Weighted Network Analysis of Biologically Relevant Chemical Spaces

Mariko I. Ito and Takaaki Ohnishi

The University of Tokyo, Bunkyo-ku Tokyo, Japan, marikoito.rnu1@gmail.com

Abstract. In cheminformatics, network representations of the space of compounds have been suggested extensively. Among these, the threshold-network consists of nodes representing molecules. In this network representation, two molecules are connected by a link if the Tanimoto coefficient, a similarity measure, between them exceeds a preset threshold. However, the topology of the threshold-network is affected significantly by the preset threshold. In this study, we collected the data of biologically relevant compounds and bioactivities. We defined the weighted network where the weight of each link between the nodes equals the Tanimoto coefficient between the bioactive compounds (nodes) without using the threshold. We investigated the relationship between the strength of the link connection and the bioactivity closeness in the weighted networks. We found that compounds with significantly high or low bioactivity have a stronger connection than those in the overall network.

 $\textbf{Keywords:} \ \ \text{compound space, chemical space networks, community structure}$

1 Introduction

The chemical space is an abstract concept but is roughly defined as a set of all possible molecules [1,2]. In cheminformatics, the central idea that structurally similar compounds tend to share the similar chemical properties is called the similarity property principle [3]. Based on this idea, the calculation of structural similarity is performed for various purposes from drug discovery to retrosynthetic analysis [4,5,6]. In drug discovery, biologically relevant chemical spaces are primarily explored. Compounds exhibit biological activity in this space. For example, ligand is a compound that binds to a receptor (the target) and inhibits biological response [7]. In these circumstances, it has been extensively investigated whether compounds with a similar structure share similar bioactivity [1].

In networks, edges represent various kinds of relationships such as interaction, social influence, and correlation between two nodes [8,9,10,11,12]. Investigating the topology of such networks affords a global view of how they are related to each other. For example, through community detection on networks, we can extract the groups of nodes each of which are densely connected. Nodes in a community can be regarded as those that are particularly interacting [13], and

having a similar feature or role [14]. 'Being similar' is a kind of relationship. Previous studies have suggested certain network representations of biologically relevant chemical spaces, where each node represents a compound and each link shows the similarity relationship between the compounds [2,3,4,5,15,16,17]. They investigated the topological features of the chemical subspace and examined how molecules with certain bioactivity are distributed among the network. In these studies, community detection was performed as well. Through community detection on such networks, we obtain groups containing nodes with a similar chemical structure.

Network representation of a chemical space was performed using the circular fingerprint technique and Tanimoto coefficient [17]. In circular fingerprint representation, a molecule is often represented by a vector called a fingerprint. In the vector, each index denotes a certain chemical substructure, and the entry denotes the count of the molecule substructure corresponding to the index. The Tanimoto coefficient is the most popular similarity measure between two molecules [18]. It takes a value from 0 to 1, and equals 1 if two molecules are the same. In previous studies regarding network representation based on the Tanimoto coefficient, two nodes were assumed to be connected by a link if the Tanimoto coefficient between them exceeded a preset threshold. In these studies, a threshold-dependent unweighted network is defined, known as the 'thresholdnetwork'. Furthermore, the value of the preset threshold was tuned such that the edge density was approximately 0.025. Consequently, a well-resolved community structure was obtained. Although the evaluation was not performed in detail, the visualized network demonstrated that compounds with similar bioactivity tend to form a community [17]. However, the topology of the threshold-network is affected significantly by the preset threshold. A point of concern is that the threshold-network constructed by an artificially preset threshold cannot capture the structure of the chemical space. While constructing the threshold-network, the structural information of the chemical subspace should be reduced significantly.

Hence, in the present study, we analyze the weighted network of biologically relevant chemical spaces as follows. Instead of applying a preset threshold to determine the existence of a link, we assume that two nodes (molecules) are connected by a link whose weight is the similarity between them. In particular, we are interested in discovering whether the weighted network topology can facilitate the investigation of compounds with high bioactivity. We evaluate the community structure on the weighted networks and discuss whether nodes that are strongly connected to each other share a similar activity.

