

RESEARCH

Open Access



Novel artificial intelligence-based identification of drug-gene-disease interaction using protein-protein interaction

Y.-h Taguchi^{1*} and Turki Turki²

*Correspondence:
tag@granular.com

¹ Department of Physics, Chuo University, 1-13-27 Kasuga, Bunkyo-ku, Tokyo 112-8551, Japan

² Department of Computer Science, King Abdulaziz University, Jeddah 21589, Saudi Arabia

Abstract

The evaluation of drug-gene-disease interactions is key for the identification of drugs effective against disease. However, at present, drugs that are effective against genes that are critical for disease are difficult to identify. Following a disease-centric approach, there is a need to identify genes critical to disease function and find drugs that are effective against them. By contrast, following a drug-centric approach comprises identifying the genes targeted by drugs, and then the diseases in which the identified genes are critical. Both of these processes are complex. Using a gene-centric approach, whereby we identify genes that are effective against the disease and can be targeted by drugs, is much easier. However, how such sets of genes can be identified without specifying either the target diseases or drugs is not known. In this study, a novel artificial intelligence-based approach that employs unsupervised methods and identifies genes without specifying neither diseases nor drugs is presented. To evaluate its feasibility, we applied tensor decomposition (TD)-based unsupervised feature extraction (FE) to perform drug repositioning from protein-protein interactions (PPI) without any other information. Proteins selected by TD-based unsupervised FE include many genes related to cancers, as well as drugs that target the selected proteins. Thus, we were able to identify cancer drugs using only PPI. Because the selected proteins had more interactions, we replaced the selected proteins with hub proteins and found that hub proteins themselves could be used for drug repositioning. In contrast to hub proteins, which can only identify cancer drugs, TD-based unsupervised FE enables the identification of drugs for other diseases. In addition, TD-based unsupervised FE can be used to identify drugs that are effective in *in vivo* experiments, which is difficult when hub proteins are used. In conclusion, TD-based unsupervised FE is a useful tool for drug repositioning using only PPI without other information.

Keywords: Protein-protein interaction, Drug repositioning, Artificial intelligence, Unsupervised learning, Tensor decomposition

Introduction

The identification of drug-gene-disease information is critical for drug repositioning. Nevertheless, due to the fact that this does not involve a simple identification of paired information but rather of triplet information, identifying drug-gene-disease information



© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

is not an easy task. The identification of drug-gene-disease information is often performed in two steps. From the perspective of the genes (i.e. drug-centric approach), the identification of drug-gene relationships is performed first, followed by the extraction of gene-disease information by identifying diseases related to the selected genes. By contrast, when starting from the perspective of the diseases (i.e. disease-centric approach), the identification of gene-disease relations is performed first, followed by the extraction of drug-gene information via the identification of drugs that target the genes selected. Since the existence of a significant relationship at the second stage is not guaranteed (i.e., gene-disease information for the drug-centric approach or drug-gene information for the disease-centric approach), the identification of drug-gene-disease information is much more difficult than the identification of only gene-disease information or drug-gene information.

Several studies have reported on such drug-gene-disease relationships. Zickenrott et al. [1] attempted to predict disease-gene-drug relationships using differential network analysis, while Wong et al. [2] searched for gene-drug-disease interactions in pharmacogenomics by using GeneDive. Yu et al. [3] predicted drugs with opposing effects on disease genes using a directed network, Sun [4] investigated gene-gene, drug-drug, and disease-disease networks and studied their relationships, and Qahwaji et al. [5] reviewed the genetic approaches to drug development and therapy. Furthermore, Iida et al. [6] investigated network-based characterization of disease-disease relationships in terms of drugs and therapeutic targets. Lastly, Quan et al. [7] considered genetic disease genes as promising sources of drug targets. While these are only some examples, all focused on diseases, drugs, or both, and no studies have focused on genes without considering drugs and diseases. Thus, the relevant literature is not free from the use of a two-stage approach.

To address this difficulty, linear algebra has often been used to identify drug-gene-disease information, since it enables the identification of drug-gene-disease information in neither a drug-centric nor disease-centric manner. Wang et al. [8] applied matrix factorization to gene expression matrices for drug and disease treatments, while Kim and Cho [9] employed tensor decomposition (TD) to extract drug-gene-disease information starting from the product of ID embedding vectors for drugs, genes, and diseases. In these approaches, because the identification of drug-gene-disease information is performed without the need for a two-stage approach, these approaches can often avoid the difficulty associated with this approach.

Another advantage of these linear-algebra-based approaches is that they are fully unsupervised. Therein, there is no need for any drug- or disease-specific information (e.g., differential expression between healthy controls and patients), and we are free to identify drug-gene-disease information in a fully data-driven manner. We recently applied a TD-based unsupervised FE [10, 11] to integrate PPI with gene expression in cancer [12]. We found that integrated analysis enhanced the coincidence with clinical labels. In this study, we also found that PPI only contained information related to cancers, even though PPI themselves are unlikely to be associated with cancers. In this paper, we attempt to determine the degree of the relationship between PPI and cancers using TD-based unsupervised FE without integrating other information (e.g. gene expression). As a result, we found that applying TD-based unsupervised FE to PPI enabled the identification of

cancer-related genes that were also used for drug repositioning. After identifying that the selected proteins were likely to be hub proteins in PPI, we replaced the proteins selected via TD-based unsupervised FE with hub proteins and found that the hub proteins could be used for drug repositioning. The distinction between hub proteins and those selected using TD-based unsupervised FE was the strength of the correlation among the interactions; that is, proteins selected using TD-based unsupervised FE were found to exhibit a greater number of shared interactions. Only proteins selected using TD-based unsupervised FE had hits in the *in vivo*-based drug database DrugMatrix. The identification of proteins that not only have more interactions (i.e., hub proteins) but also more shared interactions may be the key to identifying more promising proteins for drug repositioning.

Results

Figure 1 illustrates the analysis conducted in this study.

TD-based unsupervised FE

The first matrix to be integrated was human PPI, $n_{ii'}^1 \in \mathbb{R}^{24875 \times 24875}$ (for BioGRID) or $n_{ii'}^1 \in \mathbb{R}^{4901 \times 4901}$ (for DIP), and the second was mouse PPI, $n_{ii'}^2 \in \mathbb{R}^{16645 \times 16645}$ (for BioGRID) or $n_{ii'}^2 \in \mathbb{R}^{2342 \times 2342}$ (for DIP). Since the number of common protein, N_{common} , is 9688 (for BioGRID) or 1097 (for DIP), and the number of orthogonal protein, N_{ortho} , is 4026 (for BioGRID) or 432 (for DIP), the resulting integrated tensor is $n_{ii'k} \in \mathbb{R}^{N \times N \times 2}$ where $N = 24875 + 16645 - 9688 - 4026 = 27806$ (for BioGRID) or $N = 4901 + 2342 - 1097 - 432 = 5714$ (for DIP). After applying HOSVD to $n_{ii'k}$, we obtained Eq. (7). P_i s were attributed to i s using u_{2i} (for BioGRID and DIP) or u_{3i} (only for DIP) using Eq. (6) with replacing $u_{\ell i}$ with $u_{\ell 1i}$ and are corrected using the BH criterion [10, 11] (The reason why u_{3i} was considered only for DIP was because genes selected by

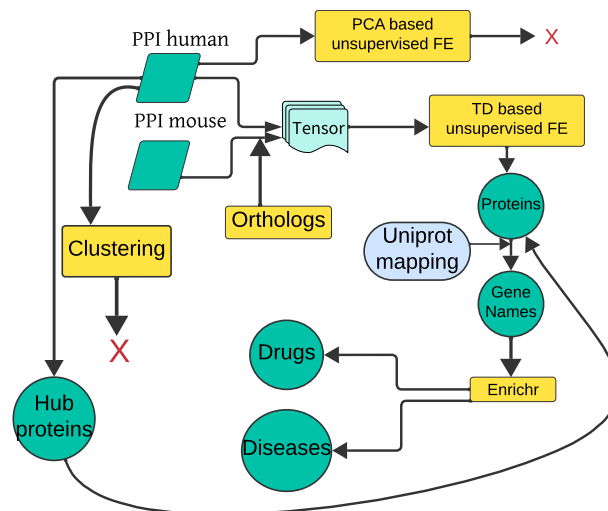


Fig. 1 Flowchart of the analyses in this study. PCA-based unsupervised FE was applied to human PPI, but failed. TD-based unsupervised FE was applied to tensor generated from human and mouse PPI. Gene names associated with the identified proteins were uploaded to Enrichr to identify associated diseases and drugs. Hub proteins and proteins selected via cluster analyses were tested and used for the identification of diseases and drugs

u_{3i} for DIP share the same enriched diseases, cancers, with those by u_{2i} , which was not the case for BioGRID. For more details, see the latter part of this paper.) Thus, 195 (using u_{2i} for BioGRID) and 196 (using u_{2i} for DIP) and 59 (using u_{3i} for DIP) proteins were associated with adjusted P -values less than 0.01. The Uniprot accession numbers associated with these proteins were converted to 217 (using u_{2i} for BioGRID), 193 (using u_{2i} for DIP), and 57 (using u_{3i} for DIP) gene names by Uniprot ID mapping (see Supplementary Information for the list of proteins and gene names). 217, 193, and 57 gene names were uploaded to Enrichr [13] for evaluation purposes.

Next, the types of diseases associated with three sets of gene names were identified. Firstly, we considered the category of “Jensen Diseases”.

For the 217 gene names selected by u_{2i} on BioGrid (Table S1 and Fig. 2), not only were there highly significant diseases, but most were cancers or tumors (“cancer,” “stomach cancer,” “adenoma,” “immune system cancer,” “ovarian cancer,” “ductal carcinoma in situ,” “esophageal carcinoma,” and “lymphoid leukemia”).

For the 193 gene names selected by u_{2i} for the DIP (see Table S2 and Fig. 3), not only were there highly significant diseases, but more than half were cancers or tumors (“lymphoid leukemia,” “immune system cancer,” “cancer,” “biliary tract cancer,” and “ovarian cancer”). Additionally, “familial adenomatous polyposis” is known to develop into cancer [14].

For the 57 gene names selected by u_{2i} for DIP (Table S3 and Fig. 4), not only were there highly significant diseases, but more than half were cancer or tumors (“ovarian cancer,” “immune system cancer,” “breast cancer,” “biliary tract cancer,” “ductal carcinoma in situ,” and “hereditary breast ovarian cancer”). In addition, “Li-Fraumeni syndrome” is known to be the cause of several cancers [15].

Next, we considered the “OMIM Diseases” category of Enrichr.

