1.DataSet 1

The first data set was introduced in 2016 by Ezzat et al[1]. Most new methods of predicting drug-target interactions have evaluated their method and reported results using this dataset[2-4]. This dataset includes three txt files of drug-target interactions, drug features, and target features, and all data in this data set are numeric.

The first data set is described in two parts; In the first part, information about drugtarget interactions and in the second part, information about the drug and target used and how they are represented as feature vectors are described.

1-1. Drug-target interaction data

Drug-target interaction file as a two-dimensional matrix $Y \in \mathbb{R}^{n \times m}$ with n drug rows and m target columns. That is, $Y_{ij} = 1$ if drug d_i and target t_j interact and $Y_{ij} = 0$ otherwise. In other words, the one indicates the presence of interaction and the zero indicates the lack of interaction between the drug d_i which is in the matrix row and the target t_j which is located in the matrix column. Also, the features of drugs are in the form of $D \in \mathbb{R}^{n \times f}$ matrix and the features of targets are in the form of $T \in \mathbb{R}^{m \times g}$ matrix, where f and g are the number of features of drugs and targets. The interaction data were obtained from the DrugBank [5]. In total, there are 12674 drug-target interactions between 5877 drugs and their 3348 protein interaction partners.

1-2. Drugs and targets information

features for drugs were calculated using the Rcpi package[6]. Examples of drug features include constitutional, topological and geometrical descriptors among other molecular properties. the target features were computed from their genomic sequences with the help of the PROFEAT web server[7]. The features that have been used to represent targets in this work are descriptors related to amino acid composition; dipeptide composition; autocorrelation; composition, transition and distribution; quasi-sequence-order; amphiphilic pseudo-amino acid composition and total amino acid properties.

The file of drug features is in the form of $D \in \mathbb{R}^{n \times f}$ matrix and the target features are in the form of $T \in \mathbb{R}^{m \times g}$ matrix, where f and g are the number of drug features and targets features. The constituents of these matrices are numbers that represent the values of the molecular descriptors of drugs and targets. The properties values of drugs and targets are decimal numbers, all of which are in the range [1,0] to avoid biasing properties in large quantities using the Min-Max normalization method. In this dataset, each drug and target is represented by a 193-dimensional and 1290-dimensional property vectors, respectively; Each drug is displayed as $d = [d_1, d_2, ..., d_p]$ and each target as $t = [t_1, t_2, ..., t_q]$. Details of this dataset are given in Table 1.

Table 1. Statistics of first dataset

Drugs	Targets	Interactions
5877	3348	12674

Formation of drug-target vectors for use in machine learning model as $[d_1, d_2, ..., d_p, t_1, t_2, ..., t_q, I]$, which label I is a binary number and indicates the presence or absence of interaction Is between the drug and the target in the vector. Label I is identified by the drug-target interaction matrix. Figures 1 and 2 show a view of the drug-target interaction file and the drug features file, respectively.

