

Study of Codon Degeneracy Based on Similarity Measure

Pankaj Hazarika¹, Pranjal Kumar Bora^{2,*}, Arun Kumar Baruah³, Papori Bora⁴

¹Department of Mathematics, Madhabdev University, Narayanpur, Assam, INDIA.

²Centre for Computer Science and Applications, Dibrugarh University, Dibrugarh, Assam, INDIA.

³Department of Mathematics, Dibrugarh University, Dibrugarh, Assam, INDIA.

⁴Department of Food Science and Nutrition, Assam Agriculture University, Jorhat, Assam, INDIA.

Submission Date: 25-06-2022; Revision Date: 16-07-2022; Accepted Date: 04-08-2022.

ABSTRACT

In genetics, Codon degeneracy is a salient feature which refers to a single amino acid being encoded by more than one codon. According to a study, Degeneracy of genetic code helps an organism to prosper on earth. Each amino acid is encoded by triplet codes of four possible (Guanine (G), Adenine (A), Cytosine (C), or Thymine (T/U)) bases. The genetic codon degeneracy occurs mainly due to the variance in third position e.g. the amino acids Glycine is encoded by four codons GGU, GGC, GGA, GGG differ only in third base. Taking part of more than one tri nucleotides sequence out of 64 triplets to encode one amino acids lead to the concept of Degeneracy. In this manuscript we formulate a new classifying technique with the help of cosine similarities to explain the degeneracy. Further we have done a comparison of our method with an existing classification technique. The consequences of our results open a new paradigm to study the genetics from a new mathematical perspective. The disassortative nature of codons networks may help us to understand the flow of genetic information in the evolution process of amino acids.

Keywords: Degeneracy, Similarity Measures, Impression.

Correspondence:

Dr. Pranjal Kumar Bora,
Centre for Computer
Science and Applications,
Dibrugarh University,
Dibrugarh-786004,
Assam, INDIA.

Email: pranjaly2k@gmail.com

INTRODUCTION

In molecular biology, transcription and translation are two processes through which transformation of genetic information from DNA to mRNA takes place. DNA has four nucleotide bases namely Adenine (A), Thymine (T), Guanine (G) and Cytosine (C), whereas RNA has A, G, C and Uracil (U) instead of T. A combination of three adjoined nucleotides (triplet) of four nucleotide bases form an amino acid.^[1] For the first time, George Gamow named a set of three adjoined nucleotides (triplet) as codon. 20 different types of amino acids which are formed by the different combinations of 61 codons play the pivotal role in protein synthesis in a living organism. Either minimum of 3 codons

or maximum of 6 codons form an amino acid. The mechanism of different codon codes for the same amino acid is called codon degeneracy. In broader sense, Degeneracy can be described as a property that elements of different structures carry out the same results. Different codons match to the same amino acid are recognized as synonymous codons. In this paper, we are trying to address the codon degeneracy using a mathematical tool.

Several attempts have been made to describe the codon degeneracy, few of them are^[1-7] etc. Since the inception of Genetic Code Table^[8] actual reason of degeneracy still remains a mystery. Sengupta *et al.*, classified 64 codons into three different classes namely transitional, weak and strong based on their strength in codon degeneracy.^[9] This classification was based on position and their nucleotides in a codon. In a study, Gonzalez *et al.*, proposed a new interpretation of degeneracy based only on symmetry principle.^[10] Unifying mathematical framework, they depicts degeneracy with integer number representation systems. Based on symmetrical properties

SCAN QR CODE TO VIEW ONLINE



www.ajbls.com

DOI: 10.5530/ajbls.2022.11.80

of codons T. Negadi^[11] has proposed degeneracy of genetic code into distinct multiplet structures. Revisiting the study of degeneracy of genetic codon Jayanta Kumar Das *et al.*^[12] introduced a new parameter “impression” to classify 20 amino acids. In their classification it has been noticed that amino acids having different degeneracy clubs into same cluster. Though in the recent few years study of genetic code gets new prospects incorporating mathematics yet a lot of avenues are still open to explore. The study of genetic code is still a subject of more or less rigorous exploration from mathematical standpoint. So, looking into these facts in our proposed work we define similarity measurement among the codons which give rise to clustering of amino acids. Further while comparing our results it is observed that the amino acid heaving same degeneracy gives the same similarity measurements. So, our clustering technique is more significant than the previous methods. Grouping different amino acids heaving same degeneracy paves the path for a potent analysis.

