# RANKING AND RELATED ASPECTS IN METABOLIC NETWORKS PROFILE

Enrolment No. - 111506

Name of Student  -  Pryanka Sharma

Name of  Supervisor - Mr. Suman Saha

Dr. Sree Krishna Chanumolu



May 2015

Submitted in partial fulfilment for the requirement of

Bachelor of Technology

DEPARTMENT OF BIOTECHNOLOGY AND

BIOINFORMATICS

Jaypee University of Information Technology

P.O.Waknaghat-173234

Himachal Pradesh (INDIA)

# **DECLARATION**

I hereby declare that the project titled "**Ranking and related aspects in metabolic networks profile***"* is submitted as a Project Work has been carried out by me at Jaypee University of Information Technology , Solan under the guidance of  Mr. suman saha and Dr. sree krishna chanumolu**.** Any further extension, continuation or use of this project has to be undertaken with prior express written consent from the Supervisor, Jaypee University of Information Technology, Solan-173234.


 I further declare that the project work or any part thereof has not been previously submitted for any degree or diploma in any university.


Signature:                                                                                Name:


Date:

# **Abstract**

Complex biological systems may be represented and analyzed as computable networks. Nodes and edges are the basic components of a network. Nodes represent units in the network, while edges represent the interactions between the units. Nodes can represent a wide-array of biological units, from individual organisms to individual neurons in the brain. Two important properties of a network are degree and betweenness. Degree is the number of edges that connect a node, while betweenness is a measure of how central a node is in a network. Nodes with high betweenness essentially serve as bridges between different portions of the network (i.e. interactions must pass through this node to reach other portions of the network). In social networks, nodes with high degree or high betweenness may play important roles in the overall composition of a network. Here metabolic networks were taken from KEGG database and were acted upon by the ranking code generated. The results were displayed and graphical representations were obtained that simplified the process of analyzing, interpreting, visualizing was resolved. The graphical representations demonstrated the node priority and thereby eased the work and effort required to analyze the network. The ranking algorithm used here counted the priority and thus displayed the important steps in metabolic networks and thereby metabolism process.

# **<u>ACKNOWLEDGEMENT</u>**

I acknowledge my sincere thanks to **Mr. Suman Saha** and **Dr. Sree Krishna Chanumolu** , *Jaypee University of Information Technology, JUIT, Solan* for giving me this great opportunity to have project work under him  and throughout this project, enlightening me on various topics and creating a congenial environment for my work. I am sure it would continue to raise research interests in young students like me.

I am grateful to the whole department of Bioinformatics at Jaypee University of Information Technology, JUIT, for extending their help throughout this project and supported me in this Project.

Date: ……….                                                                            Signature:

# TABLE OF CONTENTS

- **Introduction**
  - **Complex networks and brief introduction**
  - **Biological networks**
  - **Metabolic networks and its section profile**
  - **Various ranking algorithms**
  - **Associated aspects**

- **Objective**

- **Tools and Methodologies**

  - **Data Retrieval**

  - **Data Collection**

  - **Ranking algorithm used**

  - **Implementations**

  - **Associated work**

**Results & discussions**
  - **Ranking result in metabolic network**
  - **Plots**
  - **Adjacency matrix**
  - **Interpretations, analysis, visualizations.**
  - **Related progress**

- **References**

  - **Journals References**

  - **Database References**

# 1.0 **Introduction**

## 1.1 Complex Networks

Complex networks are loosely defined as networks with non trivial topology and dynamics, which appear as the skeleton of complex systems in the real world. Understanding complex systems often requires a bottom-up approach, breaking the system into small and elementary constituents and mapping out the interactions between these components. In many cases, the myriads of components and interactions are best characterized as networks.Networks, in general, are constituted by a set of objects and by a set of interconnections among these objects. The suitability of networks to represent many real world systems has given an impressive spur to the recent research area of complex networks. Collaboration networks, the Internet, the world-wide-web, biological networks, communication and transport networks, social networks are just some examples. An interesting property to investigate, typical to many networks, is the ranking structure that prioritizes the nodes in a network based on the no of incoming or the outgoing edges and hence the connections  as mentioned earlier. The capability of detecting the partitioning of a network in clusters can give important information and useful insights to understand how the structure of ties affects individuals and their relationships.

 The identification of high order structures in networks unveils insights into their functional organization. Complex networks can be categorized as follows:

- Social networks
- Biological networks
- Telecommunications networks
- Information networks
- Technological networks etc.

Complex networks possess many distinctive properties, of which ranking and community structure is one of the most studied. The community structure is usually as the division of networks into subsets of vertices within which intra-connections are dense while between

which inter-connections are sparse. Identifying the community structure is very helpful to obtain some important information about the relationship and interaction among nodes. Once extracted, such clusters of nodes are often interpreted as organizational units in social networks, functional units in biochemical Networks, ecological niches in food web networks, or scientific Disciplines in citation and collaboration networks. And the way to this is by the implementation of certain efficient algorithms and by the means of use of computer informatics.

The main features of complex networks are:
1. Many interacting units
2. Self organization
3. Small world
4. Scalefree heterogeneity
5. Dynamical evolution

The  first approach to capture the global properties of such systems is to model them as graphs whose nodes represents the dynamical units and whose links stand for the interactions between them .one has to cope with the structural issues such as characterization of topology of complex wiring architecture, revealing the unifying principles that at the basis of real networks , and developing models to mimic the growth of a network and reproduce its structural properties . and on the other hand , many relevant questions arise when studying complex networks dynamics such as learning how a large ensemble of dynamical systems that interact through a complex wiring topology can behave collectively. The major concepts and results recently achieved in the study of the structure and dynamics of complex networks and summarize the relevant applications of these ideas in many different disciplines ranging from nonlinear science to biology , from statistical mechanics to medicine and engineering.
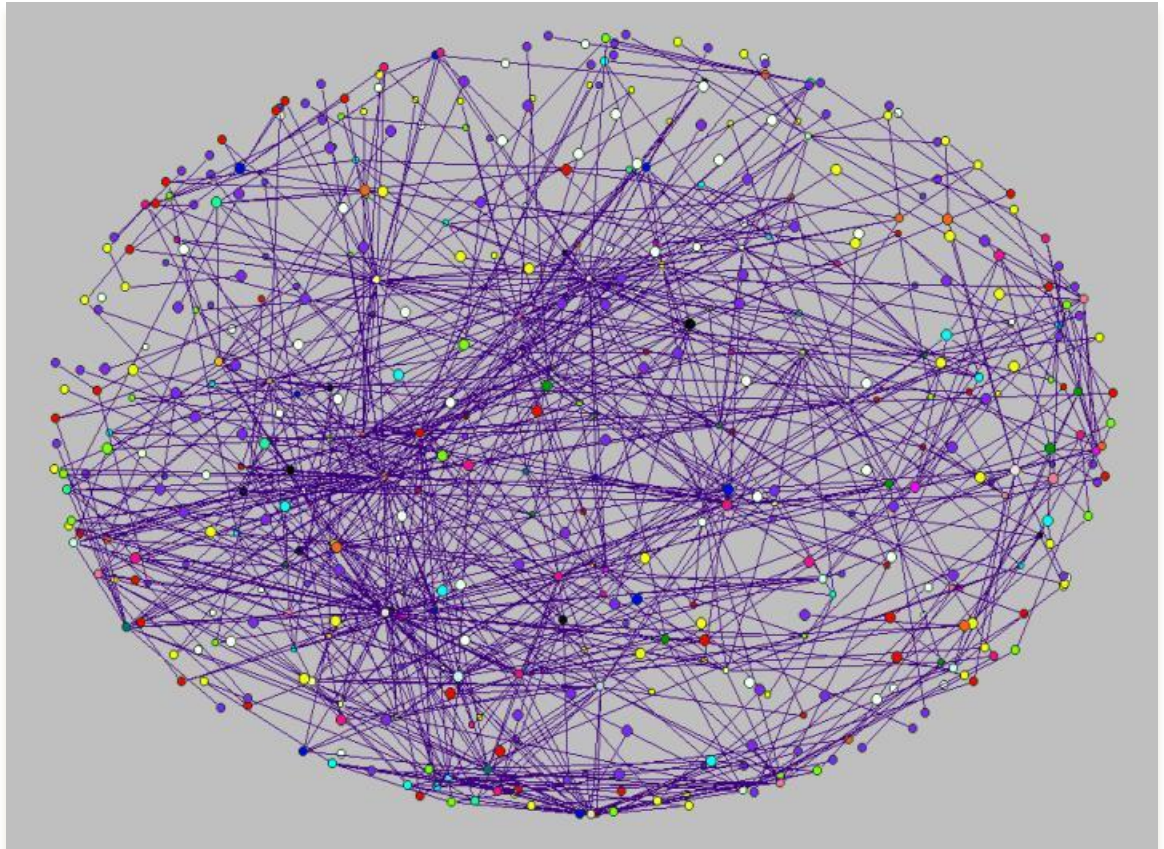
Fig: complex networks

