۱.مجموعه داده اول

مجموعه داده اول در سال ۲۰۱۶ توسط عزت و همکاران [۱] معرفی شده است. اغلب روشهای جدید پیشبینی تعاملات دارو- هدف، با استفاده از این مجموعه داده، روش خود را ارزیابی و نتایج را گزارش کردهاند [1-1]. این مجموعه داده شامل سه فایل متنی تعاملات دارو- هدف، ویژگیهای داروها و ویژگیهای اهداف میباشد و تمام دادههای موجود در این مجموعه داده، از نوع عددی هستند.

مجموعه داده اول در دو قسمت توضیح داده می شود؛ 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 می شود.

۱-۱.اطلاعات تعاملات دارو-هدف

عناصر فایل تعاملات دارو- هدف از اعداد صفر و یک تشکیل شده است؛ این فایل به صورت یک ماتریس دوبعدی d_i با d_i باشد، دارو و m ستون پروتئین هدف میباشد؛ به طوریکه اگر $I_i = 1$ باشد، دارو و $I_i = 1$ باشد، دارو و $I_i = 1$ باشد، دارو و $I_i = 1$ باشد، دارو و در غیر اینصورت $I_i = 1$ میباشد. به عبارت دیگر عدد یک نشان دهنده وجود تعامل و عدد صفر نشان دهنده عدم تعامل بین دارو $I_i = 1$ که در سطر ماتریس و هدف $I_i = 1$ که در ستون ماتریس قرار گرفته است، میباشد. تعاملات موجود در این مجموعه داده از پایگاه داده $I_i = 1$ باشد، وجود دارد. $I_i = 1$ باشد، و عناصل شناخته شده وجود دارد.

۲-۱.اطلاعات داروها و اهداف

ویژگیهای داروهای موجود در این مجموعه داده با استفاده از کتابخانه [۷] محاسبه می شود. نمونههایی از ویژگیهای داروها شامل توصیف کنندههای constitutional، توپولوژیکی داروها شامل توصیف کنندههای به دروها شامل توصیف کنند مولکولی می باشد [۱]. ویژگیهایی از اهداف که بتواند خواص پروتئینهای مختلف را به طور جامع توصیف کنند از توالی ژنومی پروتئینها با استفاده از وب سرور PROFEAT آمینه آمینه می شود؛ نمونههایی از ویژگیهای اهداف شامل توصیف کنندههای مربوط به ترکیب اسید آمینه مربوط به ترکیب اسید آمینه PROFEAT قابل مشاهده می باشد [۱].

-

¹ Topological

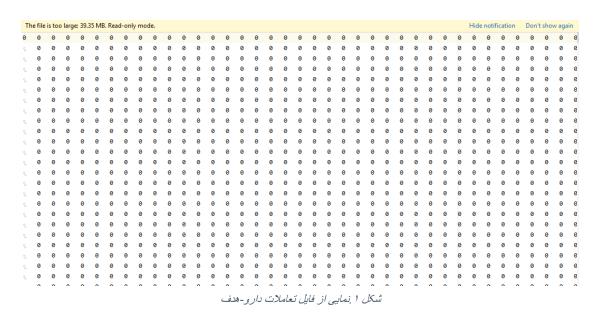
² Amino acid composition

فایل ویژگیهای داروها به صورت ماتریس $D \in \mathbb{R}^{n \times f}$ و ویژگیهای اهداف به شکل ماتریس اعدادی هستند میباشد که f و g تعداد ویژگیهای داروها و اهداف است. عناصر تشکیل دهنده این ماتریسها اعدادی هستند که مقادیر توصیف کنندههای ملکولی داروها و اهداف را نشان میدهند. مقادیر ویژگیهای داروها و اهداف اعدادی اعشاری هستند که همه آنها برای اجتناب از سوگیری ویژگیها با مقادیر بزرگ با استفاده از روش نرمال سازی Info Max در بازه $[\cdot,1]$ قرار گرفتهاند. در این مجموعه داده هر دارو و هدف به ترتیب با یک بردار ویژگی ۱۹۳ بعدی و ۱۲۹۰ بعدی نشان داده می شود؛ به این صورت که هر دارو به شکل [t,t] و و هر هدف به صورت [t,t] آورده شده است.

جدول ا .جزئيات مجموعه داده اول

تعاملات	تعداد	تعداد اهداف	تعداد داروها
1	7574	۲۳۴۸	۵۸۷۷

 $[d_1,d_2,\ldots,d_p,t_1,t_2,\ldots,t_q,I]$ تشکیل بردارهای دارو-هدف برای استفاده در مدل یادگیری ماشین به صورت I یک عدد باینری است و نشان دهنده وجود یا عدم وجود تعامل بین دارو و هدف موجود در بردار می باشد. برچسب I از طریق ماتریس تعاملات دارو-هدف مشخص می شود. شکل I و I به ترتیب نمایی از فایل تعاملات دارو-هدف و فایل ویژگیهای داروها می باشد.