The rest of the paper is organized as follows. In Section 2, we discuss basics of degeneracy, Lagerkvist’s hypothesis, Impression and similarities measures. Section 3, includes the main results of our paper followed by a comparison of our method with impression. Section 4 incorporates the concluding remarks.

Preliminaries

In this section we recall the degeneracy property of amino acids followed by Lagerkvist’s hypothesis, Impression and similarities measures.

Degeneracy

The Standard Genetic Code table clearly reveals that out of 64 codons 61 codons plays the role of forming

20 amino acids and rest of the three codons does not have any significant role in forming amino acids. Which indicates multiple codons encoded for the same amino acids and this phenomenon is called the degeneracy.^[13] The redundant codons are called synonymous codons. For ready reference redundant codons are arranged in some specific groups as shown in Table 1. Based on numbers of codons involved in forming a corresponding amino acid, degeneracy of genetic codons is divided into distinct multiplet structure^[11] From the Table 1, it is observed that nine amino acids are corresponding to groups of two codons, called twofold degenerated. Five amino acids are corresponding to groups of four codons called fourfold degenerated, and three amino acids are corresponding to groups of six codons. One amino acid is corresponding to groups of by three codons, and only two amino acids are corresponding to single codons. For example, Alanine (A) and Valine (V) are encoded by four codons called fourfold degeneracy. It has been found that some organisms have been surviving more prosperously due to degeneracy of genetic codes. It may be due to the fact that while point mutation takes place if the mutation leads to a synonymous codon then effectiveness of the mutation will not affect the functionality of the protein.

Since the emergence of the degeneracy phenomena it has been challenging many biologists to find out the reasons behind this. Crick 1968 explained this phenomena and he suggested that first two codon positions were responsible in forming synonymous codons.^[8] It is worth mentioning that genetic codon redundancy has a significant role in the mutation process. Because while substituting one nucleotide to another in the degenerated codon position gives the same amino

Table 1: Standard Genetic Code table.

	U	C	A	G	
U	UUU } Phe UUC } (F) UUA } Leu UUG } (L)	UCU UCC } Ser UCA } (S) UCG	UAU } Tyr(Y) UAC } UAA } Stop UAG } codon	UGU } Cys(C) UGC } UGA } StopCodon UGG } Try(W)	U C A G
C	CUU CUC } Leu CUA } (L) CUG	CCU CCC } Pro CCA } (P) CCG	CAU } His(H) CAC } CAA } Gln(Q) CAG }	CGU CGC } Arg CGA } (R) CGG	U C A G
A	AUU AUC } Iso AUA } (I) AUG } Met(M)	ACU ACC } Thr ACA } (T) ACG	AAU } Asn(N) AAC } AAA } Lys(L) AAG }	AGU } Ser(S) AGC } AGA } Arg(R) AGG }	U C A G
G	GUU GUC } Val GUA } (V) GUG	GCU GCC } Ala GCA } (A) GCG	GAU } Asp(D) GAC } GAA } Glu(E) GAG }	GGU } GGC } Gly(G) GGA } GGG }	U C A G

acids. For example in fourfold degenerate codon point mutation is insignificant at the third position.

Lagerkvist's Hypothesis

Followed by Crick 1968,^[8] Lagerkvist 1978^[14,15] formulated a new hypothesis which is known as hypothesis "Two out of three". In this hypothesis Lagerkvist stated that:

- If the first two base pair of codon-anticodon binds with six hydrogen bonds then the third positioned codon base becomes insignificant.
- If first two base pairs of codon-anticodon binds with four hydrogen bonds then the characteristic of the third codon base i.e., purine or pyrimidine will have a determining role towards the forming of different amino acids.
- If two first base pairs binds with 5 hydrogen bonds then rule 1 and rule 2 is relevant when the second codon base is pyrimidine and purine respectively.

