

1. Gold dataSet

The gold standard dataset is divided into four datasets; In other words, the four datasets of enzymes, ion channels, G protein-coupled receptors (GPCRs), and nuclear receptors are actually classes of proteins in the human body. This gold standard data was introduced in 2008 by Yamanishi et al[1]. The statistical quantities of existing drugs are 445, 223, 210, and 54, respectively. The numbers of known proteins are 664, 95, 204, and 26, respectively. The counts of the DTIs which have been proven are 2926, 635, 1476, and 90, respectively. Most new methods of predicting drug-target interactions have evaluated their method and reported results using this dataset [1-5]. Each data set in this gold data set includes three txt files of drug-target interactions, drug properties, and target properties, and the data in this data set are numeric.

The golden data set is explained in two parts; In the first part, information about drug-target 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

All data sets originate from the databases of DrugBank [6], SuperTarget [7], BRENDA [8], and KEGG BRITE [9]. 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.

1-2. Drugs and targets information

Note that in the enzyme class focused on the regulatory interactions between enzymes and compounds rather than the metabolic interactions, so all the ligands in the enzyme data are inhibitors or activators rather than substrates or products. Cofactors such as adenosine triphosphate (ATP) and nicotinamide adenine dinucleotide phosphate (NADPH) are also not included except when they are annotated as regulators in the BRENDA database. Also, do not use compounds whose molecular weights are <100 .

Chemical structures of the drugs were obtained from the DRUG and COMPOUND Sections in the KEGG LIGAND database. Amino acid sequences of the target proteins were obtained from the KEGG GENES database.

The files 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. Each drug is displayed as $d = [d_1, d_2, \dots, d_p]$ and each target as $t = [t_1, t_2, \dots, t_q]$. Details of this dataset are given in Table 1.

Table 1. Gold Dataset of Yamanishi et al

Dataset	Drugs	Targets	Interactions
Enzyme	445	664	2926
Ion channel	210	204	1476
GPCRs	223	95	635
Nuclear receptor	54	26	90

Formation of drug-target vectors for use in machine learning model as $[d_1, d_2, \dots, d_p, t_1, t_2, \dots, t_q, I]$, which label I is a binary number and indicates the presence or absence of interaction I s 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.

The file size (4.73 MB) exceeds the configured limit (2.56 MB). Code insight features are not available.									
1	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
2	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
3	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
4	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
5	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
6	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
7	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
8	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
9	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
10	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
11	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
12	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
13	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
14	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
15	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
16	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
17	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
18	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
19	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
20	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
21	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
22	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
23	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
24	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00