2 Materials and Methods

To investigate the structure of biologically relevant chemical spaces, we collected data from ChEMBL (version 25), an open bioactivity database [19]. We selected 19 targets based on a previous study [17], as shown in Table 1. For each target, we extracted the data of compounds whose potency has been tested against the

target by the measure of Ki, a kind of bioactivity. We regard the pChEMBL value of the compound as an indicator of its bioactivity [19,20]. The larger the pChEMBL value, the stronger is the bioactivity against the target [20]. The number of compounds corresponding to each target is shown in Table 1. Subsequently, we obtained the "Morgan fingerprints" (circular fingerprints) of these compounds using RDKit, an open-source toolkit for cheminformatics. The circular fingerprint we used in this study is Morgan fingerprint. Each fingerprint is a 2048-dimensional vector, in which the entry in each index is an integer. We calculated the Tanimoto coefficient for all pairs of compounds that correspond to each target. For two molecules that have fingerprints \mathbf{x}_i and \mathbf{x}_j , the Tanimoto coefficient T_{ij} [18] is calculated as

$$T_{ij} = \frac{\mathbf{x}_i \cdot \mathbf{x}_j}{|\mathbf{x}_i|^2 + |\mathbf{x}_j|^2 - \mathbf{x}_i \cdot \mathbf{x}_j}.$$
 (1)

Subsequently, the similarity matrix T, in which the (i, j) entry is the similarity (Tanimoto coefficient) between molecules i and j, is constructed for each target. We considered the weighted network for each target by regarding the similarity matrix as the adjacency matrix.

In weighted network analysis, not only the number of links connected to the node, but also the sum of the weight of those links should be considered [21]. The former is the degree of the node and the latter is its strength [22]. As the Tanimoto coefficient does not vanish in the case of almost all pairs of compounds, the weighted networks are almost complete and the degrees of nodes do not vary. Therefore, we examined how strength is distributed in each weighted network. To examine whether the weighted networks exhibit community structure, we applied Louvain heuristics, an algorithm used to obtain a graph partition that (locally) optimizes the modularity Q [23]. In the case of weighted networks, the modularity Q [22,14] is defined as

$$Q = \frac{1}{2m} \sum_{i,j} M_{ij} \delta(c(i), c(j)), \tag{2}$$

where

$$M_{ij} := T_{ij} - \frac{s_i s_j}{2m}. (3)$$

The strength of node i, $\sum_j T_{ij}$, is denoted by s_i . The sum of all weights $\sum_{i,j} T_{ij}$ is 2m and c(i) denotes the community to which node i belongs. Kronecker's delta is denoted by δ ; therefore, $\delta(x,y)$ equals 1 (0) when x=y ($x\neq y$). Regarding M_{ij} , the second term on the right side of Eq. (3) represents the expected strength of the link between nodes i and j in the null-network, which is random except that it has the same strength distribution as the focal network [14]. Therefore, M_{ij} represents how strongly nodes i and j are connected compared to the null-model.

3 Results

First, we show the histogram of pChEMBL value in each compound set corresponding to each target. For the three examples of networks for targets 238,

Table 1. Examined compounds. In the first column, Target ID means the ChEMBL ID assigned to each target in the ChEMBL database. The second column shows the name of the targets—Hs, Rn, and Cp represent $Homo\ sapiens$, $Cavia\ Porcellus$, and $Rattus\ norvegicus$, respectively. In the third column, the number of compounds that correspond to the target is shown. In the fourth and fifth columns, the mean and standard deviation of the node strength are shown, respectively. The sixth and seventh columns show the number of communities and the modularity Q in the resulted graph partition by the community detection, respectively.