It is worth mentioning that Crick hypothesis^[8] was revolutionary in the field of genetics to study the degeneracy property of amino acids. Crick hypothesis suggested that Wobble position determines why multiple codons can encode a single amino acid. The Wobble Hypothesis proposed by Francis Crick, states that the 3rd base in an mRNA codon can pair with 1st base of a tRNA anticodon which is a non-Watson-Crick base pairing. But the Crick hypothesis has come into questions when it was significantly failed to describe the molecular mechanism of the G-U pair formation. Later integrating H-bonding perspective of first two codon nucleotides into Crick hypothesis Lagerkvist proposed a hypothesis "two out of three"; which is later known as Lagerkvist hypothesis. Later, Lagerkvist's rules were seen in the genetic code table suggested by Rumer.^[16]

Impression

Impression is a parameter to study the underlying theme of degeneracy. Introducing impression to study the underlying theme of degeneracy Das *et al.*,^[12,17] made a paradigm shift in the field of computational biology. Based on molecular weights on amino acids they have assigned rank of the amino acids as some ternary number and finally they have calculated impression factor for individual amino acids. This impression factor leads them to incorporate mathematics to define the degeneracy of amino acids. They classified twenty amino acids into four groups based on parameter impression. Impression of amino acids (*IP*) is calculated as -

$$IP \{X_1, X_2, X_3\} = \{I_1, I_2, I_3\} \quad (1)$$

Where, $X_1, X_2, X_3 \in \{0, 1, 2\}$ are ternary number and $I_1 = |X_1 - X_2|$, $I_2 = |X_2 - X_3|$ and $I_3 = |X_3 - X_1|$.

Similarity Measures

Similarity measures are used to determine how two data objects are alike. Generally it ranges between zero (0) to one (1). Similarity of two data object is defined as 1 if they are alike and 0 otherwise. Cosine similarity is one of the most broadly used similarity measure. Cosine similarity of two vectors A and B is defined as -

$$\begin{aligned} \sin(A, B) &= \frac{AB}{\|A\| \cdot \|B\|} \\ &= \frac{\sum_i^n A_i B_i}{\sqrt{\sum_i^n A_i^2} \sqrt{\sum_i^n B_i^2}} \end{aligned} \quad (2)$$

Where,

$$A = \{A_1, A_2 \dots A_N\} \in R^n; \quad B = \{B_1, B_2 \dots B_N\} \in R^n$$

Cosine similarity is not affected by the magnitude of vectors. This is why cosine similarity is more preferable than the Euclidian distance while evaluating relationship among various features is the same. Cosine similarity is more affected by the orientation of vectors than the exact position of the vector.

A New Interpretation of Genetic Codon Degeneracy

Over the few years several researches have been carried out to describe the degeneracy of the genetic code. But still the development in this field is in the neonatal stage. The work that has been carried out here is a part of incorporating mathematical perspective to describe the degeneracy. In an approach to unifying mathematical framework Bora *et al.*,^[18] has assigned some weight on each base positions of a codon depending on positional impact, appearance of pyrimidine and purine in the second and third base position and Hydrogen bonding influence of codon bases as given by Lagerkvist's hypothesis to define a distance among the amino acids. For ready reference, consider the following example. Details of the method of giving weight on each base position of codons are discussed by Bora *et al.*, 2020.^[18]

Examples

- Alanine(A) encoded by GCU, the first two codon base position of GCU binds with six hydrogen bonds so Bora *et al.*,^[18] allocated weights on three codon base position as (2, 2, 0).
- Amino acids Asparagine (N) encoded by AAU, Lysine (K) encoded by AAA, the first two codon base position of AAU and AAA binds with four

Table 7: Codon Table for Degeneracy Four.