## 1.2 Biological Networks

Biological processes are often represented in the form of networks such as protein-protein interaction networks and metabolic pathways. In biological systems networks emerge in many disguises, from food webs in ecology to various biochemical nets in molecular biology. In particular, the wide range of interactions between genes, proteins and metabolites in a cell are best represented by various complex networks. During the last decade, genomics has produced an incredible quantity of molecular interaction data, contributing to maps of specific cellular networks. The emerging fields of transcriptomics and proteomics have the potential to join the already extensive data sources provided by the genome wide analysis of gene expression at the mRNA and protein level. The study of biological networks, their modeling, analysis, and visualization are important tasks in life science today. An understanding of these networks is essential to make biological sense of much of the complex data

that is now being generated. This increasing importance of biological networks is also evidenced by the rapid increase in publications about network-related topics and the growing number of research groups dealing with this area. Most biological networks are still far from being complete and they are usually difficult to interpret due to the complexity of the relationships and the peculiarities of the data. Network visualization is a fundamental method that helps scientists in understanding biological networks and in uncovering important properties of the underlying biochemical processes.
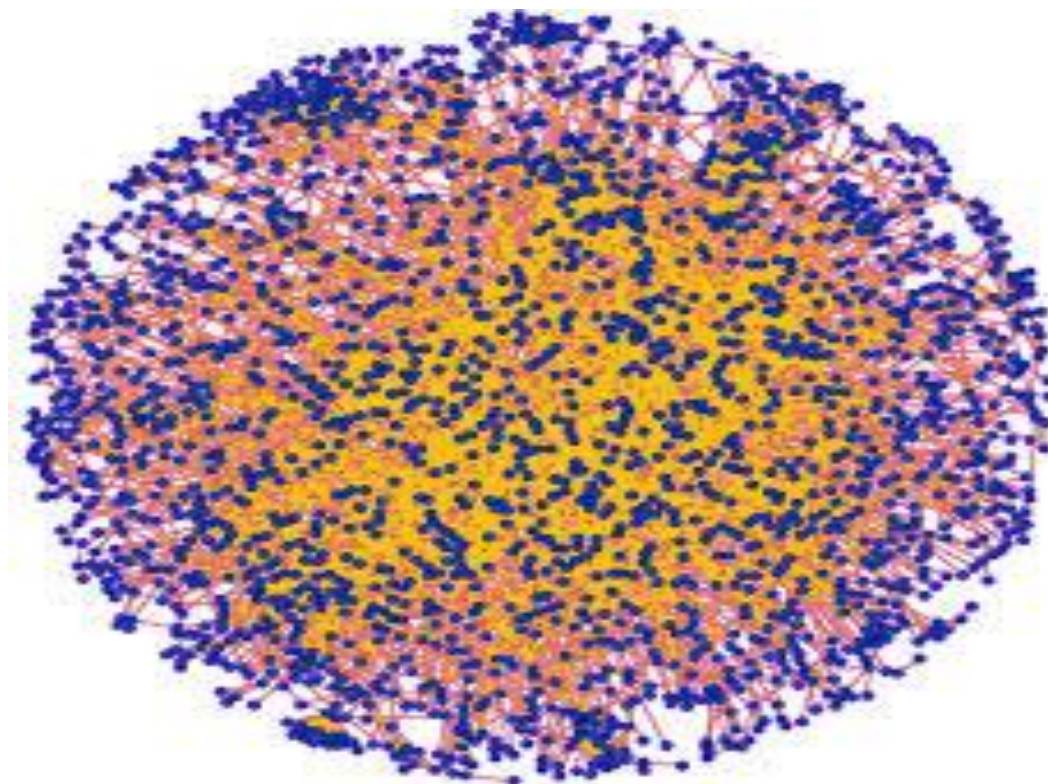


**Fig: Yeast protein protein interaction network**

Several highly important biological networks are related to molecules such as DNA, RNA, proteins and metabolites and to interactions between them. Gene regulatory and signal transduction networks describe how genes can be activated or repressed and therefore which proteins are produced in a cell at a particular time. Such regulation can be caused by regulatory proteins or external signals. Protein-protein interaction networks represent the interaction between proteins such as the building of protein complexes and the activation of one protein by another protein. Metabolic networks show how

metabolites are transformed, for example to produce energy or synthesize specific substances. We here consider phylogenetic trees, special networks or hierarchies which are often built on information from molecular biology such as DNA or protein sequences. Phylogenetic trees represent the ancestral relationships between different species. They are used to study evolution, which describes and explains the history of species, i.e., their origins, how they change, survive, or become extinct. Finally, signal transduction, gene regulatory, protein-protein interaction and metabolic networks interact with each other and build a complex network of interactions; furthermore these networks are not universal but species specific, i.e., the same network differs between different species. Metabolic networks have been studied for a long time in biology and biochemistry, and specific visualization requirements are given, e.g., by established drawing styles. We present some algorithmic extensions of the hierarchical layout approach which aim to fulfill these requirements.

Networks offer us a new way to categorize systems of very different origin under a single framework. This approach has uncovered unexpected similarities between the organization of various complex systems, indicating that the networks describing them are governed by generic organization principles and mechanisms. Understanding the driving forces which invest different networks with similar topological features enables systems biology to combine the numerous details about molecular interactions into a single framework, offering means to address the structure of the cell as a whole.

## 1.3 Metabolic Networks as profile:

Metabolic reactions are fundamental to life processes, e.g., for the production of energy and the synthesis of substances. A huge number of reactions occur at any time in living cells and the product of one reaction is usually used by another reaction, thus metabolic reactions are strongly interco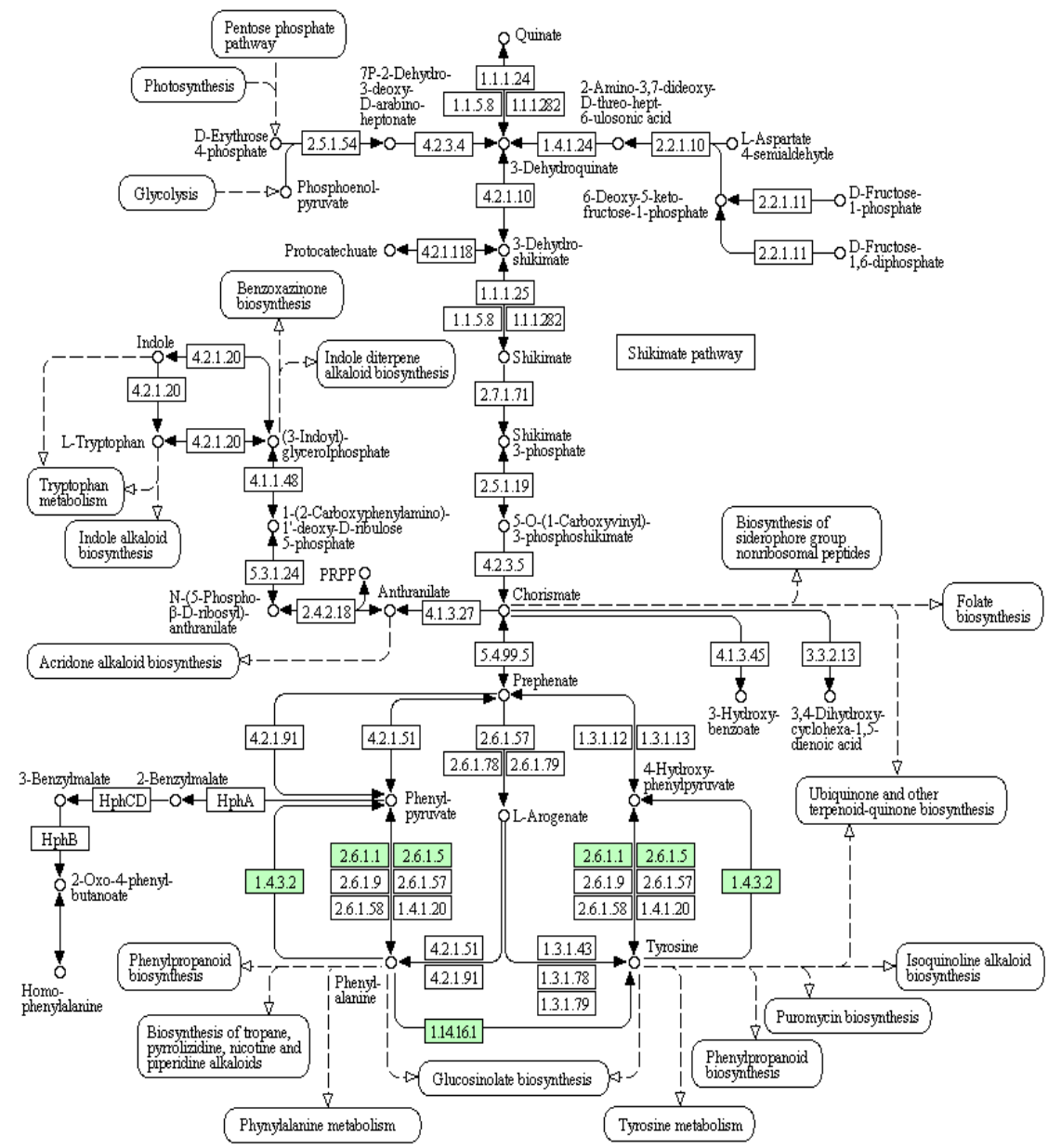nnected and form metabolic pathways and networks. Metabolic reactions are fundamental to life processes, e.g., for the production of energy and the synthesis of substances. A huge number of reactions occur at any time in living

