# Drug-Disease Association Prediction by Known Nearest Neighbour Algorithm

Gill Varghese Sajan[1] and Joby George[2]
[1,2]Mar Athanasius College of Engineering, Kothamangalam
Ernakulam, Kerala, India. Pin: 686666

**Abstract:-** **Cooperations among medications and infection give vital data to the medication disclosure. Currently, only a few number of drug-disease interactions are identified through experiments. Accordingly, the advancement of computational strategies for medication malady cooperation forecast is a dire undertaking of hypothetical intrigue and commonsense significance. Legitimate treatment plans can be made by the specialist, when he is told with every single imaginable medication that can be utilized for the treatment of every diseases. Critical data with respect to the medication disclosure and medication repositioning are gotten from medication infection affiliations computations. The worldview about medication revelation has changed from finding new medications that display restorative properties for an ailment to reusing existing medications for a fresher illness. Exploratory assurance of illnesses relationship with medications could be tedious and expensive. Computational technique can be a proficient choice to distinguish potential medications identified with every infection. Here we present a methodology which processes the missing affiliations that exists among ailments and medications utilizing known nearest neighbour algorithm.**

*Keywords:- Drug, Disease, Association, Prediction.*

## I. INTRODUCTION

Medications are synthetic compounds utilized for distinguishing and treating sicknesses. When the related illnesses are found, it very well may be utilized to anticipate the medications that can be utilized for the treatment of such ailments. Medication repositioning is the way toward recognizing new clinical applications for existing synthetic compounds. It is a basic system for the procedure called drug revelation. Drug developers are becoming increasingly innovative in discovering new functions for currently available chemicals or drugs. The process of identifying these new capabilities for currently available chemicals is called drug repositioning. Selecting the main restorative impacts for a chemical to perform clinically potential experiments is the major challenge lying in the repositioning of drugs. Blind screening and serial testing of creature models are some of the traditional repositioning efforts. Accumulation of biomedical information suggests that computational methodologies can be a viable option in anticipating the most favorable drug-disease associations for new indication experiments for currently available drugs.

## II. MATERIALS AND METHOD

Disease-disease similitudes and drug-drug similitudes are key segments in different medication ailment affiliation expectation models. Pairwise topological likenesses among sickness or medication are typically assessed using cosine similitude or Gaussian connection profile piece comparability subject to topological attributes of medications or ailments. Be that as it may, expectation prediction models utilizing these similitudes are not powerful enough.

### A. Materials

A database which is publically accessible known as comparative toxicogenomics database (CTD) [1] is utilized for putting away tentatively demonstrated drug-disease connections. The medication substructures are gathered from PubChem [2] dataset. Correspondingly compound database DrugBank [3] gives synthetic targets, medicate proteins, sedate medication joint efforts. ID transformation table from CTD Database is utilized to outline between various databases.

➢ *Disease-Disease Similarity*

Utilizing the MeSH database, each malady could be spoken to as directed acyclic graph. The MeSH dataset [4] gives a technique to ailment order and can help in the investigations of the connection between maladies. In the DAG of every sickness, every nodes speaks to infections and connections speak to node connections. Only a solitary kind of relationship exists to interface a tyke and parent hub. The relationship that exist between them are characterized utilizing 'is-a' relationship. To characterize the area of an illness in the MeSH chart, every infection have no less than one location in the DAG, called as codes. The parent hub's code is added behind the youngster's delivers and are utilized to characterize the codes of the tyke hub. For instance, there exists two conceivable locations for breast neoplasms(C04.588.180 and C17.800.090.500). The corresponding parent nodes are C17.800.090 breast diseases and C04.588 neoplasms by site as shown in Fig.1.
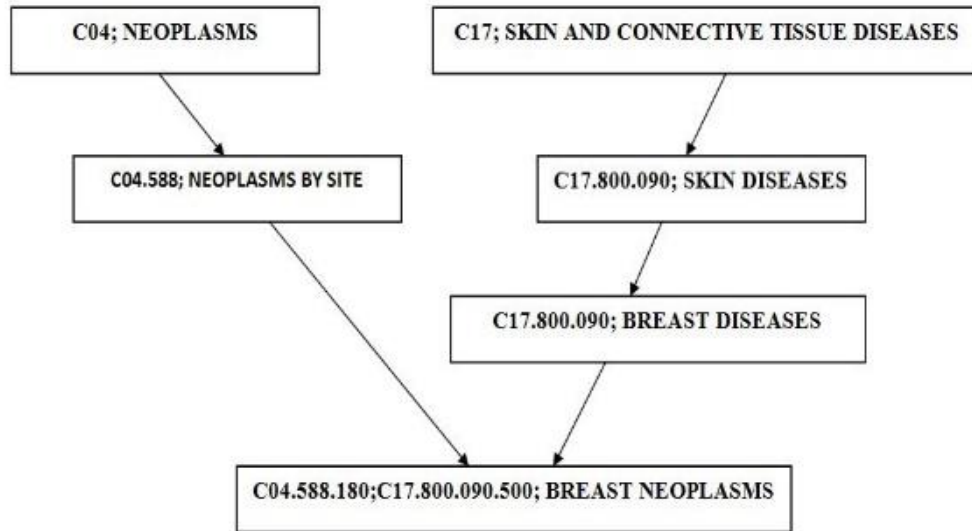
Fig 1:- DAG of disease "Breast Neoplasms"