GUU	GUC	GUA	GUG	CCU	CCA	CCG	ACU	ACC	ACA	ACG	GCU	GCC	GCA	GCG	GGU	GGC	GGA	GGG
1	1	1	1	0.9989	0.9989	0.9989	0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
1	1	1	1	0.9989	0.9989	0.9989	0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
1	1	1	1	0.9989	0.9989	0.9989	0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
1	1	1	1	0.9989	0.9989	0.9989	0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
0.9989	0.9989	0.9989	0.9989	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989
0.9989	0.9989	0.9989	0.9989	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989
0.9989	0.9989	0.9989	0.9989	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989
0.9989	0.9989	0.9989	0.9989	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989
0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	1	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	1	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	1	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
0.9998	0.9998	0.9998	0.9998	0.9989	0.9989	0.9989	1	1	1	1	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989
0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989	1	1	1	1	0.9918	0.9918	0.9918	0.9918
0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989	1	1	1	1	0.9918	0.9918	0.9918	0.9918
0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989	1	1	1	1	0.9918	0.9918	0.9918	0.9918
0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989	1	1	1	1	0.9918	0.9918	0.9918	0.9918
0.9989	0.9989	0.9989	0.9989	0.9987	0.9987	0.9987	0.9989	0.9989	0.9989	0.9989	1	1	1	1	0.9918	0.9918	0.9918	0.9918
0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9918	0.9918	0.9918	0.9918	1	1	1	1
0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9918	0.9918	0.9918	0.9918	1	1	1	1
0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9918	0.9918	0.9918	0.9918	1	1	1	1
0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9918	0.9918	0.9918	0.9918	1	1	1	1
0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9989	0.9918	0.9918	0.9918	0.9918	1	1	1	1

Table 8: Codon Table for Degeneracy Six.

	UUA	UUG	CUU	CUC	CUA	CUG	CGU	CGC	CGA	CGG	AGA	AGG	AGU	AGC	UCU	UCC	UCA	UCG	
UUA	1	1	1	1	1	1	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711
UUG	1	1	1	1	1	1	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711
CUU	1	1	1	1	1	1	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711
CUC	1	1	1	1	1	1	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711
CUA	1	1	1	1	1	1	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711
CUG	1	1	1	1	1	1	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711
CGU	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	1	1	1	1	1	1	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875
CGC	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	1	1	1	1	1	1	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875
CGA	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	1	1	1	1	1	1	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875
CGG	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	1	1	1	1	1	1	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875
AGA	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	1	1	1	1	1	1	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875
AGG	0.99946	0.99946	0.99946	0.99946	0.99946	0.99946	1	1	1	1	1	1	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875
AGU	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	1	1	1	1	1	1	1
AGC	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	1	1	1	1	1	1	1
UCU	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	1	1	1	1	1	1	1
UCC	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	1	1	1	1	1	1	1
UCA	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	1	1	1	1	1	1	1
UCG	0.99711	0.99711	0.99711	0.99711	0.99711	0.99711	0.99875	0.99875	0.99875	0.99875	0.99875	0.99875	1	1	1	1	1	1	1

Table 9: Codon Table for Degeneracy Three.

	AUU	AUA	AUC
AUU	1	1	1
AUA	1	1	1
AUC	1	1	1

Table 10: Codon Table for Degeneracy One.

	UGG	AUG
UGG	1	.99875
AUG	.99875	1

our work, as we discussed above we have assigned some weights on the base position of the codons to find out a similarity measurement among the codons based on Lagerkvist’s hypothesis i.e., based on positional impact, appearance of pyrimidine and purine in the second and third base position and Hydrogen bonding influence in nucleotides of a codon. Since the behavior of the nucleotides in all the codons shows similar in kind so we have found similar kinds of weights for each codons. Which pave the way of getting very close similarity in between the codons.