For the 217 gene names selected by u_{2i} for BioGrid (Table S4 and Fig. 5), six diseases were cancer or tumors (“breast cancer,” “ovarian cancer,” “thyroid carcinoma,” “prostate cancer,” “colorectal cancer,” “pancreatic cancer,” and “gastric cancer”), although

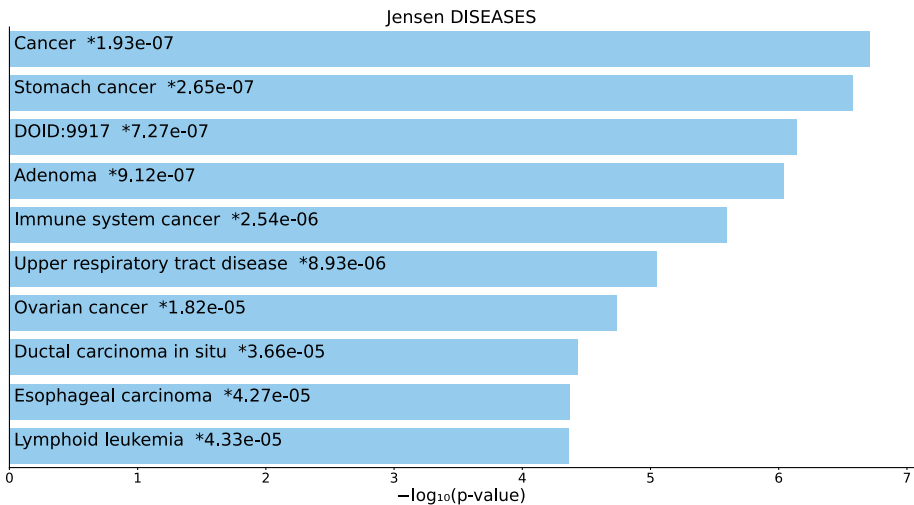


Fig. 2 Top 10 diseases in the “Jensen Diseases” category of Enrichr for 217 gene names selected by u_{2i} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

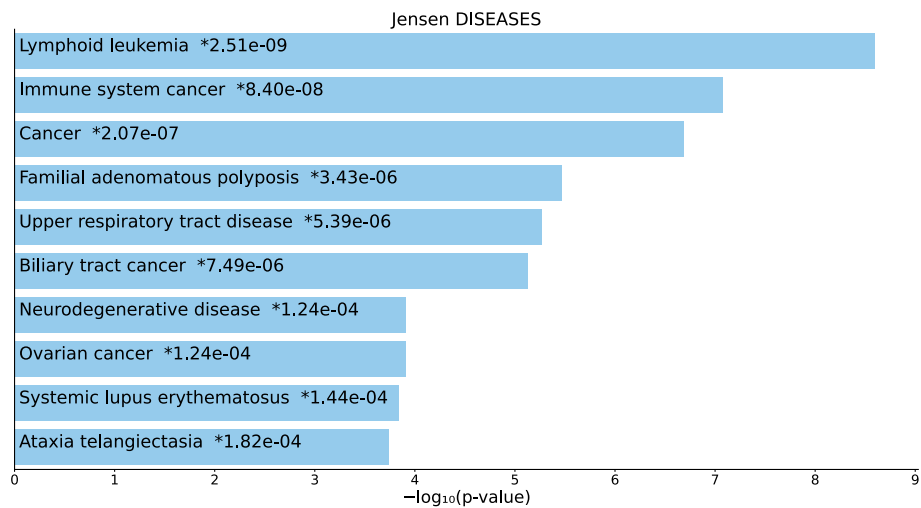


Fig. 3 Top 10 diseases in the “Jensen Diseases” category of Enrichr for 193 gene names selected by u_{2i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

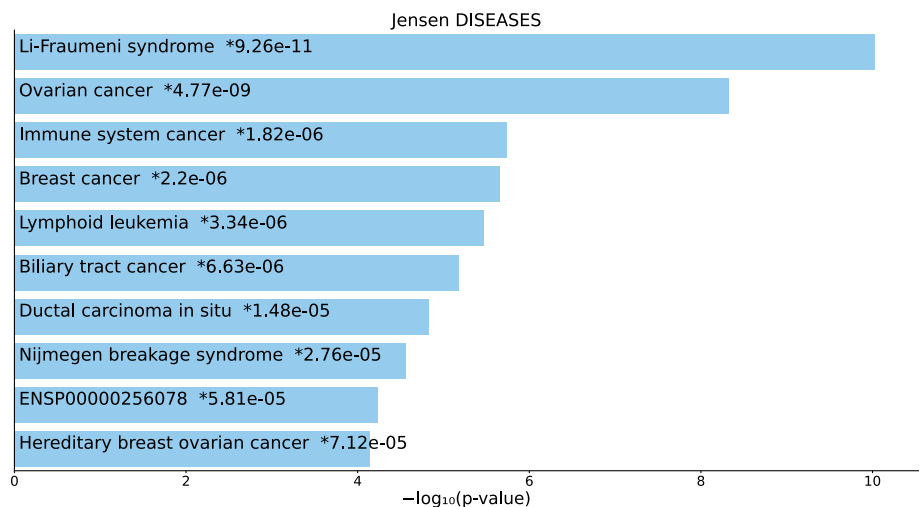


Fig. 4 Top 10 diseases in the “Jensen Diseases” category of Enrichr for 57 gene names selected by u_{3i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

statistical (i.e., associated with adjusted P -values less than 0.05) significance was only observed for the top two.

In contrast to BioGRID, which failed to identify sufficiently large significant diseases, 193 gene names selected by u_{2i} for DIP (Table S5 and Fig. 6), not only were there highly significant diseases, but four diseases were related to cancers or tumors (“colorectal cancer,” “breast cancer,” “leukemia,” and “lymphoma”). Additionally, “Fanconi anemia” is also often associated with cancer [16].

For 57 gene names selected by u_{3i} for DIP (Table S6 and Fig. 7), not only there were highly significant diseases, but also more than half of the diseases were cancers or tumors (“pancreatic cancer,” “ovarian cancer,” “colorectal cancer,” “breast cancer,” “prostate cancer,” and “melanoma”). In addition, “Fanconi anemia” was also identified.

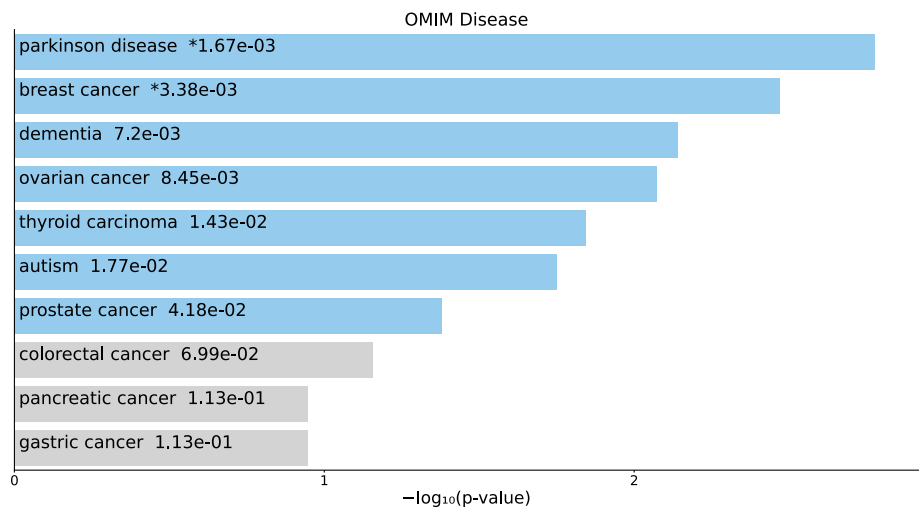


Fig. 5 Top 10 diseases in the “OMIM Diseases” category of Enrichr for 217 gene names selected by u_{2i} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

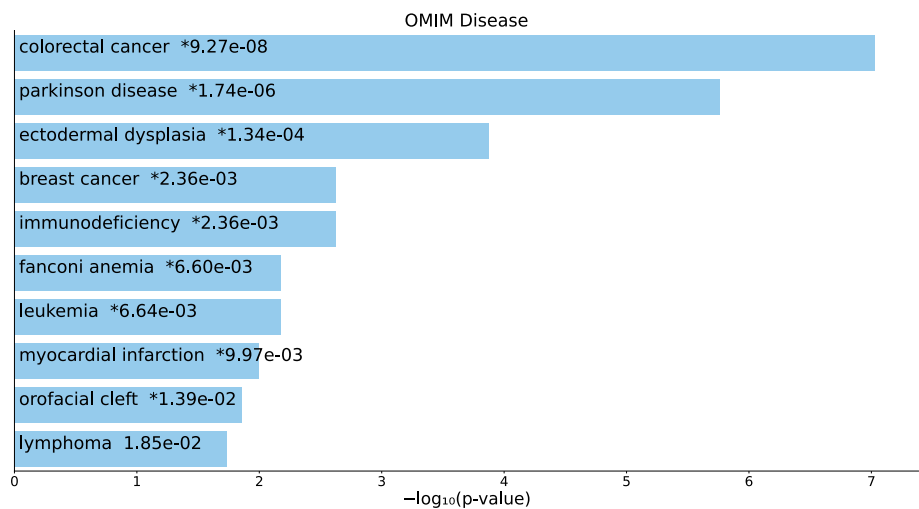


Fig. 6 Top 10 diseases in the “OMIM Diseases” category of Enrichr for 193 gene names selected by u_{2i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

In conclusion, these three sets of proteins were found to be highly related to cancers and tumors regardless of the PPI datasets used, excluding the 217 gene names selected by u_{2i} for BioGRID.

Next, we considered potential drug repositioning using these gene sets. Multiple categories can be used for drug repositioning in Enrichr: “LINCS L1000 Chem Pert Consensus Sigs,” “DSigDB,” “DrugMatrix,” “Drug Perturbations from GEO down,” and “Drug Perturbations from GEO up.” All of these categories return a list of compounds that are supposed to significantly target a set of uploaded genes. Because three sets of gene names are supposed to be deeply related to cancers and tumors, drugs that target these genes can be used to treat tumors and cancers.