Target ID	Target name	Size	Mean	Std	Com	\overline{Q}
255	Adenosine A2b receptor (Hs)	1575	723	130	3	0.060
3242	Carbonic anhydrase XII (Hs)	2392	863	262	3	0.051
269	Delta opioid receptor (Rn)	1577	719	113	4	0.081
219	Dopamine D4 receptor (Hs)	2138	1087	183	4	0.031
238	Dopamine transporter (Hs)	1406	528	86	4	0.065
65338	Dopamine transporter (Rn)	1624	723	105	3	0.074
339	Dopamine D2 receptor (Rn)	2555	1119	171	3	0.045
4124	Histamine H3 receptor (Rn)	1591	637	135	4	0.075
344	Melanin-concentrating	1430	771	68	4	0.046
	hormone receptor 1 (Hs)					
270	Mu opioid receptor (Rn)	2318	984	178	4	0.091
4354	Mu opinion receptor (Cp)	654	266	51	3	0.078
2014	Nociceptin receptor (Hs)	1105	519	97	4	0.070
2001	Purinergic receptor P2Y12 (Hs)	584	400	74	4	0.029
225	Serotonin 2c (5-HT2c) receptor (Hs)	1980	785	132	4	0.049
273	Serotonin 1a (5-HT1a) receptor (Rn)	3370	1469	249	3	0.052
322	Serotonin 2a (5-HT2a) receptor (Rn)	3076	1278	215	4	0.067
1833	Serotonin 2b (5-HT2b) receptor (Hs)	1121	421	74	5	0.058
3155	Serotonin 7 (5-HT7) receptor (Hs)	1569	775	129	3	0.035
4153	Sigma-1 receptor (Cp)	1617	717	123	3	0.048

2001, and 2014, the histograms of pChEMBL values are shown in Fig. 1(a). Few compounds have an extremely small or large value. We also observed a similar tendency in the case of other targets.

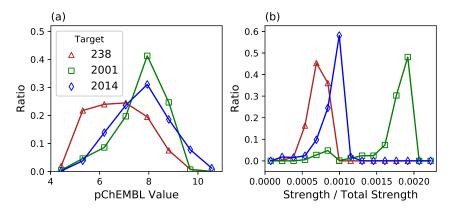


Fig. 1. Histogram of pChEMBL values and strength. (a) Histograms of pChEMBL values in networks of targets 238, 2001, and 2014. (b) Histograms of strength normalized by total strength in networks of targets 238, 2001, and 2014. Markers are same as in (a).

Subsequently, we examined the strength of the nodes in the weighted networks. In Table 1, we show the mean and standard deviation of the node strength in the weighted network for each target. Fig. 1(b) shows three histograms of the strength normalized by the total strength, $s_i/\sum_{i,j}T_{ij}$, in each network. As shown, many nodes share a similar strength, while a few nodes exhibit small strength. No node exhibited extremely large strength in all networks.

Finally, we investigated whether nodes connected with high similarity tend to share similar bioactivity. As explained in Sect. 2, we performed community detection. The number of communities and the modularity Q of the graph partition resulted from the Louvain heuristics is shown for each target in Table 1. The values of modularity are low, and the community structure in each network is weak in general.

We further inspected the community structure obtained by this community detection. Some detected communities were not connected sufficiently; as such, they could not be called 'communities'. Therefore, we extracted communities that could be regarded as connected strongly. For detected community C, we defined the extent to which the nodes are strongly connected within C, Q_C , as

$$Q_C = \frac{1}{\sum_{i,j \in C} T_{ij}} \sum_{i,j \in C} M_{ij},$$
(4)

where $\sum_{i,j\in C}$ denotes the sum of all pairs of nodes within community C. Therefore, the modularity Q equals $\sum_{C} Q_{C}$, and Q_{C} can measure how strongly links within C are connected without considering other communities. In Figs. 2(a)–(c), we show the histograms of M_{ij} for all pairs of nodes within each community and those in the overall network, for three targets. In these figures, only communities with Q_{C} exceeding 0.2 are shown. Therefore, these communities have larger values of M_{ij} than the overall network.

On the other hand, Figs. 2(d)–(f) show the histogram of pChEMBL values of nodes in each community included in Figs. 2(a)–(c) and the overall network. In the case of target 238, the histogram of pChEMBL value for community 4 shifts to the right side compared to the overall network (Fig. 2(d)). Community 5 in the network for target 2014 exhibits the same feature as well (Fig. 2(f)). Conversely, Community 1 in the network for target 2001 comprise nodes with lower pChEMBL values than those in the overall network. In Fig. 3, for all targets, we show the mean of pChEMBL values in each community that satisfies $Q_C > 0.2$ and consists of more than 20 nodes. Each error bar shows the standard deviation. In some communities, the mean pChEMBL value is located far from that of the overall network. However, in most cases, this value is within the standard deviation range of that of the overall network.

In summary, although the whole community structure is weak, we observed some communities in which the nodes are connected with a large weight. In some of them, the distribution of pChEMBL value is biased compared to that of the overall network. This suggests that certain sets of compounds are similar to each other and share stronger/weaker bioactivity against the target than the compounds in general.