Comparative Study with Impression Classification Techniques

In this section our aim is to produce a systematic comparative study between the existing work proposed by Das *et al.*, 2016^[12] and our proposed work. In order to compare these two methods, we identified some points that describe how to classify the amino acids based on degeneracy. While examining Table 11, we have noticed the following points

- I. Das *et al.*, (2016)^[12] have assigned some ternary number to each amino acid.
- II. They have mapped these ternary numbers into 10 impression values (IP).
- III. Finally, they have put these ten IP values into three groups and several subgroups whereas Total Impression IP value (TIP) same within a group.

As a consequence of these three facts different amino acids having different degeneracy clubbed into a same group. So in terms of degeneracy their method of grouping of amino acids is not significant. But, our defined similarity measure on codons put the all the amino acids having same degeneracy into one group. So our clustering method is more notable than to interpret the concept of degeneracy.

Network Representation of Codons

Every codon code a unique amino acid and the set of codons which are involved for coding a particular amino

Table 11: Classification of Amino Acids based on Impression.

Groups	Amino acid With Ternary symbol	Impression value	TIP (I1 + I2 + I3)
First group	GAP – 000, Q – 111, X – 222	(0,0,0)	0
	G – 001, D – 110 K – 112, O – 221	(0,1,1)	
Second group	L – 100, P – 011 fmet – 211, H – 122	(1,0,1)	2
	S – 010, I – 101 M – 121, Y – 212	(1,1,0)	
	A – 002, W – 220	(0, 2, 2)	
	Hyl – 200, Hyp – 022	(2, 0, 2)	
Third group	T – 020, U – 202	(2, 2, 0)	4
	E – 120, N – 102	(1, 2, 1)	
	V – 012, R – 210 F – 201, C – 021	(1, 1, 2) (2, 1, 1)	

acid are called synonymous codons. From the literature study it has been observed that the general structure of different aspects of codons such as nucleotide substitutions, degeneracy etc can be described by the methodology taken from graph theory. Akhtar *et al.* (2015)^[19] constructed three different types of undirected amino acids graph based on the point mutation in the first, second and third positions of a codon.

Figure 1 represents undirected unweighted codons graph $G(V, E)$ of 61 codons where represent set of vertices constitute all possible 61 codons and the set of edges connecting these vertices represented by E . Adjacency values between two codons of the matrix (Table 2, 3, 4 and Table 5) are greater than equal to 0.9967.

Mixing on Degree of Nodes

Definition 1

One of the important properties to study the nature of a network is mixing; it may be of assortative or disassortative in behavior. In a disassortative behavior network, higher degree has a tendency to connect with lower degree nodes. Disassortativity of a network can be measured by the Pearson correlation coefficient r of links nodes.

A positive correlation exhibits associations between the nodes of same degrees (assortativity), and the negative correlation exhibits connections between the nodes with dissimilar degrees (disassortativity).

Assessing and measuring the importance of a node in a network is of practical significance to improve the robustness, stability, and network synchronization, etc. Here, to calculate the correlation of mixing of amino acids network, equation (2) is used.

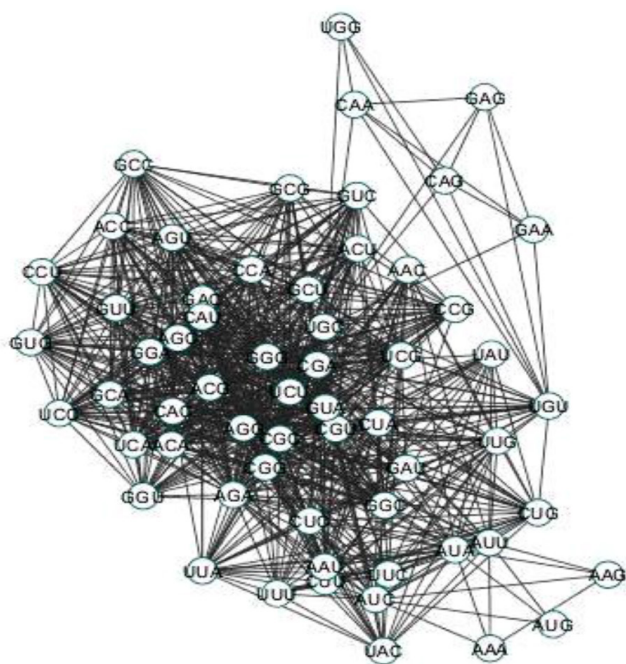


Figure 1: Codon Networks.