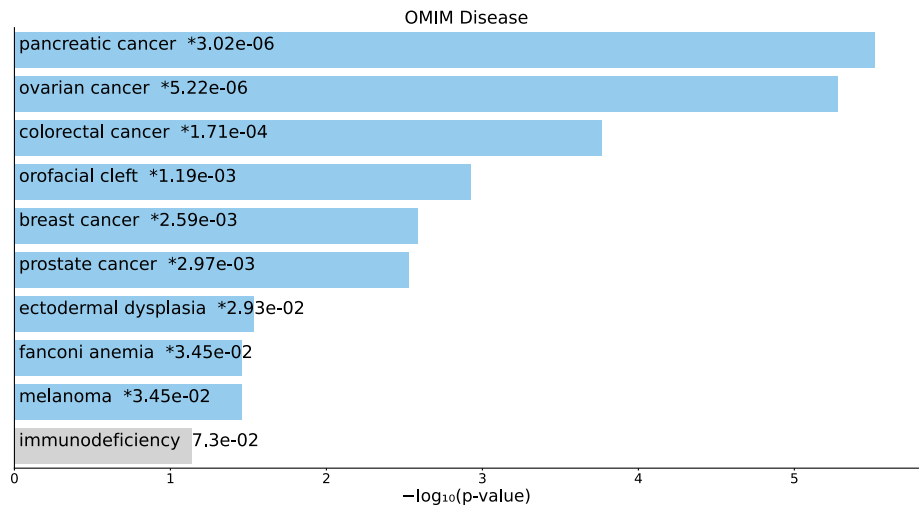


Fig. 7 Top 10 diseases in the “OMIM Diseases” category of Enrichr for 57 gene names selected by u_{3i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

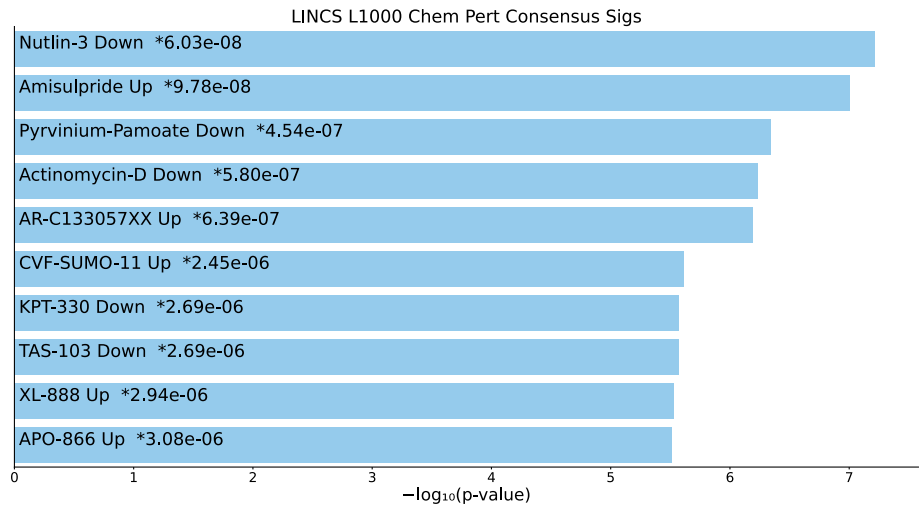


Fig. 8 Top 10 drugs in the “LINCS L1000 Chem Pert Consensus Sigs” category of Enrichr for 217 gene names selected by u_{2i} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

A list of the top 10 compounds in the “LINCS L1000 Chem Pert Consensus Sigs” category for 217 gene names selected by u_{2i} for BioGrid and 193 gene names selected by u_{2i} for DIP (no significant hit for 57 gene names selected by u_{3i} for DIP is provided in Tables S7 and S8 and Figs. 8 and 9).

Tables S9, S10, and S11 and Figs. 10, 11, and 12 list the top 10 compounds in the “DSigDB” category for 217 gene names selected by u_{2i} for BioGRID, 193 gene names selected by u_{2i} for DIP, and 57 gene names selected by u_{3i} for DIP.

Tables S12 and S13 and Figs. 13 and 14 list the top 10 compounds in the “DrugMatrix” category for 217 gene names selected by u_{2i} for BioGRID and 193 gene names selected by u_{2i} for DIP (no significant hits for 57 gene names selected by u_{3i} for DIP).

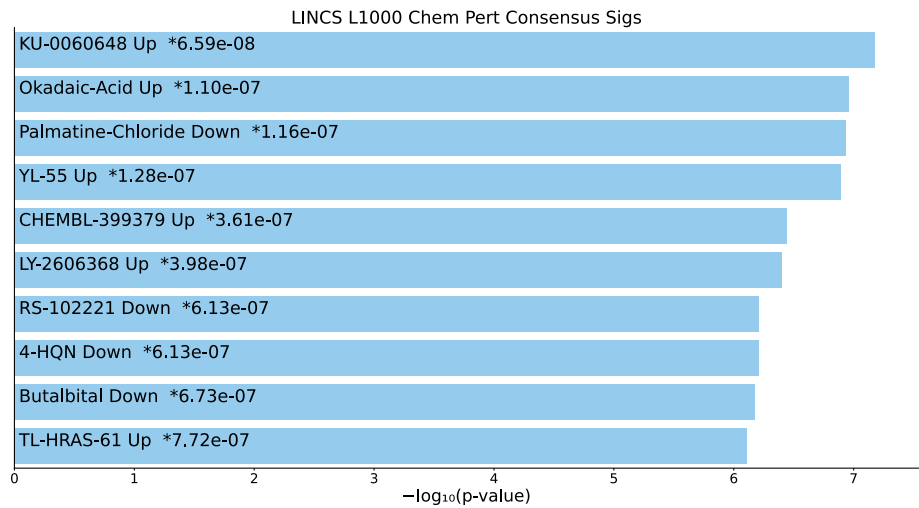


Fig. 9 Top 10 drugs in the “LINCS L1000 Chem Pert Consensus Sigs” category of Enrichr for 193 gene names selected by u_{2i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

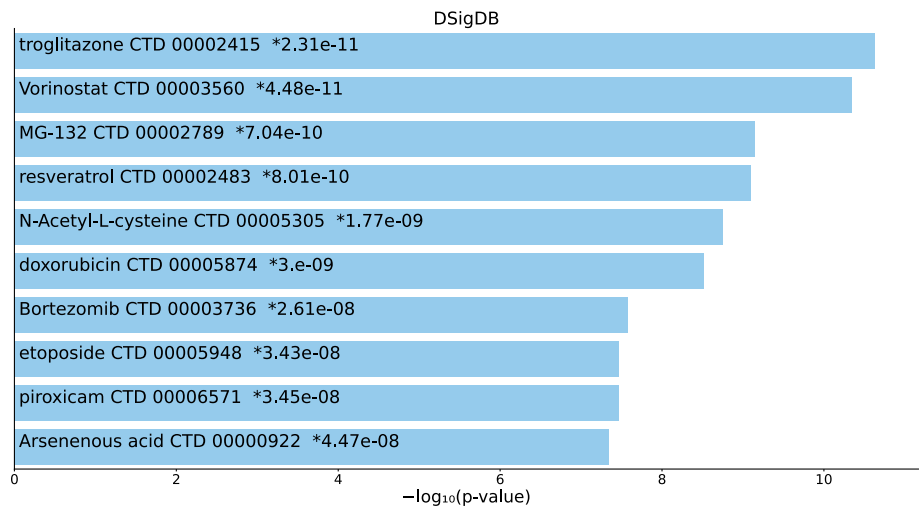


Fig. 10 Top 10 drugs in the “DSigDB” category of Enrichr for 217 gene names selected by u_{2i} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

Tables S14, S15, and S16 and Figs. 15, 16, and 17 list the top 10 compounds in the “Drug Perturbations from GEO down” category for 217 gene names selected by u_{2i} for BioGRID, 193 gene names selected by u_{2i} for DIP, and 57 gene names selected by u_{3i} for DIP.

Tables S17, S18, and S19 and Figs. 18, 19, and 20 list the top 10 compounds in the “Drug Perturbations from GEO up” category for 217 gene names selected by u_{2i} for BioGRID, 193 gene names selected by u_{2i} for DIP, and 57 gene names selected by u_{3i} for DIP.

As can be seen above, although it is not true for all categories, three sets of gene names were often associated with lists of drugs that significantly target these sets of gene names. Thus, we can conclude that our proposed gene-centric drug repositioning

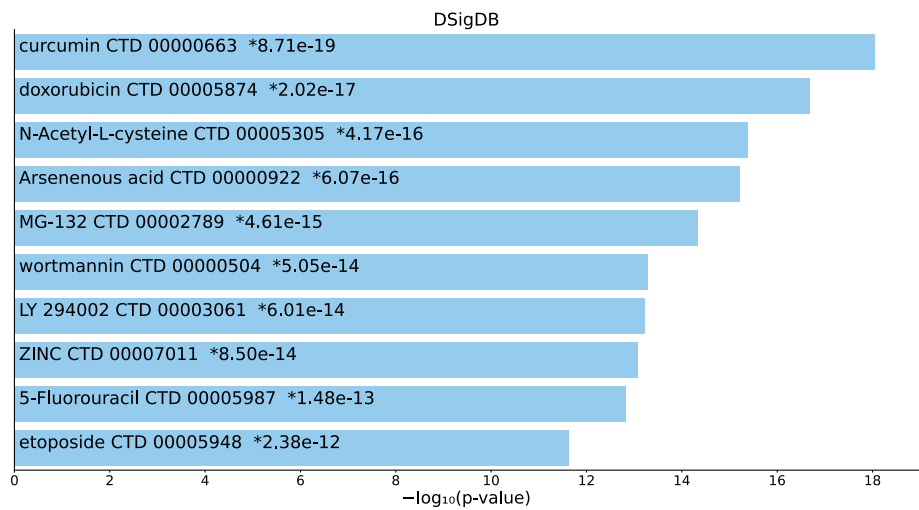


Fig. 11 Top 10 drugs in the “DSigDB” category of Enrichr for 193 gene names selected by u_{2i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

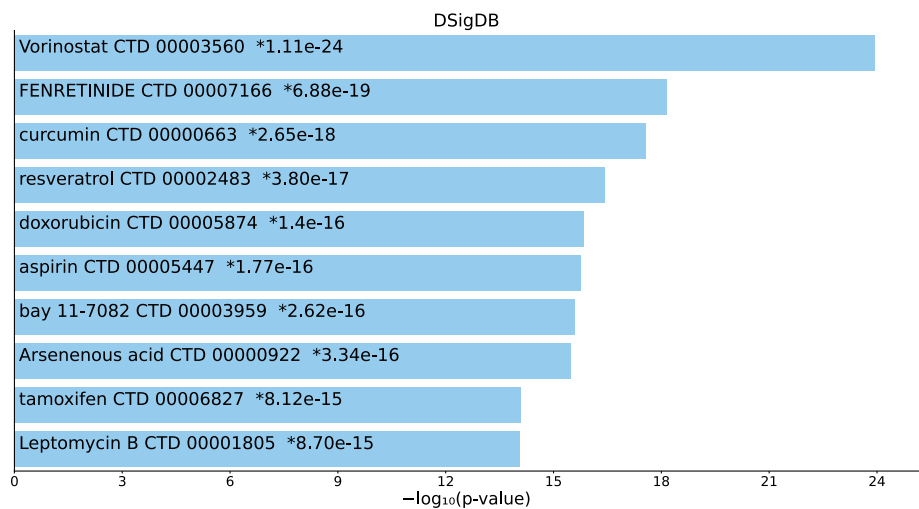


Fig. 12 Top 10 drugs in the “DSigDB” category of Enrichr for 57 gene names selected by u_{3i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

strategy can be performed employing TD-based unsupervised FE to select a set of proteins based on PPI networks, regardless of the PPI dataset employed.

