

Identification of Unique patterns, Mutation in Genetic finger printing using neural network approach

Mr.T.Saravanan¹, Head & Assistant Professor, Dr.B.Mukunthan², Associate Professor, PG & Research Department of Computer Science, Jairams Arts & Science College, Karur-639003.

Abstract – In genetic engineering, the advent of human genome project immensely increased the pressure for molecular computations and sequencing technologies dealing with data beyond the current abilities that are to be sequenced and interpreted. The automation of DNA feature extraction process achieved by applying neural network technique which has the advantage over conventional programming, in their ability to solve problem that do not have an algorithmic solution or the available solutions is too complex to be found is discussed in this paper, the above work also reduces the complication in precisely analyzing, interpreting of human DNA. The identification of exact location of occurrence of mutation in the DNA chain caused by radiation, viruses, transposons, mutagenic chemicals, as well as errors occurring during meiosis or DNA replication can be easily and exactly detected by subjecting the sample to gene sequencing process and analyzed using the above technique. In this novel approach the perfect blend made of bioinformatics and neural networks technique results in efficient DNA pattern analysis algorithm with 100% prediction accuracy, computed by number of correct identification of the target for a set of given inputs.

Key words - Neural-Fuzzy Resonance Mapping, Competitive learning, NFPR-processor, Input Generator, Preprocessor, Separator, Discriminator and Comparator, DNA profiling, DNA sequence Format, Mutation.

1. INTRODUCTION

Knowledge of DNA sequences has become indispensable for basic biological research. DNA sequencing is applied in various fields such as diagnostic, biotechnology, forensic biology and biological systematic. The DNA sequences of thousands of organisms have been decoded and stored in databases. The sequence information is analysed to determine genes that

encode polypeptides, RNA genes, regulatory sequences, structural motifs, and repetitive sequences. A comparison of genes within a species or between different species can show similarities between protein functions, or relations between species. With the growing amount of data, it became impractical to analyse DNA sequences manually.

Neural networks learn by examples so that it can be trained with known examples of a problem to gain knowledge about it so the neural network can be effective to solve unknown or untrained instances of the problem if is aptly trained. A pattern is essentially an arrangement or an ordering, in which some organization of underlying structure can be said to exist i.e. a pattern can be referred to as a quantitative or structural description of an object or some item of interest. A set of patterns that share some common properties can be regarded as pattern class in our case the unique repeated nucleotide sequence from the given Human DNA sample. The concept of applying Artificial Neural Systems (ANS) or Artificial Neural Networks (ANN) or simply Neural Networks in the field of DNA profiling is discussed in this paper.

2. ARTIFICIAL NEURAL NETWORK TECHNIQUES

Neural Networks [3] can process information in parallel, at high speed, and in a distributed manner. Neural networks which are simplified models of the biological neuron system, is a massively parallel distributed processing system made up of highly interconnected neural computing elements that have the ability to learn and thereby acquire knowledge and make it available for use. Neural Network architectures have been classified into various types based on their learning mechanisms and other features. Some classes of Neural Network refer to this learning process as training and the ability to solve a problem using the knowledge acquired as inference.

Neural Networks exhibit mapping capabilities, i.e., they can map input patterns to their associated output patterns. Neural Networks architectures can be trained with known examples of a problem before they are tested for their inference. They can, therefore, identify new objects previously untrained. Neural Networks possess the capability to generalize i.e. they can

predict new outcomes from past trends. Neural Networks are robust systems and are fault tolerant. They can therefore, recall full patterns from incomplete, partial or noisy patterns.

In Competitive Learning method those neurons which respond strongly to input stimuli have their weights updated, when an input pattern is presented, all neurons in the layer compete and the winning neuron undergoes weight adjustment. Hence it is a “Winner-takes-all” strategy.

Adaptive resonance theory employs a new principle of self organization based on competitive learning. Adaptive resonance theory nets are designed to be both stable and plastic. Neural networks suitable particularly for pattern classification problems in realistic environment is Neural- Fuzzy resonance mapping, it is a vast simplification of fuzzy resonance mapping which possess reduced computational overhead and architectural redundancy when compared to fuzzy resonance mapping.

3. DNA PROFILING AND SEQUENCING

DNA profiling also called DNA testing, DNA typing, or genetic fingerprinting, is a technique employed by forensic scientists to assist in the identification of individuals on the basis of their respective DNA profiles. DNA profiles, are encrypted sets of numbers that reflect a person's DNA makeup, which can also be used as the person's identifier. DNA sequencing theory addresses physical processes related to sequencing DNA. The term DNA sequencing [19] refers to sequencing methods for determining the order of the nucleotide bases—adenine, guanine, cytosine, thymine and uracil (rare case) in a molecule of DNA.

Single nucleotide poly-morphisms are a DNA sequence variation occurring when a single nucleotide A, T, C, or G in the genome (or other shared sequence) differs between members of a species (or between paired chromosomes in an individual). The genome [21] is the entirety of an organism's hereditary information which is encoded either in DNA or, for many types of virus, in RNA. For example, two sequenced DNA fragments from different individuals, AAGCCTA to AAGCTTA, contain a difference in a single nucleotide. Various DNA Sequence Formats available are: 1) Plain sequence format 2) EMBL format 3) GCG format 4) GCG-RSF (rich sequence format 5) Gen Bank format 6) IG format 7) FASTA format. A

sequence file in FASTA format of a given sample is used as an input to that is to be interpreted and analysed.

4. NEURAL-FUZZY PATTERN RECOGNITION PROCESSOR

4.1 Learning Input Generator

The input generator is used for input normalization and it represents the presence of particular feature in the input patterns and its absence. Various conditions for generating normalized learning input are shown below.