The coefficients of mixing of the network given by the Figure 1 is -0.0311 which reveals the fact that the codon networks disassortative mixing in nature.

Weighted Directed Graph Representation of Ordered Codons Sequences

In this section, it has been tried to represent a directed weighted graph from any random ordered codons sequences. Let $G = (V, E)$ be a directed graph where vertices $V = \{v_1, v_2, v_3, \dots, v_n\}$ are the codons and edges $E = \{e_1, e_2, e_3, \dots, e_n\}$. Individual weight has been assigned between two codons as shown in Table 2, 3, 4 and 5. This assigned weight represents the degree of similarities between the codons. To illustrate our findings let us consider the followings codons sequence. In this section we try to describe synonymous codons and non-synonymous codon with the aid of graph structure.

.....UUC GCU GCU AAA UUU UUC ACC ACC AUG.....

Now, for the above random codons sequence the graph ordered from UUC can be drawn to the next GCU, then GCU to GCU (forming a self-loop) and so on. The weight of an edge e_i is denoted as w_{e_i} and is defined as the degree of similarity of the connecting vertices.

As shown in Figure 2, GCU is the next consecutive vertex of the starting vertex UUC. So, there is a direction from UUC to GCU with a weight of 0.9972, which signifies that these two are non-synonymous codon,

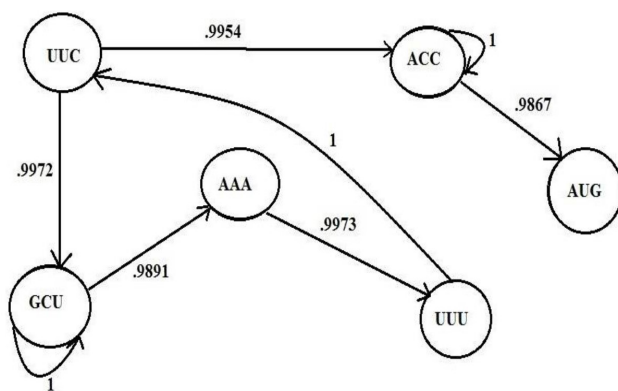


Figure 2: Weighted Directed Graph of Codons.

i.e., low degree of similarity. A self-loop from GCU to GCU signifies that they are synonymous codons with a maximum degree of similarity 1.

Similarly, to establish the other connections, it has been found that a weight 1, i.e., maximum degree of similarity between two distinct nodes UUU and UUC signifies that they are synonymous codons that represent the same amino acid.

CONCLUSION

In this paper we have interpreted the degeneracy property of genetic codons from a mathematical perspective. In our study, we demonstrated the degree of similarity among the codons using Cosine similarity. Furthermore, while analyzing the degree of similarities among 61 codons it has been noticed that some codons have degree of similarity 1, biologically which is substantiated because those codons are synonymous codons. Finally, we made a comparative study of our work with the work proposed by Das et al., 2016.^[6] In our future work we aim to analyze additional features of the genetic codon table from a mathematical perspective.

ACKNOWLEDGEMENT

We thank to Chairperson, Centre for Computer Science and Applications for providing constant guidance.

CONFLICT OF INTEREST

The authors declare no conflict of interest.

ABBREVIATIONS

DNA: Deoxyribonucleic Acid; **RNA:** Ribonucleic acid; **mRna:** Messenger RNA.

SUMMARY

In this manuscript we formulate a new classifying technique with the help of cosine similarities to explain the degeneracy. Also study the assortative and disassortative nature of codons networks which may help us to understand the flow of genetic information in the evolution process of amino acids.