The file size (6.48 MB) exceeds the configured limit (2.56 MB). Code insight features are not available.																									
1	0.74017 0.01967 0.	24839 0.2340	9 0.23077 0	0 0	0	0	0	0.2	3077	0	0	0.66	667	0.2 0.	1111	1 0	0.3	3333	0	0	0.2	175	0.10	824 @	.822 🤝
2	0.77503 0.0099534	0.20324 0.	18898 0.17949	0.166	57 0.	5 0	0.1	16667	0.1	4286	0	0.20	513	0 0	0.	33333	0	0	0	0	0	0	0.17	125 0	.088254
3	0.77162 0.010758	0.29128 0.	26601 0.25641	0.166	57 0	0	0	0	0	0.25	641	0	0	0.3333	3 0.	2 0	0	0.16	6667	0	0	0.24	1125	0.106	24 0.82
4	0.84958 0.00017337	0.25819 0.	25141 0.15385	0 0	0	0	0	0	0.1	7949	0	0	0	0.4 0	0	0	0	0.4	0.24	4375	0.0	89449)	0.823	97 0.82
5	0.76321 0.01288 0.	20109 0.1874	4 0.20513 0	0.5 0	0.	33333	0.1	L 42 86	0	0.23	3077	0	0	0.3333	3 0	0.1	1111	0	0	0	0	0.17	7	0.084	837
6	0.79337 0.0061565	0.21689 0.	18973 0.17949	0 0	0	0.3	3333	8 0	0	0.17	7949	0	0	0 0	0.	11111	. 0	0.16	6667	0.25	i	0	0.17	. 6	.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	29281	L	0.927	33 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	20701	L	0.96	5091	0.833	73 0.17
9	0.82023 0.0022332	0.086367	0.073563	0.025	541	0	0	0	0	0	0.14	4286	0.05	1282	0	0	0	0	0	0	0	0	0	0.067	5 0.04
10	0.82983 0.0013039	0.030515	0.024233	0.025	541	0	0	0	0	0	0	0.02	5641	. 0	0	0	0	0	0	0	0	0	0.02	375 0	.014919
11	0.8071 0.003908	0.070662	0.066079	0.025	541	0	0	0	0	0	0	0.02	5641	. 0	0	0	0	0	1	0	0	0	0.06	125 0	.040021
12	0.85476 5.1187e-00	5 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	14746	ó	0.979	67 0.82
13	0.86246 3.198e-006			0.025	541	0	0	0	0	0	0	0.02			0	0	0	0	0	-	0				.010025
14	0.84538 0.00032589			0 0	0	0	0	0	0	0	0	0	0	0 0	0	0	0								23 0.17
15	0.84128 0.00052035			0 0	0	0	0	0	0	0	0	-	-	0 0	0	0	0		0.0			08486			08 0.83
16	0.82662 0.0015871	0.025917		0.025		0	0	0	0	0	0	0.02		_	0	0	0	0.11			0	-			75 0.01
17	0.83371 0.00099868			0.025			.6667	-	0	0	0		0.02		0	0	0	0	0	-	0				25 0.01
18	0.83897 0.00064999			0 0	0	0	0	0	0	0	0	-	-	0 0	0	0	0	-				26036			91 0.81
19	0.81821 0.0024606	0.043922	0.0.00.2	0 0	0	0	0	0	0	0	0	•		0 0	0	0	0	0				07804			.99 0.83
20	0.84405 0.00038376			0.025		0	0	0.5	-	0	0	0.02		_	0	0	0	0	0	•	0	-			.018201
21	0.83021 0.001272	0.023591		0.025		0	0	0	0	0	-	0.02		_	0	0	0	0	0	•	0	-	0.02	-	.011076
22	0.83553 0.0008695	0.023706	0.02251 0.02		0	0	0	0		4286		0.02			0	0	0	9	0	0	0				.015912
23	0.81864 0.0024119	0.058498		0 0	0	0	0	0	0	0	0	-	•	0 0	0	0	0	-	0.0			46482			9 0.84
24	0.8597 2.0128e-00		0.063653	0 0	0	0	0	0	9	0	0	0	•	0 0	0	0	0	-				08339	_	0.980	0.83

شکل ۲ نمایی از فایل ویژگیهای دارو

- [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, "Computational prediction of drug-target interactions using chemogenomic approaches: an empirical survey," *Brief Bioinform*, vol. 8, 2018.
- [٤] A. Ezzat, M. Wu, X.-L. Li, and C.-K. Kwoh, "Drug-target interaction prediction using ensemble learning and dimensionality reduction," *Methods*, vol. 129, pp. 81-88, 2017.
- [°] 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.
- 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.
- [Y] 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.
- [^] 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.