A malady 'A' can be spoken to as a directed acyclic graph, $DAG_A$ = (A, $T_A$, $E_A$), where $T_A$ speaks to the arrangement of each precursor nodes of 'A' including node 'A' itself and $E_A$ is the arrangement of every comparing connection. The semantic contribution value speaking to the connection of an infection t to illness 'A' is spoken to by $D_A(t)$, which can be determined as pursues:

$$D_A(t) = \begin{cases} 1, & if\ t = A \\ \{0.5 \times D_A(t')|t' \in children\ of\ t\}, & if\ t \neq A \end{cases} \quad (1)$$

In the coordinated non-cyclic diagram of 'A', sickness 'An' is the most explicit ailment and hence its commitment to its own semantic esteem is taken as one though its progenitor hubs which are found more remote from hub 'An' are progressively broad groups. So these precursor hubs contribute substantially less to the semantic estimation of hub 'A'. So utilizing Equation (1), sickness A's semantic esteem can be determined as:

$$SemVal(A) = \sum_{t \in T_A} D_A(t) \quad (2)$$

By applying the basic principle that the diseases sharing larger part of their directed acyclic graphs tend to possess a higher semantic similarity [5], the similarity between different diseases can be calculated. The relative locations of two diseases in the MeSH disease directed acyclic graph are used for the calculation of the semantic similarity between them. The similarity value between any two diseases is calculated as follows:

$$S^d(A, B) = \frac{\sum_{t \in T_A \cap T_B} D_A(t) + D_B(t)}{SemVal(A) + SemVal(B)} \quad (3)$$

where $D_A(t)$ and $D_B(t)$ are the semantic values representing relationship of disease 't' to disease 'A' and disease 'B' respectively. Equation (3) is used to calculate the semantic similarity value between two diseases using the locations of those diseases in directed acyclic graphs and their relationships with their corresponding ancestor diseases.

➤ *Drug-Drug Similarity*

The fundamental standard utilized in the estimation of drug-drug likeness is comparative medications must be equipped for treating a comparable group of stars of maladies and must be comparable in component of activity. Drug-drug comparability can be useful in finding approved medication repositioning openings. Expanding on these triumphs, a medication sedate similitude estimation strategy for registering drug repositioning can be actualized like the MISIM technique which is utilized on account of microRNAs [6]. The commitment of comparative sicknesses which are related with the two synthetic substances are utilized to process the likeness esteem between two medications (drugs). The illnesses related with every compound are recovered from the CTC dataset. The semantic likeness esteem that exists between a malady and an illness bunch is determined utilizing equation (4). To calculate the comparability between two drugs or chemicals $c_1$ and $c_2$, assume $D_1$ represents the related diseases of $c_1$ and $D_2$ represents the related diseases of $c_2$. $D_1$ and $D_2$ have m and n distinct diseases respectively. While computing the similarity value between two drugs, it is necessary to consider each and every diseases in both $D_1$ and $D_2$. Hence the similarity of two drugs is computed using:

$$S^c(c_1, c_2) = \frac{\sum_{1 \leq i \leq m} S(d_{1i}, D_2) + \sum_{1 \leq i \leq n} S(d_{2j}, D_1)}{m+n} \quad (4)$$

For the better comprehension of cooperation between medications it is important to develop a solid medication useful system. Consequently a medication useful system is developed by fixing a limit. The pairwise similitude coefficients for a rundown of medications are figured. At that point a limit an incentive for similitude coefficients is chosen. Now drug pairs with similarity coefficient which is greater than or equal to the threshold value will possess a direct link between them. And finally a functional network for the drugs are constructed.