Subsequently, we investigated whether nodes with particularly high (low) pChEMBL values are connected to each other with a large weight. First, we collected the |0.01RN| nodes whose pChEMBL values exceeded the (100-R)-th percentile, where N is the number of nodes in the network. Second, we calculated the mean of M_{ij} for all pairs of the $\lfloor 0.01RN \rfloor$ nodes. Similarly, we calculated the mean of M_{ij} for nodes with low pChEMBL values (lower than the R-th percentile), and intermediate pChEMBL values (ranging from the (50 - R/2)-th to (50 + R/2)-th percentile). The results for these three cases are presented in Figs. 4(a)–(c), where the horizontal axis represents the ratio R and the vertical axis the mean of M_{ij} . The mean M_{ij} exceeded 0 when the ratio R is small in the cases of high and low pChEMBL values. The mean M_{ij} also exceeds 0 for a small ratio R in the case of intermediate pChEMBL values, but it is much lower than the means in the other cases. The mean M_{ij} decreases with the ratio and approaches the mean of the overall network, which approximately equals 0. Although Figs. 4(a)–(c) show only the targets 238, 2001, and 2014, we observed the same tendency in all other targets. Therefore, the sets of nodes with high/low pChEMBL values in particular are connected with stronger weights than the overall network. Figs. 4(a)-(c) show some consistency with Figs. 2(d)(f). For target 2001, the nodes with low pChEMBL values are connected strongly (Fig. 4(b)) and some of them are detected as those included in Community 1

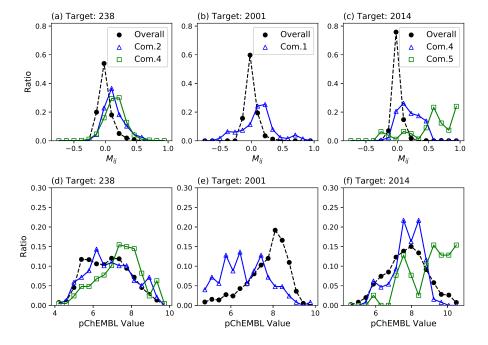


Fig. 2. Histograms of M_{ij} and pChEMBL values in communities. Histograms of M_{ij} in the overall network and in each community in the cases of targets 238 (a), 2001 (b), and 2014 (c). Histograms of pChEMBL values in the overall network and in each community in the cases of targets 238 (d), 2001 (e), and 2014 (f). In these figures, only communities with $Q_C > 0.2$ and size exceeding 20 are shown.

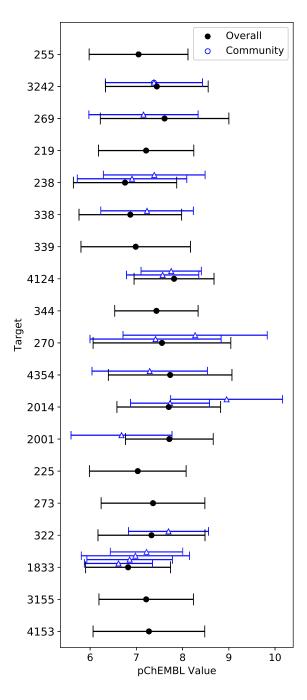


Fig. 3. Mean pChEMBL value in each community. For each target (ordinate), the mean pChEMBL value in the overall network is shown by a circle. For communities with $Q_C>0.2$ and size exceeding 20, the mean pChEMBL value in each community is shown by a triangle above that of the overall network. The error-bar represents the standard deviation of pChEMBL values.

sharing a low pChEMBL value (Fig. 2(e)). Additionally, consistency is shown between Community 5 in the network of target 2014 (Fig. 2(f)) and the set of high pChEMBL values (Fig. 4(c)).

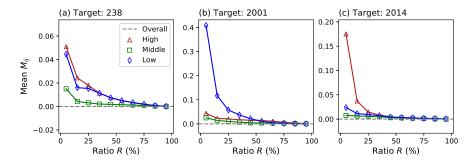


Fig. 4. Mean M_{ij} as a function of the ratio R. Mean M_{ij} of all links that connect nodes with high (more than (100 - R)-th percentile)/intermediate (from (50 - R/2)-th to (50 + R/2)-th percentile)/low (less than R-th percentile) pChEMBL values versus the ratio R, in the cases of targets 238 (a), 2001 (b), and 2014 (c).