cells and the product of one reaction is usually used by another reaction, thus metabolic reactions are strongly interconnected and form metabolic pathways and networks. . Furthermore, many reactions are regulated, i.e., they are suppressed or enhanced by other factors (allosteric control). This shifts the steady state and together with the steady supply of substances from outside and their final use, e.g., by exporting them from the cell, one can consider a main direction of a reaction. This is also expressed by the differentiation of substances into reactants and products. The metabolic network or metabolism of a particular cell or an organism is the complete network of metabolic reactions of this cell or organism. A metabolic pathway is a connected sub-network of the metabolic network either representing specific processes or defined by functional boundaries, e.g., the network between an initial and a final substance.

Since genes and proteins tend to function through interacting in networks, the interaction networks must be analyzed for the purpose of studying biological functions. It is very important to analyze the biological functions in terms of the network of interacting molecules and genes, such as genetic regulatory network, signal transduction network, protein interaction network, and metabolic network, with the aim of understanding how a biological system is organized from its individual building blocks. It represents an integration of biological knowledge from genomic information towards the understanding of the basic principles of life for biomedical applications. The analysis of metabolic networks can help to understand and utilize cellular metabolic process in order to promote the development of ferment technology and medicine industry. On the other hand, the topology of metabolic networks reflects the dynamics of their formation and evolution. A study of this realm may help to understand the evolutionary history of life.

From a formal point of view a metabolic pathway is a hyper-graph. The nodes represent the substances and the hyper-edges represent the reactions. A hyper-edge connects all substances of a reaction, is directed from reactants to products and is labeled with the enzymes that catalyze the reaction. Hyper-graphs can be represented by bipartite graphs. Additionally to the nodes representing substances, the reactions are nodes (either labeled with the enzymes or with further nodes for enzymes) and edges are binary relations connecting the substances of a reaction with the corresponding reaction node.

Fig: Metabolic network

There are several networks which are closely related to metabolic pathways or networks:

• Simplified metabolic network: A network which contains reactions, enzymes and main substances, but no co-substances.

• Metabolite network and simplified metabolite network: A network which consists only of substances (metabolites); in the simplified case only of main substances.

• Enzyme network: A network which consists only of the enzymes catalyzing the reactions.

Three different things that would be required for efficient and productive information retrieval are:

**1. Parts of reactions**: The display of substances and enzymes is application and user-specific. Usually for main substances their name, structural formula or both should be shown. Co-substances should be displayed using their name or abbreviation and enzymes should be represented by their name or EC-number [Int92].

**2. Reactions**: The reaction arrow(s) should be shown from the reactants to the products with enzymes placed on one side of the reaction arrow and co-substances on the opposite side. The reversibility of a specific reaction should be clearly visible. For co-substances their temporal order, which depends on the reaction mechanism, is important, and they should be placed according to this order.

**3. Pathways**: The main direction of reactions (e.g., from top to bottom) should be clearly visible to express the temporal order of reactions. There are important exceptions to the main direction used for the visualization of specific pathways, e.g., the citrate acid cycle or the fatty acid synthesis. The structure of these cyclic reaction chains should be emphasized. Such pathways are characterized by the continuous repetition of a reaction sequence in which the product of the sequence re-enters in the next loop as a reactant.

Now here the metabolic network file was taken from the KEGG database and were thoroughly studied. Based on the information already available from this database a profile was generated that incorporated the data for various metabolic networks ;now this

data was acted upon by the algorithm that performed ranking based on the score evaluation and hence the priority of various linkages was observed then analyzed and interpreted . The most important nodes would get a number according to the frequency of that particular linkage and hence the priorities get classified and a plot is generated that depicts the nodes and linkages getting different priorities and hence makes it simpler for the biologists to interpret and analyze the important connections and relations between various nodes respectively. Then this database could be used for community detection and thereby the clusters would be formed that will depict the functional and relevant hubs in a particular metabolic networks.

## 1.4 Various ranking algorithms:

There have been from a long time a lot of algorithms that does ranking in a variety of complex networks .but the ranking algorithm chosen here is purely a score function that counts the priorities and enables one to filter the chances of redundancy and effective information retrieval. So we get a score that would prioritize the nodes and the respective linkages; in this case the plot is generated that depicted the nodes with higher priorities and hence makes the work easy for researchers. A fundamental problem in the field of social network analysis is to rank individuals in a society according to their implicit .importance, derived from a network's underlying topology. More precisely, given a social network, the goal is to produce a (cardinal) ranking, whereby each individual is assigned a nonnegative real value, from which an ordinal ranking (an ordering of the individuals) can be extracted if desired. In this paper, we propose a solution to this problem specially geared toward social networks that possess an accompanying hierarchical structure.

A social network is typically encoded in a link graph, with individuals represented
by vertices and relationships represented by directed edges, or .links, annotated with
Weights. Given a link graph, there are multiple ways to assign meaning to the weights.
On one hand, one can view the weight on a link from i to j as expressing the distance
from i to j.a quantity inversely related to j's importance. On the other hand, one can view

each weight as the level of endorsement, or respect, i grant j.a quantity directly proportional to j's importance.

Ranking of nodes according to their centrality, or importance, in a complex network such as the Internet, the World Wide Web, and other social and biological networks, has been a hot research topic for several years in physics, mathematics, and computer science. For a comprehensive overview of the vast literature on rankings in networks, and more recently to for a thorough up-to-date mathematical classification of centrality measures.

**PageRank algorithm** is developed by Brin and Page during their Ph. D at Stanford University based on the citation analysis. PageRank algorithm is used by the famous search engine that is Google. This algorithm is the most commonly used algorithm for ranking the various pages. Working of the PageRank algorithm depends upon link structure of the web pages. The PageRank algorithm is based on the concepts that if a page contains important links towards it then the links of this page towards the other page are also to be considered as important pages. The PageRank considers the back link in deciding the rank score. If the addition of the all the ranks of the back links is large then the page then it is provided a large rank. Therefore, PageRank provides a more advanced way to compute the importance or relevance of a web page than simply counting the number of pages that are linking to it. If a back link comes from an important page, then that back link is given a higher weighting than those back links comes from non-important pages. In a simple way, link from one page to another page may be considered as a vote. However, not only the number of votes a page receives is considered important, but the importance or the relevance of the ones that cast these votes as well.

**Weighted PageRank Algorithm** is proposed by Wenpu Xing and Ali Ghorbani. Weighted PageRank algorithm (WPR) is the modification of the original PageRank algorithm. WPR decides the rank score based on the popularity of the pages by taking into consideration the importance of both the in links and out links of the pages. This algorithm provides high value of rank to the more popular pages and does not equally divide the rank of a page among it's out link pages. Every out-link page is given a rank value based on its popularity. Popularity of a page is decided by observing its number of

in links and out links. As suggested, the performance of WPR is to be tested by using different websites and future work include to calculate the rank score by utilizing more than one level of reference page list and increasing the number of human user to classify the web pages.