PCA-based unsupervised FE

One might wonder why we needed to integrate human and mouse PPI. Simply using only human PPI might result in a similar or even better performance. To address this problem, we applied PCA-based unsupervised FE [10, 11] to human PPI and attempted to obtain a set of proteins associated with adjusted P -values less than 0.01. Nevertheless, when considering the u_{2i} for BioGRID, although 158 proteins were associated with adjusted P -values less than 0.01, the “DrugMatrix” category failed to identify enriched drug for gene names associated with these proteins. Similarly, when considering u_{2i} for

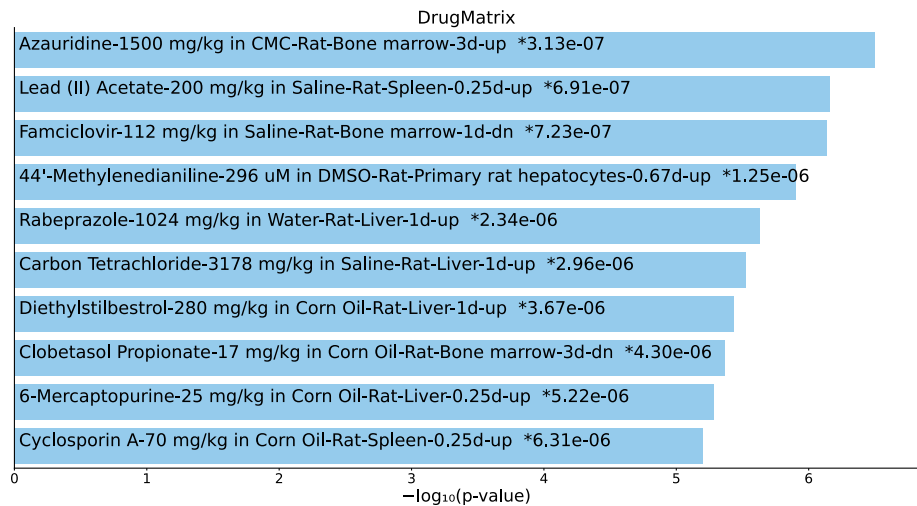


Fig. 13 Top 10 drugs in the “DrugMatrix” category of Enrichr for 217 gene names selected by u_{2i} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

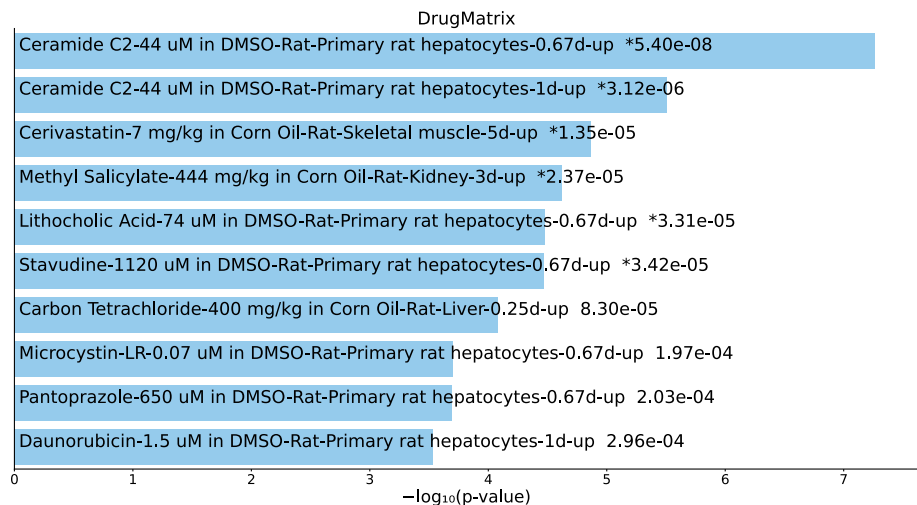


Fig. 14 Top 10 drugs in the “DrugMatrix” category of Enrichr for 193 gene names selected by u_{2i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

DIP, only one protein was found to be associated with adjusted P -values of less than 0.01. This results suggest that TD-based unsupervised FE is superior to PCA-based unsupervised FE and thereby worth testing.

Cluster analyses

However, another concern is whether TD-based unsupervised FE is required when selecting a set of proteins since several methods exist with which to select sub-clusters within large networks. Such methods can replace TD-based unsupervised FE in the selection of a set of proteins. Following Zhao et al [17], we employed three methods, multi-level [18], label-progration [19], and edge-betweenness [20], which enabled us to identify sub-clusters. When considering PPI human for BioGrid, edge-betweenness

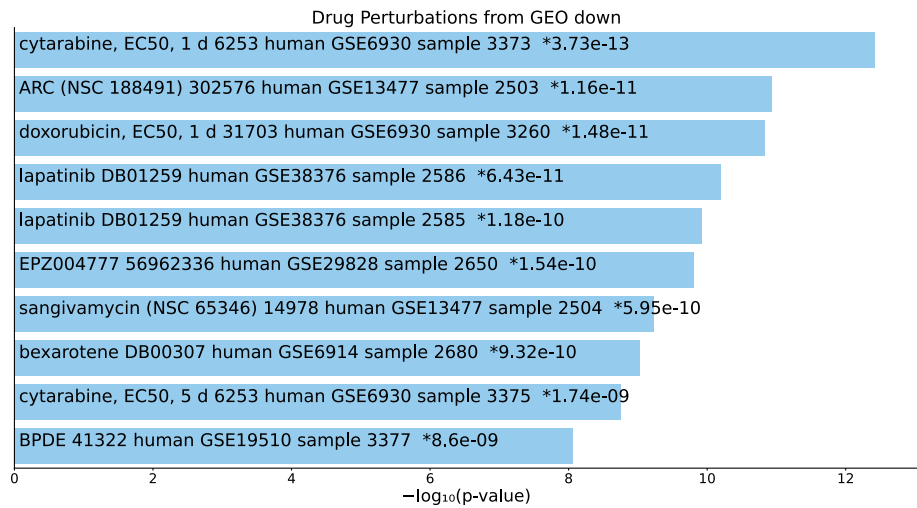


Fig. 15 Top 10 drugs in the “Drug Perturbations from GEO down” category of Enrichr for 217 gene names selected by u_{2j} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

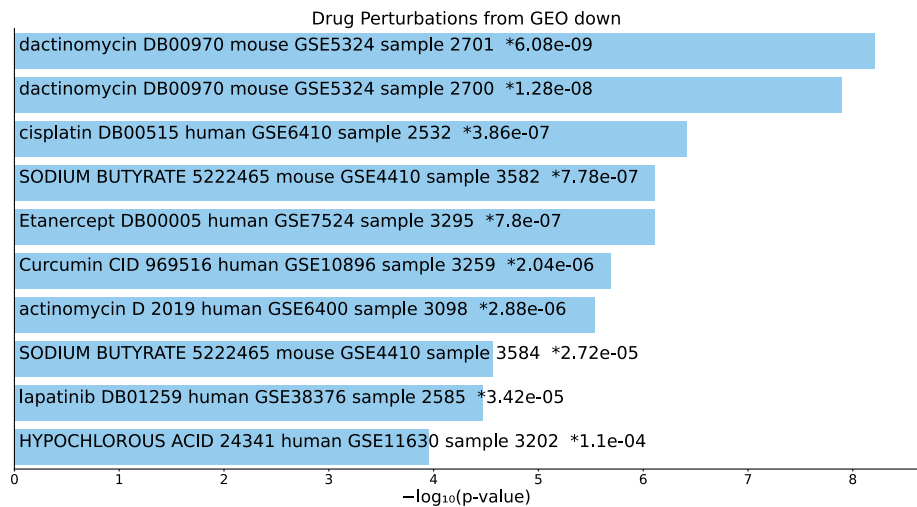


Fig. 16 Top 10 drugs in the “Drug Perturbations from GEO down” category of Enrichr for 193 gene names selected by u_{2j} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

failed to converge after waiting more than eight hours (i.e longer than the ten minutes required for TD-based unsupervised FE, which includes the time consuming process of generating tensor, $n_{ii'k} \in \mathbb{R}^{N \times N \times 2}$), and the first and the second clusters generated by multi-level were too large (6406 and 2776). The first and the second clusters generated by label-progration were 24501 and 1, which was not reasonable at all. However, when human PPI for DIP was considered, the results improved. Table 1 lists the performances for the identification of diseases and drugs significantly associated with selected proteins in individual identified sub-clusters (label-progration could not identify large enough clusters to be used for enrichment analyses). It is clear that the tensor method outperforms the clustering methods, even for DIP. In conclusion, in contrast to conventional cluster analysis in which cluster sizes varied heavily

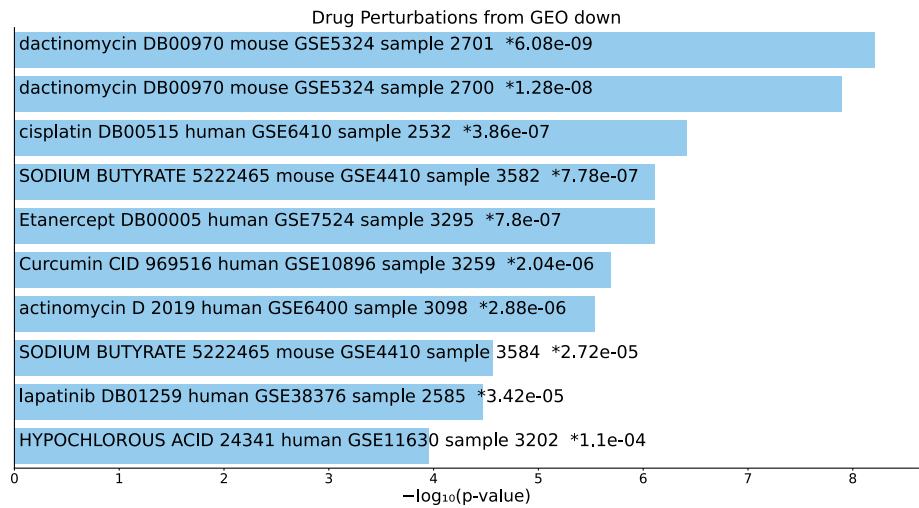


Fig. 17 Top 10 drugs in the “Drug Perturbations from GEO down” category of Enrichr for 57 gene names selected by u_{3j} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

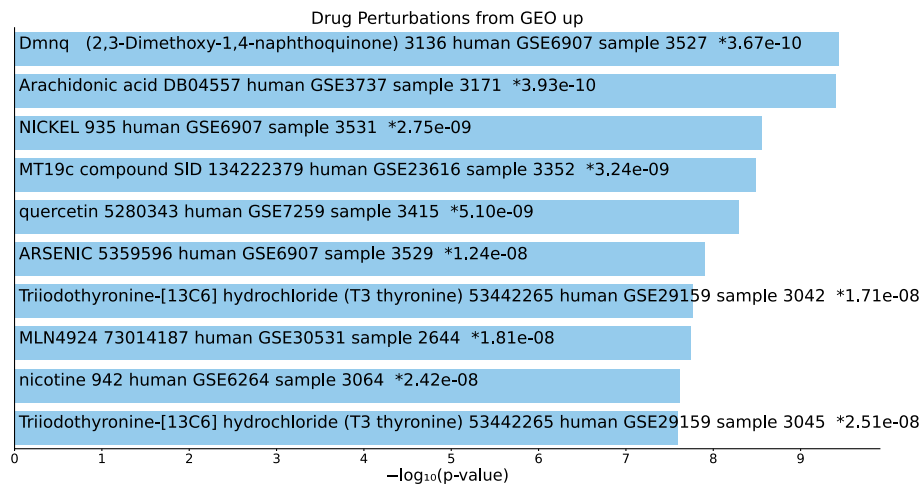


Fig. 18 Top 10 drugs in the “Drug Perturbations from GEO up” category of Enrichr for 217 gene names selected by u_{2j} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

depending on the PPI dataset used, using TD-based unsupervised FE can provide a stable and reasonable size for a set of genes. Thus, TD-based unsupervised FE appears to be superior to conventional cluster analysis.

