.g. omics data, laboratory data, imaging data) calls for the development of data fusion techniques for PSNs in order to leverage this rich heterogeneous information. In this article, we review existing methods for integrating multiple biomedical data views to construct PSNs, together with the different patient similarity measures that have been proposed. We also review methods that have appeared in the machine learning literature but have not yet been applied to PSNs, thus providing a resource to navigate the vast machine learning literature existing on this topic. In particular, we focus on methods that could be used to integrate very heterogeneous datasets, including multi-omics data as well as data derived from clinical information and medical imaging.

Key words: patient similarity networks, biomedical applications, multimodal data, data fusion

1 Introduction

- 2 In the last decades, medical research has begun to move from a population-
- 3 based perspective to a personalized one, often referred to as Precision
- 4 Medicine, where patients' biomedical characteristics are leveraged for
- 5 diagnosis, prognosis and choice of appropriate treatment [1, 2]. In this
- 6 context, it is widely accepted that if two patients share similar clinical
- 7 variables and omics profiles, their clinical outcomes should also be similar.
- 8 Pairwise similarities between patients have a natural representation as
- 9 graphs Patient Similarity Networks (PSN) where nodes represent
- 10 patients and edges represent the similarity between patients calculated
- 11 using their clinical and/or biomolecular features. In this framework
- 12 unsupervised clustering methods and supervised classification models that

leverage similarities between patients have been successfully applied to13stratify patients and to predict their phenotype or clinical outcome [3, 4, 5,146, 7, 8]. Representing data as graphs provides several advantages, including15interpretability and privacy [9], as patient specific information cannot be16recovered from the similarity measures.17

The increasing availability of high-throughput technologies able to generate high-dimensional, distributed biomedical datasets, ranging from multi-omics [8] to imaging [10], clinical and demographic data [11], calls for approaches to mine and aggregate salient information [12] with the ultimate aim of building PSNs integrating such diverse datasets. However, the majority of PSNs that have been proposed are built using only one source of information. At the same time, several methods that can integrate 24

¹AnacletoLab - Computer Science Department, Università degli Studi di Milano, Via Celoria 18, 20135, Milan, Italy, ²European Commission, Joint Research Centre (JRC), Ispra (VA), Italy, ³Department of Computer Science, Royal Holloway, University of London, Egham, TW20 0EX UK, ⁴School of Applied Mathematics (EMAp), Fundação Getúlio Vargas, Rio de Janeiro Brazil, ⁵CINI, Infolife National Laboratory, Roma, Italy, ⁶DSRC UNIMI, Data Science Research Center, Milano, 20135, Italy and ⁷ELLIS, European Laboratory for Learning and Intelligent Systems, Berlin, Germany

^{*}Corresponding author: Elena Casiraghi - email: elena.casiraghi@unimi.it

heterogenous sources of information into graph structures have appeared 25 in the past decades in the biomedical and machine learning literature. 26 27 In this article, we review existing methods for integrating multiple biomedical data views to construct patient similarity networks. Since the 28 29 type of data being integrated and the specific integration method must 30 be coupled with an appropriate choice of similarity measure, we will also 31 discuss different similarity measures. Importantly, this paper also reviews methods for integrating information into graph structures that appeared 32 33 in the machine learning literature but have not yet been used for PSNs. We believe that this will be beneficial for the reader, providing a resource 34 to navigate the vast machine learning literature existing on this topic, 35 and possibly inspire the use and development of novel techniques of data 36 integration for PSNs. Moreover, unlike earlier reviews (see e.g. [13, 14, 15, 37 38 8, 16]), we focus on methods that may be used for patients' classification and clustering that integrate not only multi-omics data, but also clinical and 39 40 image sources.

41 We propose a taxonomy that groups existing methods for building 42 PSNs into three main categories. PSN-fusion methods [3, 17, 6] build different PSNs, one for each data source, that are then fused together into a single 43 PSN. Input data-fusion methods [18, 19, 20, 21] combine the different data 44 45 sources into a single dataset that is then used for building a single PSN. Finally, Output-fusion methods [22, 23, 24] build different PSNs, one for each 46 47 data source, that are analysed separately, and results are then combined. 48 Other multimodal data-fusion surveys not specific for PSNs have been proposed in the bioinformatics field by adopting different taxonomies 49 (schematized in Fig. 1, Appendix A). Some taxonomies focus on the 50 type of multi-datasets being integrated, thus identifying horizontal 51 52 integration techniques [25] (top of Fig. 1-yellow box) and vertical integration techniques [25] (top of Fig. 1-light blue box). While the former fuse 53 54 homogeneous multisets [26], i.e. multimodal datasets where each view produces the same data type under different settings, the latter integrate 55 the classic heterogeneous [26] multimodal datasets. Vertical integration 56 57 techniques are further classified into methods applying a hierarchical (alias multi-staged [27]) integration flow, where ground knowledge about 58 the relationships between the different views is considered during the 59 integration, and methods applying a parallel (alias meta-dimensional [27]) 60 61 integration flow (bottom of Fig. 1-red-dashed box), where each view is processed in a similar but independent way. Parallel integration methods 62 are the most diffused in literature given their generalizability. For 63 this reason several reviews concentrate solely on them and introduce 64 taxonomies that distinguish, e.g., model-agnostic versus model-dependent 65 methods [28], or exploit an early-intermediate-late taxonomy [29, 30, 31, 27, 66 32, 33] (described in detail in Appendix A). 67 Anyhow, each review paper focuses on different aspects of the 68 multimodal data integration. For example, some works solely focus on 69

integrative unsupervised clustering techniques [34] or supervised multi-70 71 omics prediction models [29, 35, 33], or survey data-fusion techniques that are either applied to multi-omics data [36, 26, 27, 25, 16], or that apply 72 specific data-fusion techniques (e.g. integrative Bayesian models [37, 13] or 73 multimodal neural networks [38]). 74

Unlike previous reviews, this work specifically focuses on integrative 75 methods for PSN-based models integrating not only multi-omics data, but 76 also clinical and imaging sources. Each method is critically described to 77 highlight its main advantages and drawbacks, enabling the reader to select 78 the most appropriate approach to answer her/his scientific questions. 79

Given a set of patients and their corresponding clinical and biomolecular 80 81 features, the topology of the corresponding PSN depends crucially on how the similarity measure is calculated. Therefore, we begin describing the 82 similarity measurement methods presented in the literature. Our taxonomy 83 of existing methods for building PSNs is described in Sections PSN-fusion 84

methods and Input data-fusion and output-fusion methods. Tables 3 - 8 85 summarize the most relevant methods we surveyed.

PSN construction

The construction of the PSN is a crucial step in PSN analysis models, whose effectiveness mainly depends on the available multimodal datasets from which samples are extracted and on the choice of the measure exploited for pairwise similarity computation between samples.

Several kinds of similarity measures have been adopted in literature for PSN construction: classic distance metrics tailored to the data type [39, 40]; kernel functions [41, 42] that substitute distance metrics; "kernels on graphs" methods [43]. In the remainder, we discuss their characteristics.

The usage of classic (opportunely inverted) distances or similarity 96 metrics [39, 40] is often preferred when the data types are normalized and 97 homogeneous. As an example, PSNs on continuous, normalized, data have 98 been constructed by using the cosine similarity [44, 5], or the Euclidean [45] 99 or Mahalanobis distance [45]; PSNs on discrete data types have been built 100 by exploiting the Chi-squared distance [3, 6]; binary data has been handled 101 by using the Jaccard distance [46] or many other distance measures (see [47] 102 for a list of 76 metrics and measures specifically designed for binary data). 103

When data-blocks with heterogeneous and/or normalized variable 104 types are available, more articulated schemas [48, 6] have been proposed to 105 integrate different similarity metrics into a unique measure. As an example, 106 in [48] the authors proposed a supervised Cox regression model to initially 107 learn a weight for each variable; the learnt weights are then used to compute 108 a similarity score as a weighted sum of individual similarities obtained on 109 each feature by using standard metrics. In this way, different similarity 110 metrics can be used on the different variables based on their type, and 111 the influence of each variable to the global similarity score is weighted on 112 the prediction (e.g. survival time when using a Cox regression model). On 113 the other hand, when dealing with datasets composed by continuous non-114 normalized variable types, Pai et al. [6] propose computing the average of 115 all the normalized similarities over each variable, where the normalization 116 is essentially a min-max normalization. 117

When dealing with complex problems, literature works often rely 118 on Kernel functions [49] for PSN computation. The rationale behind this 119 choice is based on the assumption that point separability is often improved 120 after a non-linear projection of points into a higher-dimensional space. 121 Kernel functions are particularly appealing in this context since they 122 express pairwise distances in a higher-dimensional space by directly using 123 the (lower-dimensional) input samples, therefore avoiding the expensive 124 explicit computation of a non-linear higher-dimensional mapping followed 125 by pairwise similarity evaluation (using the well-known kernel-trick). Even 126 in this case the choice of the kernel function must be tailored to the 127 data type which is crucial to obtain reliable results. In this context, PSNs 128 are often computed in literature methods working on biomedical data 129 by using classic parametric normalized linear kernels [30, 50], polynomial 130 kernels, or Gaussian kernels [51, 52], whose parameters are tuned to optimize 131 performance. As an example, the prognostic approach presented in [30] 132 obtains a set of unimodal PSNs by applying normalized linear kernels on 133 each of the data-sources containing clinical and multi-omics data sets. In 134 this case, the usage of the same kernel function on different sources is 135 appropriate because they are characterized by the same data type (real-136 valued data type). 137

In a subsequent work [53], the same authors extend the data set by 138 including categorical and integer data types; therefore, they substitute 139 the linear kernels with a set of kernels tailored on each data type being 140 processed. Of note, the kernels used in [30, 53] are always normalized. 141 This is a crucial characteristic when integrating multiple kernels because 142

87 88

89

90

91

92

93

94

95



Fig. 1. Schema of the main taxonomies proposed in literature for categorizing multimodal integration methods. Considering the data integration flow, literature works identify two broad classes: horizontal integration approaches and vertical integration approaches. Horizontal integration approaches fuse multisets (i.e. datasets where each view is acquired by the same source under different conditions) by independently applying the same process on each view and then pooling the individual results. On the other hand, vertical integration approaches fuse multimodal datasets (i.e. datasets composed by semantically different views) through more complex techniques, further categorized as hierarchical-vertical integration methods and parallel-vertical integration techniques. The former fuse data views following a hierarchy driven by biological a-priori knowledge while the latter do not exploit knowledge-based dependencies between views. Parallel-vertical integration methods are the most diffused integration methods; they are further classified based on the phase when the data integration-step is performed w.r.t. the model construction (red-dashed box). Thus, methods are divided in (l) early approaches, which integrate the data types before model construction, (III) late approaches, which integrate the classified based from each view and subsequently integrated. Of note, the latter class of approaches is more dependent on the exploitel are model, and (III) intermediate approaches where intermediate models are obtained from each view and subsequently integrated. Of note, the latter class of approaches is more dependent on the exploitel dearning model, which is the reason why they have been also classified as *model-dependent* methods opposed to *model-agnostic* methods (blue-dashed boxes). We refer interested readers to Appendix A.

143 comparable kernel scales are obtained, therefore facilitating the kernel

144 integration. Moreover, in the case of kernel-aggregation systems exploiting

145 weighted averages of the unimodal kernels, normalization also improves

the interpretability of the computed integration weights, the latest being

directly related to the importance of their respective kernel [53].

148A recent advance in the field of PSN analysis is provided by149unsupervised methods that compute the PSN through the scaled exponential150Euclidean kernel [3] and its modifications [54, 55]. They essentially apply a151local normalization of the distance between a central node and any of its152neighbors, so that distances are independent from the neighborhood scales.

153 Their application in the context of unsupervised patient clustering through

PSN analysis has obtained promising results [3] (see Section SNF-based 154 methods). 155

Given its effectiveness, the scaled Euclidean distance has been extended156in [54] to deal with heterogeneous data types containing continuous157and boolean variables. More precisely, the similarity on boolean data is158measured by using the weighted Hamming distance with weights computed159by supervised approaches or pre-set based on existing knowledge. Further,160in [55] the authors propose adopting the Chebyshev distance instead of the161Euclidean distance.162

Gliozzo et al. [7] extend to PSNs a previous *kernel-based* approach 163 originally applied to the semi-supervised analysis of biomolecular networks 164 [56]. More precisely, the authors obtain promising outcome predictions on
unimodal PSNs by firstly using the filtered Pearson correlation (by setting
to zero all negative values) to measure similarities between unimodal gene
expression profiles, and then applying a random walk kernel to strengthen
high similarities while diminishing low ones. The neighborhoods identified
in the obtained PSN are then used to compute a score for each patient,

which is thresholded to obtain the desired classification. While unimodal
PSNs are exploited in [7], the works proposed in [57] and [58] exploit
random-walks to compute similarities in a multimodal setting.

To improve informativeness, Tables 1 and 2 sketch the similarity measures/ methods used for PSN construction by notable literature works exploiting multimodal datasets; for each paper we report the data types of the different data-sources exploited for the investigation, and the similarity measures/methods used for building the corresponding unimodal PSNs.

Even if a wide range of similarity computation methods has been proposed in literature, a consensus on which strategy performs better on specific data types and problems in the context of precision medicine is still lacking. Some tentative experiments have been conducted for determining the best performing strategies (see e.g. [59, 60]), but the lack of common benchmark datasets prevents an unbiased comparison of the different

185 proposed approaches.

186 **PSN-fusion methods**

PSN-fusion methods have been specifically developed to process a set 187 of unimodal PSNs and produce an integrated PSN. In Fig. 2 we sketch 188 the generic workflow of the PSN-fusion methods. They start by building 189 unimodal PSNs on each data source or data-type (Fig. 2-A). Mind that the 190 191 choice of the similarity measure/kernel function used to build each PSN (Section PSN construction) is crucial for obtaining informative unimodal 192 PSNs, which would otherwise hamper the achievement of successful 193 results. Next, the aggregation of the unimodal PSNs (Fig. 2-B) is performed 194 by either Multiple Kernel Learning methods (MKL, Section MKL-based 195 196 methods, Table 3) which run optimization algorithms inherited from the machine learning field to find the optimal weights of an additive unimodal 197 kernel aggregation, or approaches stemming from the seminal Similarity 198 Network Fusion algorithm (SNF - [3], Section SNF-based methods, Table 199 200 4), which use different strategies to diffuse the similarity information both between neighboring nodes in each unimodal PSN and between 201 corresponding nodes in different PSNs, or other network-based approaches 202 (Section Other PSN-fusion methods and Table 5). 203

The integrated PSN may be finally used as input to unsupervised clustering methods aiming at, e.g., identifying patients' subtypes, or supervised classification methods predicting, e.g., patients' risk, prognosis, or outcome (Fig. 2-C).

208 MKL-based methods

Inheriting theories and algorithms from the machine learning fields, MKL
methods [17, 65, 66, 64] view the unimodal PSNs as kernels and propose
their optimal additive combination, as a weighted sum of the available
unimodal kernels. In this context, "optimality" refers to either a *supervised*setting or an *unsupervised* one. *supervised MKL* algorithms (e.g. simpleMKL [17]) exploit a supervised

classifier model designed to work on the fused kernel. Supervision is guaranteed by the availability of a training set composed of samples whose labels are known. Such training set is used by the chosen supervised MKL method to solve a constrained optimization problem that finds the kernel weights and classifier hyper-parameters maximizing the classification accuracy on the training set. On the other side, *unsupervised MKL* methods make no use of labeled samples, but instead solve an optimization problem to find the weights that essentially lead to the maximum alignment between the integrated kernel and any of the input unimodal kernels.