**A distance rank algorithm** is proposed by Ali Mohammad Zareh Bidoki and Nasser Yazdani.This intelligent ranking algorithm based on reinforcement learning algorithm based on novel recursive method. In this algorithm, the distance between pages is considered as a distance factor to compute rank of web pages in search engine. The main goal of this ranking algorithm is computed on the basis of the shortest logarithmic distance between two pages and ranked according to them so that a page with smaller distance to assigned a higher rank. The Advantage of this algorithm is that, being less sensitive, it can find pages faster with high quality and more quickly with the use of distance based solution as compared to other algorithms. If the some algorithms provide quality output then that has some certain limitations. So the limitation for this algorithm is that the crawler should perform a large calculation to calculate the distance vector, if new page is inserted between the two pages. This Distance Rank algorithm adopts the PageRank properties i.e. the rank of each page is computed as the weighted sum of ranks of all incoming pages to that particular page. Then, a page has a high rank value if it has more incoming links on a page.

A typical search engine should use web page ranking techniques based on the specific needs of the users because the ranking algorithms provide a definite rank to resultant web pages. After going through this exhaustive analysis of algorithms for ranking of web pages against the various parameters such as methodology, input parameters, relevancy of results and importance of the results, it is concluded that existing algorithms have limitations in terms of time response, accuracy of results, importance of the results and relevancy of results.

Here in case of metabolic networks a ranking algorithm is generated that takes the essence from the algorithms mentioned above and does all the computing, about which the protocol and the algorithm are explained in the later stages.

## 1.5 Associated aspects:

**Ranking** is the primary step followed in case of metabolic networks or any other network requiring priority classification and preference evaluation. Then after ranking comes the part of **community detection** that wasn't achieved in due course of time but can then be easily implemented once done with the ranking.

Community detection is nothing but finding out the hidden clusters and hubs that are functional in a network and in case of metabolic networks the important functional hubs in the respective network that forms clusters and hence give important piece of information aspired.

**The community detection** in complex networks is an important problem in many scientific fields, from biology to sociology. This paper proposes a new algorithm; Differential Evolution based Community Detection (DECD), which employs a novel optimization algorithm, differential evolution (DE) for detecting communities in complex networks. DE uses network modularity as the fitness function to search for an optimal partition of a network. Based on the standard DE crossover operator, we design a modified binomial crossover to effectively transmit some important information about the community structure in evolution. Moreover, a biased initialization process and a clean-up operation are employed in DECD to improve the quality of individuals in the population. One of the distinct merits of DECD is that, unlike many other community detection algorithms, DECD does not require any prior knowledge about the community Structure, which is particularly useful for its application to real-world complex networks where prior knowledge is usually not available. We evaluate DECD on several artificial and real-world social and biological networks. Experimental results show that DECD has very competitive performance compared with other state-of-the-art community detection algorithms.

The community structure is usually considered as the division of networks into subsets of vertices within which intra-connections are dense while between which inter-connections are sparse. Identifying the community structure is very helpful to obtain some important information about the relationship and interaction among nodes. To detect the underlying community structure in complex networks, many successful algorithms have been proposed so far. However, the community detection in networks is a nondeterministic polynomial (NP) hard problem. Most of current community detection algorithms based on greedy algorithms perform poorly on large complex networks. Moreover, many algorithms for community detection also require some prior knowledge about the community structure, e.g., the number of the communities, which is very difficult to be obtained in real-world networks.

**Girvan and Newman** proposed the Girvan-Newman (GN) algorithm which is one of the most known algorithms proposed so far. This algorithm is a divisive method and iteratively removes the edges with the greatest betweenness value based on betweenness centrality. Newman presented an agglomerative hierarchical clustering method based on the greedy optimization of the network modularity. This method iteratively joins communities of nodes in pairs and chooses the join with the greatest increase in the network modularity at each step.

## 2.0 <u>Objective:</u>
## Ranking and related aspects in metabolic networks profile.

Ongoing researches globally, are producing huge data related to biological networks and here in this project, metabolic network with respect to the goal of efficient and useful information retrieval on daily basis in various forms, individually at different places by different scientists and scholars. In this Project, I aim to put my efforts to gather all the information that would make work easier for the biologists and also that is what clearly defines my area of interest.

## 3.0 <u>Tools and Methodologies:</u>

This involves the ranking of nodes in a network for biological networks the ranking implies as to the frequency of occurrence of a node or an established relationship in a network that periodically repeats itself and then mathematical approach applied to it and a plot is generated .This plot is between the score and the node number. The highest score would reflect the highest frequency or high connectivity thereby disclosing the probability of having some significance in the network. This adds enormously to the biological knowledge or stands for some biological significance if any. It follows a series of steps that accomplishes this which are explained in the sub parts.

Gone through databases like KEGG (Kyoto encyclopedia of genes and genome), which gave the XML format for metabolic networks, and detailed explanations through flowcharts that almost explained the whole process stepwise. And then these were complied on to a file from where various metabolic networks were implemented on to the ranking algorithm and desired results were obtained. And the front end is the graphics user interface that had ranking as a tool and the embedded code running behind the front end.

## 3.1 Data Retrieval:

Data was retrieved from the online server site **www.genome.jp/kegg/** and a whole detailed set of information were obtained. These were later collected and complied in a file and then were implemented onto by the codes made that solved the problem of ranking for the networks taken into consideration.

## 3.2 Data Collection

As mentioned above the data was obtained from the KEGG database and compiled onto a file that precisely had all the information required to run the code accordingly revealing

the important results. The GUI was created using HTML/CSS, php, javascript that had embedded php and R codes running at the back that implemented ranking in metabolic networks.

## 3.3 Ranking Algorithm Used:

Here as such no particular ranking algorithm was used instead an algorithm was generated that independently computed the ranking based on the score function that counted the frequency of the highly connected node and a plot was generated that depicted the respective ranks for each set of connectivity. The scores were generated by the adjacency matrix that and then were plotted.

## 3.4 Implementation:

This section of the report covers the protocol followed and and the implementation side of the project. The implementation is as follows:

3.4.1 The first step of the protocol that enabled me to retrieve and collect the data. The data was collected from the KEGG database and a file was generated that incorporated the XML versions of the metabolic networks file.

3.4.2 Then a GUI was created that gave the front end and had ranking as a tool in it.

3.4.3 After the GUI construction the metabolic networks file of desired choice was chosen and was implemented upon by the ranking code generated.

3.4.4 ranking code was generated that implemented the score function based on the number of connections that a particular node had and based on the maximum number of connections the priorities were decided.

3.4.5 The scoring function first identified the number of nodes in the networks file based on the entity id and displayed the total number of nodes in the file uploaded.

3.4.6 Based upon the number of nodes counted an adjacency matrix was made that automatically computed the node connections and calculated the score for individual nodes and in the end column the aggregated scores were displayed.