Other species

Although we only used human and mouse PPI, it is possible to employ other combinations because BioGRID and DIP include more species-specific PPI. To validate this, we initially examined a combination of human and rat PPI for BioGRID. However, this outcome was not very promising. Next, 271 proteins associated with adjusted P -values less than 0.01 were selected, and the gene names associated with these proteins were uploaded. The “Jensen Diseases” and “OMIM Diseases” categories identified two and zero diseases enriched with the uploaded gene names, respectively. Because the

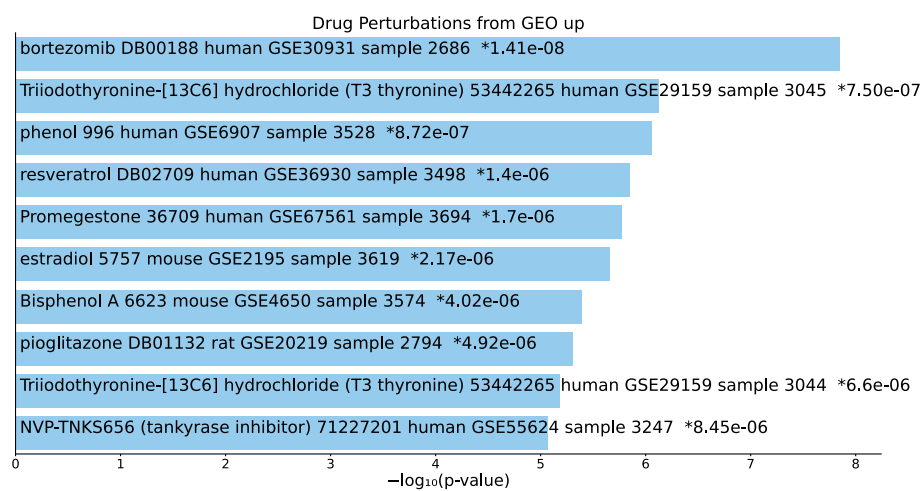


Fig. 19 Top 10 drugs in the “Drug Perturbations from GEO up” category of Enrichr for 193 gene names selected by u_{2j} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

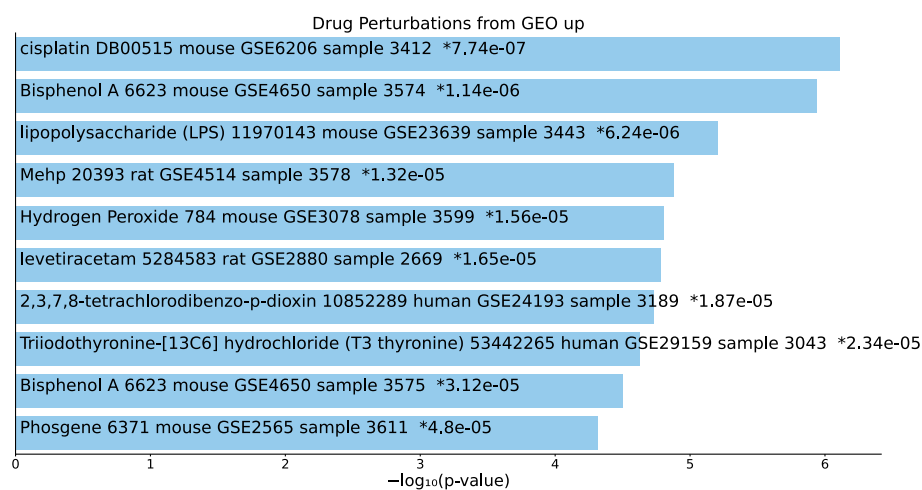


Fig. 20 Top 10 drugs in the “Drug Perturbations from GEO up” category of Enrichr for 57 gene names selected by u_{3j} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

Table 1 Performances of various protein selection methods for DIP

Methods	Tensor		multi-level		edge-betweenness	
	2nd	3rd	1st	2nd	1st	2nd
# of proteins	195	55	144	138	39	82
# of genes	193	57	128	164	30	111
Jensen Diseases	○	○	○	○	7	5
OMIM Disease	9	9	8	0	0	0
LINCS L1000 Chem Pert Consensus Sigs	○	0	4	1	0	○
DSigDB	○	○	○	○	0	○
DrugMatrix	6	0	0	0	0	0
Drug Perturbations from GEO down	○	○	0	0	0	2
Drug Perturbations from GEO up	○	○	6	0	1	0

○s denotes cases with more than or equal to 10 drugs associated with adjusted P -values less than 0.05. Otherwise, the corresponding practical numbers are presented

identification of enriched diseases is the starting point, if the process fails at this point, there are no way to proceed. The reason for this is the inadequateness of rat PPI, which includes too few PPI; integrating human and rat PPI results in the wrong conclusion, namely that missing PPI in rats indicates a lack of interactions, although this may simply mean that PPI in rats has not been thoroughly investigated. Thus, combinations other than that between human and mouse PPI will not work until a more comprehensive PPI for species close to humans can be obtained.

In conclusion, PPI can be a useful source of information for drug repositioning when used in combination with TD.

Discussion

To understand the set of genes selected in this analysis, we computed the mean number of bindings of individual proteins, defined as

$$\langle n_{ii'1} \rangle_i = \frac{1}{N} \sum_{i'=1}^N n_{ii'1} \tag{1}$$

and compared $\langle n_{ii'1} \rangle_i$ between the selected and non-selected proteins for the human PPI.

As can be seen in Table 2, $\langle n_{ii'1} \rangle_i$ for 217 proteins selected by u_{2i} for BioGRID, 195 proteins selected by u_{2i} for DIP, and 55 proteins selected by u_{3i} for DIP were always significantly larger than those of other (i.e. not selected) proteins. These results suggest that TD-based unsupervised FE selects proteins with more interactions in a fully data-driven and unsupervised manner.

Since proteins selected via TD-based unsupervised FE have more interactions, we proceeded to verify whether proteins were being selected for use in drug repositioning simply due to their greater number of interactions (hereafter, denoted as hub proteins). To verify this, we selected the top 200 and 50 proteins with more interactions (these numbers were selected to be close to the number of proteins selected using u_{2i} and u_{3i} . See Supplementary Information for a list of proteins and gene names) using human PPI.

Table 3 lists the confusion matrices for the selected proteins. Although they significantly overlapped, because the majority were not shared, distinct sets of proteins were identified.

Table 2 To compare $\langle n_{ii'1} \rangle_i$ between the selected and not selected proteins when using various $u_{\ell_1 i}$ to select proteins, P -values were computed using the t test based on the alternative hypothesis that the mean $\langle n_{ii'1} \rangle_i$ of the selected proteins was larger than that of not selected proteins

$u_{\ell_1 i}$ s used	u_{2i}		u_{3i}		u_{6i}	
	Selected	Not selected	Selected	Not selected	Selected	Not selected
BioGRID						
Mean $\langle n_{ii'1} \rangle_i$	3.63×10^{-2}	1.84×10^{-3}	1.97×10^{-2}	1.78×10^{-3}	—	—
P -value	5.32×10^{-37}		2.53×10^{-59}		—	
DIP						
Mean $\langle n_{ii'1} \rangle_i$	1.12×10^{-3}	4.17×10^{-4}	3.12×10^{-3}	4.13×10^{-4}	2.14×10^{-3}	4.26×10^{-4}
P -value	2.4×10^{-4}		5.3×10^{-5}		2.04×10^{-13}	

Table 3 Confusion matrix of selected proteins between TD-based unsupervised FE and hub proteins

TD based unsupervised FE		Hub proteins	
		Top 200	
		Not selected	Selected
BioGRID			
Using u_{2i}	Not selected	27526	80
	Selected	85	115
	Odds ratio	465.58	
	<i>P</i> -value	2.88×10^{-209}	
DIP			
Using u_{2i}	Not selected	5353	168
	Selected	165	28
	Odds ratio	5.4	
	<i>P</i> -value	4.55×10^{-11}	
		Top 50	
		Not selected	Selected
Using u_{3i}	Not selected	5625	42
	Selected	30	17
	Odds ratio	75.4	
	<i>P</i> -value	3.04×10^{-23}	

Table 4 Performances achieved by hub proteins

Methods	Tesor					Hub proteins		
	BioGRID		DIP			BioGRID	DIP	
Order	2nd	3rd	2nd	3rd	6th	Top 200	Top 200	Top 50
# of proteins	195	475	195	55	38	199	179	45
# of genes	217	503	193	57	41	229	190	59
Jensen Diseases	○	6	○	○	○	○	○	○
OMIM Disease	2	3	9	9	○	4	○	6
LINCS L1000 Chem Pert Consensus Sigs	○	○	○	0	6	○	○	0
DSigDB	○	○	○	○	○	○	○	○
DrugMatrix	○	○	6	0	0	○	0	1
Drug Perturbations from GEO down	○	○	○	○	○	○	○	○
Drug Perturbations from GEO up	○	○	○	○	○	○	○	○

○s denote case with more than or equal to 10 drugs/associated with adjusted *P*-values less than 0.05. Otherwise, the corresponding practical numbers are presented

We uploaded the associated gene names to Enrichr to evaluate the top 200 and 50 hub proteins.

Table 4 lists the performance achieved by the hub proteins. (for details on the top diseases and drugs in individual categories: see Tables S20 to S26 for the top 200 hub proteins for BioGRID, Tables S29 to S33 for the top 200 hub proteins for DIP, and Tables S34 to S40 for the top 50 hub proteins for DIP). Their performances were essentially the same as those of TD-based unsupervised FE, excluding DrugMatrix for DIP.