222

223

Recent PSN-fusion methods exploiting a supervised MKL strategy are224those presented by [30, 53, 64, 50, 67]. The work proposed in [50] designs225specific kernels for each omic type in the TCGA cancer dataset and then226computes the kernel weights by using the training set to optimize the fit of227a Cox-survival model.228

All the other works [30, 53, 64, 67] share the use of the kernelized 229 Support Vector Machine (svm) classifiers [68], opportunely modified 230 as defined in [17] and [66] to work on the kernel resulting from an 231 optimal additive sum. In particular, the works proposed by Daemen 232 et al. [30, 53] aggregate specific kernels on each clinical data type and 233 uses a classic svm optimization strategy to derive the optimal weights, 234 while the works proposed in [64] and [67] use the easyMKL algorithm 235 to optimize an svm aggregating multiple kernels defined over multimodal 236 datasets also including opportunely coded imaging sources. More precisely, 237 in [64] authors use the same Gaussian kernels to process both the real 238 CerebroSpinal Fluid (CSF) biomarkers features and the shape and texture 239 features extracted to code Magnetic Resonance Images (MRI). On the 240 other side, the work proposed in [67] improves upon the work presented 241 in [69] and defines specific kernels for the multi-omics data from the 242 TCGA cancer dataset and for the features automatically extracted from 243 histopathological images (see Table 3). The effectiveness of the simpleMKL 244 strategy is witnessed by its several extensions (easyMKL [70], SEMKL [71], 245 SpicyMKL [72]). 246

As expected, our literature search highlighted that SVMs are the most 247 widely used base-learner models in conjunction with MKL in the context 248 of biomedical predictions; however, some authors have also presented 249 MKL methods using Multiple Kernel Fisher Discriminant Analysis (MK-250 FDA [73]) or Kernel Regularized Discriminant Analysis (KRDA, [74]) as 251 base learners where the single kernel is substituted by multiple kernels. 252 Though these strategies have not been applied on patients' data, their 253 promising results on the protein sub-cellular localization prediction task 254 [73, 75] suggest they could be good options for developing a multimodal 255 PSN analysis task. 256

Unsupervised MKL approaches are described in the works of [76, 77, 257 52]. The regularized MKL with Locality Preserving Projection algorithm 258 (rMKL-LPP [76]) is an unsupervised, regularized MKL-based clustering 259 approach for the identification of cancer subtypes from multi-omics data. 260 It builds upon the MKL-DR model proposed in [78] to constrain the 261 optimization problem by handling the "small-sample-size" problems caused 262 by the high dimensionality of the input data-sources and exploits the 263 theories at the base of the locality preserving projection algorithm [79] to 264 find the integrated kernel in a lower-dimensional space that maintains the 265 local neighborhoods relationships. In other words, the model minimizes 266 a function that allows finding both the hyper-parameters of the multiple 267 kernels and their combination weights so that patients that are similar 268 according to "many" input sources (kernels) remain neighbors in the 269 integrated kernel. Further, to avoid restricting the usage of only one kernel 270 per data-source or data-type, authors add a constrained regularization that 271 avoids overfitting, so that multiple kernels can be used for each source 272 without risking to over-fit the data. Similar topological constraints are 273 used by [52] to compute kernel weights such that the resulting integrated 274 kernel maintains the neighborhood-relationship described above, and at 275 same time maximizes the alignment (similarity) to all the input kernels. 276

By contrast, Liu et al. [77] leverage the standard kernel k-means 277 clustering [80], which applies k-means in the kernel space, to a *multiple* 278 *kernel k-means clustering (MKKM)* that considers the relationships between 279 all the input kernels. The optimal clusters are found by minimizing a loss that measures the intra-class sample distance as a function of the cluster 281 assignment matrix and the kernel weights. However, differently from other 282

Reference	Data type	Data	Similarity measure/method
[46]	binary	ICD-9 diagnosis code	Jaccard similarity
[44]	continuous, categorical, discrete	clinical data	cosine similarity
[5]	continuous categorical, discrete	clinical data	cosine similarity
[61]	continuous	mRNA, PPI	Pearson correlation
[7]	continuous	mRNA	Pearson correlation
[6]	continuous	clinical variables individual gene genes in pathways/networks	mean of normalized difference normalized difference Pearson correlation
	discrete	categorical-ordinal variable (e.g. tumour stage) unbinned counts (e.g. mutation data) matrix scores	normalized difference shared incidence in a grouped unit chi-square distance
[3]	continuous discrete binary	mRNA, miRNA, DNA methylation	scaled exponential kernel of Euclidean distance chi-squared distance agreement-based measure
[54]	continuous binary	mRNA, DNA methylation somatic mutation	scaled exponential kernel of weighted Euclidean distance scaled exponential kernel of weighted Hamming distance
[62]	continuous	mRNA, miRNA, DNA methylation	scaled exponential kernel of Euclidean distance
[63]	categorical, discrete	demographic, APOE4 allele status, MRI	squared-exponential kernel
[55]	continuous	gene expression, miRNA, isoform expression	kernel of Chebyshev distance
[48]	continuous, categorical, discrete	clinical data	weighted sum of distances with weight determined by a scale Cox regression coefficient

Table 1. Similarity measures/methods used in literature to build PSNs. For notable works in literature the table reports: the reference of the literature work presenting a multimodal PSN analysis method (column *Reference*), the data type (column *Data type*) of the different sources (column *Data*) exploiting for the investigation, and the similarity measures/methods exploited for building the unimodal PSNs.

Abbreviations

ICD-9: International Classification of Diseases Version 9; CNV: Copy Number Variation; miRNA: micro RNA; MRI: Magnetic Resonance Imaging; mRNA: messenger RNA; PPI: Protein-Protein Interaction;

multiple kernel clustering models, the MKKM loss function includes a termthat promotes the choice of higher weights for uncorrelated kernels.

285 SNF-based methods

- 286 PSNs are similarity graphs by definition; therefore, recent promising works
- 287 apply graph-based algorithms and theories to integrate them. In particular,
- 288 some authors simply integrate the information from different similarity
- graphs by using graph kernels [57], or by averaging [58, 81].

On the other side, Similarity Network Fusion (SNF [3]) exploits a 290 nonlinear message-passing algorithm [82] that diffuses the information 291 between all the unimodal PSNs constructed on each data-block until they 292 converge to the integrated PSN. The diffusion process is designed so that 293 the similarity between any two points computed over a specific source is 294 updated and diffused if the two points are neighbors or share common 295 neighbours in the other modalities. SNF has proven to be successful 296 when compared to relevant PSN-fusion methods [83] in the unsupervised 297 clustering task on three real, complex, multi-omics datasets (murine liver -298

Reference	Data type	Data	Similarity measure/method
[30, 53]	continuous categorical, discrete, binary	mRNA, clinical clinical	normalized linear kernel
[64]	discrete continuous	MRI CSF	gaussian kernel
[51]	continuous discrete	mRNA, miRNA, CNV, DNA methylation, clinical	gaussian kernel
[50]	continuous, binary, discrete	mRNA, miRNA, CNV, DNA methylation, RPPA, somatic mutations, clinical data	linear kernel
[65]	continuous	mRNA, CNV, DNA methylation	normalized linear kernel, normalized polynomial kernel, normalized gaussian kernel
[53]	continuous, categorical (ordinal) categorical (nominal)	clinical variables	absolute difference of values/ranks of two subjects compared and rescaled using variable range kernel defined using Kronecker delta function
[57]	continuous, binary	mRNA, RPPA, somatic mutation	novel graph kernel called SmSPK

Table 2. Similarity measures/methods used in literature to build PSNs. For notable works in literature the table reports: the reference of the literature work presenting a multimodal PSN analysis method (column *Reference*), the data type (column *Data type*) of the different sources (column *Data*) exploiting for the investigation, and the similarity measures/methods exploited for building the unimodal PSNs.

Abbreviations

CSF: CerebroSpinal Fluid; CNV: Copy Number Variation; miRNA: micro RNA; MRI: Magnetic Resonance Imaging; mRNA: messenger RNA; RPPA: Reverse-Phase Protein Arrays.

BXD [84], platelet reactivity [85], and Breast Cancer dataset from TCGA BRCA [86]).

Several works extended SNF in different ways, thus creating a group of 301 302 algorithms (called SNF-based methods). As an example, Affinity Network Fusion (ANF) [87] has been developed to diminish the computational costs 303 of SNF, by reducing the iterative integration strategy of SNF to a unique 304 step. To this aim, authors design a multigraph where each layer corresponds 305 to a source-specific PSN, and then apply the one-step random walk kernel, 306 307 where user-defined parameters are the transition probabilities between 308 different layers, and the PSN for a specific layer represents the transition probabilities between nodes in that layer. When tested on multiple TCGA 309 datasets, AFN outperforms SNF both in terms of clustering efficacy and 310 311 computational costs.

By taking into account that the Euclidean distance metric employed 312 in SNF suffers the curse of dimensionality [88] and may affect the results, 313 [89] presented HSNF (hierarchical SNF), which essentially runs SNF several 314 times, where each iteration uses a set of unimodal PSNs, generated on each 315 316 data-block by using a randomly sampled feature set. At each iteration, the computed PSNs are fused with the integrated network computed in the 317 precedent steps through SNF. The method is evaluated by its capacity to 318 identify cancer subtypes by applying spectral clustering on the integrated 319 matrix. Though outperforming SNF on several cancer datasets, HSNF has 320 321 a higher computational cost because of the iteration of SNF. To reduce noise in the integrated network, the Similarity Kernel Fusion 322 323

iterative update function is added to control the amount of information 326 to be retained from the integrated kernel at the preceding step. When 327 compared to SNF and to a simple average fusion of different kernels, SKF 328 obtains comparable or even better performance in the discovery of cancer 329 subtypes from real cancer datasets. 330

The association-signal-annotation boosted similarity network fusion 331 (ab-SNF) method [54] tries to improve SNF by considering a weighted 332 version of distance measures with the goal to up-weight signal features 333 and down-weight noisy ones. In this work, the weight for continuous 334 variables consists in a p-value computed by the univariate t-test to assess 335 the feature significance in predicting the outcome variable; the weights 336 for binary features, such as mutation data, are obtained by considering 337 prior knowledge from databases (e.g. 1 for features related to cancer and 0 338 otherwise). Given the computed weights, the unimodal PSNs are obtained 339 by using the scaled exponential kernel [3], where the Euclidean distance is 340 substituted by the weighted Euclidean distance, for continuous variables, or 341 the weighted Hamming distance, for binary variables. The use of feature-342 level weights leads to superior performance in clustering accuracy with 343 respect to SNF on both simulated and real data, while subtypes captured 344 by ab-SNF are significant in terms of patient survival on real cancer data. 345

Other PSN-fusion methods

To reduce noise in the integrated network, the Similarity Kernel Fusion algorithm (SKF) [90] multiplies the PSN built by using SNF with a matrix of weights, where the weight is higher if two samples are included in each other neighbourhood. Moreover, different from SNF, a term in the

NetDx [6] fuses unimodal PSNs by a simple weighted network sum, where347the weights for each network are identified by ridge regression to a target348network constructed on the training patients in order to enforce higher349

346

PSN-fusion methods



Fig. 2. High-level representation of PSN-fusion methods. (A) Given a set of matrices, each representing the patients vectors acquired from one source, proper similarity measures or kernel functions are used to build a set of unimodal PSNs (one PSN per data-source or data-type); (B) all the PSNs are then fused through either MKL methods, SNF methods or other PSN-fusion approaches; (C) the integrated PSN is processed either by unsupervised clustering algorithms for solving, e.g., patients' subtype prediction tasks, or by supervised classifier models for, e.g., patients' outcome prediction.

similarities between positive nodes and lower similarities between nodesbelonging to different classes.

Some recent integration methods propose integrating the different 352 353 PSNs by using a graph-based construction, and then compute integrated similarities by visiting the graph through random walk kernels. As an 354 example, [58] propose computing similarities over a multiplex graph 355 356 composed by a collection of PSNs (layers) each built on an individual 357 data-block. The different layers share the same set of nodes [91], and 358 corresponding nodes in different layers are connected to guarantee connectivity across multiple layers, but are considered as different entities 359 to avoid disrupting the difference between the multiple views available 360 for each node/sample. Then authors use the random walk kernel with 361 362 restart [92] to express the similarities as the probabilities of reaching a node

in a specific layer when another node in the same or in another layer is 363 used as the starting point of the walk. To account for multimodality, that is 364 with the presence of multiple layers, the probability of "jumping" to another 365 layer during the walk is weighted by a parameter λ . The probabilities 366 are computed by an iterative process that continues until a stationary 367 point is reached. RWRNF [58] is an extension of this method that allows 368 connecting multiple layers by also using edges between neighbourhoods of 369 corresponding nodes. The use of many random walks, starting from all the 370 nodes in each layer, adjusts the weights of the multiplex network taking 371 into account its global topology. Finally, an integrated similarity network 372 is computed by averaging corresponding weights across different layers of 373 the network. 374

Name	Matched Samples	Dataset	Sample Cardinality	Data type	Integration approach	Task	Code and Language
PAMOGK ¹ [57]	x	TCGA KIRC NCI-PID at NDEXBio	361	somatic mutation mRNA RPPA	MKKM [77]	Unsupervised Clustering (Patient subtype identification)	MATLAB, Python code
[64]	x	ADNI	120	CSF features MRI	MKL [17]	Supervised Classification (HC vs MCI patients)	
[69]	×	TCGA	585	Histopathological images clinical data	simpleMKL [17]	Supervised Classification (Patient's Prognosis)	
				mRNA methy RPPA			
[67]	×	TCGA GBM	125	Histopathological images CNV mRNA miRNA	simpleMKL [17]	Supervised Classification (Patient's Prognosis)	
MK-FDA [73] [75]	x	Protein dataset	Not provided	protein sequences	MKL	Supervised Multiclass Classification (Protein subcellular localization)	
[50]	×	14 TCGA datasets	3382	gernline variants somatic mutation CNV mRNA miRNA methy	MKL	Supervised Classification (Patient's Survival)	
[52]	х	TCGA from mixOmics	989	mRNA miRNA methy	MKL	Unsupervised Clustering (Patients' subtype identification)	
rMKL-LPP [76]	×	TCGA GBM TCGA BIC TCGA KIRC TCGA LUSC TCGA COAD	213 105 1122 106 92	mRNA miRNA methy	MKL	Unsupervised Clustering (Patient subtype identification)	
Abbreviations ADNI: Alzheime DNA methylation	r's Disease N ; miRNA : mi	leuroimaging Initiative; CN icro RNA; MKKM: Multiple	V: Copy Number kernel k-means cl	Variation; CSF: CerebroSpina ustering; MKL: Multiple Kern	l Fluid; HC: He el Learning: MR	althy Control; MCI: Mild Cognitive Impa J. Magnetic Resonance Imaging; mRNA : 1	airment; methy : messenger RNA;

Table 3. MKL-based PSN-fusion methods. For each method, the table reports: the name/acronym with the corresponding reference paper; whether it requires the same set of patients across all data modalities (i.e. "Matched Samples"); the dataset used to develop and evaluate the approach in the reference paper and the corresponding sample cardinality and data types composing the dataset; the exploited integration method; the application task and the code availability (with link to the repository and programming languages for which the code is available).

The efficacy provided by the use of similarities computed across local 375 neighborhoods is proven by its use in simpler unsupervised PSN analysis 376 methods. As an example, NEMO (NEighborhood based Multi-Omics 377 clustering, [93, 94]) is an unsupervised clustering approach where authors 378 use a scaled normalized euclidean kernel to compute similarities, which are 379 then made symmetric in a way very similar to SNF and are designed to have 380 values equal to zero for nodes that are not neighbors. Extensive experiments 381 on simulated and real datasets showed the competitive effectiveness and 382

efficiency of NEMO with respect to 9 state-of-the-art methods among 383 which one MKL-based method, a spectral clustering method, the classic 384 k-means clustering approach, and 6 clustering methods exploiting an 385 input data-fusion approach (Section Input data-fusion and output-fusion 386 methods). 387

Finally, a noteworthy PSN-fusion method applied for unsupervised 388 patient subtype identification in the TCGA dataset is Multi-view Spectral 389 Clustering Based on Multi-smooth Representation Fusion (MRF-MSC) 390

NCI-PID: National Cancer Institute-Pathway Interaction Database; RPPA: Reverse-Phase Protein Arrays; TCGA+cancer code: The Cancer Genome Atlas+ link to complete cancer

391 [95]. MRF-MSC starts by individually processing each data-block to

392 obtain a smoothed similarity matrix where strong/weak similarities 393 are strengthened/eliminated; this is obtained by solving a regularized

are strengthened/eliminated; this is obtained by solving a regularized optimization problem that computes the similarity matrix in a feature

space that minimizes the point-reconstruction error while strengthening

the point groupings. Next, a fused similarity matrix that minimizes the

weighted distance from all the smoothed source-similarity matrices is

obtained by integrating a self-weighting method [96] into the distance

³⁹⁹ minimization problem. Finally, the clusters in fused similarity networks

- 400 are strengthened by applying the constrained Laplacian rank method and
- 401 Spectral clustering is then applied to solve the clustering problem.

402 Input data-fusion and output-fusion methods

403 Opposite to PSN-fusion models, the input data-fusion and the output404 fusion techniques reviewed in this section integrate the information
405 available either in the multimodal input data (*input data-fusion* methods 406 Fig. 3) or in the output computed by a set of individual unimodal PSN407 analysis models (*output-fusion* methods - Fig. 4).

408 Input data-fusion methods are schematized in Fig. 3. These approaches are based on the assumption that the input samples originally lied in a latent 409 (eventually orthogonal) space from which the multiple source-views have 410 been generated by unknown projections. This results in data-blocks being 411 expressed into separate source-specific spaces that are characterized by: 1) 412 an individual source-specific structure generating an individual variability 413 within each data-block; 2) a joint sample-specific structure [18] resulting 414 in shared variance (collinearities) between data blocks. Therefore, input 415 data-fusion methods estimate the embedding that back-projects the input 416 data-blocks into a shared latent space minimizing redundancy between 417 the data-blocks while maximizing the individual data-block variability. In 418 other words, all the methods find the joint components (Fig. 3) allowing 419 to capture the greatest amount of shared variance; most of the methods 420 also define ways to identify the individual components capturing the source-421

422 specific variability (Fig. 3).