3.4.7 Column wise scores were added for the individual nodes and hence an aggregated score was obtained.

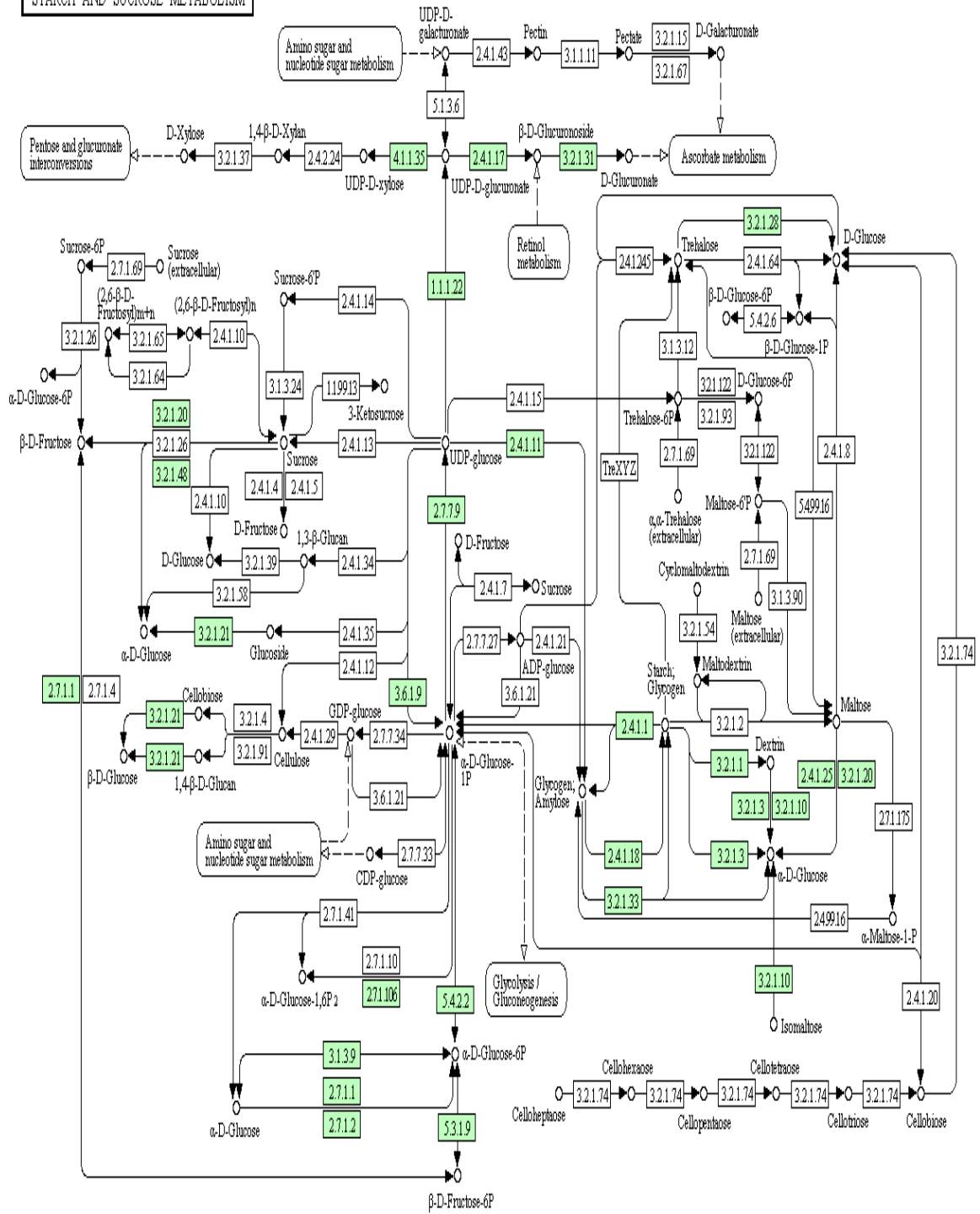3.4.8The adjacency matrix was created that had N×N rows and columns for n number of nodes.

3.4.9 Now the resulting scores were stored in a document file and were made to run at the background R code that displayed the desired plots signifying the effective and important information. The plot depicted the node score on the one end and the node number on the other end.

The step wise implementation encompasses the following:

```
<?xml version="1.0"?>
<!DOCTYPE pathway SYSTEM "http://www.kegg.jp/kegg/xml/KGML_v0.7.1_.dtd">
<!-- Creation date: Oct 20, 2014 18:02:37 +0900 (GMT+09:00) -->
- <pathway title="Histidine metabolism" link="http://www.kegg.jp/kegg-bin/show_pathway?hsa00340" image="http://www.kegg.jp/kegg/pathway/hsa/hsa00340.png"
  number="00340" org="hsa" name="path:hsa00340">
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?K00492" name="ko:K00492" reaction="rn:R00069" type="ortholog" id="37">
      <graphics name="K00492" type="rectangle" height="17" width="46" y="487" x="785" bgcolor="#FFFFFF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?K05603" name="ko:K05603" reaction="rn:R02286" type="ortholog" id="38">
      <graphics name="K05603" type="rectangle" height="17" width="46" y="560" x="885" bgcolor="#FFFFFF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?K01458" name="ko:K01458" reaction="rn:R00525" type="ortholog" id="39">
      <graphics name="K01458" type="rectangle" height="17" width="46" y="599" x="944" bgcolor="#FFFFFF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?K01479" name="ko:K01479" reaction="rn:R02285" type="ortholog" id="40">
      <graphics name="K01479" type="rectangle" height="17" width="46" y="609" x="806" bgcolor="#FFFFFF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?hsa:10841" name="hsa:10841" reaction="rn:R02287" type="gene" id="41">
      <graphics name="FTCD, LCHC1" type="rectangle" height="17" width="46" y="590" x="849" bgcolor="#BFFFBF" fgcolor="#000000"/>
    </entry>                                                              Histidine metabolism
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?hsa:144193" name="hsa:144193" reaction="rn:R02288" type="gene" id="42">
      <graphics name="AMDHD1" type="rectangle" height="17" width="46" y="518" x="829" bgcolor="#BFFFBF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?hsa:131669" name="hsa:131669" reaction="rn:R02914" type="gene" id="43">
      <graphics name="UROC1, HMFN0320" type="rectangle" height="17" width="46" y="445" x="829" bgcolor="#BFFFBF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?hsa:3034" name="hsa:3034" reaction="rn:R01168" type="gene" id="44">
      <graphics name="HAL, HIS, HSTD" type="rectangle" height="17" width="46" y="373" x="829" bgcolor="#BFFFBF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?hsa:55748" name="hsa:55748" reaction="rn:R01166" type="gene" id="47">
      <graphics name="CNDP2, CN2, CPGL, HEL-S-13, HsT2298, PEPA" type="rectangle" height="17" width="46" y="297" x="659" bgcolor="#BFFFBF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?hsa:57571" name="hsa:57571" reaction="rn:R01164" type="gene" id="48">
      <graphics name="CARNS1, ATPGD1" type="rectangle" height="17" width="46" y="267" x="758" bgcolor="#BFFFBF" fgcolor="#000000"/>
    </entry>
  - <entry link="http://www.kegg.jp/dbget-bin/www_bget?hsa:1644" name="hsa:1644" reaction="rn:R01167" type="gene" id="49">
      <graphics name="DDC, AADC" type="rectangle" height="17" width="46" y="351" x="657" bgcolor="#BFFFBF" fgcolor="#000000"/>
```

**Fig: XML file of metabolic network**

HISTIDINE METABOLISM

PRPP  2.4.2.17  Phosphoribosyl-ATP  3.6.1.31  Phosphoribosyl-AMP  3.5.4.19  Phosphoribosyl-formimino-AICAR-P  5.3.1.16  Phosphoribulosyl-formimino-AICAR-P  HisF HisH  4.2.1.19  Imidazole-acetol-P  2.6.1.9  L-Histidinol-P  3.1.3.15  L-Histidinol

Pentose phosphate pathway

AICAR

Imidazole-glycerol-3P

Purine metabolism

1.1.1.23

1-Methyl-L-histidine  L-Histidinal

Anserine  3.4.13.5  6.3.2.11  2.1.1.22  2.1.1.-  1.1.1.23

Carnosine  6.3.2.11

4-(β-Acetylaminoethyl)-imidazole

N-Formyl-L-aspartate  3.5.3.5  Imidazolone acetate  3.5.2.-  Imidazole-4-acetate  1.14.13.5  Imidazole acetaldehyde  1.2.1.3  2.3.1.-  1.4.3.22  3.4.13.18  3.4.13.20  4.1.1.22  Hercynine  Thiourocanic acid

3.5.1.15  3.5.1.8  N-Formimino-L-aspartate  6.3.4.8  Histamine  4.1.1.28  L-Histidine  Ergothioneine

Aspartate  2.1.1.8  4.3.1.3

1-(5-Phosphoribosyl)-imidazole-4-acetate  N-Methylhistamine  2.6.1.38  Urocanate  1.3.99.33  Dihydrourocanate

Alanine, aspartate and glutamate metabolism  1.4.3.4  4.2.1.49

(1-Ribosylimidazole)-4-acetate  Methylimidazole acetaldehyde  Hydantoin-5-propionate  1.14.13.-  4-Imidazolone-5-propanoate  4-Oxoglutaramate  2-Oxoglutarate

1.2.1.5  3.5.2.7  Formylisoglutamine  Isoglutamine

Methylimidazole-acetic acid  N-Formimino-L-glutamate  3.5.3.13  N-Formyl-L-glutamate

Imidazole-pyruvate  2.1.2.5  3.5.1.68

3.5.3.8

Imidazole-lactate  N-Carbamyl-L-glutamate  L-Glutamate  Alanine, aspartate and glutamate metabolism