Learning Inputs
$$LIN_{i,n} = I_1, I_2, \dots, I_p$$

(1)

Where $0.1 \leq i \leq 0.5, 0.1 \leq n \leq 0.5$
and $p = 4$

TABLE 1

CONDITIONS FOR LEARNING INPUT NORMALIZATION

	Condition	Learning Input	Category
Case 1	$i \neq n$ or $i=n=0.1$ and $n \leq 0.5$	$LIN_{i,n} = i, n, 1-i, 1-n$ e.g. $LIN_{0.1, 0.1} = 0.1, 0.1, (1-0.1), (1-0.1)$ $LIN_{0.1, 0.1} = 0.1, 0.1, 0.9, 0.9$ $LIN_{0.2, 0.5} = 0.2, 0.5, (1-0.2), (1-0.5)$ $LIN_{0.2, 0.5} = 0.2, 0.5, 0.8, 0.5$	Category=L(logical)
Case 2	$i = n$ and $0.1 > i, n < 0.5$	$LIN_{i,n} = i, 1-i, 1-n, n$ e.g. $LIN_{0.2, 0.2} = 0.2, (1-0.2), (1-0.2), 0.2$ $LIN_{0.2, 0.2} = 0.2, 0.8, 0.8,$	Category=ILL(illogical)

		0.2 LIN _{0.3, 0.3} = 0.3, (1-0.3), (1-0.3), 0.3 LIN _{0.3, 0.3} = 0.3, 0.7, 0.7, 0.3	
Case 3	$i \neq n$ and $n > 0.5$	LIN _{i, n} = $i, n, 1-i, 1-n$ e.g. LIN _{0.5, 0.6} = 0.5, 0.6, (1-0.5), (1-0.4) LIN _{0.5, 0.6} = 0.5, 0.6, 0.5, 0.4	Category=ILL (illogical)

4.2 Activation Function Generator

When coded input patterns from input generator are presented to NFPR-Processor all output nodes become active to varying degrees. The output activation denoted by ACF_j referred to as the activation function for the jth output node. Where LIN is the learning input and LIW_j is the corresponding learning input weights.

$$ACF_j = \frac{|LIN \wedge LIW_j|}{\alpha + |LIW_j|} \quad (2)$$

Here α is kept as a small value close to 0 it's about 0.0000001. The node which registers the highest activation function is deemed Winner node i.e.

$$\text{Winner node} = \max(ACF_j) \quad (3)$$

In the event of more than one node emerging as the winner owing to the same activation function value some mechanism such as choosing a node with the smallest index may be devised to break the tie.

4.3 Match Function Generator

The match function which helps to determine whether the network must adjust its learning parameters is given by

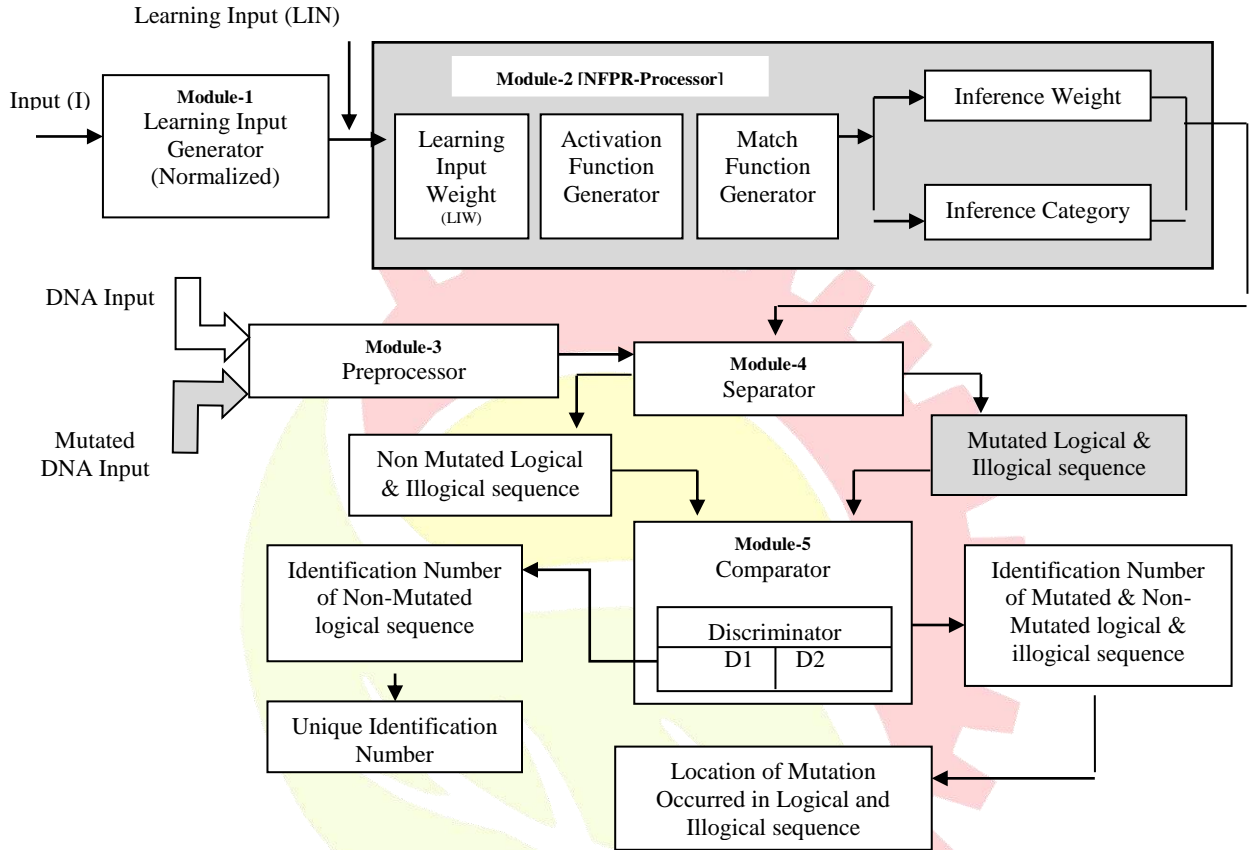


Figure 1 Block diagram of Neural-Fuzzy Pattern Recognition System

$$MAF_j = \frac{|LIN \wedge LIW_j|}{|LIN|} \quad (4)$$

The match function in association with the vigilance parameter decides on whether a particular output node is good enough to encode a given input pattern or whether a new output node should be opened to encode the same. The network is said to be in a state of resonance, if the match function value exceeds vigilance parameter. However, for a node to exhibit resonance, it is essential that it not only encodes the given input pattern but should also represent the same category as that of the input pattern.

