

۱. مجموعه داده طلایی

مجموعه داده طلایی به چهار مجموعه داده تقسیم می‌شود؛ به عبارت دیگر چهار مجموعه داده آنزیم‌ها^۱، کانال‌های یونی^۲، GPCRها^۳ و گیرنده‌های هسته‌ای^۴ در حقیقت کلاس‌های پروتئین در داخل بدن انسان هستند. این مجموعه داده طلایی در سال ۲۰۰۸ توسط یامانیشی^۵ و همکاران [۱] معرفی شده است. این مجموعه داده به ترتیب شامل ۹۰، ۶۳۵، ۱۴۷۶ و ۲۰۲۹ تعامل شناخته‌شده میان ۵۴، ۲۲۳، ۲۱۰ و ۴۴۵ دارو و ۲۶، ۹۵، ۲۰۴ و ۶۶۴ هدف است. عموم روش‌های پیش‌بینی تعاملات دارو-هدف، روش خود را بر پایه این مجموعه داده‌ها پیاده‌سازی و نتایج را گزارش کرده‌اند [۲-۴]. هر مجموعه داده موجود در این مجموعه طلایی شامل سه فایل متنی تعاملات دارو-هدف، ویژگی‌های داروها و ویژگی‌های اهداف می‌باشد و داده‌های موجود در این مجموعه داده، از نوع عددی هستند.

مجموعه داده طلایی در دو قسمت توضیح داده می‌شود؛ در قسمت اول، اطلاعاتی درباره تعاملات دارو-هدف و در قسمت دوم اطلاعات دارو و هدفی که مورد استفاده قرار گرفته شده است و نحوه نمایش آن‌ها به صورت بردارهای ویژگی، تشریح می‌شود.

۱-۱. اطلاعات تعاملات دارو-هدف

تعاملات موجود در این چهار مجموعه داده با استفاده از چند منبع معتبر KEGG BRITE [۵]، BRENDA [۶]، SuperTarget [۷] و DrugBank [۸] جمع‌آوری شده است. عناصر فایل‌های تعاملات دارو-هدف از اعداد صفر و یک تشکیل شده است؛ این فایل‌ها به صورت یک ماتریس دوبعدی $Y \in \mathbb{R}^{n \times m}$ دارای n سطر دارو و m ستون پروتئین هدف می‌باشد؛ به طوریکه اگر $Y_{ij} = 1$ باشد، دارو d_i با هدف t_j تعامل دارد و در غیر اینصورت $Y_{ij} = 0$ می‌باشد. به عبارت دیگر عدد یک نشان‌دهنده وجود تعامل و عدد صفر نشان‌دهنده عدم تعامل بین دارو d_i که در سطر ماتریس و هدف t_j که در ستون ماتریس قرار گرفته است، می‌باشد.

¹ Enzymes (E)

² Ion channels (IC)

³ G protein-coupled receptors (GPCRs)

⁴ Nuclear receptors (NR)

⁵ Yamanishi

۲-۱. اطلاعات داروها و اهداف

مورد توجه است که در مجموعه داده آنزیم‌ها بر روی فعل و انفعالات بین آنزیم‌ها و ترکیبات به جای فعل و انفعالات متابولیکی، بنابراین تمام لیگاندهای موجود در داده‌های آنزیمی به جای سوبستراها یا محصولات، بازدارنده یا فعال کننده هستند. کوفاکتورهایی مانند آدنوزین تری فسفات (ATP) و نیکوتین آمید آدنین دی نوکلئوتید فسفات (NADPH) نیز شامل نمی شوند، مگر زمانی که به عنوان تنظیم کننده در پایگاه داده BRENDA مشروح شده باشند. همچنین، از ترکیباتی که وزن مولکولی آنها کمتر از ۱۰۰ است استفاده نمی‌شود.