00340 10/20/14
(c) Kanehisa Laboratories

**Fig: Diagrammatic representation of metabolic network**

STARCH AND SUCROSE METABOLISM

00500 11/28/14
(c) Kanehisa Laboratories

## The code for adjacency matrix:

```php
<?php
    $file="D:/r/hsa00340.xml";
    $f=fopen($file,"r");
    $l=fgets($f);
    $l=fgets($f);
    $l=fgets($f);
    $l=fgets($f);
    $l=fgets($f);
    $l=fgets($f);
    $l=fgets($f);
    $l=fgets($f);
    $data=explode("'",$l);
    $start=$data[1];
    $l3=array();
    $i=0;
    while(!strstr($l,"<relation "))
    {
        $l=fgets($f);
        $l3[$i++]=$l;
        if($i==7)
        $i=0;
    }
    $data=explode("'",$l3[geti($i)]);
    $end=$data[1];
    $matrix=array(array());
    for($i=$start;$i<=$end;$i++)
    for($j=$start;$j<=$end;$j++)
    $matrix[$i][$j]=0;
```

```
$data=explode("",$l);
$i=$data[1];
$j=$data[3];
$matrix[$i][$j]=1;
while($l=fgets($f))
{
        if(strstr($l,"<relation entry1"))
        {
                $data=explode("",$l);
                $i=$data[1];
                $j=$data[3];
                $matrix[$i][$j]=1;
        }
}
fclose($f);
$f=fopen("D:/r/result.txt","w");
$f2=fopen("D:/r/adjacencymatrix.txt","w");
fwrite($f,"node score\n");
for($i=$start;$i<=$end;$i++)
{
        for($j=$start;$j<=$end;$j++)
        {
                echo $matrix[$i][$j].' ';
                fwrite($f2,$matrix[$i][$j]." ");
        }
        echo "\n";
        fwrite($f2,"\n");
}
fclose($f2);
echo "Total Nodes: ".($end-$start+1)."\n";
for($i=$start;$i<=$end;$i++)
```

```php
        {
                $score=0;
                for($j=$start;$j<=$end;$j++)
                {
                        //echo $matrix[$i][$j];
                        if($matrix[$i][$j]==1)
                        {
                                $score++;
                                echo "Relation found between ".$i." and ".$j."\n";
                        }
                }
                if($score!=0)
                fwrite($f,$i." ".$score."\n");
                //echo "\n";
        }
        fclose($f);
function geti($i)
{
        if($i<6)
        return $i+1;
        else
        return 0;
}
/*
test <- read.table("d:/r/result.txt",header=T,fill=T)
plot(score ~ node,data=test,pch=16)
*/
?>
```

**Code in R for plot generation and ranking:**

```
k <- read.table("d:/r/result.txt", header = T, fill = T)
plot(k, data=k, pch=16)
```

The ranking plots were obtained and are shown in the results section.

## 3.5 Associated work:

This section includes the related work and serves as an extension to the project. Due to less time the other goal of community detection couldn't be achieved fully but the research part is complete just the implementation part is yet to be achieved. Whole research exhibited the different ways and algorithms that could implement the community detection for the respective metabolic profile. Also this part of the project required a lot of research as there are many ways to accomplish this task. There were quite a couple of algorithms used for clustering and gave efficient results. Also in RStudios package there is an inbuilt library that makes communities and hence sorts the problem of community detection .

Also on the contrary one can compare various networks at the same time that would unveil the required set of information effectively and efficiently.

## 4.0 <u>Results and Discussions:</u>

The results for data retrieval and the related step wise progress in achieving the goal of ranking are shown below:
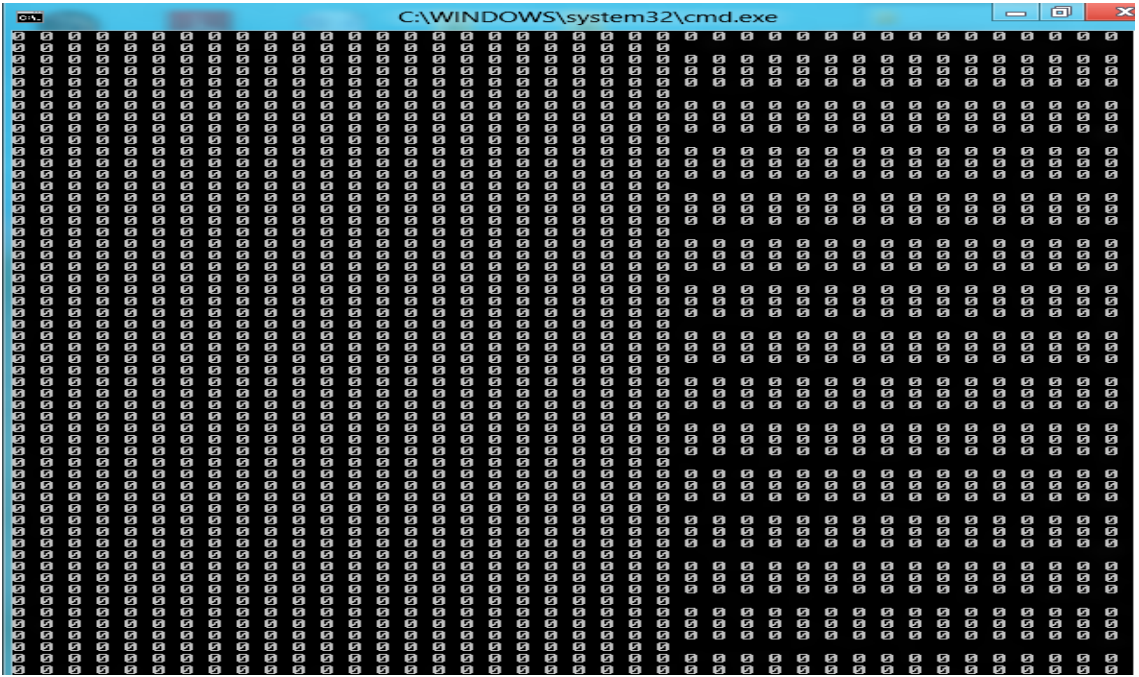


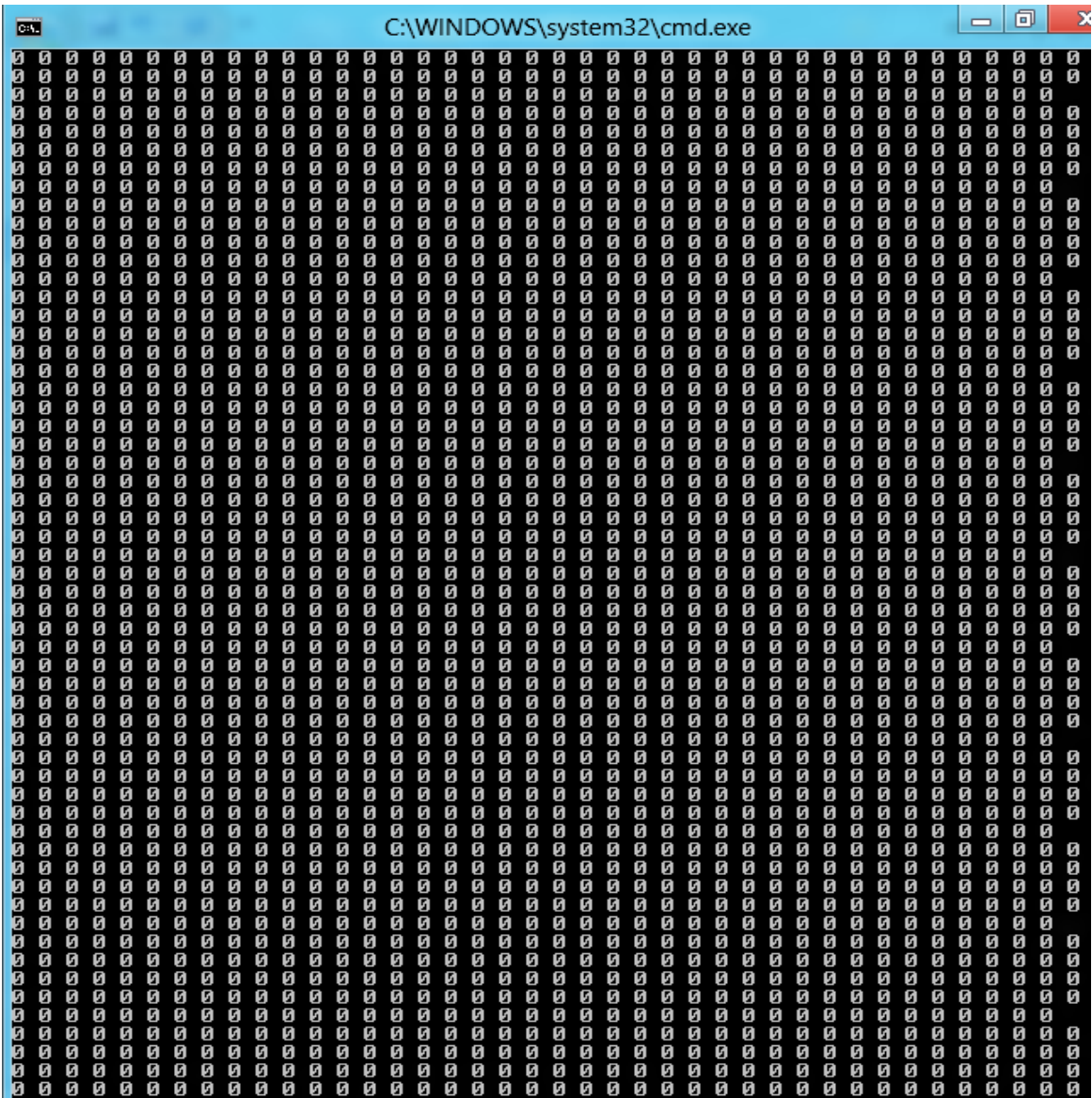Fig: **Adjacency matrix for histadine metabolism**