Comparing the performance of TD-based unsupervised FE with hub proteins, to the best of our knowledge, there are no other studies that use hub proteins derived from only PPI for drug repositioning without using other information. This is because hub proteins derived from PPI themselves, without any other information, cannot be used for drug repositioning for cancer. This was clear from the fact that proteins selected by TD-based unsupervised FE were hub proteins. Without TD-based unsupervised FE, we could not identify hub proteins derived only from PPI that could be used for drug repositioning for cancer. Secondly, the hub proteins for DIP failed to identify hits in DrugMatrix. DrugMatrix is an *in vivo* specific database. Thus, the hub proteins for DIP do not have the ability to identify drugs that may be useful in *in vivo* experiments. Thus, even if hub proteins can achieve a performance similar to that of the *in vitro* experiments, TD-based unsupervised FE is still useful, at least for DIP. Third, using hub proteins, we cannot target diseases other than cancers because hub proteins are enriched mainly in cancers. For the top 200 hub proteins for BioGRID, seven cancers were identified (“stomach cancer,” “adenoma,” “cancer,” “immune system cancer,” “esophageal carcinoma,” “intestinal benign neoplasm,” and “lymphoid leukemia”) within the top ranked diseases in the “Jensen Diseases” category and as many as five cancers (“ovarian cancer,” “breast cancer,” “colorectal cancer,” “thyroid carcinoma,” and “prostate cancer”) within the top ranked diseases in the “OMIM Diseases” category (Table S20 and S21). For the top 200 hub proteins for DIP, as many as six cancers were identified (“cancer,” “lymphoid leukemia,” “intestinal benign neoplasm,” “immune system cancer,” “biliary tract cancer,” and “stomach cancer”) within the top ranked diseases in the “Jensen Diseases” category and as many as seven (“colorectal cancer,” “ovarian cancer,” “breast cancer,” “leukemia,” “lung cancer,” “melanoma,” and “pancreatic cancer”) within the top ranked diseases in the “OMIM Diseases” category (Table S27 and S28). For the top 50 hub proteins for DIP, as many as five cancers (“cancer,” “DOID:9917” (pleural cancer), “thyroid cancer,” “esophageal carcinoma,” and “intestinal benign neoplasms”) were identified within the top ranked diseases in the “Jensen Diseases” category and as many as five (“colorectal cancer,” “ovarian cancer,” “breast cancer,” “lung cancer,” and “pancreatic cancer”) within the top ranked diseases in the “OMIM Diseases” category (Table S34 and S35). Nevertheless, because TD-based unsupervised FE has the freedom to select $u_{\ell_1 i}$, it can provide a set of proteins enriched in diseases other than cancer. Tables S41 and S42, and Figs. 21 and 22 list the diseases enriched in proteins selected by u_{3i} for BioGRID, while Tables S43 and S44, and Figs. 23 and 24 list the diseases enriched in proteins selected by u_{6i} for DIP (see Supplementary Information for the list of proteins and gene names); the diseases are mainly distinct from cancers since only one cancer (“stomach cancer”) for BioGRID and only one cancer (“cancer”) for DIP within the top 10 diseases in the “Jensen Diseases” category (Tables S41 and S43 and Figs. 21 and 23) and three cancers (“thyroid carcinoma,” “gastric cancer,” and “ovarian cancer”) for BioGRID and three cancers (“pancreatic cancer,” “melanoma,” and “lung cancer”) for DIP within the top 10 diseases in the “OMIM Diseases” category (Tables S42 and S44 and Figs. 22 and 24) were identified. Also for these diseases other than cancers, sufficient drugs exist (for detailed drug names, see Table S22–S31). Thus, the results of TD-based

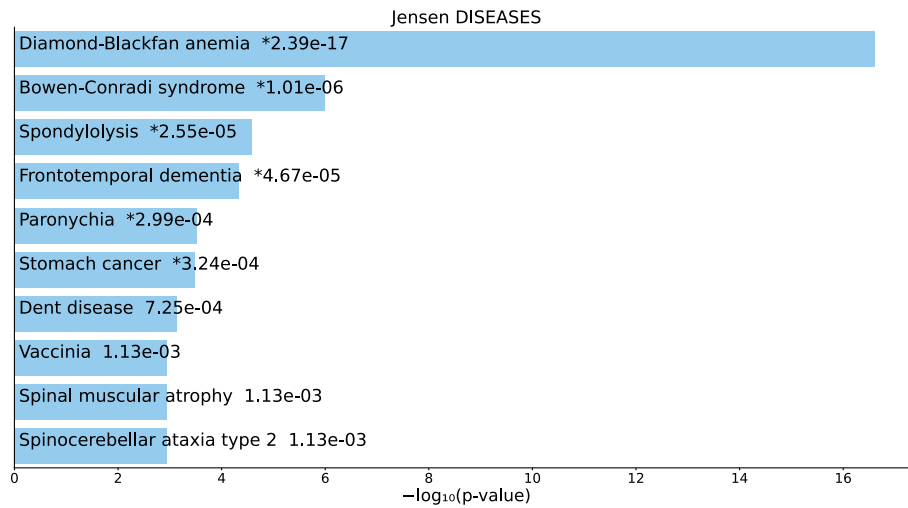


Fig. 21 Top 10 diseases in the “Jensen Diseases” category of Enrichr for 502 gene names selected by u_{3i} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

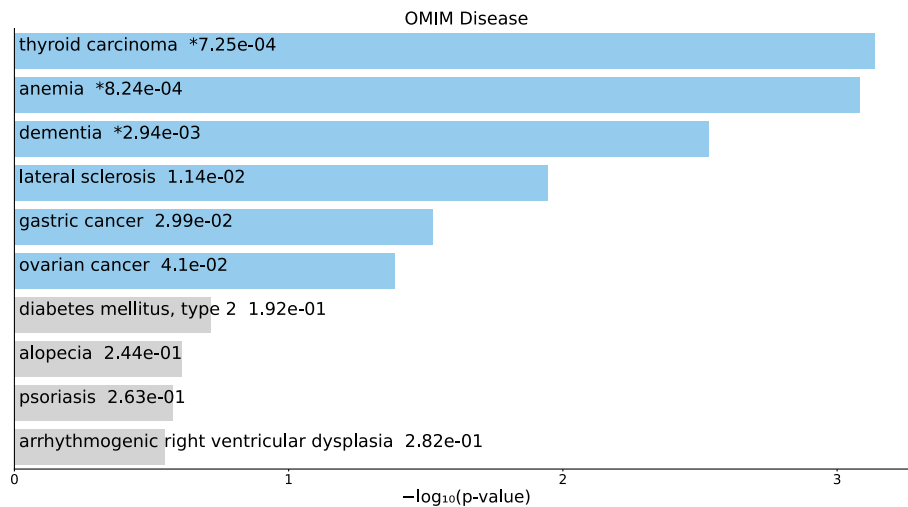


Fig. 22 Top 10 diseases in the “OMIM Diseases” category of Enrichr for 502 gene names selected by u_{3i} for BioGrid. Blue: $P < 0.05$, *: adjusted $P < 0.05$

unsupervised FE are useful to find drugs effective toward the diseases other than cancers, even if hub proteins can be used to identify drugs for cancers.

It is worth considering the differences between hub proteins and the proteins selected by TD-based unsupervised FE if they are not identical. To determine this, the mean correlation coefficients between proteins

$$\langle r \rangle = \frac{2}{N(N-1)} \sum_{i'' \neq i'} r(n_{ii'1}, n_{ii''1}) \quad (2)$$

where $r(n_{ii'1}, n_{ii''1})$ is Pearson's correlation coefficient between $\{n_{ii'1} | 1 \leq i \leq N\}$ and $\{n_{ii''1} | 1 \leq i \leq N\}$, were computed with human PPI.

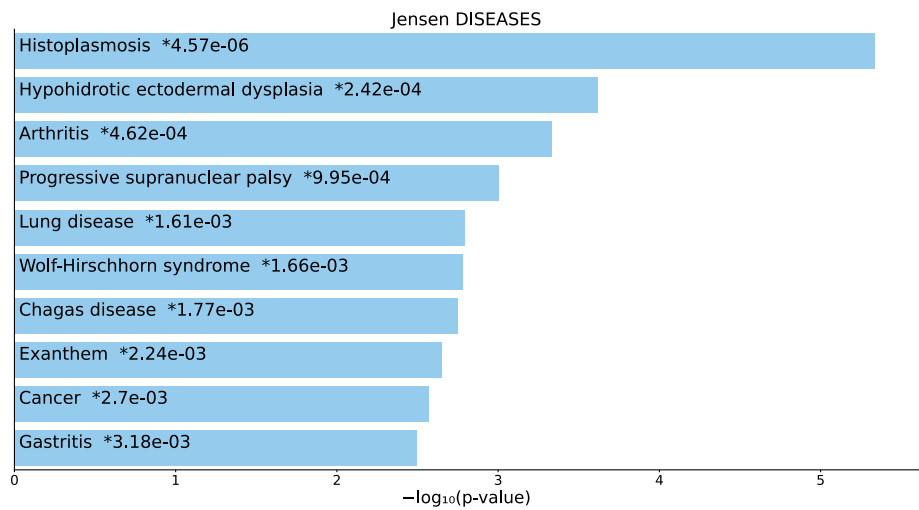


Fig. 23 Top 10 diseases in the “Jensen Diseases” category of Enrichr for 41 gene names selected by u_{6i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

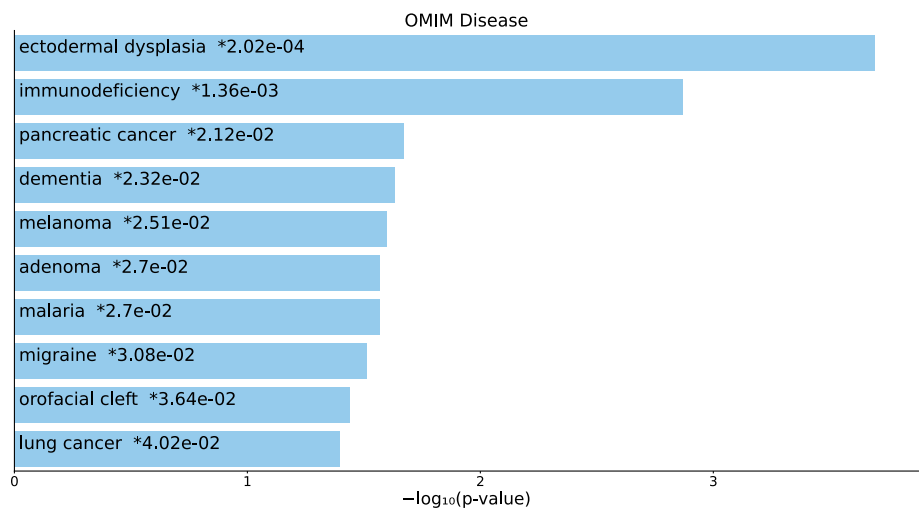


Fig. 24 Top 10 diseases in the “OMIM Diseases” category of Enrichr for 41 gene names selected by u_{6i} for DIP. Blue: $P < 0.05$, *: adjusted $P < 0.05$

Table 5 $\langle r \rangle$ s and P -values computed using a one way t test

	$\langle r \rangle$	P -values
BioGRID		
Proteins selected using u_{2i}	1.97×10^{-1}	0
Proteins selected using u_{3i}	1.60×10^{-1}	0
Top 200 proteins	2.02×10^{-1}	0
DIP		
Proteins selected using u_{2i}	3.29×10^{-1}	0
Proteins selected using u_{3i}	9.64×10^{-2}	7.70×10^{-211}
Proteins selected using u_{6i}	8.14×10^{-2}	5.84×10^{-70}
Top 200 proteins	7.49×10^{-3}	3.97×10^{-128}
Top 50 proteins	1.19×10^{-3}	6.84×10^{-18}

Table 5 lists the values of $\langle r \rangle$. Because proteins selected by TD-based unsupervised FE for DIP had a substantially larger $\langle r \rangle$, the distinctions between the proteins selected by TD-based unsupervised FE and the hub proteins were that between correlations. The fact that a larger $\langle r \rangle$ means that proteins share the proteins with which they interact suggests that proteins that share interactions may be the key to identifying effective drugs in *in vivo* experiments. In reality, the top 200 hub proteins for BioGRID have a larger $\langle r \rangle$ and many enrichments in the “DrugMatrix” category (Table S12 and Fig. 13). Thus, we can guarantee a high correlation for hub proteins without consulting TD-based unsupervised FE, which allows us to identify genes that can be used to identify effective drugs, even in *in vivo* experiments.