ساختارهای شیمیایی داروها از بخش DRUG و COMPOUND در پایگاه داده KEGG LIGAND و توالی‌های اسید آمینه پروتئین‌های مورد نظر از پایگاه داده KEGG GENES به دست آمده است؛ در واقع در این مطالعه بر روی پروتئین‌های موجود در بدن انسان تمرکز شده است.

فایل‌های ویژگی‌های داروها به صورت ماتریس $D \in \mathbb{R}^{n \times f}$ و ویژگی‌های اهداف به شکل ماتریس $T \in \mathbb{R}^{m \times g}$ می‌باشد که f و g تعداد ویژگی‌های داروها و اهداف است؛ به این صورت که هر دارو به شکل $d = [d_1, d_2, \dots, d_p]$ و هر هدف به صورت $t = [t_1, t_2, \dots, t_q]$ نمایش داده می‌شود. جزئیات این مجموعه داده در جدول ۱ آورده شده است.

جدول ۱. جزئیات مجموعه داده طلایی یامانیشی و همکاران

مجموعه داده	تعداد داروها	تعداد اهداف	تعداد تعاملات
آنزیم‌ها	۴۴۵	۶۶۴	۲۹۲۶
کانال‌های یونی	۲۱۰	۲۰۴	۱۴۷۶
GPCRs	۲۲۳	۹۵	۶۳۵
گیرنده‌های هسته‌ای	۵۴	۲۶	۹۰

تشکیل بردارهای دارو-هدف برای استفاده در مدل یادگیری ماشین در هر مجموعه داده به صورت $[d_1, d_2, \dots, d_p, t_1, t_2, \dots, t_q, I]$ می‌باشد که برچسب I یک عدد باینری است و نشان دهنده وجود یا عدم وجود تعامل بین دارو و هدف موجود در بردار می‌باشد. برچسب I از طریق ماتریس تعاملات دارو-هدف مشخص می‌شود. شکل ۱ و ۲ به ترتیب نمایی از فایل تعاملات دارو-هدف آنزیم‌ها و فایل ویژگی‌های داروهای موجود در آن می‌باشد.

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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	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	0.000000e+00
25	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	0.000000e+00

شکل ۱. نمایی از فایل تعاملات دارو-هدف آنزیم‌ها

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	7.692298e-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	7.692298e-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	1.777780e-01	8.771900e-02	5.555599e-02	7.407400e-02	0.000000e+00	4.000000e-01	2.142860e-01	0.000000e+00	2.500000e-01	
12	5.555599e-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	3.333329e-01	
14	1.960800e-02	5.172400e-02	3.846150e-01	2.173910e-01	5.263200e-02	3.333329e-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	3.333329e-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	3.333329e-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	3.333329e-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	2.777780e-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	3.333329e-01	
24	2.459020e-01	2.318840e-01	2.727270e-01	3.333329e-01	4.761900e-02	2.571429e-01	1.428570e-01	1.707320e-01	2.285710e-01	

شکل ۲. نمایشی از فایل ویژگی‌های داروهای موجود در مجموعه داده آنزیم

- [١] 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.
- [٢] 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.
- [٣] Y. Ding, J. Tang, and F. Guo, "Identification of drug–target interactions via fuzzy bipartite local model," *Neural Computing and Applications*, vol. 32, no. 14, pp. 10303-10319, 2020.
- [٤] C. Chen *et al.*, "DNN-DTIs: Improved drug-target interactions prediction using XGBoost feature selection and deep neural network," *Computers in Biology and Medicine*, vol. 136, p. 104676, 2021.
- [٥] M. Kanehisa, M. Furumichi, M. Tanabe, Y. Sato, and K. Morishima, "KEGG: new perspectives on genomes, pathways, diseases and drugs," *Nucleic acids research*, vol. 45, no. D1, pp. D353-D361, 2017.
- [٦] 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.
- [٧] 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.
- [٨] D. S. Wishart *et al.*, "DrugBank: a knowledgebase for drugs, drug actions and drug targets," *Nucleic acids research*, vol. 36, no. suppl_1, pp. D901-D906, 2008.