The network is said to be in state of mismatch reset if the vigilance parameter exceeds match function, Such a state only means that the particular output node is not fit enough to learn the given input pattern and thereby cannot update its weights even though the category of the output node may be the same as that of the input pattern. This is so, since the output node has fallen short of the expected encoding granularity indicated by the vigilance parameter.

If match function is greater than vigilance parameter and category of input pattern is not same with the learning input, the vigilance parameter is updated and is given by

$$\rho = \text{MAF} + \delta \quad (\delta = 0.001) \quad (5)$$

TABLE 2
 GENERATING WEIGHTS FOR INFERENCE, CATEGORY FOR
 INFERENCE FROM LEARNING INPUTS

Nucleotide Pair	A,A	A,U	T,A	T,T	T,U	G,A	G,G	G,U	C,A	C,C	C,U	U,A	U,U
Category	L	L	L	ILL	L	L	ILL	L	L	ILL	L	L	ILL
Fuzzy Equivalent	0.1, 0.1*	0.1, 0.5*	0.2, 0.1*	0.2, 0.8*	0.2, 0.5*	0.3, 0.1*	0.3, 0.7*	0.3, 0.5*	0.4, 0.1*	0.4, 0.6*	0.4, 0.5*	0.5, 0.1*	0.5, 0.6***
Complement of Learning Input	0.9, 0.9	0.9, 0.5	0.8, 0.9	0.8, 0.2	0.8, 0.5	0.7, 0.9	0.7, 0.3	0.7, 0.5	0.6, 0.9	0.6, 0.4	0.6, 0.5	0.5, 0.9	0.5, 0.4
Augmented Input / Learning Input(LI)	0.1, 0.1, 0.9, 0.9	0.1, 0.5, 0.9, 0.5	0.2, 0.1, 0.8, 0.9	0.2, 0.8, 0.8, 0.2	0.2, 0.5, 0.8, 0.5	0.3, 0.1, 0.7, 0.9	0.3, 0.7, 0.7, 0.3	0.3, 0.5, 0.7, 0.5	0.4, 0.1, 0.6, 0.9	0.4, 0.6, 0.6, 0.4	0.4, 0.5, 0.6, 0.5	0.5, 0.1, 0.5, 0.9	0.5, 0.6, 0.5, 0.4

δ		0.00 1	0.00 1	0.00 1	0.00 1	0.0 01	0.00 1	0.00 1	0.00 1	0.00 1	0.00 1	0.00 1	0.00 1	0.001
ρ		0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5, 0.60 0+ δ	0.60 1	0.60 1	0.601
β		1	1	1	1	1	1	1	1	1	1	1	1	1
Acti vatio n Func tion	AC F(1)	0.99 99	0.79 99	0.93 75	0.79 99	0.9 99 9	0.93 33	0.87 51	0.99 99	0.92 85	0.92 30	0.99 99	0.92 30	0.8461
	AC F(2)	~	~	~	~	0.8 49 9	0.59 99	0.89 99	0.88 88	0.61 11	0.88 88	0.93 75	0.62 49	0.9999
	AC F(3)	~	~	~	~	~	~	~	~	~	~	~	~	0.7499
Highest Activation Function		AC F(1)	AC F(1)	AC F(1)	AC F(2)	AC F(1))	AC F(1)	AC F(2)	AC F(1)	AC F(1)	AC F(2)	AC F(1)	AC F(1) , AC F(2)	ACF(2))
Mat ch Func tion	MA F(1)	1.00 00	0.80 00	0.75 00	0.60 00	0.7 50 0	0.70 00	0.60 00	0.70 00	0.65 00	0.60 00	0.65 00	0.60 00	0.4500
	MA F(2)	~	~	~	~	0.8 50 0	0.60 00	0.90 00	0.80 00	0.55 00	0.80 00	0.75 00	0.50 00	0.7500
	MA F(3)	~	~	~	~	~	~	~	~	~	~	~	~	0.700
Category Match / Mismatch		Mat ch	Mat ch	Mat ch		Ma tch	Mat ch		Mat ch	Mat ch	Mis mat ch	Mat ch	Mat ch	
					Mat ch			Mat ch				Mat ch		Mat ch
Lear ning	LI W(=LI	Up dat	Up dat	Not Up	Up dat	Up dat	Not Up	Up dat	Up dat	Not Up	Up dat	Not Up	Not Updat

Input Weights Updated/ Not Updated/ Added	1)	ed	ed	dated	ed	ed	dated	ed	ed	dated	ed	dated	ed	
	LI W(2)	~	~	~	=LI	Not Updated	Not Updated	Updated	Not Updated	Not Updated	Updated	Not Updated	Not Updated (<ρ)	
	LI W(3)	~	~	~	~	~	~	~	~	~	~	~	Added	
Learning Input Weights & Category	LI W(1) L	0.1, 0.1, 0.9, 0.9	0.1, 0.1, 0.9, 0.5	0.1, 0.1, 0.8, 0.5	0.1, 0.1, 0.8, 0.5	0.1, 0.1, 0.8, 0.5	0.1, 0.1, 0.7, 0.5	0.1, 0.1, 0.7, 0.5	0.1, 0.1, 0.7, 0.5	0.1, 0.1, 0.6, 0.5	0.1, 0.1, 0.6, 0.5	0.1, 0.1, 0.6, 0.5	0.1, 0.1, 0.6, 0.5	WFI(1)=0.1, 0.1, 0.6, 0.5 CFI(1)=L
	LI W(2) ILL	~	~	~	0.2, 0.8, 0.8, 0.2	0.2, 0.8, 0.8, 0.2	0.2, 0.8, 0.8, 0.2	0.2, 0.7, 0.7, 0.2	0.2, 0.7, 0.7, 0.2	0.2, 0.7, 0.7, 0.2	0.2, 0.6, 0.6, 0.2	0.2, 0.6, 0.6, 0.2	0.2, 0.6, 0.6, 0.2	WFI(2)=0.2, 0.5, 0.5, 0.2 CFI(2)=ILL
	LI W(3) L	~	~	~	~	~	~	~	~	~	~	~	~	0.5, 0.1, 0.5, 0.9