➢ *Heterogeneous Network Construction*

The drug-drug closeness system and disease-disease closeness system which are constructed utilizing the way toward fixing a limit for the comparability esteems in the past areas goes about as the subnets for the calculation of the heterogeneous net-works. After the computation of drug-drug similitude and disease-disease similitude, the subsequent stage is to join them utilizing the realized relationship to frame bipartite graph. The graphs is drug-disease heterogeneous network. The two subnets in particular, the ailment similitude system and medication comparability organize, could now be associated utilizing a tentatively demonstrated medication malady collaborations to shape a heterogeneous system of medications and ailments. CTD dataset is utilized to get the tentatively confirmed associations between drugs and diseases.

Suppose that $S^d = S^d(i,j)_{i=1,j=1}^{n,n}$ and $S^c = S^c(i,j)_{i=1,j=1}^{l,l}$ are adjacency matrices representing the disease similarity network and drug similarity network respectively. Similarly $Z = Z(i,j)_{i=1,j=1}^{l,n}$ represents the drug-disease interaction network, where n and l are the numbers of disease and drug entities. An example of the drug-disease heterogeneous network is given in Fig.2.
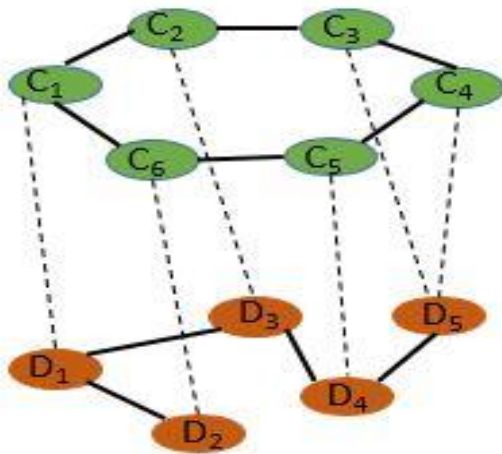


Fig 2:- Drug-Disease Heterogeneous Network

*B. Unknown Association Prediction*

➢ *Drug-Disease Association*

Let $C = C_1, C_2, …, C_l$ denote a set of drugs and $D = D_1, D_2, …, D_n$ denote a set of diseases. The aim is unobserved drug-disease association to prediction by making use of known associations. As given in the matrix in Fig.3 drugs(chemicals), diseases and their associations could be denoted as an adjacency matrix Z.ie., $Z_{ij} = 1$ if drug $C_i$ is associated with $D_j$; otherwise, $Z_{ij} = 0$. Based on this matrix, the interaction profiles for diseases and drugs are introduced. As shown in the matrix of Fig.3, the row vectors and column vectors of the adjacency matrix Z represents the drug interaction profiles and disease interaction profiles respectively.

| | $D_1$ | $D_2$ | $D_3$ | | $D_n$ |
|---|---|---|---|---|---|
| $C_1$ | 1 | 0 | 0 | … | 0 |
| $C_2$ | 1 | 1 | 0 | … | 0 |
| $C_3$ | 0 | 1 | 0 | … | 0 |
| … | … | … | … | … | … |
| $C_l$ | 0 | 0 | 1 | … | 1 |

Fig 3:- Interaction Profiles

Most of the non interactions or 0's in matrices Z are unknown cases which can actually be true interactions.i.e. they can be false negative values. The non interacting drug-disease pairs in Z are actually missing edges. Hence a matrix updation process called Known Nearest Neighbours(KNN) could be utilized to calculate the interaction likelihood score for these non interacting pairs.ie., the procedure tries to replace the values of the $Z_{ij}$ equals 0, by a value from 0 to 1 using the algorithm given below:

➢ *Known Nearest Neighbour Algorithm*
Input: Adjacency Matrices $Z \in R^{l \times n}$, $S^c \in R^{l \times l}$ and $S^d \in R^{m \times m}$ & decay term 'T'.
Output: Modified matrices Z.
Algorithm:
1: For p=1 to l do
1.1: cn = knownNeighbour(p, $S^c$)
1.2: k = length(cn)
1.3: For i=1 to K do
1.3.1: $w_i = T^{i-1} S^c$ (p,$cn_i$)
1.4: end for
1.5: $Q_c = \sum_{i=1}^{k} S^c(p, cn_i)$
1.6: $Z_c(p) = (\sum_{i=1}^{k} w_i Z(cn_i))/Q_c$
2: end for
3: For q=1 to m do
3.1: dn = knownNeighbour(q, $S^d$ )
3.2: k = length(dn)
3.3: For j=1 to K do