4 Discussion

In this study, we investigated the structure of biologically relevant compounds that share the same target. Previous studies have suggested the network representations of the structure. In the network representation of these compounds, a node represents a compound, and a link between two nodes is drawn if the similarity between them exceeds a preset threshold. The topology of these thresholdnetworks greatly depends on the preset threshold. Therefore, to understand the true nature of the structure, we considered the weighted network, where the weight of each link is the similarity between connecting compounds.

For each target, the corresponding weighted network showed a homogeneous structure, which comprises a rare node exhibiting extremely strong bioactivity, or that connecting to other nodes with extreme strength. This homogeneity is attributable to the sample bias—the compounds in each network are those sharing the same target.

In cheminformatics, the question of whether compounds that are structurally similar share similar chemical properties needs to be elucidated [1]. We performed community detection on the weighted networks to investigate whether strongly connected nodes exhibit similar bioactivity. We found that, in general, the community structure was weak in all the weighted networks. However, we observed that nodes with high/low bioactivity against the target were connected strongly to each other compared to the nodes of the overall network. Some detected communities reflected this tendency and each of their nodes exhibited a

different range of pChEMBL values compared to those in the overall network. As a practical application, such communities can help us predict whether a novel compound exhibits high bioactivity. If the novel compound is structurally similar to compounds in a community sharing high bioactivity, we can expect it to exhibit high bioactivity. Such a prediction is useful in drug discovery.

In a previous study concerning the threshold-network, community detection was performed using modularity as the quality function of graph partition [17]. The values of modularity were much greater than those in our study and a well-resolved community structures were obtained. In fact, the threshold was set to obtain high modularity without resulting in an extremely sparse community structure. The threshold was set based on network visualization, which appears to be intuitive. The weight of a link should be determined mathematically considering observations from a community structure; this is aimed to be explored in future works. For example, a definition of the weight link may be represented as a sigmoid curve

$$f(x) = \left[\frac{1}{1 + \left(\frac{1}{x^{\alpha}} - 1\right)^{\beta}}\right]^{\gamma},\tag{5}$$

where $x (\in [0,1])$ is the similarity. As α approaches infinity, f(x) approaches 1 (0) when $x>2^{-1/\alpha}$ ($x<2^{-1/\alpha}$). This limit of $\beta\to\infty$ corresponds to the construction of the threshold-network, in which the preset threshold is $2^{-1/\alpha}$. On the other hand, setting the weight of each link in our study corresponds to a limiting function f(x)=x, which is obtained when α , β and γ equal 1. In future works, the optimization of parameters α , β and γ should be investigated. Accordingly, it will be challenging to evaluate the efficacy of the detected community structure considering the structural similarity, bioactivity distribution, and application, for example, to drug discovery.

Finally, the weighted network representation with other definitions of similarity should be considered. Although the Tanimoto coefficient is a popular similarity measure, it measures the global similarity of compounds, which is sometimes disadvantageous. The bioactivity of a compound often depends on the structure of a certain part of the compound. In some studies, networks were considered and evaluated based on other types of similarity [2,5,15]. Weighted networks with those similarity measures have not been investigated yet. We expect that weighted network analysis with similarity measures other than the Tanimoto coefficient can promote a better understanding of the structure–activity relationship in a biologically relevant chemical space.

Acknowledgement

This work was supported by the JSPS Grant-in-Aid for Scientific Research on Innovative Areas: 17H06468.