REFERENCES

1. Antoneli F, Forger M. Symmetry breaking in the genetic code: Finite groups. *Math Comput Modell.* 2011;53(7-8):1469-88. doi: 10.1016/j.mcm.2010.03.050.
2. Hornos JE, Hornos YM. Algebraic model for the evolution of the genetic code. *Phys Rev Lett.* 1993;71(26):4401-4. doi: 10.1103/PhysRevLett.71.4401, PMID 10055237.
3. Kwon I, Kirshenbaum K, Tirrell DA. Breaking the degeneracy of the genetic code. *J Am Chem Soc.* 2003;125(25):7512-3. doi: 10.1021/ja0350076, PMID 12812480.
4. Origin of genetic code. *Nature.* 1966;212(5069):1397-. doi: 10.1038/2121397a0.
5. Sharma S, Bora PK, Ali T, Baruah AK. A survey on genetic code degeneracy in the Multiple structure. *J Comput Math Sci.* 2019;10(5):1052-60. doi: 10.29055/jcms/1098.
6. Lehmann J, Libchaber A. Degeneracy of the genetic code and stability of the base pair at the second position of the anticodon. *RNA.* 2008;14(7):1264-9. doi: 10.1261/ma.1029808, PMID 18495942.
7. Lenstra R. Evolution of the genetic code through progressive symmetry breaking. *J Theor Biol.* 2014;347:95-108. doi: 10.1016/j.jtbi.2014.01.002, PMID 24434741.
8. Crick FH. Codon–Anticodon pairing: The wobble hypothesis. *J Mol Biol.* 1966;19(2):548-55. doi: 10.1016/s0022-2836(66)80022-0, PMID 5969078.
9. Sengupta S, Higgs PG. A Unified Model of codon reassignment in alternative genetic codes. *Genetics.* 2005;170(2):831-40. doi: 10.1534/genetics.104.037887, PMID 15781705.
10. Gonzalez DL, Giannerini S, Rosa R. On the origin of degeneracy in the genetic code. *Interface Focus.* 2019;9(6):20190038. doi: 10.1098/rsfs.2019.0038, PMID 31641429.
11. Négadi T. The multiplet structure of the genetic code, from one and small number. *Neuro Quantology.* 2011;9(4). doi: 10.14704/nq.2011.9.4.379.
12. Das J, Majumder A, Choudhury P. Understanding of genetic code degeneracy and new way of classifying of protein family: A Mathematical Approach; 2016.
13. Ikehara K. Degeneracy of the Genetic Code has Played an Important Role in Evolution of Organisms. *Int J Genet Sci.* 2016;3(1):1-3. doi: 10.15226/2377-4274/3/1/00111.
14. Lagerkvist U. Unconventional methods in codon reading. *Bio Essays.* 1986;4(5):223-6. doi: 10.1002/bies.950040509, PMID 3790123.
15. Lagerkvist U. 'Two out of three': An alternative method for codon reading. *Proc Natl Acad Sci U S A.* 1978;75(4):1759-62. doi: 10.1073/pnas.75.4.1759, PMID 273907.
16. Rumer YB. Translation of 'Systematization of Codons in the Genetic Code [III]' by Yu. B. Rumer (1968). *Phil Trans R Soc A.* 2016;374(2063) (2063):20150447. doi: 10.1098/rsta.2015.0447.
17. Ohno S. *Evolution by Gene Duplication* (Softcover reprint of the original 1st ed. 1970 ed.). Springer; 2014.
18. Bora PK, Hazarika P, Baruah AK. Distance based amino acids network analysis. *Gene Rep.* 2020;21:100933. doi: 10.1016/j.genrep.2020.100933
19. Akhtar A, Ali T. Networks in amino acids based on mutation. *Stud Microeconomics.* 2015;3(2):89-100. doi: 10.1177/2321022215588863.

Cite this article: Hazarika P, Bora PK, Baruah AK, Bora P. Study of Codon Degeneracy Based on Similarity Measure. *Asian J Biol Life Sci.* 2022;11(2):594-604.