Depending on the technique used to project the data into the 423 shared latent space, we can distinguish input data-fusion methods into 424 PCA-based techniques (Table 6) or Matrix Factorization (MF) or Blind 425 Source Separation based methods (Table 7). One advantage of solving 426 the information-fusion in a pre-processing phase, i.e. preceding the 427 construction of an integrated PSN, is that a standard unimodal PSN-428 analysis model can be subsequently applied (Fig. 3-B) to deal with clustering 429 or supervised classification problems (Fig. 3-C). In particular, the input 430 data-fusion methods make the choice of the similarity measure to be used 431 for PSN construction particularly easy, since they compute normalized, a-432 dimensional, integrated point representations, whose pairwise similarities 433 could be handled by classic measures such as the cosine similarity or 434 the inverted euclidean distance. Moreover, a side-effect of the estimated 435 embedding is that the estimated component loadings or factors may 436 be analyzed for uncovering hidden relationships between variables (data 437 analysis task in Tables 6 and 7 and in Fig. 3). 438

The strategy applied by *output-fusion* methods is sketched in Fig. 4 and their experimental design is summarized in Table 8. They apply individual PSN pipelines on each data source to obtain individual clustering or supervised prediction results (Fig. 4-A and 4-B). All the obtained results are then fused by aggregation strategies that, acting as judges, compute a final decision by considering all the individual decisions taken by each unimodal pipeline.

Input data-fusion via PCA-based and CCA-based methods

In the bioinformatics field, Consensus PCA (CPCA [99]), hierarchical PCA447(HPCA [106]), and Multiple Factor Analysis (MFA [107]), are some of448the most used PCA-based integrative methods. They achieved interesting449results on multimodal datasets including different types of patient data,450from omics [18] to images [108, 109, 110].451

Their effectiveness is due to their ability to project the data-blocks 452 into a lower dimensional space spanned by not-correlated axis (principal 453 components) maximizing the within-block variances and between-block 454 covariances [111, 112]. By stretching the data along those axis, they induce 455 a natural separability that improves the performance of the downstream 456 algorithms, which are mostly devoted to data-exploration and unsupervised 457 clustering, though some exceptions using supervised clustering exist [20] 458 (Table 6). 459

The difference between the three approaches relies on the way the latent 460 space is found. Indeed, while CPCA solves an optimization problem by 461 an iterative algorithm in the set of nonlinear iterative partial least squares 462 methods (NIPALS [113]). HPCA [106] and MFA [107] consecutively apply 463 PCA on respectively: a) each block separately to derive lower-dimensional 464 "stretched" block representations maximizing the within-block variance: 465 b) the concatenation of the obtained block representations to derive a 466 stretched latent space maximizing the between-block covariance. 467

A notable generalization of PCA for multimodal data is *IIVE* (Joint and 468 Individual Variation Explained, [18]), which explicitly models each data-469 block X_i as the sum of a matrix representing the joint structure associated 470 with X_i and shared with other sources, and a matrix representing the 471 source-specific structure characterizing X_i , and residual noise. Given 472 this formulation, authors apply an iterative estimation procedure that 473 minimizes the reconstruction error, while constraining the axis of the joint 474 and individual structures to be orthogonal (that is, the joint and individual 475 structures must be uncorrelated). In practice the estimation iterates over 476 the following two steps: 1) having removed the individual structure, 477 apply a sparse Singular Value Decomposition (SVD) to estimate a lower-478 dimensional joint structure; 2) having removed the joint structure, apply a 479 sparse SVD to find a lower-dimensional individual structure. Interestingly, 480 JIVE also provides a permutation test to select the optimal ranks for the 481 estimated structures. When experimented on multi-omics data from the 482 glioblastoma multiforme (TCGA-GBM) dataset [18], JIVE showed its ability 483 to effectively uncover the individual and joint data structures, thus leading 484 to a better interpretation of interactions among data types and improving 485 unsupervised classification results. Since the computational complexity of 486 JIVE hampers its applicability, it has been recently reformulated (Angle 487 Based JIVE - aJIVE [97]) by using a hierarchical strategy similar to 488 HPCA, which also produces more intuitive interpretations of the obtained 489 decomposition, especially in the presence of strong collinearities. The 490 effectiveness of aJIVE is witnessed by the promising results obtained when 491 applied to an extract of the TCGA breast cancer dataset from [101] for 492 the (supervised) task of tumor subtype prediction [114]. In particular the 493 estimated joint components and the first five individual components for 494 each data block are used to compose the integrated sample views to train 495 Random Forest classifiers [115]. 496

Opposite to PCA-based integrative models, Canonical Correlation 497 Analysis-based (CCA-based) integrative models, e.g. Regularized Generalized 498 CCA (RGCCA) [104, 105] and its sparse counterpart Sparse Generalized 499 CCA (SGCCA) [19, 105], find the latent space maximizing the correlation 500 within and between the different data-blocks. They are generally used 501 for exploratory variable analysis since they try to bring all the data 502 blocks to a unique distribution, therefore uncovering hidden relationships 503 between different sources. However, DIABLO [20] has shown that SGCCA 504 is also effective in the context of supervised clustering for patients' 505

446

Name	Matched Samples	Dataset	Sample Cardinality	Data type	Integration approach	Task	Code and Language
SNF [3]	×	TCGA GBM	215	mRNA miRNA methy	SNF	Unsupervised Clustering (Patient subtype identification)	MATLAB code, , R code
ANF [87]	×	TCGA LUSC TCGA Adrenal TCGA Gland TCGA KIRC TCGA Uterus	2193	mRNA miRNA methy	SNF	Unsupervised Clustering (Patient subtype identification)	R code
HSNF [89]	×	TCGA BIC TCGA BIC TCGA GBM TCGA KIRC TCGA LUSC TCGA COAD	105 215 122 106 92	mRNA miRNA methy	SNF	Unsupervised Clustering (Patient subtype identification)	
SKF [90]	×	TCGA BIC TCGA COAD TCGA COAD TCGA KIRC TCGA LUSC TCGA Stomach	1071 426 868 981 377	mRNA miRNA isoform level	SNF	Unsupervised Clustering (Patient subtype identification)	MATLAB code
ab-SNF [54]	×	TCGA LIHC TCGA KIRP TCGA BIC	Not provided	somatic mutation mRNA methy	SNF	Unsupervised Clustering (Patient subtype identification)	R code
NEMO [93, 94]		TCGA AML TCGA BIC TCGA BIC TCGA COAD TCGA COAD TCGA COAD TCGA COAD TCGA COAD TCGA LUSC TCGA SKCM TCGA SKC TCGA SARC	3168 across all datasets	mRNA miRNA methy	SNF	Unsupervised Clustering (Patient subtype identification)	R code
Abbreviations methy: DNA meth	ylation; miR	NA: micro RNA; mR	UNA: messenger Rl	NA; SNF: Similarity N	etwork Fusion; TCGA+ca	<i>ncer code</i> : The Cancer Genome Atlas	s+ link to complete cancer codes.

Table 4. (I) SNF-based PSN-fusion methods. For each method, the table reports: the name/acronym with the corresponding reference paper; whether it requires the same set of patients across all data modalities (i.e. "Matched Samples"); the dataset used to develop and evaluate the approach in the reference paper and the corresponding sample cardinality and data types composing the dataset; the exploited integration method; the application task and the code availability (with link to the repository and programming languages for which the code is available).

subtype prediction. In practice, given a multimodal dataset containing ${\cal N}$ 506 507 samples organized into \boldsymbol{Y} classes, DIABLO firstly creates an extra dummy 508 (supervising) data-block where each column is an indicator variable for 509 the point-class $(1\ldots Y).$ Next, it uses SGCCA to maximize the covariance between all the data-blocks, including the supervising data-block. Given 510 this representation, supervised clusters may be identified either 1) by 511 512 averaging the components across data-blocks, to obtain an integrated patient representation that is then used by any supervised clustering 513 algorithm (such as the Maximum Centroids algorithm [116]); 2) by applying the Maximum Centroids algorithm on each projected data-block to obtain individual clustering results, subsequently aggregated via a majority voting algorithm.

Though effective in several applications, all the aforementioned PCA-518 based methods suffer from two main limitations: sensitiveness to outliers 519

Name	Matched Samples	Dataset	Sample Cardinality	Data type	Integration approach	Task	Code and Language
netDx [6]	x	TCGA KIRC TCGA OV TCGA GBM TCGA LUSC	150 252 77	mRNA miRNA methy CNV RPPA	average score	Supervised Classification (Patient's Survival)	R code
RWRF, RWRNF [58]	×	TCGA ACC TCGA BLCA TCGA HNSC TCGA UVM TCGA PAAD TCGA THCA	76 396 469 80 492	mRNA mRNA miRNA methy	RWR	Unsupervised Clustering (Patient subtype identification)	R code
MRF-MSC [95]	×	TCGA COAD TCGA GBM TCGA BRCA TCGA KIRC TCGA LSCC	92 215 105 122 106	mRNA miRNA methy	maximization of alignment to all the unimodal PSNs	Unsupervised Clustering (Patient subtype identification)	
Abbreviations CNV: Copy Number Va Array, RWR: Random Wé	riation; met l alk Kernel wi	ay: DNA methylati (th Restart; TCGA -	ion; miRNA : n +cancer code: Tl	nicro RNA; mR J he Cancer Geno	NA: messenger RNA; PSN me Atlas+ link to complete o	: Patient Similarity Network; RPPA 1ncer codes.	r: Reverse Phase Protein

520 and inability of handling missing data. Generalized Integrative PCA (GIPCA) [100] has been recently proposed as an extension of Consensus 521 PCA for dealing with missingness of some values and of entire views. To 522

523 this aim, eigenvectors are used to explain the intra/inter-block variance by

524 neglecting those samples/views with missing values/views.

Input data-fusion via Matrix Factorization-based methods 525

Matrix Factorization (MF) methods [117] embed the points into a latent 526 space that minimizes the reconstruction error and whose components 527

(factors) are not constrained to be orthogonal (as in PCA) [118, 31, 119]. 528 The most effective and used MF method applied on unimodal data is 529 Non-negative MF (NMF, [120]); it constrains both the component and 530 loading matrices to be non-negative, which makes the approximation 531 purely additive. 532

Given its effectiveness, several works proposed methods where NMF 533 is extended to the integration of multimodal datasets (see Table 7). The 534 most relevant example is joint NMF (jNMF [121]) where multiple NMF 535 problems are solved subject to a shared factor matrix that contains the 536 basis vectors of the shared latent space. However, jNMF is sensitive to 537

Table 5. (II) Other PSN-fusion methods. For each method, the table reports: the name/acronym with the corresponding reference paper; whether it requires the same set of patients across all data modalities (i.e. "Matched Samples"); the dataset used to develop and evaluate the approach in the reference paper and the corresponding sample cardinality and data types composing the dataset; the exploited integration method; the application task and the code availability (with link to the repository and programming languages for which the code is available).



Fig. 3. Input data-fusion. (A) During the pre-processing phase the data are integrated by either a PCA-based integrative model or a MF-based model. They estimate a shared latent space where the integrated, normalized point representations express the joint structure underlying all the data blocks plus, eventually, the individual structures characterizing each data block (e.g. JIVE [18], aJIVE [97], iNMF [98]); (B) a PSN model is then constructed on the integrated profiles by using a classic similarity measure; (C) a clustering or supervised classification model is applied to the computed PSN.

random noise and confounding effects [98] that are specific to each source, 538 and that cannot be detected if a unique shared factor matrix is estimated. 539 This affects the accuracy of the common structure estimation computed 540 by jNMF [98]. Therefore, integrative Non-negative Matrix Factorization 541 542 (iNMF [98]) uses an approach similar to JIVE, where the factor matrices to be estimated are composed both by a shared and a source-specific structure. 543 Unsupervised clustering experiments on the TCGA dataset [98, 122] have 544 proven the superiority of iNMF with respect to jNMF [121], NMF [123], 545

⁵⁴⁶ and to integrative Bayesian methods [124, 125].

547 Integrative Graph Regularized Non-Negative Matrix Factorization
 548 for Network Analysis (iGMFNA [126, 127]) proposes improving the

minimization of the reconstruction error, typical of NMF, by exploiting549a graph-view on each data block. Thanks to such representation, the550designed iterative optimization minimizes the reconstruction error while551maintaining the topology of the graph-views. When compared to jNMF552and iNMF to prioritize genes associated with cancer in two TCGA datasets553by an unsupervised clustering approach, iGMFNA showed its superior554performance.555

The popular Penalized Non-negative Matrix Tri-Factorization (NMTF,556[128, 31]) starts from a relational matrix $R_{1,2}$ containing non-negative557elements that represent the strengths of the relationships between objects of558two different types, ϵ_1 and ϵ_2 , whose respective characteristics are defined559

subject to constraints θ_1 and θ_2 , such that: $R_{1,2} \approx G_1 S_{1,2} G_2^T$ so that 561 G_1 and G_2 are the low-dimensional representations of objects with types, 562

respectively, ϵ_1 and ϵ_2 , and $S_{1,2}$ is the backbone matrix linking the two 563 564 types.

565 NMTF is exploited by [31] in DFMF, where the reliability of the 566 integrated low dimensional estimates computed over a multimodal dataset is improved by considering all the relational matrices (and corresponding 567 568 constraints) linking the different sources between each other and with the patient data. Given all the relational matrices, $R_{i,j}$, and respective 569 constraints, each $R_{i,j}$ is decomposed so that each backbone matrix 570 represents the latent structure between two data types, the generic low-571 dimensional data representations of objects with a specific type, G_i , is 572 573 bound to be used in the reconstruction of every relational matrix involving that type. Thanks to the abundance of information, the proposed model can 574 also handle missing data and treat sparse relational matrices. Furthermore, 575 576 it does not make any assumption about the structural properties of relations, 577 which can also be asymmetric. DFMF can also be used in a semi-supervised setting. During training, the model parameters (i.e. the factorization ranks) 578 are learnt, and are then used in a matrix completion problem, where 579 unobserved entries in the target matrix $R_{i,j}$ are reconstructed for elements 580 that were not present in the training set. 581

DFMF has been successfully used in the Matrix trifactorization for 582 Discovery of Data similarity and Association (MaDDA) algorithm proposed 583 by [129] to construct PSNs for unsupervised clustering. In particular, 584 given n patients to be partitioned into k clusters, the low-rank matrix 585 $G \in \mathbb{R}^{n \times k}$ estimated through DFMF is viewed as a membership 586 587 matrix relating each patient to the k ranks/groups. After repeating the factorization multiple times with different initialization parameters, a final 588 consensus matrix is obtained by element-wise averaging all membership 589 matrices and then composing a PSN where the similarity between two 590 patients (weight of the edge connecting them) represents how many times 591 they ended up in the same group. 592

Multi-Omics Factor Analysis+ (MOFA+ [130]) is an integrative method 593 exploiting Bayesian group factor analysis [51] with regularization to 594 impose: (i) a view-wise and factor-wise sparsity, which shrinks to zero the 595 loading for the m-th modality and the k-th factor if the latest does not 596 explain any variability of the m-th view; (ii) a feature-wise sparsity, which 597 sets to zero loading on individual features from active factors so that only 598 a small number of features "actively" contribute to each factor. MOFA+ can 599 handle missing values as well as entirely missing views for some samples; 600 moreover, it can cope with heterogeneous data types, which is exactly what 601 is needed when dealing with multimodal datasets containing multi-omics, 602 clinical, and imaging data. 603

Given the successful results of MF-based integrative techniques, 604 some authors have included them as a pre-processing step in their 605 clustering/classification algorithms. As an example, iCluster+ [123, 131] 606 uses NMF to fuse the heterogeneous data-blocks and then clusters the 607 integrated views. It also exploits the obtained factor loadings to identify the 608 relevant features in the cluster generation. 609

Input data-fusion via Blind Source Separation 610

In their original formulation, Blind Source Separation (BSS) models were 611

defined as an extension of NMF techniques for "recovering unobservable 612 source signals s from measurements x (i.e., data), with no knowledge of the

613 parameters θ of the generative system $x = f(s, \theta)^{"}$ [132].