3.3.1: $w_j = T^{j-1} S^d(\text{p}, dn_j)$
3.4: end for
3.5: $Q_d = \sum_{i=1}^{k} S^d(p, dn_j)$
3.6: $Z_d(p) = (\sum_{j=1}^{k} w_j Z(dn_j))/Q_d$
4: end for
5: $Z_{cd} = (Z_c + Z_d)/2$
6: $Z = \max(Z, Z_{cd})$
7: return Z

The refreshed Z framework goes about as new association profile,ie., closest neighbor cooperation profile for medications and ailments. They are utilized to build the expectation display. The framework Z has its lines as medications and sections as infections and the qualities speaks to sedate sickness associations.ie., it indicates how likely is the medication in that line would be related with the maladies in every segment. A large number of the zero qualities in the framework (shaped utilizing known affiliations) have been supplanted with qualities inside the range 0 and 1,ie., more infections will be related with every malady.

## III. EXPERIMENTS AND RESULTS

To assess the execution of the proposed technique in anticipating drug-sickness affiliations, the methodology got forget one cross approval was embraced in our analyses. At first we figured the different illnesses that are related with every medication. At that point to check the prescient exhibitions of our technique we connected forget one cross validation.ie., for every medication in the dataset, it was considered as a test sedate once and its affiliation data was erased. The rest of the medications were taken as the preparation dataset. The anticipated signs for the test sedate were positioned by the last outcome got after KNN. That is we registered the expectation effectiveness of our framework by contrasting them and tentatively demonstrated affiliations. For every particular positioning edge, if the heaviness of a gathered medication sickness affiliation was over the edge, it was viewed as a genuine positive. Else, it was viewed as a bogus positive. True positive rate (TPR) and false positive rate (FPR) were determined by fluctuating edge esteems. The forecast capacity of the strategy was spoken to utilizing the area under curve (AUC) esteem. The forget one cross approval methodology was utilized to test the forecast execution of this framework and trial results show that when forget one medication cross approvals were actualized, a normal AUC estimation of 0.68 was gotten when KNN calculation was connected. These outcomes demonstrated that dependable medication illness affiliation expectation results could be accomplished by our strategy.

## IV. CONCLUSION

The way toward repositioning the accessible synthetic compounds for new purposes could be useful to both pharmaceutical associations similarly as individuals. The forecasts of medication illness relationship utilizing computational techniques are gainful ways to deal with give synthetic substances with the most sign contender for further biomedical tests. Henceforth, determining compelling calculations to assemble quiet contamination affiliations is of high centrality and limitless undertakings are being made to this field. In our work, we use a method named Known Nearest Neighbour algorithm to foresee unobserved drug-disease associations by utilizing known drug-disease affiliations. Initially we calculated drug-drug and disease-disease similarities and formed two adjacency matrices. KNN algorithm is used to make the matrices more suited for predictions. Known Nearest Neighbour algorithm is used to consolidate these outcomes to foresee unknown associations. In the computational trials, our method can produce good performances.

## REFERENCES

[1]. P. Davis et al., Text mining effectively scores and ranks the literature for improving chemical-gene-disease curation at the comparative toxicogenomics database, PLoS ONE, vol. 8, no. 4, p. e58201, 2013.

[2]. Y. Wang et al., PubChem: A public information system for analyzing bioactivities of small molecules, Nucleic Acids Res., vol. 37, pp. W623W633, Jul. 2009.

[3]. Knox et al., DrugBank 3.0: A comprehensive resource for Omics research on drugs, Nucleic Acids Res., vol. 39, pp. D1035D1041, Jan. 2010.

[4]. I.Lee,U.M.Blom,P.I.Wang and J.E.Shim, Prioritizing candidate disease genes by networkbased boosting of genome-wide association data, Genome Res,vol. 21, pp.1109-1121.

[5]. L.Cheng,J.Li,P.Ju,J.Peng,Y.Wang, SemFunSim:a new method for measuring disease similarity by integrating semantic and gene functional association, PLoS One,vol. 9,no. 6,2014.

[6]. D.Wang, J.Wang,M.Lu,F.Song and Q.Cui, "Inferring the human miRNA functional similarity and functional network based on miRNA-associated diseases,Bionformatics,vol. 26,pp . 1644-1650,2010.