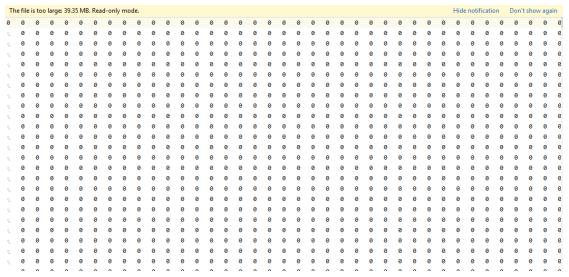


Figure 1. view of the drug-target interaction file

The file size	(6.48 MB) e	exceeds the con	figured limit (2	2.56 MB). Code ins	ight	featu	res ar	e not	avail	able.																			
1	0.74017	0.01967 0.24	4839 0.2340	9 0.23077 0	0	0	0	0	0	0.2	3077	7 0	0	0.	66667	0.2	0.1	1111	0	0.3	3333	0	0	0.2	2175	0.	1082	4 0.8	822 🛷
2	0.77503	0.0099534	0.20324 0.	18898 0.17949	0.1	6667	0.5	0	0.1	6667	0.1	14286	0	0.	20513	0	0	0.3	3333	0	0	0	0	0	0	0.	1712	5 0.0	088254
3	0.77162	0.010758	0.29128 0.	26601 0.25641	0.1	6667	0	0	0	0	0	0.25	641	0	0	0.3	3333	0.2	0	0	0.1	6667	0	0	0.2	2412	5 0.:	10624	4 0.82
4	0.84958	0.00017337	0.25819 0.	25141 0.15385	0	0	0	0	0	0	0.1	17949	0	0	0	0.4	0	0	0	0	0.4	0.2	4375	0.0	38944	19	0.8	82397	7 0.82
5	0.76321	0.01288 0.20	0109 0.1874	4 0.20513 0	0.5	0	0.3	3333	0.1	4286	0	0.2	8077	0	0	0.3	3333	0	0.1	1111	0	0	0	0	0.1	١7	0.0	08483	37
6	0.79337	0.0061565	0.21689 0.	18973 0.17949	0	0	0	0.3	3333	0	0	0.17	7949	0	0	0	0	0.1	1111	. 0	0.1	6667	0.2	5	0	0.	17	0.0	069826
7	0.83887	0.00065599	0.043929	0.035165	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.0	3125	0.0	32928	31	0.9	9273	3 0.85
8	0.83527	0.00088729	0.26454 0.	25528 0 0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.2	2625	0.0	2070:	1	0.9)609	1 0.	8337	3 0.17
9	0.82023	0.0022332	0.086367	0.073563	0.0	2564	1	0	0	0	0	0	0.14	128	6 0.05	5128	2	0	0	0	0	0	0	0	0	0	0.0	0675	0.04
10	0.82983	0.0013039	0.030515	0.024233	0.0	2564	1	0	0	0	0	0	0	0.	025641	l	0	0	0	0	0	0	0	0	0	0.0	0237	5 0.0	014919
11	0.8071	0.003908	0.070662	0.066079	0.0	2564	1	0	0	0	0	0	0	0.	025641	L	0	0	0	0	0	1	0	0	0	0.0	0612	5 0.0	040021
12	0.85476	5.1187e-005	0.014376	0.010469	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.0	1125	0.0	1474	16	0.9	97967	7 0.82
13	0.86246	3.198e-006	0.039448	0.029701	0.0	2564	1	0	0	0	0	0	0	0.	025641	L	0	0	0	0	0	0	0	0	0	0.0	0275	0.0	010025
14			0.042468	0.04196 0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.0								3 0.17
15		0.00052035		0.022464	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.0		0.0	90848	166			8 0.83
16		0.0015871	0.025917			2564		0	0	0	0	0			025641		0	0	0	0		1111		0	0	0	-		5 0.01
17		0.00099868	0.029167			2564	_		6667	0	0	0	0	0	0.02		_	0	0	0	0	0	0	0	0	0			5 0.01
18			0.029135	0.023293	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		2375						1 0.81
19		0.0024606	0.043922	0.046372	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		4875						9 0.83
20			0.021203			2564	_	0	0	0.5	0	0	-		025641		0	0	0	0	0	0	0	0	0				018201
21		0.001272	0.023591	0.023169		2564	_	0	0	0	0	0	-		025641		0	0	0	0	0	0	0	0	0		025		011076
22		0.0008695	0.023706	0.02251 0.02		_	0	0	0	0	0.1	14286	0		025641		0	0	0	0	0	0	0	0	0				015912
23		0.0024119	0.058498	0.049273	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.0			34648				0.84
24	0.8597	2.0128e-006	0.064148	0.063653	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.0	6125	0.0	10833	192	0.9	9806	0.83

Figure 2. view of the drug features file

References:

- [1] A. Ezzat, M. Wu, X.-L. Li, and C.-K. Kwoh, "Drug-target interaction prediction via class imbalance-aware ensemble learning," *BMC bioinformatics*, vol. 17, no. 19, p. 509, 2016.
- [Y] A. Sharma and R. Rani, "BE-DTI': ensemble framework for drug target interaction prediction using dimensionality reduction and active learning," *Computer methods and programs in biomedicine*, vol. 165, pp. 151-162, 2018.
- [*] A. Ezzat, M. Wu, X.-L. Li, and C.-K. Kwoh, "Drug-target interaction prediction via class imbalance-aware ensemble learning," *BMC bioinformatics*, vol. 17, no. 19, pp. 267-276, 2016.
- [٤] A. Ezzat, M. Wu, X. Li, and C.-K. Kwoh, "Computational prediction of drug-target interactions via ensemble learning," in *Computational methods for drug repurposing*: Springer, Young, pp. 239-254.
- [°] C. Knox *et al.*, "DrugBank 3.0: a comprehensive resource for 'omics' research on drugs," *Nucleic acids research*, vol. 39, no. suppl_1, pp. D1035-D1041, 2010.
- [1] D.-S. Cao, N. Xiao, Q.-S. Xu, and A. F. Chen, "Rcpi: R/Bioconductor package to generate various descriptors of proteins, compounds and their interactions," *Bioinformatics*, vol. 31, no. 2, pp. 279-281, 2014.
- [Y] Z.-R. Li, H. H. Lin, L. Han, L. Jiang, X. Chen, and Y. Z. Chen, "PROFEAT: a web server for computing structural and physicochemical features of proteins and peptides from amino acid sequence," *Nucleic Acids Research*, vol. 34, no. suppl_2, pp. W32-W37, 2006.