References

- 1. G. Opassi, A. Gesù, and A. Massarotti, "The hitchhikers guide to the chemical-biological galaxy," *Drug discovery today*, vol. 23, no. 3, pp. 565–574, 2018.
- M. Vogt, D. Stumpfe, G. M. Maggiora, and J. Bajorath, "Lessons learned from the design of chemical space networks and opportunities for new applications," *Journal of computer-aided molecular design*, vol. 30, no. 3, pp. 191–208, 2016.
- 3. M. Vogt, "Progress with modeling activity landscapes in drug discovery," *Expert opinion on drug discovery*, vol. 13, no. 7, pp. 605–615, 2018.
- 4. R. Kunimoto, M. Vogt, and J. Bajorath, "Tracing compound pathways using chemical space networks," *MedChemComm*, vol. 8, no. 2, pp. 376–384, 2017.
- R. Kunimoto and J. Bajorath, "Combining similarity searching and network analysis for the identification of active compounds," ACS omega, vol. 3, no. 4, pp. 3768–3777, 2018.
- C. W. Coley, L. Rogers, W. H. Green, and K. F. Jensen, "Computer-assisted retrosynthesis based on molecular similarity," ACS central science, vol. 3, no. 12, pp. 1237–1245, 2017.
- 7. C. M. Dobson, "Chemical space and biology," 2004.
- 8. M. E. Newman, "The structure and function of complex networks," *SIAM review*, vol. 45, no. 2, pp. 167–256, 2003.
- 9. R. Albert and A.-L. Barabási, "Statistical mechanics of complex networks," *Reviews of modern physics*, vol. 74, no. 1, p. 47, 2002.
- T. Ohnishi, H. Takayasu, and M. Takayasu, "Network motifs in an inter-firm network," *Journal of Economic Interaction and Coordination*, vol. 5, no. 2, pp. 171– 180, 2010.
- M. I. Ito, H. Ohtsuki, and A. Sasaki, "Emergence of opinion leaders in reference networks," *PloS one*, vol. 13, no. 3, p. e0193983, 2018.
- M. Bazzi, M. A. Porter, S. Williams, M. McDonald, D. J. Fenn, and S. D. Howison, "Community detection in temporal multilayer networks, with an application to correlation networks," *Multiscale Modeling & Simulation*, vol. 14, no. 1, pp. 1–41, 2016.
- 13. C. Mizokami and T. Ohnishi, "Revealing persistent structure of international trade by nonnegative matrix factorization," in *International Conference on Complex Networks and their Applications*, pp. 1088–1099, Springer, 2017.
- S. Fortunato, "Community detection in graphs," *Physics reports*, vol. 486, no. 3-5, pp. 75–174, 2010.
- 15. M. Wu, M. Vogt, G. M. Maggiora, and J. Bajorath, "Design of chemical space networks on the basis of tversky similarity," *Journal of computer-aided molecular design*, vol. 30, no. 1, pp. 1–12, 2016.
- B. Zhang, M. Vogt, G. M. Maggiora, and J. Bajorath, "Design of chemical space networks using a tanimoto similarity variant based upon maximum common substructures," *Journal of computer-aided molecular design*, vol. 29, no. 10, pp. 937– 950, 2015.
- 17. M. Zwierzyna, M. Vogt, G. M. Maggiora, and J. Bajorath, "Design and characterization of chemical space networks for different compound data sets," *Journal of computer-aided molecular design*, vol. 29, no. 2, pp. 113–125, 2015.
- 18. A. H. Lipkus, "A proof of the triangle inequality for the tanimoto distance," *Journal of Mathematical Chemistry*, vol. 26, no. 1-3, pp. 263–265, 1999.
- A. P. Bento, A. Gaulton, A. Hersey, L. J. Bellis, J. Chambers, M. Davies, F. A. Krüger, Y. Light, L. Mak, S. McGlinchey, et al., "The chembl bioactivity database: an update," Nucleic acids research, vol. 42, no. D1, pp. D1083–D1090, 2014.

- 20. F. P. Steinmetz, C. L. Mellor, T. Meinl, and M. T. Cronin, "Screening chemicals for receptor-mediated toxicological and pharmacological endpoints: Using public data to build screening tools within a knime workflow," *Molecular informatics*, vol. 34, no. 2-3, pp. 171–178, 2015.
- 21. A. Barrat, M. Barthelemy, R. Pastor-Satorras, and A. Vespignani, "The architecture of complex weighted networks," *Proceedings of the national academy of sciences*, vol. 101, no. 11, pp. 3747–3752, 2004.
- 22. L. Lü, D. Chen, X.-L. Ren, Q.-M. Zhang, Y.-C. Zhang, and T. Zhou, "Vital nodes identification in complex networks," *Physics Reports*, vol. 650, pp. 1–63, 2016.
- 23. V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, p. P10008, 2008.