A-ADENINE, T-THYMINE, G-GUANINE, C-CYTOSINE, U-URACIL, ACF=ACTIVATION FUNCTION, MAF=MATCH FUNCTION, CFI=CATEGORY FOR INFERENCE, WFI= WEIGHT FOR INFERENCE , L=LOGICAL, ILL=ILLOGICAL ρ=VIGILANCE PARAMETER *- CASE1,-CASE2,***CASE3,=LI-EQUAL TO LEARNING INPUT, <ρ=LESS THAN RHO**

DNA INPUTS

The weight updating equation of an output node j when it proceeds to learn the given input pattern I is given by

$$WFI_j^{new} = \beta (LIN \wedge WFI_j^{old}) + (1 - \beta)WFI_j^{old} \quad (6)$$

where $0 \leq \beta \leq 1$ ($\beta=1$)

Once the network has been trained, the inference of patterns, logical or illogical i.e. the categories to which the patterns belong may be easily computed. This is accomplished by passing the input pattern into the preprocessor and then to the input layer. All the output nodes compute the activation functions with respect to the input. The winner, node with the highest activation function, is chosen. The category to which output node belongs is the one to which given input pattern is classified by the network.

$$CIF_j = \frac{|PPO \wedge WFI_j|}{|WFI_j|} \quad (7)$$

If CIF (1) or CIF (3) is greater than CIF (2) the inferred category is logical else if CIF (2) is greater than CIF (1) and CIF (3) then inferred category is illogical. For the DNA inputs of fast a format whose category is logical the corresponding seven consecutive nucleotide base in the DNA sample is chosen as single logical sequence and DNA inputs whose category is illogical, two consecutive nucleotide base is considered as an illogical sequence with base pair thirty two.

TABLE 3
GENERATING WEIGHTS FOR INFERENCE, CATEGORY INFERENCE
FUNCTION FROM LEARNING INPUTS

		0.1,0	0.2,0	0.4,0	0.4,0	0.4,0	0.4,0	0.2,0	0.1,0	0.2,0	0.2,0	0.2,0
		.1	.3	.4	.4	.4	.2	.4	.1	.2	.2	.2
		A,A	T,G	C,C	C,C	C,C	C,T	T,C	A,A	T,T	T,T	T,T
		0.1,0	0.2,0	0.4,0	0.4,0	0.4,0	0.4,0	0.2,0	0.1,0	0.2,0	0.2,0	0.2,0
		.1,	.3,	.6,	.6,	.6,	.2,	.4,	.1,	.8,	.8,	.8,
		0.9,0	0.8,0	0.6,0	0.6,0	0.6,0	0.6,0	0.8,0	0.9,0	0.8,0	0.8,0	0.8,0
		.9	.7	.4	.4	.4	.8	.6	.9	.2	.2	.2
		0.1,0	0.1,0	0.1,0	0.1,0	0.1,0	0.1,0	0.1,0	0.1,0	0.1,0	0.1,0	0.1,0
		.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,
		0.6,0	0.6,0	0.6,0	0.6,0	0.6,0	0.6,0	0.6,0	0.6,0	0.6,0	0.6,0	0.6,0
		.5	.5	.5	.5	.5	.5	.5	.5	.5	.5	.5
		0.2,0	0.2,0	0.2,0	0.2,0	0.2,0	0.2,0	0.2,0	0.2,0	0.2,0	0.2,0	0.2,0
		.5,	.5,	.5,	.5,	.5,	.5,	.5,	.5,	.5,	.5,	.5,
		0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0
		.2	.2	.2	.2	.2	.2	.2	.2	.2	.2	.2
		0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0
		.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,	.1,
		0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0	0.5,0
		.9	.9	.9	.9	.9	.9	.9	.9	.9	.9	.9
		1.00	1.00	0.92	0.92	0.92	1.00	1.00	1.00	0.76	0.76	0.76
		00	00	30	30	30	00	00	00	92	92	92
		0.64	0.85	1.00	1.00	1.00	0.78	0.92	0.64	1.00	1.00	1.00
		28	71	00	00	00	51	85	28	00	00	00
		0.80	0.75	0.70	0.70	0.70	0.90	0.70	0.80	0.85	0.85	0.85

)	00	00	00	00	00	00	00	00	00	00	00	
GIC	CIF(CIF	CIF	CIF	CIF	CIF	CIF	CIF	CIF	CIF	CIF	CIF(
	1)	(1)	(2)	(2)	(2)	(1)	(1)	(1)	(2)	(2)	(2)	1)	
LOGI	L	L				L	L	L				L	
ICAL													
IC	ILLO												
	GICAL		ILL	ILL	ILL				ILL	ILL	ILL		
	L												
		0.1,0	0.2,0			0.4,0	0.2,0	0.1,0				0.2,0	
		.1,0.	.3,0.			.2,0.	.4,0.	.1,0.				.1,0.	
SOP		2,0.3	2,0.3	0.4,0	0.4,0	0.4,0	4,0.1	2,0.4	2,0.3	0.2,0	0.2,0	0.2,0	4,0.1
		,0.2,	,0.2,	.4	.4	.4	,0.1,	,0.2,	,0.2,	.2	.2	.2	,0.4,
		0.3,0	0.3,				0.1,	0.4,	0.3,				0.2,
		.2	0.1				0.1	0.1	0.2				0.4
<p>A=ADENINE,T=THYMINE,G=GUANINE,C=CYTOSINE, PPO=PREPROCESSOR OUTPUT,WFI=WEIGHT FOR INFERENCE, CIF=CATEGORY INFERENCE FUNCTION, GIC=GREATEST INFERRED CATEGORY,IC=INFERRED CATEGORY,L=LOGICAL,ILL=ILLOGICAL,SOP=SEPARATOR OUTPUT</p>													