614

Given their documented ability [133, 134, 132] of uncovering hidden 615 structures underlying the observed unimodal signals, several BSS models 616 have been extended to handle multimodal datasets comprising also 617

multisets² (Table 7), by a further step that estimates the mixing matrix that 618 recombines all the estimated latent sources so as to compute an integrated, 619 more informative signal with no redundancies [135, 136, 132, 21]. 620

Given the lack of information about the mixing process and the 621 source signals, BSS models often differ for the constraints they impose to 622 counter the ill-conditioned problem and obtain essentially unique source 623 estimates [137, 132, 134]. As an example, the well known Independent 624 Component Analysis model (ICA [138]), and its extensions to multimodal 625 data (joint ICA - jICA [139, 140]), to multisets (Independent Vector Analysis 626 - IVA [141]), and to multidimensional sources (Independent Subspace 627 Analysis - ISA [142]), assume a linear (additive) mixture with mutually 628 independent sources and a non-Gaussian distribution of each independent 629 component in the latent space. 630

All the BSS models base their computations on the existence of 631 collinearities between the observed multimodal data components, so that 632 unreliable results may be obtained when this assumption is not satisfied. 633 Some authors [135] circumvent this problem by pre-processing the data 634 with CCA (or its multimodal extension), to obtain a projected data 635 representation along correlated components. 636

The most representative BSS-based multimodal data integration 637 technique is Multidataset Independent Subspace Analysis (MISA [132, 638 21]), which was recently proposed to generalize all the BSS models 639 to the fusion of any kind of multimodal-multisets. Motivated by the 640 definition of multiset, MISA is driven by statistical independence between 641 latent subspaces while assuming correspondence within the subspaces 642 underlying the input multisets. In practice, it firstly removes redundancies 643 by estimating non-orthogonal demixing matrices, projecting each multiset 644 into a respective (intermediate) lower-dimensional space spanned by 645 independent components. The sources from all the computed latent spaces 646 are then combined through another demixing matrix that brings all the 647 data-blocks into a unique shared latent space, resulting in an integrated 648 patient view. The de-mixing matrices are estimated by minimizing the 649 mutual information in the final space, while maximizing the mutual 650 information in the intermediate spaces, so as to capture as much correlation 651 as possible. When applied to the integration of the information extracted 652 from Functional Multi-Resonance Imaging (fMRI), Structural Multi-653 Resonance Imaging (sMRI), and Electroencephalogram (EEG) data, MISA 654 has proven its robustness with respect to high Signal to Noise Ratios 655 (SNR) as well as its ability to produce effective data fusion in different ICA 656 contexts. 657

Output-fusion methods

658 659

Following Fig. 4, in the context of multimodal PSN analysis the ouputfusion methods described in this section may be applied to combine 660 the (unsupervised clustering or supervised classification) results (Fig. 4-661 B) computed by individual PSN analyses applied on each data block (see 662 Fig. 4-A). In Fig. 4-C, the combination of the unimodal results is performed 663 either by some heuristics, or by majority voting, or by using a meta-664 model that learns from the predictions performed by each unimodal PSN 665 analysis. Output-fusion techniques have been proposed for clustering 666 samples (mainly from the TCGA datasets, see Table 8) to identify patients' 667 subtypes [22, 144, 23] and for patients' classification [145, 63] (see Table 8). 668

² Multisets are multimodal datasets containing multiple views acquired by the same source under different acquisition conditions (e.g. observation times, experiments, tasks, machines). They are therefore homogeneous [26] in semantic, type, and dimensionality. Multimodalmultisets are multimodal datasets acquired by different sources, among which sources producing multisets.



Output-fusion methods

Fig. 4. Output-fusion. (A) Unimodal PSNs are constructed for each data-type or data-source, and (B) each one is individually processed to identify clusters or to classify unknown samples; subsequently, (C) a simple aggregation technique or a meta-model is used to obtain the fused/consensus clustering/classification result.

As an example, in Cluster-of-Cluster-Assignments (COCA [22]), 669 authors combine the clustering results individually obtained by NMF [146] 670 on each of the six data types of the TCGA datasets. To this aim, the samples 671 are coded into vectors composed of indicator variables representing the 672 clusters they have been assigned in each modality, so that they can be 673 reclustered according to those vectors by Consensus Clustering Plus $\left[147\right]^4$ 674 PINSPlus [144, 23] similarly exploits Consensus Clustering [148] for 675 reaching the final partition. In practice, PINS (Perturbation clustering for data 676 INtegration and disease Subtyping) starts by applying any classic unsupervised 677 clustering algorithm (e.g., k-means) individually on each of the M-th 678 679 datasets. If n is the number of patients, for the m-th dataset ($m \in$ $1, \ldots, M)$ the clustering result is expressed by a square matrix $C_m~\in$ 680 $\mathbb{R}^{n \, \times \, n}$, such that $C_m(i,j) \; = \; 1$ if samples i and j fall in the same 681 cluster, and $C_m(i,j)\,=\,0$ otherwise. All the resulting matrices are then 682 averaged to obtain a consensus matrix $S = \frac{\sum_{m=1}^{M} C_m}{M}$. Even though 683 matrix S may highlight that some points do not reach a strong agreement, 684 authors consider that S itself may be used as a pairwise similarity matrix 685 686 (since S = 1 for points for which there is a strong agreement, viewed 687 as similarity, across all the dataset, and S = 0 otherwise) that is suitable

 $^4\,$ Given the number of clusters k, Consensus Clustering Plus works on a consensus matrix (CM_k) representing "the proportion of clustering runs in which two items are [grouped] together" [148]. Given CM_k an agglomerative hierarchical consensus clustering using distance of 1-consensus values is completed and pruned to k groups that are returned as consensus clusters.

for similarity/distance-based clustering algorithms such as any Hierarchical688Clustering algorithm [149], Partitioning Around Medoids [150], or dynamic689tree cut [151]. In their work, authors propose testing different clustering690algorithms and then choose the partition that agrees the most with the691partitioning of individual data types.692

Consensus clustering has also been successfully applied by the recently 693 published SUMO [62], an integrative clustering algorithm that starts 694 by computing several unimodal PSNs by using a scaled-normalized 695 Euclidean kernel similar to the one exploited by SNF [3]. SUMO 696 then formulates a constrained NMTF (see Section Input data-fusion 697 via Matrix Factorization-based methods) to find a sparse shared 698 representation of all the samples in the cluster subspace by accounting for 699 the adjacencies observed in all the data types. The NMTF optimization 700 problem is solved by an iterative procedure that is applied several times 701 on several sample subsets to ensure robustness with respect to the initial 702 conditions and to the input data; consensus clustering is then exploited 703 to pool together the clustering results. When compared to the most 704 promising integrative clustering methods (e.g., iCluster [152], MCCA [102], 705 NEMO [94], SNF [3], PINSPlus [23]) SUMO obtained impressive results. 706

The Fuzzy-Hierarchical CLUSTering - FH-Clust method [24] interestingly 707 proposes to use fuzzy logic for identifying patients' prognostic subgroups from multiomics data, resting on the fact that in nature there is often no clear cut between subtypes. Unimodal data are separately analyzed using a fuzzy-based hierarchical clustering approach exploiting a Lukasiewicz 711 valued fuzzy similarity and individual results are then fused through a consensus matrix. Extensive experiments on 10 cancer datasets from TCGA 713

Clust is competitive with state-of-the-art methods (i.e. k-means, Spectral 715 Clustering, LRACluster, PINS, SNF, MCCA). 716 Interesting output-fusion approaches aimed at patients' classification 717 718 are described in [145, 63]. In [145] the authors obtain effective cancergrade and patient-survival classifications for cancer patients represented 719 720 in the TCGA renal (TCGA KIRC) and TCGA ovarian (TCGA OV) datasets by using all the data types included in TCGA, including hematoxylin and 721 722 eosin (H&E) stained whole-slide images of tissue samples that are processed by digital image processing techniques to extract more that 400 features 723 per sample. In practice authors firstly individually process each data block 724 to apply an internal cross-validation approach to choose 1) the number 725 of informative features to be extracted by the minimum Redundancy 726 727 Maximum Relevance (mRMR) method [153] and 2) the best performing 5fold cross classifier among SVM, logistic regression, K-nearest neighbors, 728 and Linear Discriminant Analysis. To compose all the predictions from 729 730 the different modalities authors compare the stacked generalization model 731 [154], which essentially trains a logistic regression classifier on the obtained predictions, to the majority vote strategy. The best results are obtained by 732 the stacked prediction model, which leverages the results obtained by any 733 of the multimodal predictions, independent from the classifier that is used 734 for producing them. 735 In [63] authors simply use the average to integrate the different 736

(considering gene expression, miRNA, methylation data) show that FH-

714

prognostic classifications computed over multimodal profiles of suspected 737 Alzheimer Disease (AD) patients, with the aim of identifying patients 738 who are vulnerable to conversion from mild-cognitive impairment to AD. 739 In particular, the squared-exponential kernels are firstly used to build 740 741 unimodal PSNs, and, for each unimodal network, a Gaussian process is then 742 exploited to assign labels to unknown points based on the nearest known 743 points. Finally, the unknown patients' condition is computed as the average over all the unimodal predictions. 744

Discussion and Conclusion

745

In the context of Precision Medicine, PSNs are gaining momentum given 746 their ability to uncover and exploit relationships among patients when 747 applied to clustering and classification tasks [9]. According to the state-of-748 the-art surveys describing the application of PSNs for precision medicine 749 or health-data processing [9, 45, 157, 158], PSN-based models benefit from 750 several advantages; they are: (I) easy to understand, (II) interpretable by 751 design, (III) privacy preserving, (IV) competitive or even superior to state-752 of-the-art clustering/classification methods, (V) potentially able to integrate 753 different data views. In particular, the possibility of using PSN models 754 in a multimodal setting is especially relevant in light of the increasing 755 availability of digital technologies by means of which huge amount of 756 multimodal data can be collected that describe each patient/sample by 757 considering different biological/medical views. Moreover, in the past 758 few years the increasing availability of cloud technologies allowed us 759 to distribute data processing across multiple local servers belonging to, 760 e.g. different institutions. In this context, the development of promising 761 information integration models would allow the application of a Federated 762 763 Learning strategy [159], where a central server collects, further integrates, and eventually processes, the (already) integrated data, or the individual 764 765 PSNs, or the predictions individually computed by local servers located in 766 the institutions where the data belong. In this way, the initial processing of 767 the sensitive data would be demanded to the local institutions to protect patient privacy, and the central server would have access only to pre-768 processed information, thus hiding explicit sensitive data. 769 Though in the biomedical context several multimodal approaches 770

have already shown their ability to integrate multimodal data to improve 771

the results obtained from a single view (unimodal data) [114], and the 772 survey literature about data integration methods for multimodal data is 773 wide [16, 13, 15], in the field of PSN analysis only few methods have 774 already investigated the usage of multimodal data, by building integrated 775 PSNs that exploit both the joint and the individual information from all 776 the available sources. Moreover, no state-of-the-art survey has focused 777 on the role of PSN as a cornerstone for data fusion. In this survey, we 778 aim at filling this gap with the goal of providing interested readers with a 779 comprehensive collection of integrative methods that may be exploited to 780 build PSN approaches efficiently handling multimodal data. 781

Besides an extensive literature search, the integration approaches have 782 been organized into three broad classes on the basis of the type of data 783 that is fused: PSN-fusion, Input data-fusion and Output-fusion methods. 784 More precisely, PSN-fusion methods may be split into the three sub-classes 785 of MKL, SNF-based and other methods while Input data-fusion approaches 786 comprehend algorithms PCA-based, CCA-based and MF-based. 787

The survey has highlighted the promising results and advantages that characterize the methods belonging to the three classes of our proposed taxonomy.

Methods based on PSN-fusion techniques are particularly useful in 791 Network Medicine applications [160], that study human diseases through 792 "systemic" approaches in which diseases are interpreted as perturbations 793 in complex biomolecular networks. In this context, transductive strategies 794 working on individual PSN models [7] would benefit from the application 795 of PSN-fusion approaches, as shown by recent promising results [3, 6]. 796

Methods based on input-data fusion techniques rely on factor analysis 797 models for the removal of data collinearities and the simultaneous 798 enhancement of the individual structure characterizing each view. For this 799 reason, we believe such techniques are particularly useful when dealing with 800 multiview data involving follow-up examinations, where the multiple views 801 likely contain correlated information. 802

Output-fusion techniques should be used when the differences between 803 the multimodal views impose the usage of peculiar and specific unimodal 804 PSN models for obtaining individual inferences. This is the case, for 805 example when we need to combine data having substantially different 806 structures, ranging from vectorial to sequence and graph-structured data. 807

Though being effective, our thorough review also evidenced difficulties 808 and drawbacks that harbour from the data-fusion strategy. In particular, 809 PSN-fusion models require to build an individual PSNs on each data-type. 810 This raises the crucial, still open, and often overlooked problem of choosing 811 proper individual similarity measures for building each unimodal PSNs. 812 Indeed, only few methods [60, 161, 162] reported exhaustive comparative 813 evaluations among few distance metrics applied to genetic data. Bv 814 considering that several problems in Precision Medicine are characterized 815 by non-linearly separable omics data, and given the experimental results 816 we have collected during our literature search, we recommend computing 817 PSNs by exploiting a kernel function. In this context, though several 818 functions have been successfully proposed and used in literature, when 819 dealing with continuous data, we suggest using the scaled exponential 820 kernel of Euclidean distance [3, 62], due to its ability to adapt to different 821 neighborhood sizes. This allows dealing with datasets distributed on 822 complex manifolds where datapoints are not evenly distributed in space, as 823 it often happens in real-world problems. On the other hand, when dealing 824 with simpler data-types with lower dimensionalities and complexities 825 (e.g. clinical data), simpler normalized similarities may be sufficient to 826 appropriately capture the data structure. Clinical datasets usually contain 827 categorical variables, often mixed with numeric features. The former 828 situation can be appropriately addressed by averaging the normalized 829 similarities individually computed on each variable [6], while Chi-squared 830 distances are the most suitable for categorical data [3, 6]. Of note, the 831

788

789

790

subsequent application of a Random Walk kernel, as proposed by Gliozzoet al [7], is a promising step to refine the obtained PSN.

On the other side, input data-fusion techniques integrate the input data by projecting them into a shared space with lower dimensionality, thus making these approaches strongly dependent on the chosen final dimensionality *d*.

838 While classic approaches have been proposed to automatically set d [163, 164, 165, 100], this value is often user-defined after observation of 839 840 the scree plot. However, observing that the optimal latent vector space is the one that allows to capture the intrinsic data structure, we instead 841 suggest setting d to the intrinsic data dimensionality (id) [166], which is 842 the minimum number of parameters needed to represent the data without 843 844 information loss. 845 Finally, output-data fusion methods are often too generic or use very

simple output-aggregation strategies, e.g. average or majority voting, that may produce sub-optimal results.

Generally speaking, our survey evidenced some important open issues
 in the context of data integration methods for PSN that call for the future
 research directions summarized in the following subsection.

851 Future Research Directions

While conducting our survey we noted the need of investigating methods 852 for data pre-processing, with the aim of, e.g., detecting and eliminating 853 noise with heterogeneous characteristics, collinearities between different 854 views, and confoundings that could bias the final results (as per [27]). Indeed, 855 only few recently proposed preliminary attempts were able to explicitly 856 consider the presence of noise with heterogeneous characteristics [98, 122]. 857 Moreover, future research should be devoted to the investigation of 858 novel multimodal feature-selection algorithms. Indeed, the few methods 859 applying a feature selection step exploit either classic univariate statistics, 860 or algorithms, such as mRMR [153], that analyze group of features by 861 neglecting their multimodal characteristics. 862 On the other side, missing data imputation needs deeper investigation 863 to handle two types of biomedical data-missingness: 1) missingness of some 864 data values in some views; 2) missingness of entire views for some samples. 865 While missingness is becoming a common problem in different fields, 866 in the bio-medical field few approaches present thorough missing data 867 imputation studies [11]. Besides, among the approaches we have surveyed, 868 only GIPCA [100] specifically addressed both these types of missingness. 869 Finally, given the big-data produced by high-throughput technologies, 870 scalability is becoming an important and often overlooked issue, nowadays 871 hampering the applicability of several promising tools. 872 Though the aforementioned issues are still open, all the surveyed 873 strategies have reported promising results that might improve knowledge 874 in then field of Precision Medicine. Unfortunately, different similarity 875 metrics, experimental setups, and evaluation measures are used for model 876 assessment: this hampers an objective comparison between the different 877 integration techniques and data analysis models. Furthermore, we found 878 no evidence about data integration approaches that should be preferred 879 over the others. Instead, the type and semantic of the available data 880 type and the specific biomedical question to address should guide the 881 882 choice. An additional open problem regards the identification of the most appropriate similarity/distance measure for each biological data modality. 883 884 To the best of our knowledge, only few works tried to investigate this issue 885 by comparing different metrics for specific data views and most of them 886 are focused on gene expression data [60, 161, 162]. Comprehensive studies comparing the usage of different similarity measures in different contexts 887 (e.g. when applied to different biological data types and in supervised 888

and unsupervised prediction contexts) would provide fruitful insights to guide the scientific community towards effective PSN construction. We also remark that, though some algorithms are already available as open891source packages/repositories (mostly coded using R, Python and Matlab)892[16], many others are not, thus slowing down their diffusion and testing by
the community.893

Another interesting research line that should be given attention is represented by the development of Web applications extending, e.g., those presented in [167, 168], for the visual analysis of PSN models. Indeed, the graphical tools can enable the visual comparison of different PSN models realized according to any of the methods discussed in this survey. This in turn can improve the explainability of the computed results and would allow the user to choose the approach mostly suited to her/his needs. 901