Figure 1. view of the drug-target interaction file

The file size (3.17 MB) exceeds the configured limit (2.56 MB). Code insight features are not available.									
1	1.000000e+00	5.156250e-01	3.846200e-02	8.474600e-02	9.803900e-02	1.200000e-01	8.333300e-02	9.090899e-02	1.000000e-01
2	4.696970e-01	1.000000e+00	3.278699e-02	7.352899e-02	8.333300e-02	8.333300e-02	1.090910e-01	9.523800e-02	8.474600e-02
3	3.846200e-02	3.278699e-02	1.000000e+00	4.285709e-01	1.000000e-01	3.750000e-01	0.000000e+00	2.380950e-01	4.000000e-01
4	8.474600e-02	7.352899e-02	4.285709e-01	1.000000e+00	6.666699e-02	2.307689e-01	0.000000e+00	2.000000e-01	2.399999e-01
5	9.803900e-02	8.333300e-02	1.000000e-01	6.666699e-02	1.000000e+00	9.090899e-02	0.000000e+00	7.692299e-02	9.523800e-02
6	1.200000e-01	8.333300e-02	3.750000e-01	2.307689e-01	9.090899e-02	1.000000e+00	1.764709e-01	1.666670e-01	6.428570e-01
7	8.333300e-02	1.090910e-01	0.000000e+00	0.000000e+00	0.000000e+00	1.764709e-01	1.000000e+00	4.347800e-02	1.875000e-01
8	9.090899e-02	9.523800e-02	2.380950e-01	2.000000e-01	7.692299e-02	1.666670e-01	4.347800e-02	1.000000e+00	1.739130e-01
9	1.000000e-01	8.474600e-02	4.000000e-01	2.399999e-01	9.523800e-02	6.428570e-01	1.875000e-01	1.739130e-01	1.000000e+00
10	1.960800e-02	1.666699e-02	2.000000e-01	1.200000e-01	5.263200e-02	1.764709e-01	0.000000e+00	2.631579e-01	1.875000e-01
11	1.777800e-01	8.771900e-02	5.555999e-02	7.407400e-02	0.000000e+00	4.000000e-01	2.142860e-01	0.000000e+00	2.500000e-01
12	5.555999e-02	4.761900e-02	4.375000e-01	3.199999e-01	8.695700e-02	1.904760e-01	0.000000e+00	2.083330e-01	2.000000e-01
13	1.041670e-01	5.084700e-02	1.875000e-01	1.153850e-01	0.000000e+00	6.153850e-01	2.142860e-01	1.363640e-01	3.333299e-01
14	1.960800e-02	5.172400e-02	3.846150e-01	2.173910e-01	5.263200e-02	3.333299e-01	0.000000e+00	1.428570e-01	3.571430e-01
15	3.999999e-02	5.172400e-02	3.846150e-01	2.173910e-01	1.111110e-01	3.333299e-01	0.000000e+00	1.428570e-01	3.571430e-01
16	5.769199e-02	1.034480e-01	5.000000e-02	6.896600e-02	4.545500e-02	4.545500e-02	0.000000e+00	0.000000e+00	4.761900e-02
17	4.000000e-01	3.333299e-01	7.407400e-02	5.405399e-02	1.923079e-01	1.481480e-01	1.250000e-01	1.290320e-01	2.000000e-01
18	1.041670e-01	6.896600e-02	1.875000e-01	1.153850e-01	0.000000e+00	5.000000e-01	2.142860e-01	1.363640e-01	4.285709e-01
19	2.241380e-01	4.814809e-01	5.714299e-02	9.302300e-02	1.818179e-01	1.142860e-01	2.068970e-01	7.500000e-02	1.176470e-01
20	3.225800e-02	4.285699e-02	3.043479e-01	3.333299e-01	3.225800e-02	1.428570e-01	0.000000e+00	2.000000e-01	1.481480e-01
21	3.404260e-01	2.857140e-01	3.571400e-02	2.631600e-02	1.923079e-01	1.071430e-01	2.272730e-01	1.666670e-01	1.538460e-01
22	5.769199e-02	1.034480e-01	1.666670e-01	1.071430e-01	0.000000e+00	2.777800e-01	5.833330e-01	1.739130e-01	2.941179e-01
23	3.921600e-02	5.084700e-02	4.615379e-01	2.608700e-01	1.052630e-01	3.125000e-01	0.000000e+00	1.904760e-01	3.333299e-01
24	2.459020e-01	2.318840e-01	2.727270e-01	3.333299e-01	4.761900e-02	2.571429e-01	1.428570e-01	1.707320e-01	2.285710e-01

Figure 2. view of the drug features file

References:

- [1] Y. Yamanishi, M. Araki, A. Gutteridge, W. Honda, and M. Kanehisa, "Prediction of drug–target interaction networks from the integration of chemical and genomic spaces," *Bioinformatics*, vol. 24, no. 13, pp. i232-i240, 2008.
- [2] Z.-Y. Zhao *et al.*, "An Ensemble Learning-Based Method for Inferring Drug-Target Interactions Combining Protein Sequences and Drug Fingerprints," *BioMed Research International*, vol. 2021, 2021.
- [3] Y. Tabei, E. Pauwels, V. Stoven, K. Takemoto, and Y. Yamanishi, "Identification of chemogenomic features from drug–target interaction networks using interpretable classifiers," *Bioinformatics*, vol. 28, no. 18, pp. i487-i494, 2012.
- [4] 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.
- [5] C. Réda, E. Kaufmann, and A. Delahaye-Duriez, "Machine learning applications in drug development," *Computational and structural biotechnology journal*, vol. 18, pp. 241-252, 2020.
- [6] D. S. Wishart *et al.*, "DrugBank: a comprehensive resource for in silico drug discovery and exploration," *Nucleic acids research*, vol. 34, no. suppl_1, pp. D668-D672, 2006.
- [7] S. Günther *et al.*, "SuperTarget and Matador: resources for exploring drug-target relationships," *Nucleic acids research*, vol. 36, no. suppl_1, pp. D919-D922, 2007.
- [8] I. Schomburg *et al.*, "BRENDA, the enzyme database: updates and major new developments," *Nucleic acids research*, vol. 32, no. suppl_1, pp. D431-D433, 2004.
- [9] M. Kanehisa *et al.*, "KEGG for linking genomes to life and the environment," *Nucleic acids research*, vol. 36, no. suppl_1, pp. D480-D484, 2007.