Fig : Matrix for hsa00400

Fig : Relations between nodes

Fig: Matrix for hsa00360

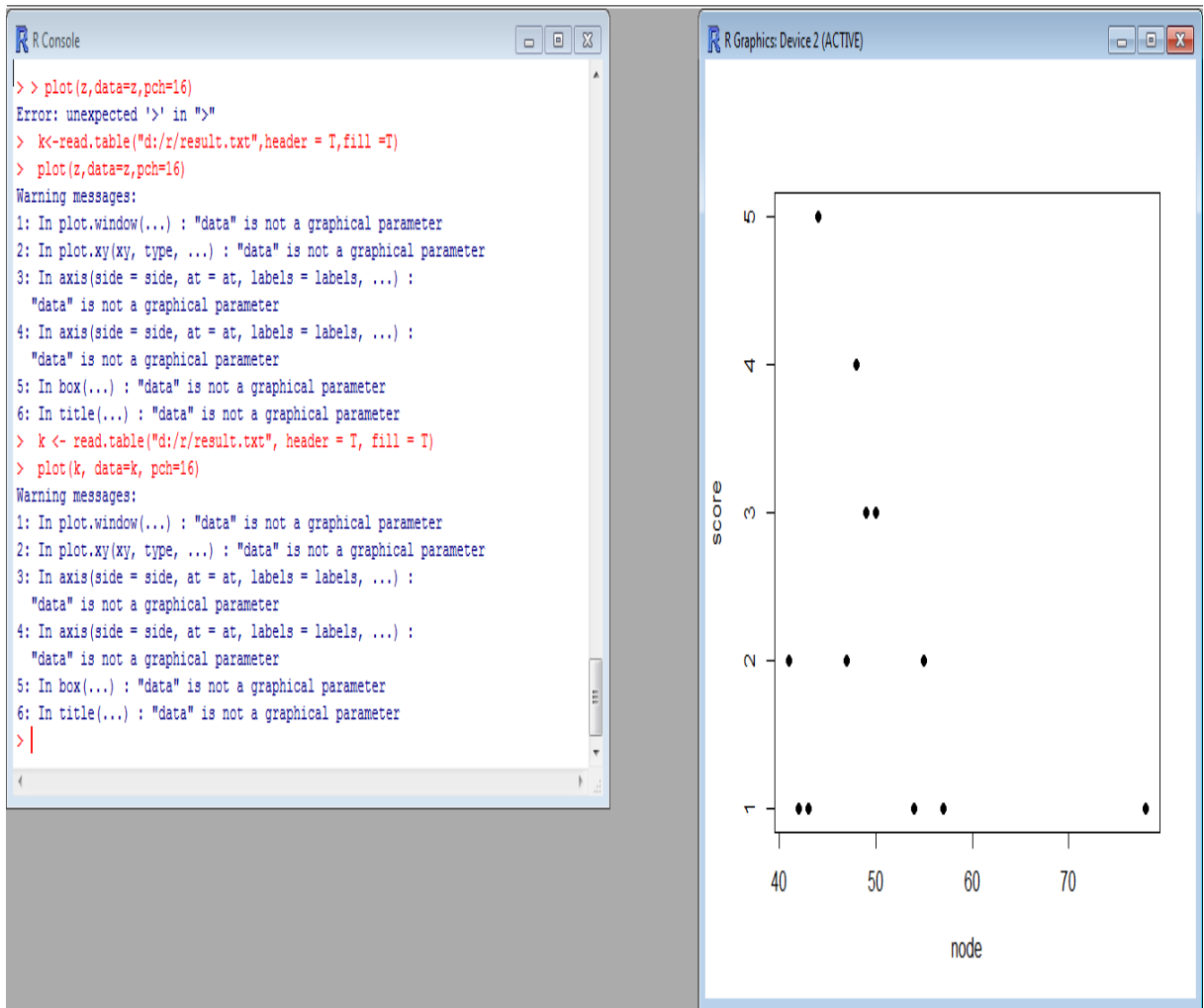Fig: Relations between nodes

Fig: Relation in sucrose metabolism

Fig: ranking result for hsa00340
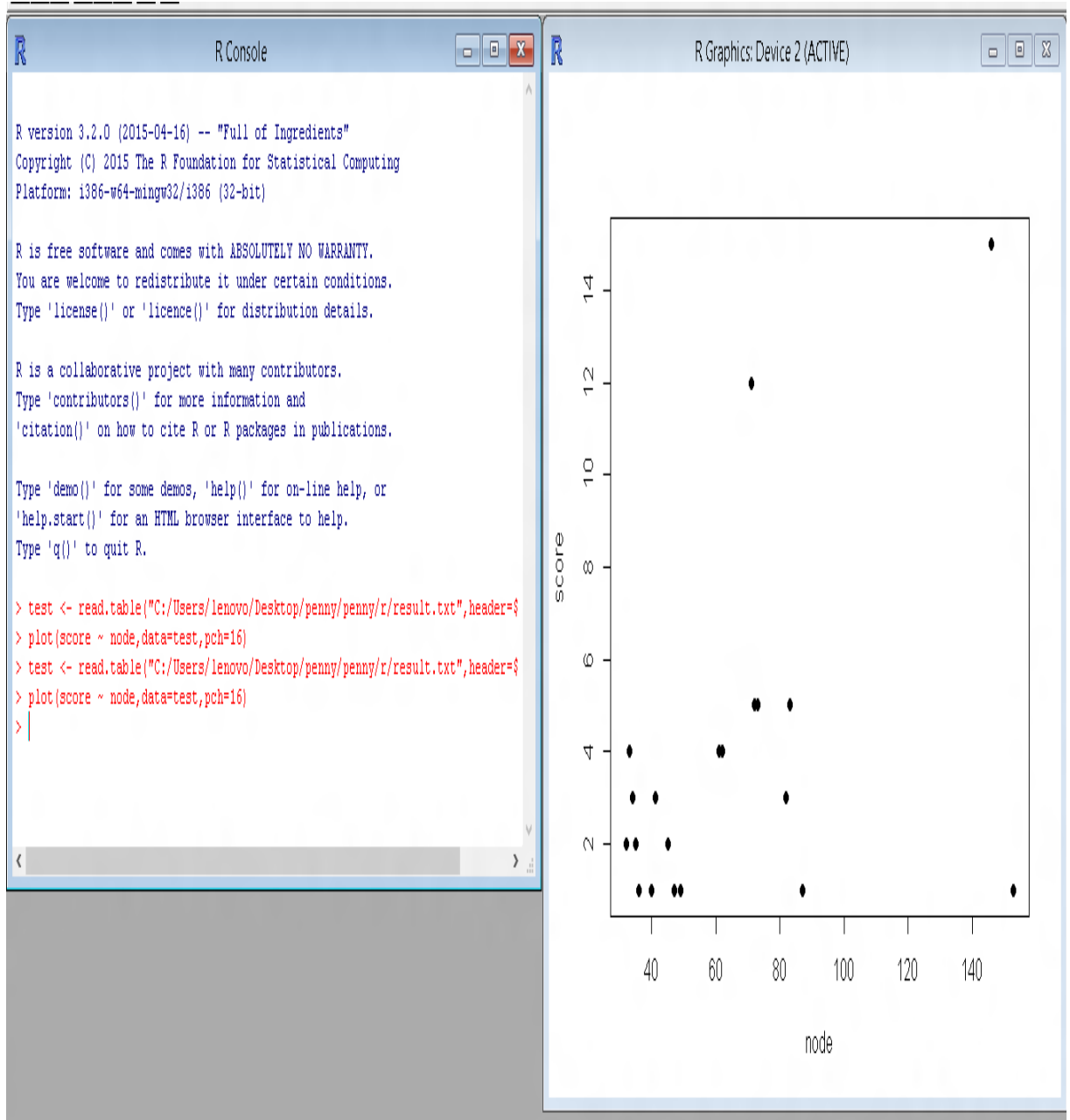
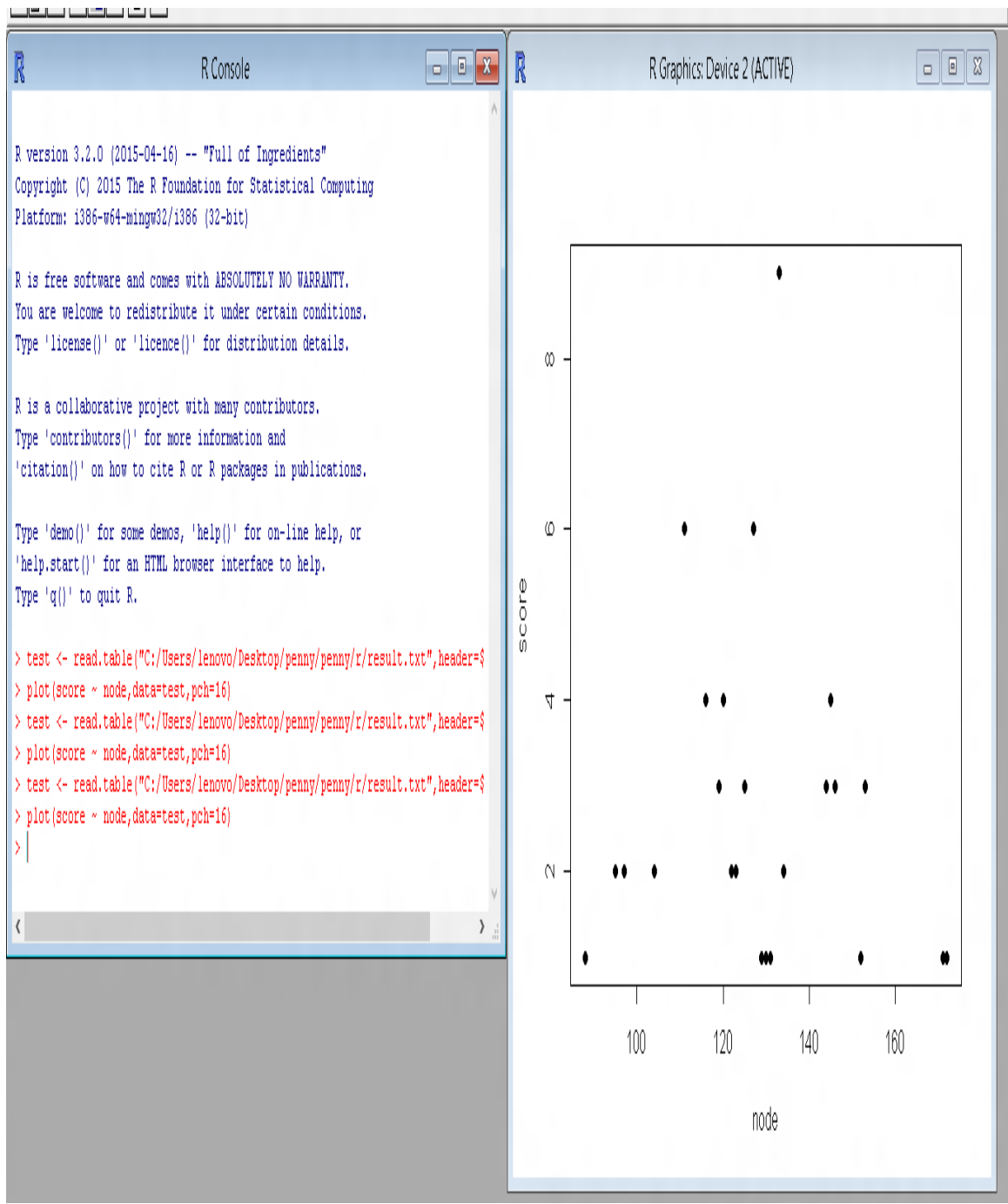Fig: ranking result for hsa00360

Fig: ranking result for hsa00400

Fig: ranking for sucrose metabolism

# CONCLUSION

The above followed protocol gave us a brief idea about how ranking can be implemented in complex networks or biologically called as metabolic networks. Due to the high complexity of the networks, the existing algorithms cannot be applied straightaway to biological systems especially metabolic networks. They require special algorithms maintaining the integrity of the data. This algorithm is designed to do so keeping in mind the high efficiency, accuracy and cost effectiveness.

This algorithm gave appropriate ranking results, thereby creating a platform to explore more in this field and prove beneficial to the biologists. This can firther be extended to different area like community detection and can be easily implemented in "R" programming language that have in built community algorithms.

# 5.0 References

## 5.1 Journal References

5.1.1Ma, Jian James; Zeng, Daniel; and Huff, Richard A. "Complex Network Analysis," *Journal of International Technology and Information Management, 2013.*

5.1.2Da Fontura Costa, L., Oliveira Jr., O.N., Travieso, G., Rodrigues, r.A., Villas Boas, P.R., Antiqueira, L., Viana, M.P., da Rocha, L.E.C.: Analyzing and Modeling Real-World Phenomena with Complex Networks: A Survey of Applications. arXiv physics.soc-ph, 0711.3199 ,2008.

5.1.3Maron, M.E. and J.L. Kuhn's. On Relevance, Probabilistic Indexing and Information Retrieval. JACM,7, 1960.

5.1.4Robertson, S.E. The Probability Ranking Principle in IR. Journal of Documentation, 33:4, 1977.