902

Abbreviations

ab-SNF - Association-signal-annotation Boosted Similarity Network	903
Fusion	904
AD - Alzheimer's Disease	905
ADNI - Alzheimer's Disease Neuroimaging Initiative	906
ANF - Affinity Network Fusion	907
aJIVE - Angle based Joint and Individual Variation Explained	908
BSS - Blind Source Separation	909
CCA - Canonical Correlation Analysis	910
CGH - Comparative Genomic Hybridization	911
CNV - DNA Copy Number Variation	912
COCA - Cluster of Cluster Assignment	913
CPCA - Consensus Principal Component Analysis	914
CSF - CerebroSpinal Fluid	915
DFMF - Data Fusion by Matrix Factorization	916
DILI - Drug-Induced Liver Injury	917
DLBCL - Diffuse Large B-Cell Lymphoma	918
EEG - Electroencephalography	919
fMRI - functional Magnetic Resonance Imaging	920
FH-Clust - Fuzzy-Hierarchical CLUSTering	921
GIPCA - Generalized Integrative PCA	922
GNMF - Graph Regularized Non-negative Matrix Factorization	923
GO - Gene Ontology	924
HPCA - Hierarchical Principal Component Analysis	925
HC - Healthy Control	926
HSNF - Hierarchical Similarity Network Fusion	927
ICA (and DS-ICA) - Independent Component Analysis (for Disjoint	928
Subspaces)	929
<i>id</i> - intrinsic data dimensionality	930
<i>iGMFNA</i> - integrated Graph Regularized Non-negative Matrix Factorization	931
for Network Analysis	932
iNMF - integrative Non-negative Matrix Factorization	933
ISA - Independent Subspace Analysis	934
IVA - Independent Vector Analysis	935
jICA - joint Independent Component Analysis	936
JIVE - Joint and Individual Variation Explained	937
jNMF - joint Non-negative Matrix Factorization	938
KEGG - Kyoto Encyclopedia of Genes and Genomes	939
LPP - Locality Preserving Projections	940
LRAcluster - Low-Rank Approximation based multi-omics data clustering	941
MaDDA - Matrix trifactorization for Discovery of Data similarity and	942
Association	943
MCCA - Multiple Canonical Correlation Analysis	944
MCI - Mild Cognitive Impairment	945
methy - DNA methylation	946
MFA - Multiple Factor Analysis	947
MF - Matrix Factorization	948

- 950 MISA Multidataset Independent Subspace Analysis
- 951 MK-FDA Multiple Kernel Fisher Discriminant Analysis
- 952 MKKM Multiple kernel k-means clustering
- 953 MKL Multiple Kernel Learning
- 954 MOFA Multi-Omics Factor Analysis
- 955 MRF-MSC Multi-view Spectral Clustering Based on Multi-smooth
- 956 Representation Fusion
- 957 MRI Magnetic Resonance Imaging
- 958 mRMR maximum Relevance Minimum Redundancy
- 959 *mRNA* messenger RNA
- 960 MTF Matrix Tri-Factorization
- 961 MultiNMF Multi-View Non-negative Matrix Factorization
- 962 NCI-PID National Cancer Institute—Pathway Interaction Database
- 963 NEMO NEighborhood based Multi-Omics clustering
- 964 NIPALS Nonlinear Iterative Partial Least Squares
- 965 NMF Non-negative Matrix Factorization
- 966 NMTF Penalized Non-negative Matrix Tri-Factorization
- 967 PAMOGK Pathway-based MultiOmic Graph Kernel clustering
- 968 PCA Principal Component Analysis
- 969 PINS Perturbation clustering for data INtegration and disease Subtyping
- 970 PPI Protein-Protein Interaction
- 971 *PSN* Patient Similarity Network
- 972 RDA Regularized Discriminant Analysis
- 973 RGCCA Regularized Generalized Canonical Correlation Analysis
- 974 rMKL Regularized Multiple Kernel Learning
- 975 *rMKL-LPP* Regularized Multiple Kernel Learning Locality Preserving
 976 Projections
- 977 RPPA Reverse-Phase Protein Arrays
- 978 RW Random Walk
- 979 RWR Random Walk with Restart
- 980 RWRF Random Walk with Restart for multi-dimensional data Fusion
- 981 RWRNF Random Walk with Restart and Neighbor information-based
- 982 multi-dimensional data Fusion
- 983 SCA SpinoCerebellar Ataxia
- 984 SGCCA Sparse Generalized Canonical Correlation Analysis
- 985 SKF Similarity Kernel Fusion
- 986 sMRI structural Magnetic Resonance Imaging
- 987 SNF Similarity Network Fusion
- 988 SNR Signal to Noise Ratio
- 989 SVD Singular Value Decomposition
- 990 SVM Support vector Machine
- 991 TCGA study name The Cancer Genome Atlas link to official study
- 992 abbreviations

993

Appendix A - Data integration in Medicine: previous surveys and taxonomies

The abundance of multimodal data integration approaches developed in the past decade in the biomedical context has motivated many relevant surveys [29, 36, 37, 34, 35, 16], which proposed different definitions and taxonomies, schematized in Fig. 1.

In the context of Precision Medicine, multimodal sets are composed of multiple views (or data-blocks) for the same set of patients. They are either *multisets*, or *multimodal datasets* [142]. *Multisets* (top of Fig. 1yellow box) contain multiple views acquired by the same source under different acquisition conditions (e.g. observation times, experiments, tasks, machines), and are therefore *homogeneous* [26] in semantic, type, and dimensionality. Conversely, *multimodal datasets* (alias *heterogeneous* sets [26],

 Fig. 1-light blue box) contain data-blocks acquired by different sources,
 1007

 characterized by different semantics, type and dimensionalities. Among the
 1008

 latter, multimodal-multiset are datasets acquired by different sources, some
 1009

 of which are used to produce multisets.
 1010

Given their different characteristics, multisets and multimodal datasets 1011 are generally fused by following different integration flows. Horizontal 1012 integration methods [169, 25] are usually used for multisets because they 1013 equally process all the data-blocks and then pool the obtained results by 1014 e.g. summary statistics [170]. By contrast, vertical integration methods are 1015 used for processing multimodal datasets, which are more articulated and 1016 are usually grouped in the hierarchical-vertical class and the parallel-vertical 1017 class. 1018

Hierarchical-vertical [171, 172, 173] or multi-staged analysis methods [27] 1019 consider omics data being interrelated by regulatory mechanisms and 1020 exploit such prior knowledge during the integration procedure. Since 1021 these methods are tailored for the treatment of specific data types and 1022 applications and cannot be generalized to different research contexts, 1023 they will not further considered in this survey. Parallel-vertical integration 1024 techniques, alias meta-dimensional analysis methods [27], are the most 1025 diffused and generalizable ones because dependencies between data-blocks 1026 injected by prior information are not considered. To categorize parallel-1027 vertical approaches several interrelated taxonomies have been defined. The 1028 categorization reported in the red-dashed box in Fig. 1 is the one adopted 1029 by several authors [29, 30, 31, 27, 32, 33] that relies on the processing stage 1030 (early, intermediate, late) in which the data fusion happens, which also 1031 influences the kind of information that is fused. 1032

Early integration techniques (also called concatenation-based models [27, 1033 33]) are applied on the input data-blocks in an early stage to compose 1034 the integrated input vectors subsequently used in the analysis by either 1035 using a simple data-concatenation, or by exploiting joint latent space 1036 estimation models. The evident advantage of early methods relies on their 1037 ability to uncover the individual information characterizing each of the 1038 different sources as well as the hidden relationships between them. Another 1039 advantage is brought by the fact that early methods solve the integration 1040 problem in the first stage, so that any unimodal analysis process may be 1041 subsequently applied. 1042

Intermediate integration approaches (also named transformation-based 1043 models) [27, 33] individually transform the data-blocks into intermediate 1044 (unimodal) models that are subsequently integrated to produce a unique 1045 fused model to be analyzed. In the taxonomy proposed by [28] (blue-1046 dashed box in Fig. 1), these methods have been classified as model-dependent 1047 approaches for highlighting their dependency from the data analysis model, 1048 which guarantees the ability to retain the original data structure by 1049 explicitly addressing the fusion task in the construction of the predictive 1050 model itself. 1051

Late integration approaches (also named model-based approaches in [27, 1052 33]) separately analyze each of the incoming data-blocks to produce 1053 individual results, subsequently integrated in a late phase by some meta-1054 learner acting as the final judge or by simple techniques such as majority 1055 voting. These approaches along with the early integration ones are classified 1056 as model-agnostic in the taxonomy proposed by [28] (blue-dashed box in Fig. 1057 1) and are contrasted with the model-independent approaches previously 1058 discussed. They are named "agnostic" because they are independent from 1059 the specific algorithm applied in the preceding unimodal analysis, which 1060 can be therefore tailored to the processed type. 1061

Even if the aforementioned early/intermediate/late taxonomy is the most diffused in literature, other taxonomies have been defined in the context of integrative (multi-omics) methods for Precision Medicine. As an example, [174, 25, 175] consider three classes: 1) *statistical-based methods*, most of which can be considered instances of the class of early integrative methods; 2) *unsupervised methods* neglecting the outcome variable during the integration phase, that may be applied in any (early, intermediate,

1069 or late) phase and are mainly devoted to unsupervised clustering; 3) 1070 supervised integration methods fusing the available information to maximize

the outcome prediction performance by mainly using an intermediate

1072 Multiple Kernel Learning (MKL) integration approach or a late fusion

1073 approach.

1074 Other taxonomies [14, 32, 8, 16] consider the specific algorithm used for the integration; they recognize network-based approaches (among 1075 1076 which deep-network based approaches, not treated in this survey), feature transformation models mainly applying an early integration approach (e.g., 1077 Principal Component Analysis - PCA, Canonical Correlation Analysis -1078 CCA), integrative models exploiting Matrix Factorization (MF) techniques 1079 in an early integrative fashion, MKL models belonging to the class of 1080 1081 intermediate methods, and Bayesian techniques applied in an early phase. Note that Bayesian models are not considered in this work since they have 1082 been exhaustively described in a dedicated survey [13]. 1083

Finally, the relevant survey by [16] is focused on the description of publicly available multimodal datasets in the context of multi-omics and in the critical analysis of open source integrative models. After a thorough study, the authors conclude that an objective comparison between different models is difficult, and highlight the lack of an easy-to-use multiomics data fusion model providing a "biologist-friendly" visualization and interpretation.

1091 Competing interests

1092 There is NO Competing Interest.

1093 Author contributions statement

J.G., G.V. and E.C. conceived the work, J.G. collected the literature papers,
 J.G. and E.C. studied the literature, selected the most relevant works and

1096 drafted them; J.G, M.M., M.N., A.Pac., G.V., E.C. wrote the paper; all the

1097 authors validated the work.

1098 Acknowledgments

1099 Funding

This work was supported by "Piano di sostegno alla ricerca" PSR2015-17
 funded by the University of Milan and by the Transition Grant "UNIMI
 partenariati H2020" PSR2015-1720GVALE_01.

1103 Key Points

- Patients similarity networks (PSN) are explainable and privacy preserving representations of patients that leverage the similarity of their clinical/biomolecular profiles to construct graphs of patients.
- Network Medicine algorithms on PSNs for patient stratification, phenotype and outcome prediction and disease risk assessment represent novel tools for Genomic and Precision Medicine
- The combination of clinical, omics and imaging bio-medical data can lead to novel PSNs able to leverage the synergy of multiple views of the patients.
- Several reviews about data integration methods in Bioinformatics and bio-medical applications have been proposed but no specific reviews about the emerging field of heterogeneous data integration methods for patient similarity networks are actually available.

- We provide a thorough review and propose a taxonomy of heterogeneous
 1117
 data integration methods for PSNs, together with the different patient
 similarity measures proposed in literature.
- We also review methods that have appeared in the machine learning literature but have not yet been applied to PSNs, thus providing a resource to navigate the vast machine learning literature existing on this topic.
- Strengths and limitations of the proposed methods are discussed to 1124 both assist researchers in the design and development of novel methods 1125 and to guide the selection of PSN integration methods for specific 1126 applications, focusing on methods that could be used to integrate very diverse datasets, including multi-omics data as well as data derived 1128 from clinical information and medical imaging. 1129

References

1.

- Inke R Koenig, Oliver Fuchs, Gesine Hansen, Erika von Mutius, and
 1131

 Matthias V Kopp. What is precision medicine? European respiratory
 1132

 journal, 50(4), 2017.
 1133
- Samuel J Aronson and Heidi L Rehm. Building the foundation for genomics in precision medicine. *Nature*, 526(7573):336–342, 2015. 1135
- Bo Wang, Aziz M Mezlini, Feyyaz Demir, Marc Fiume, Zhuowen 1136 Tu, Michael Brudno, Benjamin Haibe-Kains, and Anna Goldenberg. 1137 Similarity network fusion for aggregating data types on a genomic 1138 scale. Nature methods, 11(3):333, 2014. 1139
- Dokyoon Kim, Je-Gun Joung, Kyung-Ah Sohn, Hyunjung Shin, 1140 Yu Rang Park, Marylyn D Ritchie, and Ju Han Kim. Knowledge 1141 boosting: a graph-based integration approach with multi-omics data 1142 and genomic knowledge for cancer clinical outcome prediction. 1143 *Journal of the American Medical Informatics Association*, 22(1):109–120, 1144 2015.
- Li Li, Wei-Yi Cheng, Benjamin S Glicksberg, Omri Gottesman, 1146 Ronald Tamler, Rong Chen, Erwin P Bottinger, and Joel T 1147 Dudley. Identification of type 2 diabetes subgroups through 1148 topological analysis of patient similarity. *Science translational medicine*, 1149 7(311):311ra174–311ra174, 2015. 1150
- Shraddha Pai, Shirley Hui, Ruth Isserlin, Muhammad A Shah, Hussam
 Kaka, and Gary D Bader. netdx: interpretable patient classification
 using integrated patient similarity networks. *Molecular Systems Biology*, 15(3):e8497, 2019.
- Jessica Gliozzo, Paolo Perlasca, Marco Mesiti, Elena Casiraghi, Viviana Vallacchi, Elisabetta Vergani, Marco Frasca, Giuliano Grossi, Alessandro Petrini, Matteo Re, et al. Network modeling of patients' biomolecular profiles for clinical phenotype/outcome prediction.
 Scientific reports, 10(1):1–15, 2020.
- Giovanna Nicora, Francesca Vitali, Arianna Dagliati, Nophar 1160 Geifman, and Riccardo Bellazzi. Integrated multi-omics analyses in 1161 oncology: A review of machine learning methods and tools. *Frontiers* 1162 *in Oncology*, 10:1030, 2020. 1163
- Shraddha Pai and Gary D Bader. Patient similarity networks for precision medicine. *Journal of molecular biology*, 430(18):2924–2938, 1165 2018. 1166
- Noël Malod-Dognin, Julia Petschnigg, and Nataša Pržulj. Precision 1167 medicine—a promising, yet challenging road lies ahead. Current 1168 Opinion in Systems Biology, 7:1–7, 2018. 1169
- Elena Casiraghi, Dario Malchiodi, Gabriella Trucco, Marco Frasca, Luca Cappelletti, Tommaso Fontana, Alessandro Andrea Esposito, Emanuele Avola, Alessandro Jachetti, Justin Reese, et al. Explainable machine learning for early assessment of covid-19 risk prediction in emergency departments. *IEEE Access*, 8:196299–196325, 2020.