Logical sequence (LS): $LS_{p,s,k} = Lseq_{p,s,1}, Lseq_{p,s,2}, \dots, Lseq_{p,s,k}$
 (8)

Where $p,s = 1$ to ∞
 and $k = 1$ to 7

The separator outputs which are logical in their category are fed to the discriminator (D1) where identification number is computed using the equation below to generate unique identification number

$$D1_{p,s} = \sum_{k=1}^7 k(Lseq_{p,s,k})^k$$

$p,s = 1 \text{ to } \infty$

(9)

Illogical sequence (IS): $IS_{p,s} = ILseq_s, ILseq_s, \dots, ILseq_\infty$

(10)

The separator outputs which are illogical in their category are fed to the discriminator (D2) where the discriminator output defined by

$$D2_{p,s} = ILseq_s^m$$

where $p, s, m = 1 \text{ to } \infty$

(11)

$m =$ Number of times nucleotide base is repeated

The comparator unit compares the identification numbers of all logical sequences of mutated and non-mutated DNA inputs from Discriminator (D1) and illogical sequences of mutated and non-mutated sequences from (D2) to identify the location of mutation in the given sample.

DNA SAMPLE: HUMAN-1 [BASE PAIR =32, SEQUENCE =25]
 AATGTGTTGTGTGACCCCTCAAATCTCTCAAATGTGTTTTTACAC
 TCCGTTGGTAATATGGAATGTGTTAAAGTTGCTACCCGGGGTTTT
 TTAATGTGTCTCT

TABLE 4
 DISCRIMINATOR (D1) OUTPUTS FOR NON-MUTATED LOGICAL SEQUENCE OF HUMAN-1

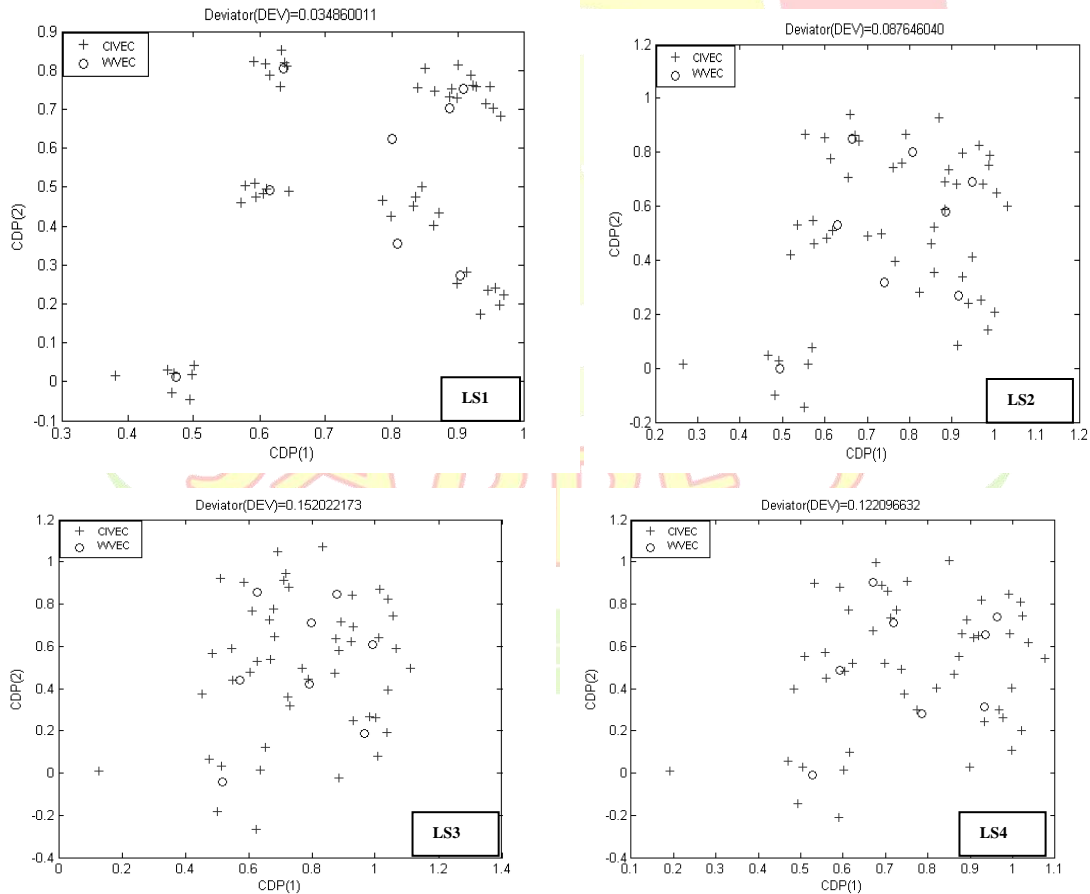
Sequ	Discr	Logical Sequence(LS _{p,s,k})	Identific	Uniq
------	-------	----------------------------------------	-----------	------