5.1.5Neelam Duhan, A. K. Sharma and Komal Kumar Bhatia, "Page Ranking Algorithms: A Survey", In proceedings of the IEEE International Advanced Computing Conference (IACC), 2009.

5.1.6Newman, M.E.J., Girvan, M.: Finding and Evaluating Community Structure in Networks. Phys Rev E 69, 026113 ,2004.

5.1.7Girvan, M., Newman, M.E.J.: Community Structure in Social and Biological Networks. PNAS 99, 7821-7826, 2002.

5.1.8Stanley Wasserman and Katherine Faust. *Social Network Analysis: Methods and Applications. Cambridge University* Press, Cambridge, 1994.

5.1.9R. Andersen, F. Chung, and K. Lang. Local graph partitioning using PageRank vectors. In *FOCS '06: Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science, pages 475–486, 2006*

5.1.10M. Gaertler. Clustering. In U. Brandes and T. Erlebach, editors, *Network Analysis: Methodological Foundations,* pages 178–215. Springer, 2005.

5.1.11R. Guimerà, M. Sales- Pardo , and L. Amaral. Modularity from fluctuations in random graphs and complex networks. *Physical Review E, 70:025101, 2004.*

5.1.12Weston, J., Eliseeff, A., Zhou, D., Leslie, C., & Noble, W. S. Protein ranking: from local to global structure in the protein similarity network. Proceedings of the National Academy of Science, 101, 6559–6563 (2004).

5.1.13da Fontura Costa, L., Oliveira Jr., O.N., Travieso, G., Rodrigues, r.A., Villas Boas, P.R., Antiqueira, L., Viana, M.P., da Rocha, L.E.C.: Analyzing and Modeling Real-World Phenomena with Complex Networks: A Survey of Applications. arXiv physics.soc-ph, 0711.3199 (2008).

5.1.14Guimerà, R., Amaral, L.A.N.: Functional Cartography of Complex Metabolic Networks. Nature 433, 895-900 (2005).

5.1.15R. Andersen, F. Chung, and K. Lang. Local graph partitioning using PageRank vectors. In *Proceedings of FOCS2006*, pages 475–486, 2006.

5.1.16Stanley Wasserman and Katherine Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, Cambridge, 1994.

## 5.2Database references:

- **http://www.ncbi.nlm.nih.gov/**
- **www.genome.jp/kegg/**