- Thirunavukarasu Ramkumar, Shanmugasundaram Hariharan, and
 Shanmugam Selvamuthukumaran. A survey on mining multiple data
 sources. Wiley Interdisciplinary Reviews: Data Mining and Knowledge
 Discovery, 3(1):1–11, 2013.
- Matteo Bersanelli, Ettore Mosca, Daniel Remondini, Enrico Giampieri, Claudia Sala, Gastone Castellani, and Luciano Milanesi.
 Methods for the integration of multi-omics data: mathematical aspects. *BMC bioinformatics*, 17(2):167–177, 2016.
- 1183 14. Vladimir Gligorijević, Noël Malod-Dognin, and Nataša Pržulj.
 1184 Integrative methods for analyzing big data in precision medicine.
 1185 Proteomics, 16(5):741–758, 2016.
- Chen Meng, Oana A Zeleznik, Gerhard G Thallinger, Bernhard
 Kuster, Amin M Gholami, and Aedín C Culhane. Dimension
 reduction techniques for the integrative analysis of multi-omics data.
 Briefings in bioinformatics, 17(4):628–641, 2016.
- Indhupriya Subramanian, Srikant Verma, Shiva Kumar, Abhay
 Jere, and Krishanpal Anamika. Multi-omics data integration,
 interpretation, and its application. *Bioinformatics and biology insights*,
 14:1177932219899051, 2020.
- 1194 17. Alain Rakotomamonjy, Francis R Bach, Stéphane Canu, and Yves
 1195 Grandvalet. Simplemkl. *Journal of Machine Learning Research*,
 1196 9(Nov):2491–2521, 2008.
- 1197
 18. Eric F Lock, Katherine A Hoadley, James Stephen Marron, and
 1198
 Andrew B Nobel. Joint and individual variation explained (jive)
 1199
 for integrated analysis of multiple data types. *The annals of applied* 1200
 statistics, 7(1):523, 2013.
- Arthur Tenenhaus, Cathy Philippe, Vincent Guillemot, Kim-Anh
 Le Cao, Jacques Grill, and Vincent Frouin. Variable selection for
 generalized canonical correlation analysis. *Biostatistics*, 15(3):569–
 583, 2014.
- Amrit Singh, Casey P Shannon, Benoît Gautier, Florian Rohart, Michaël Vacher, Scott J Tebbutt, and Kim-Anh Lê Cao. Diablo:
 an integrative approach for identifying key molecular drivers from multi-omics assays. *Bioinformatics*, 35(17):3055–3062, 2019.
- Rogers F Silva, Sergey M Plis, Tülay Adalı, Marios S Pattichis, and Vince D Calhoun. Multidataset independent subspace analysis
 with application to multimodal fusion. *IEEE Transactions on Image Processing*, 30:588–602, 2020.
- 22. Katherine A Hoadley, Christina Yau, Denise M Wolf, Andrew D
 Cherniack, David Tamborero, Sam Ng, Max DM Leiserson, Beifang
 Niu, Michael D McLellan, Vladislav Uzunangelov, et al. Multiplatform
 analysis of 12 cancer types reveals molecular classification within and
 across tissues of origin. *Cell*, 158(4):929–944, 2014.
- Hung Nguyen, Sangam Shrestha, Sorin Draghici, and Tin Nguyen.
 Pinsplus: a tool for tumor subtype discovery in integrated genomic
 data. *Bioinformatics*, 35(16):2843–2846, 2019.
- 1221 24. Angelo Ciaramella, Davide Nardone, and Antonino Staiano. Data
 1222 integration by fuzzy similarity-based hierarchical clustering. *BMC* 1223 *bioinformatics*, 21(10):1–15, 2020.
- Cen Wu, Fei Zhou, Jie Ren, Xiaoxi Li, Yu Jiang, and Shuangge Ma. A
 selective review of multi-level omics data integration using variable
 selection. *High-throughput*, 8(1):4, 2019.
- 1227 26. Vladimir Gligorijević and Nataša Pržulj. Methods for biological data
 1228 integration: perspectives and challenges. *Journal of the Royal Society* 1229 *Interface*, 12(112):20150571, 2015.
- 27. Marylyn D Ritchie, Emily R Holzinger, Ruowang Li, Sarah A
 Pendergrass, and Dokyoon Kim. Methods of integrating data to
 uncover genotype-phenotype interactions. *Nature Reviews Genetics*,
 16(2):85–97, 2015.
- 1234 28. Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency.
 1235 Multimodal machine learning: A survey and taxonomy. *IEEE*

transactions on pattern analysis and machine intelligence, 41(2):423–443, 1236 2018. 1237

- Paul Pavlidis, Jason Weston, Jinsong Cai, and William Stafford Noble. 1238
 Learning gene functional classifications from multiple data types. 1239
 Journal of computational biology, 9(2):401–411, 2002. 1240
- Anneleen Daemen, Olivier Gevaert, and Bart De Moor. Integration of clinical and microarray data with kernel methods. In 2007 29th Annual International Conference of the IEEE Engineering in Medicine and Biology Society, pages 5411–5415. IEEE, 2007. 1244
- Marinka Žitnik and Blaž Zupan. Data fusion by matrix factorization. 1245 IEEE transactions on pattern analysis and machine intelligence, 37(1):41– 1246 53, 2014. 1247
- Yifeng Li, Fang-Xiang Wu, and Alioune Ngom. A review 1248 on machine learning principles for multi-view biological data 1249 integration. *Briefings in bioinformatics*, 19(2):325–340, 2018. 1250
- Zahra Momeni, Esmail Hassanzadeh, Mohammad Saniee Abadeh, and Riccardo Bellazzi. A survey on single and multi omics data mining methods in cancer data classification. *Journal of Biomedical Informatics*, page 103466, 2020.
- Wei Tang, Zhengdong Lu, and Inderjit S Dhillon. Clustering with 1255 multiple graphs. In 2009 Ninth IEEE International Conference on Data 1256 Mining, pages 1016–1021. IEEE, 2009. 1257
- Martin H Van Vliet, Hugo M Horlings, Marc J Van De Vijver, Marcel JT 1258
 Reinders, and Lodewyk FA Wessels. Integration of clinical and gene expression data has a synergetic effect on predicting breast cancer outcome. *PloS one*, 7(7):e40358, 2012. 1261
- Gert RG Lanckriet, Tijl De Bie, Nello Cristianini, Michael I Jordan, and William Stafford Noble. A statistical framework for genomic data fusion. *Bioinformatics*, 20(16):2626–2635, 2004.
- Olivier Gevaert, Frank De Smet, Dirk Timmerman, Yves Moreau, and Bart De Moor. Predicting the prognosis of breast cancer by integrating clinical and microarray data with bayesian networks. *Bioinformatics*, 22(14):e184–e190, 2006.
- Hossein Sharifi-Noghabi, Olga Zolotareva, Colin C Collins, and Martin Ester. Moli: multi-omics late integration with deep neural networks for drug response prediction. *Bioinformatics*, 35(14):i501– i509, 2019.
- Shihyen Chen, Bin Ma, and Kaizhong Zhang. On the similarity 1273 metric and the distance metric. *Theoretical Computer Science*, 410(24-25):2365–2376, 2009. 1275
- Lluís Belanche and Jorge Orozco. Things to know about a (dis) similarity measure. In International Conference on Knowledge-Based and Intelligent Information and Engineering Systems, pages 100–109.
 Springer, 2011.
- 41. Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller.
 1280

 Kernel Principal Component Analysis. In International Conference on Artificial Neural Networks, pages 583–588. Springer, 1997.
 1281
- 42. Nils M Kriege, Fredrik D Johansson, and Christopher Morris. A 1283 survey on graph kernels. *Applied Network Science*, 5(1):1–42, 2020. 1284
- François Fouss, Kevin Francoisse, Luh Yen, Alain Pirotte, and Marco Saerens. An experimental investigation of kernels on graphs for collaborative recommendation and semisupervised classification. *Neural networks*, 31:53–72, 2012.
- Joon Lee, David M Maslove, and Joel A Dubin. Personalized mortality prediction driven by electronic medical data and a patient similarity metric. *PloS one*, 10(5):e0127428, 2015.
- Anis Sharafoddini, Joel A Dubin, and Joon Lee. Patient similarity 1292 in prediction models based on health data: a scoping review. JMIR 1293 medical informatics, 5(1):e7, 2017. 1294
- Ping Zhang, Fei Wang, Jianying Hu, and Robert Sorrentino. 1295 Towards personalized medicine: leveraging patient similarity and 1296

- 1297drug similarity analytics.AMIA Summits on Translational Science1298Proceedings, 2014:132, 2014.129947.Seung-Seok Choi, Sung-Hyuk Cha, and Charles C Tappert. A survey
- of binary similarity and distance measures. Journal of Systemics,
 Cybernetics and Informatics, 8(1):43–48, 2010.

Sebastian Klenk, Jürgen Dippon, Peter Fritz, and Gunther Heidemann.
 Determining patient similarity in medical social networks. In
 Proceedings of the First International Workshop on Web Science and
 Information Exchange in the Medical Web, pages 6–14, 2010.

- Bernhard Schölkopf. The kernel trick for distances. In Advances in neural information processing systems, pages 301–307, 2001.
- 50. Bin Zhu, Nan Song, Ronglai Shen, Arshi Arora, Mitchell J Machiela,
 Lei Song, Maria Teresa Landi, Debashis Ghosh, Nilanjan Chatterjee,
 Veera Baladandayuthapani, et al. Integrating clinical and multiple
 omics data for prognostic assessment across human cancers. *Scientific reports*, 7(1):1–13, 2017.
- 1313 51. Ya Zhang, Ao Li, Chen Peng, and Minghui Wang. Improve
 1314 glioblastoma multiforme prognosis prediction by using feature
 1315 selection and multiple kernel learning. *IEEE/ACM transactions on* 1316 computational biology and bioinformatics, 13(5):825–835, 2016.
- 1317 52. Jérôme Mariette and Nathalie Villa-Vialaneix. Unsupervised multiple
 1318 kernel learning for heterogeneous data integration. *Bioinformatics*,
 1319 34(6):1009–1015, 2018.
- Anneleen Daemen, Dirk Timmerman, Thierry Van den Bosch, Cecilia
 Bottomley, Emma Kirk, Caroline Van Holsbeke, Lil Valentin, Tom
 Bourne, and Bart De Moor. Improved modeling of clinical data with
 kernel methods. Artificial intelligence in medicine, 54(2):103–114, 2012.
- 54. Peifeng Ruan, Ya Wang, Ronglai Shen, and Shuang Wang. Using
 association signal annotations to boost similarity network fusion.
 Bioinformatics, 35(19):3718-3726, 2019.
- 1327 55. Shuhao Li, Limin Jiang, Jijun Tang, Nan Gao, and Fei Guo. Kernel
 1328 fusion method for detecting cancer subtypes via selecting relevant
 1329 expression data. *Frontiers in Genetics*, 11, 2020.
- 1330 56. Giorgio Valentini, Giuliano Armano, Marco Frasca, Jianyi Lin, Marco
 1331 Mesiti, and Matteo Re. RANKS: a flexible tool for node label ranking
 1332 and classification in biological networks. *Bioinformatics*, 32(18):2872–
 1333 2874, 06 2016.
- 1334 57. Yasin Ilkagan Tepeli, Ali Burak Ünal, Furkan Mustafa Akdemir, and
 Oznur Tastan. Pamogk: A pathway graph kernel based multi-omics
 1336 approach for patient clustering. *Bioinformatics*, 2020.
- 1337 58. Yuqi Wen, Xinyu Song, Bowei Yan, Xiaoxi Yang, Lianlian Wu, Dongjin
 1338 Leng, Song He, and Xiaochen Bo. Multi-dimensional data integration
 1339 algorithm based on random walk with restart. *BMC bioinformatics*,
 1340 22(1):1–22, 2021.
- 1341 59. Pablo A Jaskowiak, Ricardo JGB Campello, and Ivan G Costa.
 1342 Proximity measures for clustering gene expression microarray data:
 1343 a validation methodology and a comparative analysis. *IEEE/ACM*1344 *transactions on computational biology and bioinformatics*, 10(4):845–857,
 1345 2013.
- Pablo A Jaskowiak, Ricardo JGB Campello, and Ivan G Costa. On the
 selection of appropriate distances for gene expression data clustering.
 BMC bioinformatics, 15(2):1–17, 2014.
- 1349 61. Chihyun Park, Jaegyoon Ahn, Hyunjin Kim, and Sanghyun Park.
 1350 Integrative gene network construction to analyze cancer recurrence
 1351 using semi-supervised learning. *PLOS ONE*, 9(1):1–9, 01 2014.
- Karolina Sienkiewicz, Jinyu Chen, Ajay Chatrath, John T Lawson,
 Nathan C Sheffield, Louxin Zhang, and Aakrosh Ratan. Detecting
 molecular subtypes from multi-omics datasets using sumo. *Cell Reports Methods*, page 100152, 2022.
- 1356 63. Hongjiu Zhang, Fan Zhu, Hiroko H Dodge, Gerald A
 1357 Higgins, Gilbert S Omenn, Yuanfang Guan, and Alzheimer's

Disease Neuroimaging Initiative.A similarity-based approach to1358leverage multi-cohort medical data on the diagnosis and prognosis of1359alzheimer's disease.GigaScience, 7(7):giy085, 2018.1360

- Fayao Liu, Luping Zhou, Chunhua Shen, and Jianping Yin. Multiple 1361
 kernel learning in the primal for multimodal alzheimer's disease classification. *IEEE journal of biomedical and health informatics*, 1363
 18(3):984–990, 2013. 1364
- Mingxin Tao, Tianci Song, Wei Du, Siyu Han, Chunman Zuo, Ying Li, Yan Wang, and Zekun Yang. Classifying breast cancer subtypes using multiple kernel learning based on omics data. *Genes*, 10(3):200, 2019. 1367
- Mehmet Gönen and Ethem Alpaydın. Multiple kernel learning 1368 algorithms. *The Journal of Machine Learning Research*, 12:2211–2268, 1369 2011. 1370
- Ya Zhang, Ao Li, Jie He, and Minghui Wang. A novel mkl method 1371 for gbm prognosis prediction by integrating histopathological image and multi-omics data. *IEEE journal of biomedical and health informatics*, 24(1):171–179, 2019. 1374
- Nello Cristianini and Bernhard Scholkopf. Support vector machines 1375 and kernel methods: the new generation of learning machines. Ai 1376 Magazine, 23(3):31–31, 2002. 1377
- Dongdong Sun, Ao Li, Bo Tang, and Minghui Wang. Integrating genomic data and pathological images to effectively predict breast cancer clinical outcome. *Computer methods and programs in biomedicine*, 161:45–53, 2018.
 1381
- Fabio Aiolli and Michele Donini. Easymkl: a scalable multiple kernel learning algorithm. *Neurocomputing*, 169:215–224, 2015.
- Zenglin Xu, Rong Jin, Haiqin Yang, Irwin King, and Michael R
 Lyu. Simple and efficient multiple kernel learning by group lasso.
 In Proceedings of the 27th international conference on machine learning (ICML-10), pages 1175–1182. Citeseer, 2010.
- 72. Taiji Suzuki and Ryota Tomioka. Spicymkl: a fast algorithm for 1388 multiple kernel learning with thousands of kernels. *Machine learning*, 1389 85(1-2):77–108, 2011.
 1390
- Fei Yan, Josef Kittler, Krystian Mikolajczyk, and Atif Tahir. Nonsparse multiple kernel fisher discriminant analysis. *The Journal of Machine Learning Research*, 13(1):607–642, 2012.
 1393
- Gaston Baudat and Fatiha Anouar. Generalized discriminant analysis 1394 using a kernel approach. *Neural computation*, 12(10):2385–2404, 2000. 1395
- Cheng Soon Ong and Alexander Zien. An automated combination of kernels for predicting protein subcellular localization. In *International Workshop on Algorithms in Bioinformatics*, pages 186–197. Springer, 2008.
- Nora K Speicher and Nico Pfeifer. Integrating different data types by regularized unsupervised multiple kernel learning with application to cancer subtype discovery. *Bioinformatics*, 31(12):i268–i275, 2015.
- Xinwang Liu, Yong Dou, Jianping Yin, Lei Wang, and En Zhu. Multiple
 kernel k-means clustering with matrix-induced regularization. In
 Proceedings of the AAAI Conference on Artificial Intelligence, 2016.
- Yen-Yu Lin, Tyng-Luh Liu, and Chiou-Shann Fuh. Multiple kernel learning for dimensionality reduction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(6):1147–1160, 2011. 1408
- Xiaofei He and Partha Niyogi. Locality preserving projections. 1409 Advances in neural information processing systems, 16(16):153–160, 1410 2004. 1411
- Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. 1412 Nonlinear component analysis as a kernel eigenvalue problem. *Neural* 1413 *computation*, 10(5):1299–1319, 1998. 1414
- Giorgio Valentini, Alberto Paccanaro, Horacio Caniza, Alfonso E. 1415
 Romero, and Matteo Re. An extensive analysis of disease-gene associations using network integration and fast kernel-based gene prioritization methods. *Artificial Intelligence in Medicine*, 61(2):63–78, 1418

1496

1497

	2014
1419	2014.

- 1420 82. Judea Pearl. Probabilistic reasoning in intelligent systems: networks of
 plausible inference. Elsevier, 2014.
- 1422 83. Giulia Tini, Luca Marchetti, Corrado Priami, and Marie-Pier Scott1423 Boyer. Multi-omics integration—a comparison of unsupervised
 1424 clustering methodologies. *Briefings in bioinformatics*, 20(4):1269–
 1425 1279, 2019.
- 1426 84. Evan G Williams, Yibo Wu, Pooja Jha, Sébastien Dubuis, Peter
 1427 Blattmann, Carmen A Argmann, Sander M Houten, Tiffany Amariuta,
 1428 Witold Wolski, Nicola Zamboni, et al. Systems proteomics of liver
 1429 mitochondria function. *Science*, 352(6291), 2016.
- 1430 85. Anne Zufferey, Mark Ibberson, Jean-Luc Reny, Séverine Nolli,
 Domitille Schvartz, Mylène Docquier, Ioannis Xenarios, Jean-Charles
 1432 Sanchez, and Pierre Fontana. New molecular insights into modulation
 of platelet reactivity in aspirin-treated patients using a network-based
 1434 approach. Human genetics, 135(4):403–414, 2016.
- 86. DCFR Koboldt, Robert Fulton, Michael McLellan, Heather Schmidt,
 Joelle Kalicki-Veizer, Joshua McMichael, Lucinda Fulton, David
 Dooling, Li Ding, Elaine Mardis, et al. Comprehensive molecular
 portraits of human breast tumours. *Nature*, 490(7418):61–70, 2012.
- 1439 87. Tianle Ma and Aidong Zhang. Integrate multi-omic data using affinity network fusion (anf) for cancer patient clustering. In 2017 1441 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pages 398–403. IEEE, 2017.
- 1443 88. Alessandro Rozza, Gabriele Lombardi, Elena Casiraghi, and Paola
 1444 Campadelli. Novel fisher discriminant classifiers. *Pattern Recognition*,
 1445 45(10):3725–3737, 2012.
- 1446 89. Shuhui Liu and Xuequn Shang. Hierarchical similarity network
 1447 fusion for discovering cancer subtypes. In *International Symposium*1448 on *Bioinformatics Research and Applications*, pages 125–136. Springer,
 1449 2018.
- Limin Jiang, Yongkang Xiao, Yijie Ding, Jijun Tang, and Fei Guo.
 Discovering cancer subtypes via an accurate fusion strategy on multiple profile data. *Frontiers in genetics*, 10:20, 2019.
- Mikko Kivelä, Alex Arenas, Marc Barthelemy, James P Gleeson, Yamir
 Moreno, and Mason A Porter. Multilayer networks. *Journal of complex networks*, 2(3):203–271, 2014.
- 1456 92. Alberto Valdeolivas, Laurent Tichit, Claire Navarro, Sophie Perrin,
 1457 Gaelle Odelin, Nicolas Levy, Pierre Cau, Elisabeth Remy, and Anaïs
 1458 Baudot. Random walk with restart on multiplex and heterogeneous
 1459 biological networks. *Bioinformatics*, 35(3):497–505, 2019.
- Nimrod Rappoport and Ron Shamir. Multi-omic and multi-view
 clustering algorithms: review and cancer benchmark. *Nucleic acids research*, 46(20):10546–10562, 2018.
- 1463 94. Nimrod Rappoport and Ron Shamir. Nemo: Cancer subtyping by
 1464 integration of partial multi-omic data. *Bioinformatics*, 35(18):3348–
 1465 3356, 2019.
- Jian Liu, Shuguang Ge, Yuhu Cheng, and Xuesong Wang. Multi-view
 spectral clustering based on multi-smooth representation fusion for
 cancer subtype prediction. *Frontiers in Genetics*, page 1574, 2021.
- Peiping Nie, Jing Li, Xuelong Li, et al. Self-weighted multiview
 clustering with multiple graphs. In *IJCAI*, pages 2564–2570, 2017.
- 1471 97. Qing Feng, Meilei Jiang, Jan Hannig, and JS Marron. Angle-based joint
 1472 and individual variation explained. *Journal of multivariate analysis*,
 166:241–265, 2018.
- 1474 98. Zi Yang and George Michailidis. A non-negative matrix factorization
 1475 method for detecting modules in heterogeneous omics multi-modal
 1476 data. *Bioinformatics*, 32(1):1–8, 2016.
- 1477 99. Johan A Westerhuis, Theodora Kourti, and John F MacGregor.1478 Analysis of multiblock and hierarchical pca and pls models. *Journal*
- 1479 of chemometrics, 12(5):301–321, 1998.