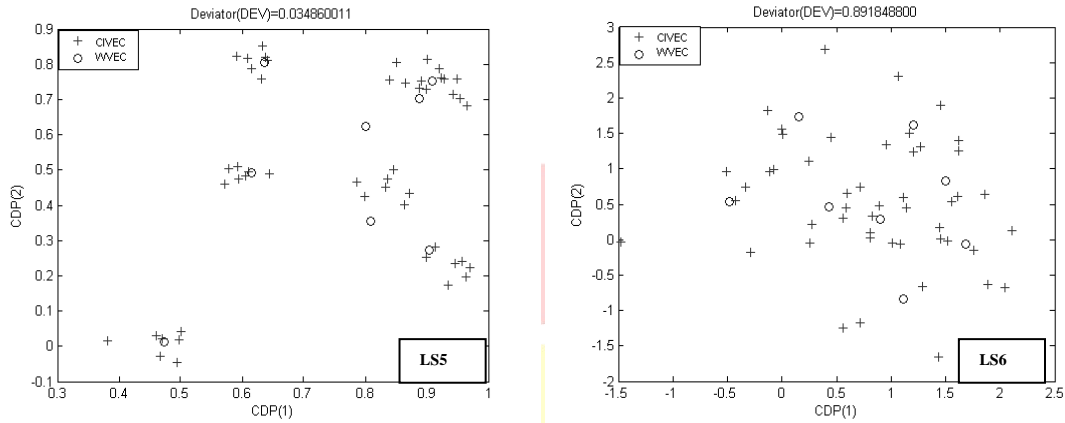
Logi cal Sequ ence	Hu man (p)	ence (s)	imin- ator(D1) Input s ($l_{p,s,k}$)	Lseq p,s,k (k=1)	Lseq p,s,k (k=2)	Lseq p,s,k (k=3)	Lseq p,s,k (k=4)	Lseq p,s,k (k=5)	Lseq p,s,k (k=6)	Lseq p,s,k (k=7)	ation Number ($D1_{p,s}$)	ue Identi ficat -ion numb er (UIN_p)
LS1	1	1	$L_{1,1,k}$	0.1	0.1	0.2	0.3	0.2	0.3	0.2	0.18246 4	0.182 464 (REP EAT ED PAT TER N)
LS2	1	2	$L_{1,2,k}$	0.2	0.3	0.2	0.3	0.2	0.3	0.1	0.44237 5	
LS3	1	3	$L_{1,3,k}$	0.4	0.2	0.4	0.1	0.1	0.1	0.1	0.67245 7	
LS4	1	4	$L_{1,4,k}$	0.2	0.4	0.2	0.4	0.2	0.4	0.1	0.67113 7	
LS5	1	5	$L_{1,5,k}$	0.1	0.1	0.2	0.3	0.2	0.3	0.2	0.18246 4	
LS6	1	6	$L_{1,6,k}$	0.2	0.1	0.4	0.1	0.4	0.2	0.4	0.47545 3	
LS7	1	7	$L_{1,7,k}$	0.4	0.3	0.2	0.2	0.3	0.3	0.2	0.62701 4	
LS8	1	8	$L_{1,8,k}$	0.1	0.1	0.2	0.1	0.2	0.3	0.3	0.15046 5	
LS9	1	9	$L_{1,9,k}$	0.1	0.1	0.2	0.3	0.2	0.3	0.2	0.18246 4	
LS1 0	1	10	$L_{1,10,k}$	0.2	0.1	0.1	0.1	0.3	0.2	0.2	0.23948 0	
LS1 1	1	11	$L_{1,11,k}$	0.3	0.4	0.2	0.1	0.4	0.4	0.4	0.73164 5	
LS1 2	1	12	$L_{1,12,k}$	0.3	0.2	0.2	0.2	0.2	0.2	0.2	0.41103 3	
LS1 3	1	13	$L_{1,13,k}$	0.1	0.1	0.2	0.3	0.2	0.3	0.2	0.18246 4	

TABLE 5
PERFORMANCE OF PROPOSED SYSTEM FOR VARIOUS NUMBERS
OF EPOCHS

Vigilance Parameter (ρ):0.5				
S. no	Learning Vector (Number of Epochs)	Number of Learning Inputs	Learning Time (Seconds)	Accuracy%
1	25	25	62.49	100%
2	13	13	34	100%

(No. of Epochs=25(For All Possible Combinations) and No. of Epochs=13)





(CDP=Clustered Data Points, CIVEC=Cluster of Input

Figure 3 MATLS Vectors, WVEC=Weight Vectors) input (LS₍₁₎-LS₍₆₎)
 Showing LS₍₁₎ and LS₍₅₎ are unique

5. IDENTIFICATION OF MUTATION IN THE SAMPLE

Mutation is a change of DNA sequence within a gene or chromosome of an organism resulting in the creation of a new character or trait not found in the parental type [22]. The mutation results when a change occurs in a chromosome, either through an alteration in the nucleotide sequence of the DNA coding for a gene or through a change in the physical arrangement[24] [25] of a chromosome.

There are many different types of mutations; a point mutation (base pair substitution) is a simple change in one base of the gene sequence. In this case, the entire meaning of the sentence has been altered with a one letter change. In neutral or silent mutation, another one letter point mutation has occurred. However, the meaning of the sentence has not been altered.

In a frame shift mutation, one or more bases are inserted or deleted into the sequence of the gene, the equivalent of adding or removing letters in a sentence, adding or removing one letter changes each subsequent word. This type of mutation can make the DNA meaningless and often results in shortened and functionless protein. Mutations that result in missing DNA are called deletions. These can be small, or longer deletions that affect a large number of genes on the chromosome. Deletions can also cause frame-shift mutations. Mutations that result in the addition of extra DNA are called insertions. Insertions can also cause frame-shift mutations, and generally result in a nonfunctional protein.

In an inversion mutation, an entire section of DNA is reversed. A small inversion may involve only a few bases within a gene, while longer inversions involve large regions of a chromosome containing several genes.

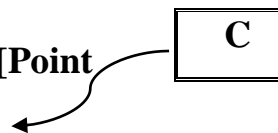
5.1 Various Types of Mutation identification in Human-1 sample

Before Mutation:

LS1/RS	LS2	IS1	IS1	IS1	LS3	LS4
LS5/RS	IS2	IS2				
AATGTGT	TGTGTGA	C	C	C	CTCAAAA	
TCTCTCA	AATGTGT	T	T			
IS2	LS6	LS7	LS8	LS9/RS	LS10	
LS11	IS3	IS3				
T	TACACTC	CGTTGGT	AATATGG	AATGTGT	TAAAGTT	
GCTACCC	G	G				
IS3	LS12	LS13/RS	LS14			
G	GTTTTTT	AATGTGT	CTCTXXX			

Case 1:-

After Mutation in Logical sequence: [Point C mutation]



LS1/RS LS2 IS1 IS1 IS1 LS3 LS4
 LS5/RS IS2 IS2
AATGTGT TGTGTGA C C C CTCA C AA
TCTCTCA AATGTGT T T
 IS2 LS6 LS7 LS8 LS9/RS LS10
 LS11 IS3 IS3
T TACACTC CGTTGGT AATATGG AATGTGT TAAAGTT
GCTACCC G G
 IS3 LS12 LS13/RS LS14
G GTTTTTT AATGTGT CTCTXXX

In case 1 the point mutation occurred in logical sequence (LS_{1,3}) by the mutant C can be identified with the change in identification number of LS_{1,3} where identification number of illogical sequence remains unaltered.

Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number
LS _{1,1}	RS	0.182464
LS _{1,2}		0.442375
LS _{1,3}		0.672457
LS _{1,4}		0.671137
LS _{1,5}	RS	0.182464
LS _{1,6}		0.475453
LS _{1,7}		0.627014
LS _{1,8}		0.150465
LS _{1,9}	RS	0.182464
LS _{1,10}		0.239480
LS _{1,11}		0.731645
LS _{1,12}		0.411033

LS _{1,13}	RS	0.182464
--------------------	----	----------

Before Mutation

After [Point Mutation] in Logical Sequence

<u>Before Mutation</u>		Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number
<u>After Mutation in logical Sequence (Not Altered)</u>		LS _{1,1}	RS	0.182464
<u>Illogical Sequence (IS)</u>	<u>Identification Number</u>	IS _{1,1}		0.442375
		IS _{1,2}		0.723607
		IS _{1,3}		0.671137
		IS _{1,1}	RS	0.182464
IS _{1,2}			0.475453	
IS _{1,3}			0.627014	
		LS _{1,8}		0.150465
		LS _{1,9}	RS	0.182464
		LS _{1,10}		0.239480
		LS _{1,11}		0.731645
		LS _{1,12}		0.411033
		LS _{1,13}	RS	0.182464

[Result: Change in polypeptide sequence might change the shape or function of the protein, depending on where in the sequence occurs]

Case 2:-

After Mutation in Logical sequence: [Frame shift mutation-Insertion]

LS1/RS	LS2	IS1	IS1	IS1	LS3	LS4
LS5/RS	IS2	IS2				
AATGTGT	TGTGTGA		C	C	C	CTCAAAA
TCTCTCA	AATGTGT	T	T			
IS2	LS6	LS7		LS8	LS9/RS	LS10
LS11	IS3	LS12				

T TACACTC CGTTGGT AATATGG AATGTGT TAAAGTT
 GCTACCC G CGGGTTT

IS4 1S4 LS13 LS14

T T TAATGTG TCTCTXX

C

In case 2 the frame shift mutation (insertion) occurred in one of the IS_{1,3} by the mutant C which alters both the logical sequence (LS_{1,12}) and illogical sequence (IS_{1,3}) that can be identified by the change in identification number of both logical sequence (LS_{1,12}) and illogical sequence (IS_{1,3}) after mutation.

Before Mutation
Illogical Sequence

After Mutation in

Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number	Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number
LS _{1,1}	RS	0.182464	LS _{1,1}	RS	0.182464
LS _{1,2}		0.442375	LS _{1,2}		0.442375
LS _{1,3}		0.672457	LS _{1,3}		0.672457
LS _{1,4}		0.671137	LS _{1,4}		0.671137
LS _{1,5}	RS	0.182464	LS _{1,5}	RS	0.182464
LS _{1,6}		0.475453	LS _{1,6}		0.475453
LS _{1,7}		0.627014	LS _{1,7}		0.627014
LS _{1,8}		0.150465	LS _{1,8}		0.150465
LS _{1,9}	RS	0.182464	LS _{1,9}	RS	0.182464
LS _{1,10}		0.239480	LS _{1,10}		0.239480
LS _{1,11}		0.731645	LS _{1,11}		0.731645
LS _{1,12}		0.411033	LS _{1,12}		0.985633
LS _{1,13}	RS	0.182464	LS _{1,13}		0.243464

Mutation occurred after LS_{1,12}

Before Mutation
Illogical Sequence

After Mutation in

Illogical Sequence (IS)	Identification Number
IS _{1,1}	0.064000
IS _{1,2}	0.008000
IS _{1,3}	0.027000

Illogical Sequence (IS)	Identification Number
IS _{1,1}	0.064000
IS _{1,2}	0.008000
IS _{1,3}	-0.300000
IS _{1,4}	0.040000

Mutation occurred in IS_{1,3}

[Result:

Change in polypeptide sequence might change the shape or function of the protein, depending on where in the sequence occurs.]

Case 3:-

After Mutation in Illogical sequence: [Point mutation (Neutral or Silent)]

LS1/RS	LS2	IS1	IS1	IS1	LS3	LS4
LS5/RS	IS2	IS2	IS2			
AATGTGT	TGTGTGA	C	C	C	CTCAAAA	TCTCTCA
AATGTGT	T	T	T			
LS6	LS7	LS8	LS9/RS	LS10		
LS11	IS3	IS3				
TACACTC	CGTTGGT	AATATGG	AATGTGT	TAAAGTT		
GCTACCC	G	G				
IS3	IS3	LS12	LS13/RS	LS14		
G	G	GTTTTTT	AATGTGT	CTCTXXX		



Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number
LS _{1,1}	RS	0.182464

<u>Before Mutation</u>		<u>After Mutation</u>		
		LS _{1,2}		0.442375
		LS _{1,3}		0.672457
		LS _{1,4}		0.672577
		LS _{1,5}	RS	0.182464
		LS _{1,6}		0.475453
		LS _{1,7}		0.627014
		LS _{1,8}		0.151905
		LS _{1,9}	RS	0.182464
		LS _{1,10}		0.236024
Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number		
LS _{1,1}	RS	0.182464	LS _{1,11}	0.731645
LS _{1,2}		0.442375	LS _{1,12}	0.412474
LS _{1,3}		0.672457	LS _{1,13}	RS 0.182464
LS _{1,4}		0.671137		
LS _{1,5}	RS	0.182464		
LS _{1,6}		0.475453		
LS _{1,7}		0.627014		
LS _{1,8}		0.150465		
LS _{1,9}	RS	0.182464		
LS _{1,10}		0.239480		
LS _{1,11}		0.731645		
LS _{1,12}		0.411033		
LS _{1,13}	RS	0.182464		