Because we employed a gene-centric strategy, the drugs and diseases identified were associated with common genes. Nevertheless, one might wonder whether this always guarantees the effectiveness of the identified drugs against the diseases. The so-called docking simulation is not a good idea for validation because Enrichr is used to relate genes to drugs or diseases in order to identify relationships based on gene expression. Thus, this does not always mean that the selected genes are direct targets of drugs, but simply those whose expression is altered by drug treatment, since the genes are in the downstream pathway. However, to the best of our knowledge, search engines specifically for studies on drug-disease pairs within the proposed list of drugs and diseases do not exist. Alternatively, to achieve this, we employed a large language model (LLM). We sought to identify pairs of drugs and diseases for which research has been reported. To avoid the wrong relationships being reported by LLM due to hallucinations, we verified whether other papers existed in support of the relationship reported by LLM.

Table 6 lists the reported and validated combination of drugs and diseases between the “Jensen Diseases” category by u_{3i} for BioGRID (Table S41 and Fig. 21) and the “DSigDB” category or “DrugMatrix” category by u_{3i} for BioGRID (Table S23 or S24). The reason why we employ this is simply because it is difficult to validate the correspondence between diseases and drugs if diseases are mostly composed of cancers because many cancers share effective drugs. Therefore, it is preferable to use a set of diseases other than cancer.

It is evident that the combinations of drugs and diseases reported by TD-based unsupervised FE are associated with previously published studies, despite the diverse types of

Table 6 Diseases and drugs whose relation is reported by LLM. Their corresponding ranks in categories appear in parentheses

DSigDB	Jensen disease	References
Verteporfin (1st)	Stomach Cancer (6th)	[21]
Clindamycin (2nd)	Paronychia (5th)	[22]
Captopril (3rd)	Frontotemporal Dementia (4th)	[23]
Glibenclamide (7th)	Spinal Muscular Atrophy (9th)	[24]
Puromycin (8th)	Spinocerebellar Ataxia Type 2 (10th)	[25]
DrugMatrix		
Mitomycin C (8th)	Stomach Cancer (6th)	[26]
Pravastatin (7th)	Frontotemporal Dementia (4th)	[27]
Miconazole (2nd)	Paronychia (6th)	[28]
NN-Dimethylformamide (1st)	Various Cancers (6th)	[29]

diseases. Thus, we concluded that the pairs of diseases and drugs identified in this study are likely to include promising pairs.

The limitation of our present methods is that we cannot always get enough number of PPI information for most species. As can be seen above, when we considered rat, we could not get any good results because of the lack of enough PPI information. We expect that more advanced methods for PPI identification will be developed.

Methods

Protein-protein interactions (PPI)

BioGRID

The PPI data were downloaded from BioGRID [30]. PPI classified by species, BIOGRID-ORGANISM-4.4.236.tab3.zip, were downloaded. Three species-specific files, BIOGRID-ORGANISM-Homo_sapiens-4.4.236.tab3.txt, BIOGRID-ORGANISM-Mus_musculus-4.4.236.tab3.txt, and BIOGRID-ORGANISM-Rattus_norvegicus-4.4.236.tab3.txt were extracted for analysis.

DIP

The PPI data were also downloaded from DIP [31]. The datasets used were species-specific sets for *Homo Sapiens*, Hsapi20170205, and *Mus musculus*, Mmusc20170205. “full [gz]” were downloaded as tab-limited files.

Tensor format

The PPI files were further loaded into R using the `read_csv` command as $n_{ii'}^k \in \mathbb{R}^{N_k \times N_k}$ ($k = 1$: human, $k = 2$: mouse, and $k = 3$: rat (only for BioGRID)), which takes 1 when the i th and i' th proteins interact with each other; otherwise, 0. As these matrices were sparse, they were stored in a sparse matrix format using the Matrix [32] package.

Principal component analysis (PCA)-based unsupervised feature extraction (FE)

Singular value decomposition was applied to the `irlba` function in the `irlba` package [33] for human PPI $n_{ii'}^1$. As a result, we obtained

$$n_{ii'}^1 = \sum_{\ell=1}^L \lambda_{\ell} u_{\ell i} u_{\ell i'} \quad (3)$$

where λ_{ℓ} is a singular value, and $u_{\ell i} \in \mathbb{R}^{N \times N}$ is the singular value vector and orthogonal vector.

P -values are attributed to i th protein as

$$\langle u_{\ell i} \rangle = \frac{1}{N} \sum_{i=1}^N u_{\ell i} \quad (4)$$

$$\sigma_{\ell} = \sqrt{\frac{1}{N} \sum_{i=1}^N (u_{\ell i} - \langle u_{\ell i} \rangle)^2} \quad (5)$$

$$P_i = P_{\chi^2} \left[> \left(\frac{u_{\ell i} - \langle u_{\ell i} \rangle}{\sigma_{\ell}} \right)^2 \right] \quad (6)$$

where $P_{\chi^2}[> x]$ is the cumulative χ^2 distribution, the argument is larger than x and $\langle u_{\ell i} \rangle$, and σ_{ℓ} are the mean and standard deviation, respectively. Thus, we assumed the following null hypothesis: $u_{\ell i}$ follows a Gaussian distribution. P -values were corrected using the BH criterion, and proteins associated with adjusted P -values less than 0.01 were selected.

Tensor decomposition (TD)-based unsupervised FE

To apply TD-based unsupervised FE to PPI, an integrated tensor that stores multiple PPIs must be constructed. Suppose we have two PPI matrices, $n_{ii'}^1 \in \mathbb{R}^{N_1 \times N_1}$ and $n_{ii'}^2 \in \mathbb{R}^{N_2 \times N_2}$. It is also assumed that there are N_{common} common proteins and N_{ortho} orthogonal proteins between the two datasets. Then, we merge these datasets into one tensor, $n_{ii'2} \in \mathbb{R}^{N \times N \times 2}$ where $N = N_1 + N_2 - N_{\text{common}} - N_{\text{ortho}}$ as shown in (Fig. 25).

$n_{ii'}^k$, ($1 \leq i, i' \leq N_{\text{common}} + N_{\text{ortho}}$) are placed in $n_{ii'k}$ as is (yellow and red regions in Fig. 25). $n_{ii'}^1$, ($N_{\text{common}} + N_{\text{ortho}} < i, i' \leq N_1$) are placed in $n_{ii'1}$ (the blue region in Fig. 25). $n_{ii'}^2$, ($N_{\text{common}} + N_{\text{ortho}} < i, i' \leq N_2$) are placed in $n_{ii'2}$ in a fragmented manner (green region in Fig. 25). As a result, the blue region of $n_{ii'2}$ and green region of $n_{ii'1}$ are blank. The shaded region is blank in $n_{ii'k}$.

HOSVD was applied to $n_{ii'k}$, resulting in

$$n_{ii'k} = \sum_{\ell_1=1}^N \sum_{\ell_2=1}^N \sum_{\ell_3=1}^2 G(\ell_1 \ell_2 \ell_3) u_{\ell_1 i} u_{\ell_2 i'} u_{\ell_3 k} \quad (7)$$

where G is a core tensor representing the weight (contribution) of $u_{\ell_1 i} u_{\ell_2 i'} u_{\ell_3 k}$ to $n_{ii'2}$ and $u_{\ell_1 i} = u_{\ell_2 i'} \in \mathbb{R}^{N \times N}$ and $u_{\ell_3 k} \in \mathbb{R}^{2 \times 2}$ are singular value and orthogonal matrices, respectively. The attribution of P_i s to i s and the selection of proteins were performed by replacing $u_{\ell i}$ where $u_{\ell i}$ in eq. (6).

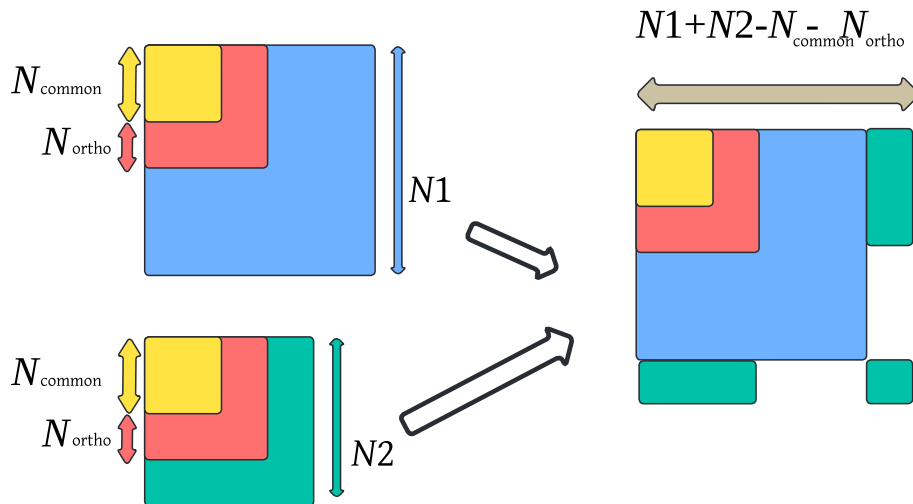


Fig. 25 How to merge two PPI matrices into one tensor

HOSVD was performed using the `+irlba+` function by applying SVD to unfolded matrices because the usual R function that can perform HOSVD does not accept a sparse matrix as input.

Identification of orthogonal proteins

Orthologs between human and mouse were retrieved by querying “((organism_id:10090) OR ((organism_id:9606) AND (reviewed:true))) AND (database:orthodb)” in uniprot search (<https://www.uniprot.org/>). Human Uniprot accessions were converted to gene name as “XXX_HUMAN” using the above retrieved information. Then, after seeking “XXX_MOUSE” in the above retrieved information (using the first hit if multiple hits are found), “XXX_MOUSE” was converted into the corresponding mouse Uniprot accession numbers in the information retrieved above. This results in the corresponding table of Uniprot ortholog (the above retrieved information is provided as Supplementary Material).

Gene ID conversion

The obtained Uniprot accession numbers attributed to proteins were converted to gene names using Uniprot ID mapping (<https://www.uniprot.org/id-mapping>).