- Huichen Zhu, Gen Li, and Eric F Lock. Generalized integrative 1480 principal component analysis for multi-type data with block-wise missing structure. *Biostatistics*, 21(2):302–318, 2020. 1482
- 101. Giovanni Ciriello, Michael L Gatza, Andrew H Beck, Matthew D
 1483
 Wilkerson, Suhn K Rhie, Alessandro Pastore, Hailei Zhang, Michael
 1484
 McLellan, Christina Yau, Cyriac Kandoth, et al. Comprehensive
 molecular portraits of invasive lobular breast cancer. *Cell*, 163(2):506–
 519, 2015.
 1487
- 102. Daniela M Witten and Robert J Tibshirani. Extensions of sparse 1488 canonical correlation analysis with applications to genomic data. 1489 Statistical applications in genetics and molecular biology, 8(1), 2009. 1490
- 103. Georg Lenz, George W Wright, NC Tolga Emre, Holger Kohlhammer, 1491
 Sandeep S Dave, R Eric Davis, Shannon Carty, Lloyd T Lam, 1492
 AL Shaffer, Wenming Xiao, et al. Molecular subtypes of diffuse large b-cell lymphoma arise by distinct genetic pathways. *Proceedings of the National Academy of Sciences*, 105(36):13520–13525, 2008. 1495
- 104. Arthur Tenenhaus and Michel Tenenhaus. Regularized generalized canonical correlation analysis. *Psychometrika*, 76(2):257, 2011.
- 105. Imene Garali, Isaac M Adanyeguh, Farid Ichou, Vincent Perlbarg, 1498
 Alexandre Seyer, Benoit Colsch, Ivan Moszer, Vincent Guillemot, 1499
 Alexandra Durr, Fanny Mochel, et al. A strategy for multimodal 1500
 data integration: application to biomarkers identification in 1501
 spinocerebellar ataxia. *Briefings in bioinformatics*, 19(6):1356–1369, 1502
 2018.
- 106. Age K Smilde, Johan A Westerhuis, and Sijmen de Jong. A 1504 framework for sequential multiblock component methods. Journal 1505 of Chemometrics: A Journal of the Chemometrics Society, 17(6):323–337, 1506 2003.
- 107. M. De Tayrac, L[^] E.S. Aubry, M. Mosser, and J.F. Husson.
 Simultaneous analysis of distinct omics data sets with integration
 of biological knowledge: Multiple factor analysis approach. BMC
 Genomics, 10:32, 2009.
- 108. NE Kucukboyaci, N Kemmotsu, KM Leyden, HM Girard, ES Tecoma,
 VJ Iragui, and CR McDonald. Integration of multimodal mri data via
 pca to explain language performance. *NeuroImage: Clinical*, 5:197–
 1514
 207, 2014.
- Maxime Chamberland, Erika P Raven, Sila Genc, Kate Duffy, Maxime
 Descoteaux, Greg D Parker, Chantal MW Tax, and Derek K Jones.
 Dimensionality reduction of diffusion mri measures for improved tractometry of the human brain. *NeuroImage*, 200:89–100, 2019.
- 110. Bryce Landon Geeraert, Maxime Chamberland, R Marc Lebel, and
 Catherine Lebel. Multimodal principal component analysis to identify
 major features of white matter structure and links to reading. *bioRxiv*,
 2020.
- 111. Bradley Worley and Robert Powers. A sequential algorithm for 1524 multiblock orthogonal projections to latent structures. *Chemometrics* 1525 *and Intelligent Laboratory Systems*, 149:33–39, 2015. 1526
- 112. Li Zhang, Chenkai Lv, Yaqiong Jin, Ganqi Cheng, Yibao Fu, 1527
 Dongsheng Yuan, Yiran Tao, Yongli Guo, Xin Ni, and Tieliu Shi. Deep learning-based multi-omics data integration reveals two prognostic 1529
 subtypes in high-risk neuroblastoma. *Frontiers in genetics*, 9:477, 2018. 1530
- 113. Svante Wold, Michael Sjöström, and Lennart Eriksson. Plsregression: a basic tool of chemometrics. Chemometrics and intelligent laboratory systems, 58(2):109–130, 2001.
- 114. Erica Ponzi, Magne Thoresen, Therese Haugdahl Nøst, and Kajsa
 Møllersen. Integrative, multi-omics, analysis of blood samples
 improves model predictions: applications to cancer. *bioRxiv*, pages
 2020–10, 2021.
- 115. Leo Breiman. Random forests. Machine learning, 45(1):5–32, 2001. 1538
- Florian Rohart, Benoit Gautier, Amrit Singh, and Kim-Anh Lê Cao.
 mixomics: An r package for 'omics feature selection and multiple data
 1540

- 1541 integration. PLoS computational biology, 13(11):e1005752, 2017.
- 117. Sini Isokääntä, Eetu Kari, Angela Buchholz, Liqing Hao, Siegfried Schobesberger, Annele Virtanen, and Santtu Mikkonen.
 Comparison of dimension reduction techniques in the analysis of mass spectrometry data. *Atmospheric Measurement Techniques*, 13(6):2995–3022, 2020.
- 118. Nicolas Gillis. Sparse and unique nonnegative matrix factorization
 through data preprocessing. *The Journal of Machine Learning Research*,
 13(1):3349–3386, 2012.
- 119. Yifeng Li. Advances in multi-view matrix factorizations. In 2016
 International Joint Conference on Neural Networks (IJCNN), pages 3793– 3800. IEEE, 2016.
- 120. Patrik O Hoyer. Non-negative matrix factorization with sparseness
 constraints. *Journal of machine learning research*, 5(9), 2004.
- 121. Shihua Zhang, Chun-Chi Liu, Wenyuan Li, Hui Shen, Peter W Laird,
 and Xianghong Jasmine Zhou. Discovery of multi-dimensional
 modules by integrative analysis of cancer genomic data. *Nucleic acids research*, 40(19):9379–9391, 2012.
- 122. Cécile Chauvel, Alexei Novoloaca, Pierre Veyre, Frédéric Reynier, and
 Jérémie Becker. Evaluation of integrative clustering methods for the
 analysis of multi-omics data. *Briefings in Bioinformatics*, 21(2):541–
 552, 2020.
- 123. Ronglai Shen, Adam B Olshen, and Marc Ladanyi. Integrative clustering of multiple genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis.
 Bioinformatics, 25(22):2906–2912, 2009.
- 1567 124. Eric F Lock and David B Dunson. Bayesian consensus clustering.
 1568 Bioinformatics, 29(20):2610–2616, 2013.
- 125. Paul Kirk, Jim E Griffin, Richard S Savage, Zoubin Ghahramani, and
 David L Wild. Bayesian correlated clustering to integrate multiple
 datasets. *Bioinformatics*, 28(24):3290–3297, 2012.
- 126. Ying-Lian Gao, Mi-Xiao Hou, Jin-Xing Liu, and Xiang-Zhen Kong. An
 integrated graph regularized non-negative matrix factorization model
 for gene co-expression network analysis. *IEEE Access*, 7:126594–
 126602, 2019.
- 127. Deng Cai, Xiaofei He, Jiawei Han, and Thomas S Huang. Graph
 regularized nonnegative matrix factorization for data representation.
 IEEE transactions on pattern analysis and machine intelligence,
 33(8):1548–1560, 2010.
- 128. Fei Wang, Tao Li, and Changshui Zhang. Semi-supervised clustering
 via matrix factorization. In *Proc. SIAM Int. Conf. on Data Mining*, 2008.
- 129. F Vitali, S Marini, D Pala, A Demartini, S Montoli, A Zambelli,
 and R Bellazzi. Patient similarity by joint matrix trifactorization to
 identify subgroups in acute myeloid leukemia. *JAMIA open*, 1(1):75–
 86, 2018.
- 130. Ricard Argelaguet, Damien Arnol, Danila Bredikhin, Yonatan Deloro,
 Britta Velten, John C Marioni, and Oliver Stegle. Mofa+: a statistical
 framework for comprehensive integration of multi-modal single-cell
 data. *Genome Biology*, 21(1):1–17, 2020.
- 131. Qianxing Mo, Sijian Wang, Venkatraman E Seshan, Adam B Olshen,
 Nikolaus Schultz, Chris Sander, R Scott Powers, Marc Ladanyi, and
 Ronglai Shen. Pattern discovery and cancer gene identification in
 integrated cancer genomic data. *Proceedings of the National Academy* of Sciences, 110(11):4245–4250, 2013.
- 132. Rogers F Silva, Sergey M Plis, Jing Sui, Marios S Pattichis, Tülay
 Adalı, and Vince D Calhoun. Blind source separation for unimodal
 and multimodal brain networks: A unifying framework for subspace
 modeling. *IEEE journal of selected topics in signal processing*, 10(7):1134–
 1149, 2016.
- 133. Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.

- 134. Guoxu Zhou, Qibin Zhao, Yu Zhang, Tülay Adalı, Shengli Xie, and
 Andrzej Cichocki. Linked component analysis from matrices to high order tensors: Applications to biomedical data. *Proceedings of the IEEE*,
 104(2):310–331, 2016.
- 135. Tulay Adali, Matthew Anderson, and Geng-Shen Fu. Diversity 1606
 in independent component and vector analyses: Identifiability, 1607
 algorithms, and applications in medical imaging. *IEEE Signal* 1608
 Processing Magazine, 31(3):18–33, 2014. 1609
- 136. Dana Lahat, Tülay Adali, and Christian Jutten. Multimodal data 1610 fusion: an overview of methods, challenges, and prospects. *Proceedings* 1611 of the IEEE, 103(9):1449–1477, 2015. 1612
- Pierre Comon and Christian Jutten. Handbook of Blind Source 1613 Separation: Independent component analysis and applications. Academic 1614 press, 2010. 1615
- Aapo Hyvärinen and Erkki Oja. Independent component analysis: 1616 algorithms and applications. Neural networks, 13(4-5):411–430, 2000. 1617
- 139. V Calhoun, Tulay Adali, and Jingyu Liu. A feature-based approach to combine functional mri, structural mri and eeg brain imaging data. In 1619
 2006 International Conference of the IEEE Engineering in Medicine and Biology Society, pages 3672–3675. IEEE, 2006. 1621
- 140. Matthias Moosmann, Tom Eichele, Helge Nordby, Kenneth Hugdahl,
 and Vince D Calhoun. Joint independent component analysis for
 simultaneous eeg-fmri: principle and simulation. *International Journal* of Psychophysiology, 67(3):212–221, 2008.
- 141. Taesu Kim, Torbjørn Eltoft, and Te-Won Lee. Independent vector analysis: An extension of ica to multivariate components. In International conference on independent component analysis and signal separation, pages 165–172. Springer, 2006.
- 142. Tülay Adali, MABS Akhonda, and Vince D Calhoun. Ica and iva for data fusion: An overview and a new approach based on disjoint subspaces. *IEEE sensors letters*, 3(1):1–4, 2018.
- 143. Marinka Žitnik and Blaž Zupan. Matrix factorization-based data fusion for drug-induced liver injury prediction. Systems Biomedicine, 2(1):16–22, 2014.
 1635
- 144. Tin Nguyen, Rebecca Tagett, Diana Diaz, and Sorin Draghici. A novel approach for data integration and disease subtyping. *Genome research*, 27(12):2025–2039, 2017.
- 145. John H Phan, Ryan Hoffman, Sonal Kothari, Po-Yen Wu, and May D
 Wang. Integration of multi-modal biomedical data to predict
 cancer grade and patient survival. In 2016 IEEE-EMBS International
 Conference on Biomedical and Health Informatics (BHI), pages 577–580.
 IEEE, 2016.
- 146. Jean-Philippe Brunet, Pablo Tamayo, Todd R. Golub, and Jill P.
 Mesirov. Metagenes and molecular pattern discovery using matrix factorization. Proceedings of the National Academy of Sciences, 101(12):4164-4169, 2004.
- 147. Matthew D Wilkerson and D Neil Hayes. Consensusclusterplus: a 1648 class discovery tool with confidence assessments and item tracking. *Bioinformatics*, 26(12):1572–1573, 2010. 1650
- 148. Stefano Monti, Pablo Tamayo, Jill Mesirov, and Todd Golub.
 Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data. *Machine learning*, 52(1):91–118, 2003.
- 149. Fionn Murtagh and Pedro Contreras. Algorithms for hierarchical 1655 clustering: an overview. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 2(1):86–97, 2012. 1657
- 150. Leonard Kaufmann and Peter Rousseeuw. Clustering by means of medoids. Data Analysis based on the L1-Norm and Related Methods, pages 405–416, 01 1987.
 1660
- 151. Peter Langfelder, Bin Zhang, and Steve Horvath. Defining clusters 1661 from a hierarchical cluster tree: the dynamic tree cut package for r. 1662

- 1663 Bioinformatics, 24(5):719–720, 2008.
- 1664 152. Dingming Wu, Dongfang Wang, Michael Q Zhang, and Jin Gu.
- Fast dimension reduction and integrative clustering of multi-omics
 data using low-rank approximation: application to cancer molecular
 classification. *BMC genomics*, 16(1):1022, 2015.
- 1668 153. Chris Ding and Hanchuan Peng. Minimum redundancy feature 1669 selection from microarray gene expression data. *Journal of*

bioinformatics and computational biology, 3(02):185–205, 2005.

- 1671 154. David H Wolpert. Stacked generalization. *Neural networks*, 5(2):241–
 259, 1992.
- 155. Clifford R Jack Jr, Matt A Bernstein, Nick C Fox, Paul Thompson,
 Gene Alexander, Danielle Harvey, Bret Borowski, Paula J Britson,
 Jennifer L. Whitwell, Chadwick Ward, et al. The alzheimer's disease
 neuroimaging initiative (adni): Mri methods. Journal of Magnetic
 Resonance Imaging: An Official Journal of the International Society for
 Magnetic Resonance in Medicine, 27(4):685-691, 2008.
- 156. Simon Lovestone, Paul Francis, Iwona Kloszewska, Patrizia Mecocci,
 Andrew Simmons, Hilkka Soininen, Christian Spenger, Magda
 Tsolaki, Bruno Vellas, Lars-Olof Wahlund, et al. Addneuromed—the
 european collaboration for the discovery of novel biomarkers for
 alzheimer's disease. Annals of the New York Academy of Sciences,
 1180(1):36–46, 2009.
- 157. Sherry-Ann Brown. Patient similarity: emerging concepts in systems
 and precision medicine. *Frontiers in physiology*, 7:561, 2016.
- 158. Leyu Dai, He Zhu, and Dianbo Liu. Patient similarity: methods and
 applications. arXiv preprint arXiv:2012.01976, 2020.
- 159. Jie Xu, Benjamin S Glicksberg, Chang Su, Peter Walker, Jiang Bian,
 and Fei Wang. Federated learning for healthcare informatics. *Journal of Healthcare Informatics Research*, 5(1):1–19, 2021.
- 1692 160. Albert-László Barabási, Natali Gulbahce, and Joseph Loscalzo.
 1693 Network medicine: a network-based approach to human disease.
 1694 Nature reviews genetics, 12(1):56–68, 2011.
- 161. Raffaele Giancarlo, Giosuè Lo Bosco, and Luca Pinello. Distance
 functions, clustering algorithms and microarray data analysis. In
 International Conference on Learning and Intelligent Optimization, pages
 1698 125–138. Springer, 2010.
- 1699 162. Taiyun Kim, Irene Rui Chen, Yingxin Lin, Andy Yi-Yang Wang,
 Jean Yee Hwa Yang, and Pengyi Yang. Impact of similarity metrics
 1701 on single-cell rna-seq data clustering. *Briefings in bioinformatics*,
 1702 20(6):2316-2326, 2019.
- 163. Ian T Jolliffe. Principal component analysis: a beginner's guide—ii.
 pitfalls, myths and extensions. *Weather*, 48(8):246–253, 1993.
- 1705 164. Yunjin Choi, Jonathan Taylor, and Robert Tibshirani. Selecting the
 number of principal components: Estimation of the true rank of a
 1707 noisy matrix. *The Annals of Statistics*, pages 2590–2617, 2017.
- 165. Gen Li and Irina Gaynanova. A general framework for association
 analysis of heterogeneous data. *The Annals of Applied Statistics*,
 12(3):1700–1726, 2018.
- 166. P Campadelli, E Casiraghi, C Ceruti, and A Rozza. Intrinsic dimension
 estimation: Relevant techniques and a benchmark framework.
 Mathematical Problems in Engineering, 2015, 2015.
- 167. Paolo Perlasca, Marco Frasca, Cheick Ba, Marco Notaro, Alessandro
 Petrini, Elena Casiraghi, Giuliano Grossi, Jessica Gliozzo, Giorgio
 Valentini, and Marco Mesiti. Unipred-web: a web tool for the
 integration and visualization of biomolecular networks for protein
 function prediction *BMC Bioinformatics* 20, 12 2019
- 168. Paolo Perlasca, Marco Frasca, Cheick Tidiane Ba, Jessica Gliozzo,
 Marco Notaro, Mario Pennacchioni, Giorgio Valentini, and Marco
 Mesiti. Multi-resolution visualization and analysis of biomolecular
 networks through hierarchical community detection and web-based
 graphical tools. *PLOS ONE*, 15(12):1–28, 12 2020.

- 169. Sylvia Richardson, George C Tseng, and Wei Sun. Statistical methods in integrative genomics. Annual review of statistics and its application, 3:181–209, 2016. 1726
- 170. Qing Zhao, Xingjie Shi, Jian Huang, Jin Liu, Yang Li, and Shuangge
 Ma. Integrative analysis of '-omics' data using penalty functions. Wiley
 Interdisciplinary Reviews: Computational Statistics, 7(1):99–108, 2015.
 1729
- 171. Wenting Wang, Veerabhadran Baladandayuthapani, Jeffrey S Morris, 1730
 Bradley M Broom, Ganiraju Manyam, and Kim-Anh Do. ibag: 1731
 integrative bayesian analysis of high-dimensional multiplatform 1732
 genomics data. *Bioinformatics*, 29(2):149–159, 2013. 1733
- 172. Ruoqing Zhu, Qing Zhao, Hongyu Zhao, and Shuangge Ma. 1734
 Integrating multidimensional omics data for cancer outcome. 1735
 Biostatistics, 17(4):605–618, 2016. 1736
- 173. Cen Wu, Qingzhao Zhang, Yu Jiang, and Shuangge Ma. Robust 1737 network-based analysis of the associations between (epi) genetic 1738 measurements. *Journal of multivariate analysis*, 168:119–130, 2018. 1739
- 174. Sijia Huang, Kumardeep Chaudhary, and Lana X Garmire. More
 is better: recent progress in multi-omics data integration methods.
 Frontiers in genetics, 8:84, 2017.
 1740
- 175. Elad Noor, Sarah Cherkaoui, and Uwe Sauer. Biological insights
 1743 through omics data integration. *Current Opinion in Systems Biology*, 1744
 15:39–47, 2019.



Jessica Gliozzo. Jessica Gliozzo is a PhD 1746 student in Computer Science enrolled in the 1747 Collaborative Doctoral Partnership program 1748 between the University of Milan and the Joint 1749 Research Center of European Commission. 1750 Her latest research works comprehend the 1751 development of a multi-modal semi-supervised 1752 method based on patient similarity networks 1753 for patients' outcome prediction; the application 1754

of deep neural networks to predict the tissue-specific activity status of
cis-regulatory regions in the genome (i.e. promoters and enhancers);
the use of compression methods to obtain compact representations
of convolution neural networks in the biological domain (e.g. ki67
and TIL-index prediction), showing their adavantages when limited
computational resources are available. She is author of a few works
in the fields of Machine Learning and Bioinformatics.1755



Marco Mesiti. He is associate professor at 1762 the Department of Computer Science Giovanni 1763 degli Antoni, Università degli Studi di Milano. 1764 He has got a master and phD degree from 1765 the University of Genova in 1998 and 2003. 1766 His research interest are in the integration. 1767 querying and visualization of different kinds 1768 of information (structured and semi-structured) 1769 according to different data models (relational, 1770

graph, and nosql). Moreover, he has involved in different projects for1771protein network integration, protein function prediction, and protein1772networks visualization. On these topics is ha published more than 1001773articles in international conferences and journals. He is associate editor1774for the Springer Data Science and Engineering Journal and MDPI1775Applied Sciences.1776



Marco Notaro. He is a postdoctoral fellow at
the Computer Science Department of Milan
University. His research interests touch the
fields of Bioinformatics, Computational Biology,1780
Biological Network and Machine Learning.
His main expertise is the analysis and construction
rese
of complex biomolecular networks and the
design and implementation of output-structured1784

learning algorithms to discover novel gene-1785 disease associations or to predict novel protein function. His Ph.D. 1786 paper was awarded by International Medical Informatics Association 1787 as one of the best five papers of 2017 in the field of Medical Informatics. 1788



Alessandro Petrini. He is a postdoctoral researcher at the department of Computer Science of Università degli Studi di Milano. He is currently a member of the laboratory of Bioinformatics and Computational Biology - AnacletoLab - and his main research is focused on High Performance Computing and Machine Learning. He is author of more than 30 articles in international journals and

conferences. He designed and developed parallel and accelerated ML 1798 algorithms for image and video processing / encoding / compression, 1799 omics analisys, graph modeling and analysis, data visualization, neural 1800 network compression, MRI volumes processing and analysis. 1801

1804

1807

1821 1822

1824

1825

1826

1827 1828

1829

1831

1832

1833

1834

1835

1836

1837

1838

1839



Alex Patak. Alex Patak, PhD MD, graduated in Medicine and Surgery at the School of Medicine at "Universidad Autónoma de Barcelona", Barcelona (Spain) and holds a Master in Medical Bioengineering from the "Universidad Politécnica de Cataluña". At the Instituto Municipal de Investigación Médica (Barcelona) he has been working on Expert Systems for medical diagnostic and did

his PhD on computer-assisted medical education after a stage at 1811 Dartmouth Medical College in Vermont (USA). Since 1994 works 1812 at the Joint Research Centre in Ispra (Italy) where he has been 1813 working on Three-dimensional medical imaging, and from 2003 to 1814 2017 was responsible for the Bioinformatics team at the Molecular 1815 Biology and Genomics Unit of the Institute for Health and Consumer 1816 Protection in Ispra. He is now a team leader at Knowledge for Health 1817 1818 & Consumer Safety and is responsible for the Collaborative Doctoral 1819 Partnership programme in Genomics and Bioinformatics, working on the application of artificial intelligence to omics data and Microbiome. 1820



and Machine Learning. 1830



Antonio Puertas-Gallardo. Antonio Puertas Gallardo is an IT Project Manager at the Joint Research Center (JRC) of the European Commission. He provides High Performance Computing (HPC) support to bioinformatics members of the Knowledge for Health and Consumer Safety Unit at JRC, and he has recently begun to collaborate with the unit's data scientists on Natural Language Processing

Alberto Paccanaro. Alberto Paccanaro is full Professor in Machine Learning and Computational Biology at the School of Applied Mathematics of the Fundação Getúlio Vargas in Rio de Janeiro and at the Department of Computer Science at Royal Holloway University of London, where he is also Director of the Centre for Systems and Synthetic Biology. He completed his

undergraduate studies in Computer Science at the University of Milan 1840 1841 and received his PhD from the University of Toronto in 2002. His research interests are in applying and developing machine learning 1842 1843 algorithms for solving problems in molecular biology, medicine and pharmacology and he has led a number of international research 1844 1845 projects in this area.



Giorgio Valentini. Giorgio Valentini is a full 1846 Professor at the Department of Computer 1847 Science, University of Milan (UNIMI). 1848 Director for UNIMI of the European doctorate 1849 in Genomics and Bioinformatics in collaboration1850 with the Joint Research Center of the European 1851 Union. Director of AnacletoLab, Computational 1852 Biology and Bioinformatics Laboratory of 1853 the Department of Computer Science of 1854

the University of Milano. He has been PI in several national 1855 and international research projects funded by public and private 1856 institutions in the area of Bioinformatics, Machine learning and Big 1857 Data analytics. He is author of over 150 scientific publications with 1858 peer-review in collaboration with several research groups in Europe 1859 and America in the field of Bioinformatics, Computational Biology and 1860 Machine Learning. 1861



Elena Casiraghi. She is associate professor 1862 at the Department of Computer Science 1863 Giovanni degli Antoni, Università degli Studi 1864 di Milano. She is co-lead of the AnacletoLab, 1865 Computational Biology and Bioinformatics 1866 Laboratory of the Department of Computer 1867 Science of the University of Milano. 1868 Her research interests are focused on the 1869

development of applications that process 1870 biomedical data (images, multi-omics, clinical) for causal inference, 1871 inductive inference, machine learning, and pattern recognition. She 1872 currently cooperates with several hospitals and reserach centers in 1873 Italy (Istituto Nazionale dei Tumori, Ospedale Policlinico e Regina 1874 Margherita, Ospedale San Raffaele, Humanitas), Europe (Berlin 1875 Institute of Health at Charité, grupo espanol de investigacion de 1876 sarcomas, groupe sarcomes francais), and United States (Berkeley 1877 Laboratories, Jackson Laboratories) and she is the scientific supervisor 1878 for UNIMI of all the researches in cooperation with the N3C Enclave 1879 (USA, https://covid.cd2h.org/) created and funded by the NIH. 1880

CPCA [99] x		Caramany		approach		
	Simulated	Not provided	Numeric	PCA	Data Analysis	
CPCA for missing data [100] x	Human Mortality Database (Italy + Switzerland)	143	exposure-to-risk	PCA	Data Analysis	
JIVE[18] x	TCGA BIC	348	mRNA miRNA methy RPPA	PCA	Unsupervised Clustering (Patient subtype identification)	R code
ajIVE [97] x	TCGA extract from [101]	616	mRNA miRNA somatic mutation CNV RPPA	PCA	Data Analysis and Unsupervised Clustering (Patient subtype identification)	R code
MCCA [102] x	DLBCL Dataset [103]	203	mRNA array CGH measurements	CCA	Data analysis	R code
RGCCA [104] x SGCCA [19] [105]	SCA Dataset Private	67 SCA + 35 Healthy	pons volume metabolic features	CCA	Data Analysis	RGCCA/SGCCA R code
DIABLO [20] x	TCGA COAD TCGA KIRC TCGA GBM TCGA LUSC TCGA BRCA	92 122 213 106 989	mRNA miRNA methy	(SG)CCA	Data Analysis and Supervised Clustering (Patient's Survival)	R code

Table 6. PCA-based and CCA-based input data-fusion methods. For each method, the table reports: the name/acronym with the corresponding reference paper; whether it requires the same set of patients across all data modalities (i.e. "Matched Samples"); the dataset used to develop and evaluate the approach in the reference paper and the corresponding sample cardinality and data types composing the dataset; the exploited integration method; the application task and the code availability (with link to the repository and programming languages for which the code is available).

RNA; **mRNA**: messenger RNA; **PCA**: Principal Component Analysis; **RPPA**. **TCGA+cancer code**: The Cancer Genome Atlas+ link to complete cancer codes.

Name	Matched Samples	Dataset	Sample Cardinality	Data type	Integration approach	Task	Code and Language
MFA [107]		Brain Cancer Dataset Private	Not provided	multi-omics	MF	Data Analysis	R code
jNMF [121]	×	TCGA OV	385	mRNA miRNA methy	NMF	Data Analysis	ی موج د ک ک ک
iNMF [98]	×	TCGA OV	592	mRNA miRNA methy	NMF	Unsupervised Clustering (Patient subtype identification)	Python code
iGMFNA [126]	×	TCGA CHOL TCGA PAAD	45 180	mRNA methy CNV	NMF	Data Analysis	uuc,
MOFA+ [130]		Private	Not provided	multi-omics	NMF	Data Analysis	Python and R code
iCluster [123] and iCluster+ [131]	×	TCGA CRC	189	exome sequence mRNA methy CNV	NMF	Unsupervised Clustering (Patient subtype identification)	R code
DFMF ³ [31] [143]				GO terms GO annotations drugs tissue samples DILI potentials	MTF	Unsupervised Clustering (hepatotoxic risk associated with individual drugs)	Python code
MaDDa [129]		TCGA BRCA, BioGRID KEGG Disease Ontology DisGeNET	200	gene-gene interactions, gene-pathway associations disease-disease relationships, disease-gene associations, disease-pathway relations	MTF	Unsupervised Clustering (Patient subtype identification)	Matlab code
DS-ICA [142]	x	Private	38	features from EEG and fMRI images	ICA	Data Analysis	
MISA [21]	х	Private	1001	EEG images sMRI and fMRI images	BSS	Data Analysis	MATLAB code
Abbreviations BSS: Blind Sourc ICA: Independent Matrix Tri-Factor	ce Separation; CNV: C t Component Analysis; ization; NMF: Non-ne	opy Number Variation KEGG: Kyoto Encyclo igative Matrix Factorize	; DILI: Drug-Ind pedia of Genes ar ation; sMRI : stru	luced Liver Injury; EEG : Electroen d Genomes; methy : DNA methyls ctural Magnetic Resonance Imagin	teephalography; fMR I: fur ation; miRNA : micro RN/ g; TCGA+ <i>cancer code</i> : Th	ictional Magnetic Resonance Imagin y, MF: Matrix Factorization; mRNA . e Cancer Genome Atlas+ link to comp	ış; GO: Gene Ontology; : messenger RNA; MTF: blete cancer codes.

Table 7. MF-based input data-fusion methods. For each method, the table reports: the name/acronym with the corresponding reference paper; whether it requires the same set of patients across all data modalities (i.e. "Matched Samples"); the dataset used to develop and evaluate the approach in the reference paper and the corresponding sample cardinality and data types composing the dataset; the exploited integration method; the application task and the code availability (with link to the repository and programming languages for which the code is available).

Name	Dataset	Sample Cardinality	Data type	Integration approach	Task	Code and Language
COCA [22]	TCGA AML TCGA BIC TCGA COAD TCGA READ TCGA GBM TCGA GBM TCGA LUSC TCGA UUSC TCGA UUSC TCGA UUSC TCGA LUAD TCGA HNSC	161 834 182 73 475 238 329 329 329 230 270	DNA sequence mRNA miRNA methy CNV RPPA	Consensus Clustering	Unsupervised Clustering (Patient subtype identification)	ConsensusClusterPlus R code
PINS/PINSPlus [23]	34 TCGA datasets 2 Metabric datasets	12158	mRNA miRNA methy	Consensus Clustering	Unsupervised Clustering (Patient subtype identification)	R code
SUMO [62]	TCGA extract from NEMO [94] (Table 4)	3168 across	mRNA miRNA methy	Consensus Clustering	Unsupervised Clustering (Patient subtype identification)	Python code
FH-Clust [24]	TCGA AML TCGA BIC TCGA COAD TCGA GBM TCGA GBM TCGA LIHC TCGA LIHC TCGA LUSC TCGA SKCM TCGA OV TCGA SARC	170 621 220 274 183 367 287 287 287	mRNA miRNA methy	Consensus Clustering	Unsupervised Clustering (Patients' clusters related to known Survival)	R code
[145]	TCGA KIRC TCGA OV TCGA KIRC TCGA OV	418 250 220 160	histopathological image RNA-seq data	Stacked Generalization: linear regression of unimodal classifiers	Supervised Classification (Cancer grade < 3 vs Cancer grade >= 3) Supervised Classification (Known Survival < 5 vs Known Survival >= 5)	
[63]	ADNI Phase 1 [155] AddNeuroMed study [156]	628 for training 94 for validation 88 for testing	demographic data APOE e4 allele information anatomical brain features from 1.5T MRI scans	Average	Supervised Multiclass classification (HC vs MCI vs AD)	Python code
Abbreviations						

Table 8. Output-fusion methods. For each method, the table reports: the name/acronym with the corresponding reference paper; the dataset used to develop and evaluate the approach in the reference paper and the corresponding sample cardinality and data types composing the dataset; the exploited integration method; the application task and the code availability (with link to the repository and programming languages for which the code is available).