Relating genes to drugs and diseases through enrichment analysis

Enrichr [13] enables the validation of overlaps between two sets of genes using statistical tests that measure the probability that the observed overlap occurs by chance (the so-called *P*-value). If *P*-values are sufficiently small, the two sets of genes are significantly related. In Enrichr, one set of genes was provided by the researchers, and its overlap with the prepared sets of genes was evaluated. If the uploaded set of genes significantly overlaps with one of the prepared sets of genes, the uploaded set of genes can be said to be associated with the properties of the overlapping set of genes. In this analysis, we considered two sets of genes, diseases and drugs. For the disease gene sets, the uploaded set of genes was evaluated if it overlapped with a set of genes known to be related to diseases. The “Jensen Diseases” and “OMIM disease” are independent categories that collect the genes related to diseases. Thus, if the uploaded gene sets have significant overlap with sets of genes in either “Jensen Diseases” or “OMIM disease,” we can regard the uploaded gene sets as being related to diseases. For drug gene sets, the uploaded set of genes was evaluated if it overlapped with a set of genes whose expression was altered by drug treatment. If the uploaded set of genes has a significant overlap with the gene sets whose expression is known to be altered by some drugs, we can consider that the expression of the uploaded gene set is also altered by drug treatment. Thus, if the uploaded gene set significantly overlaps with the drug and disease gene sets simultaneously, we can expect that drug treatment can significantly alter the expression of genes whose expression is known to be altered by drug treatment. Consequently, drugs related to diseases through genes are potential drug compounds, although they are not always effective against specific diseases.

Query of drug-disease relation by LLM

The actual queries were performed using Microsoft Copilot. The basic structure of the prompt is provided in Supplementary Information.

The actual query prompts and LLMs replies are accessible through the following two URLs: <https://copilot.microsoft.com/sl/dCHesf9R01c> and <https://copilot.microsoft.com/sl/jCQxrKJzZoO>.

Conclusions

In this study, we demonstrated the usefulness of drug repositioning TD-based unsupervised FE applied to PPI. Although it is unlikely that only the information retrieved from PPI will be useful for disease-specific drug repositioning, our findings indicate that it works in a practical sense. Thus, TD-based unsupervised FE applied to PPI is likely to be useful for drug repositioning. Further studies will be needed to understand the extent to which this strategy is effective.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-024-06009-9>.

Supplementary file 1 (zip 11054 KB)

Author Contributions

Y.-H.T. planned the study and performed the analyses. Y.-H.T. and T.T. evaluated the results, discussions, and outcomes and wrote and reviewed the manuscript. All the authors have read and agreed to the published version of this manuscript.

Data and availability

All the data used in this study can be downloaded from BioGRID [30] and DIP [31]. Sample code for data processing is in <https://github.com/tagtag/TDbasedUFEPPi>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Conflict of interest

The authors declare no Conflict of interest.

Received: 1 June 2024 Accepted: 5 December 2024

Published online: 18 December 2024

References

1. Zickenrott S, Angarica VE, Upadhyaya BB, Del Sol A. Prediction of disease-gene-drug relationships following a differential network analysis. *Cell Death Disease*. 2016;7(1):2040–2040. <https://doi.org/10.1038/cddis.2015.393>.
2. Wong M, Previde P, Cole J, Thomas B, Laxmeshwar N, Mallory E, Lever J, Petkovic D, Altman RB, Kulkarni A. Search and visualization of gene-drug-disease interactions for pharmacogenomics and precision medicine research using genedive. *J Biomed Inform*. 2021;117: 103732. <https://doi.org/10.1016/j.jbi.2021.103732>.
3. Yu H, Choo S, Park J, Jung J, Kang Y, Lee D. Prediction of drugs having opposite effects on disease genes in a directed network. *BMC Syst Biol*. 2016;10(Suppl 1):2. <https://doi.org/10.1186/s12918-015-0243-2>.
4. Sun PG. The human drug-disease-gene network. *Inf Sci*. 2015;306:70–80. <https://doi.org/10.1016/j.ins.2015.01.036>.
5. Qahwaji R, Ashankyty I, Sannan NS, Hazzazi MS, Basabrain AA, Mobashir M. Pharmacogenomics: A genetic approach to drug development and therapy. *Pharmaceuticals*. 2024. <https://doi.org/10.3390/ph17070940>.
6. Iida M, Iwata M, Yamanishi Y. Network-based characterization of disease-disease relationships in terms of drugs and therapeutic targets. *Bioinformatics*. 2020;36:516–24. <https://doi.org/10.1093/bioinformatics/btaa439>.
7. Quan Y, Luo Z-H, Yang Q-Y, Li J, Zhu Q, Liu Y-M, Lv B-M, Cui Z-J, Qin X, Xu Y-H, Zhu L-D, Zhang H-Y. Systems chemical genetics-based drug discovery: Prioritizing agents targeting multiple/reliable disease-associated genes as drug candidates. *Front Genetics*. 2019. <https://doi.org/10.3389/fgene.2019.00474>.

8. Wang L, Wang Y, Hu Q, Li S. Systematic analysis of new drug indications by drug-gene-disease coherent subnetworks. *CPT Pharmacometrics Syst Pharmacol*. 2014;3(11):146. <https://doi.org/10.1038/psp.2014.44>.
9. Kim Y, Cho Y-R. Predicting drug-gene-disease associations by tensor decomposition for network-based computational drug repositioning. *Biomedicines*. 2023. <https://doi.org/10.3390/biomedicines11071998>.
10. Taguchi Y-h. Unsupervised Feature Extraction Applied to Bioinformatics: A PCA Based and TD Based Approach. 1st ed. Berlin: Springer; 2020. <https://doi.org/10.1007/978-3-030-22456-1>.
11. Taguchi Y-h. Unsupervised Feature Extraction Applied to Bioinformatics: A PCA Based and TD Based Approach. 2nd ed. Berlin: Springer; 2024. <https://doi.org/10.1007/978-3-031-60982-4>.
12. Taguchi Y-H, Turki T. Integrated analysis of gene expression and protein-protein interaction with tensor decomposition. *Mathematics*. 2023. <https://doi.org/10.3390/math11173655>.
13. Xie Z, Bailey A, Kuleshov MV, Clarke DJB, Evangelista JE, Jenkins SL, Lachmann A, Wojciechowski ML, Kropiwnicki E, Jagodnik KM, Jeon M, Ma'ayan A. Gene set knowledge discovery with enrichr. *Current Protocols*. 2021;1(3):90. <https://doi.org/10.1002/cpz1.90>.
14. Carr S, Kasi A (2024) Familial Adenomatous Polyposis. StatPearls Publishing, Treasure Island (FL) . Updated 2023 Feb 25. <https://www.ncbi.nlm.nih.gov/books/NBK538233/>
15. Aedma SK, Kasi A (2024) Li-Fraumeni Syndrome. StatPearls Publishing, Treasure Island (FL) . PMID: 30335319. <https://pubmed.ncbi.nlm.nih.gov/30335319/>
16. Bhandari J, Thada PK, Puckett Y, Fanconi Anemia. StatPearls Publishing, Treasure Island (FL) (2024). Last Update: August 10, 2022. <https://www.ncbi.nlm.nih.gov/books/NBK559133/>
17. Zhao Q, Zhang Y, Shao S, Sun Y, Lin Z. Identification of hub genes and biological pathways in hepatocellular carcinoma by integrated bioinformatics analysis. *PeerJ*. 2021;9:10594. <https://doi.org/10.7717/peerj.10594>.
18. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *J Stat Mech: Theory Exp*. 2008;2008(10):10008. <https://doi.org/10.1088/1742-5468/2008/10/P10008>.
19. Raghavan UN, Albert R, Kumara S. Near linear time algorithm to detect community structures in large-scale networks. *Phys Rev E*. 2007;76: 036106. <https://doi.org/10.1103/PhysRevE.76.036106>.
20. Newman MEJ, Girvan M. Finding and evaluating community structure in networks. *Phys Rev E*. 2004;69: 026113. <https://doi.org/10.1103/PhysRevE.69.026113>.
21. Kang M-H, Seok Jeong G, Smoot DT, Ashktorab H, Mo Hwang C, Sik Kim B, Sung Kim H, Park Y-Y. Verteporfin inhibits gastric cancer cell growth by suppressing adhesion molecule fat1. *Oncotarget*. 2017;8(58):98887–97. <https://doi.org/10.18632/oncotarget.21946>.
22. Rigopoulos D, Larios G, Gregoriou S, Alevizos A. Acute and chronic paronychia. *Am Fam Physician*. 2008;77(3):339–46.
23. Arjmand Abbassi Y, Mohammadi MT, Sarami Foroshani M, Raouf Sarshoori J. Captopril and valsartan may improve cognitive function through potentiation of the brain antioxidant defense system and attenuation of oxidative/nitrosative damage in stz-induced dementia in rat. *Adv Pharm Bull*. 2016;6(4):531–9. <https://doi.org/10.15171/apb.2016.067>.
24. Michela C, Antonietta M, Domenico T. Effects of the antidiabetic drugs on the age-related atrophy and sarcopenia associated with diabetes type ii. *Curr Diabetes Rev*. 2014;10(4):231–7. <https://doi.org/10.2174/1573399810666140918121022>.
25. Ohno T, Nakane T, Akase T, Kurasawa H, Aizawa Y. Development of an isogenic human cell trio that models polyglutamine disease. *Genes Genetic Syst*. 2023;98(4):179–89. <https://doi.org/10.1266/ggs.22-00030>.
26. Giuliani F, Molica S, Maiello E, Battaglia C, Gebbia V, Bisceglie MD, Vinciarelli G, Gebbia N, Colucci G (2005) Irinotecan (cpt-11) and mitomycin-c (mmc) as second-line therapy in advanced gastric cancer: A phase ii study of the gruppo oncologico dell'italia meridionale (prot. 2106). *American Journal of Clinical Oncology*. 28(6), 581–585 <https://doi.org/10.1097/01.coc.0000190398.52142.7f>
27. Zhu X-C, Dai W-Z, Ma T. Overview the effect of statin therapy on dementia risk, cognitive changes and its pathologic change: a systematic review and meta-analysis. *Ann Transl Med*. 2018;6(22):435.
28. Billingsley E, Vidimos A, Paronychia Treatment & Management. *Medscape*, 2022 ;1106062. <https://emedicine.medscape.com/article/1106062-treatment?form=fpf>
29. Yoon J-H, Yoo C-I, Ahn Y-S. N, n-dimethylformamide: evidence of carcinogenicity from national representative cohort study in south korea. *Scandinavian J Work Environ Health*. 2019;4:396–401. <https://doi.org/10.5271/sjweh.3802>.
30. Bughtred R, Rust J, Chang C, Breitkreutz B-J, Stark C, Willems A, Boucher L, Leung G, Kolas N, Zhang F, Dolma S, Coulombe-Huntington J, Chatr-aryamontri A, Dolinski K, Tyers M. The BioGRID database: a comprehensive biomedical resource of curated protein, genetic, and chemical interactions. *Protein Sci*. 2021;30(1):187–200. <https://doi.org/10.1002/pro.3978>.
31. Xenarios I, Rice DW, Salwinski L, Baron MK, Marcotte EM, Eisenberg D. DIP: the Database of Interacting Proteins. *Nucleic Acids Res*. 2000;28(1):289–91. <https://doi.org/10.1093/nar/28.1.289>.
32. Bates D, Maechler M, Jagan M (2024) Matrix: Sparse and Dense Matrix Classes and Methods. R package version 1.7-0. <https://CRAN.R-project.org/package=Matrix>
33. Baglama J, Reichel L, Lewis BW (2022) Irlba: Fast Truncated Singular Value Decomposition and Principal Components Analysis for Large Dense and Sparse Matrices. R package version 2.3.5.1. <https://CRAN.R-project.org/package=irlba>

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.