[Result: No change in polypeptide sequence, possible consequence for the organism =none]

In case 3 the point mutation is occurred in same IS_{1,3} as case 2 but with mutant G that only alters the illogical sequence (IS_{1,3}) and not any of the logical sequences that can be identified only using

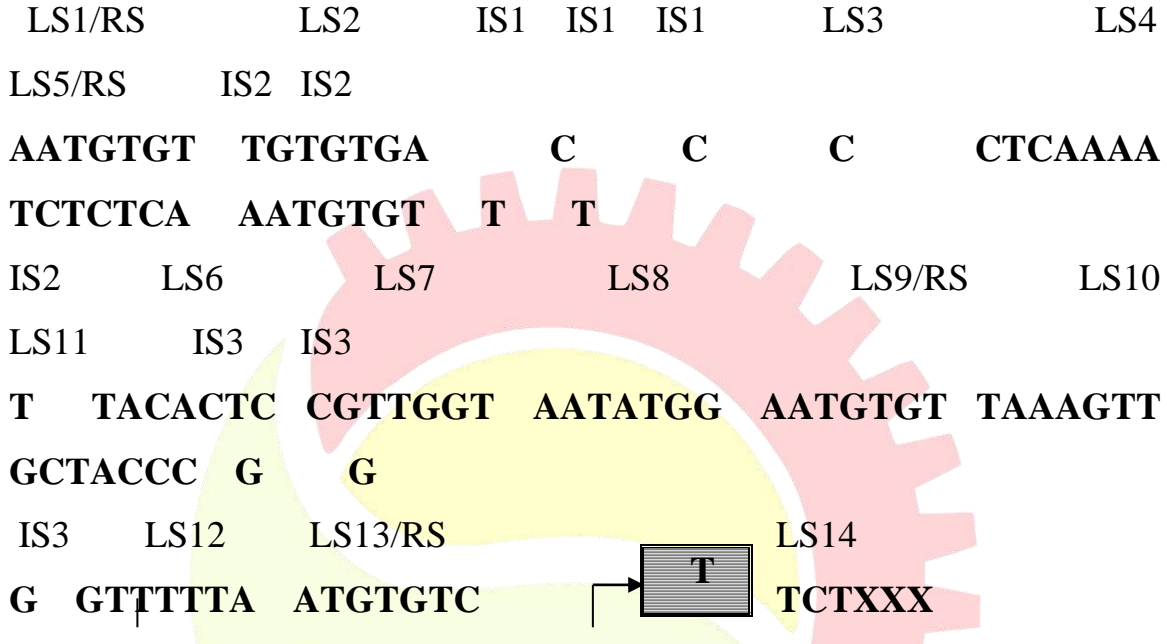
the change in identification number of illogical sequence (IS_{1,3}).

<u>Before Mutation</u>		<u>After Mutation in</u>	
<u>Illogical Sequence</u>		<u>Illogical Sequence</u>	
Illogical Sequence (IS)	Identification Number	Illogical Sequence (IS)	Identification Number
IS _{1,1}	0.064000	IS _{1,1}	0.064000
IS _{1,2}	0.008000	IS _{1,2}	0.008000
IS _{1,3}	0.027000	IS _{1,3}	0.008100

Mutation occurred in IS_{1,3}

Case 4:-

After Mutation in Logical sequence: [Frame shift mutation-Deletion]



In case 4 the frame mutation [deletion] occurred in logical sequence (LS_{1,12}) by the removal of mutant T and can be identified with the change in identification number of logical sequence (LS_{1,12}) with no alteration in any of the illogical sequence .

<u>Before Mutation</u>		<u>After Mutation in Logical</u>
Logical Sequence (LS)	Repeated Sequence (RS)	Sequence Identification Number
LS _{1,1}	RS	0.182464
LS _{1,2}		0.442375
LS _{1,3}		0.672457
LS _{1,4}		0.671137
LS _{1,5}	RS	0.182464
LS _{1,6}		0.475453
LS _{1,7}		0.627014
LS _{1,8}		0.150465
LS _{1,9}	RS	0.182464

Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number
LS _{1,10}		0.239480
LS _{1,11}		0.731645
LS _{1,12}		0.411033
LS _{1,13}	RS	0.182464
Before Mutation		
After Mutation in Illogical Sequence (Not Altered)		
LS _{1,1}		0.182464
LS _{1,2}		0.442375
LS _{1,3}		0.672457
LS _{1,4}		0.671113
LS _{1,5}	RS	0.182464
LS _{1,6}		0.475453
LS _{1,7}		0.627014
LS _{1,8}	Illogical Sequence (IS)	0.150465
LS _{1,9}	RS	0.182464
LS _{1,10}	IS	0.239480
LS _{1,11}	IS _{1,1}	0.064000
LS _{1,11}	IS _{1,1}	0.731645
LS _{1,12}	IS _{1,2}	0.008000
LS _{1,12}	IS _{1,2}	0.411033
LS _{1,13}	IS _{1,3}	0.027000
LS _{1,13}	IS _{1,3}	0.291402

Mutation occurred after LS_{1,12}

Change in polypeptide sequence might change the shape or function of the protein, depending on where in the sequence occurs.]

In case 5 below the inversion mutation occurred in logical sequence (LS_{1, 10}) by replacing TAAAGTT with mutant TTGAAAT that can be identified with the change in identification number of logical sequence (LS_{1, 10}) alone with no alteration in any of the illogical sequence.

Case 5:-

After Mutation in Logical sequence: [Inversion mutation]

LS1/RS	LS2	IS1	IS1	IS1	LS3	LS4
LS5/RS	IS2	IS2				
AATGTGT	TGTGTGA	C	C	C	CTCAAAA	
TCTCTCA	AATGTGT	T	T			

IS2 LS6 LS7 LS8 LS9/RS LS10
 LS11 IS3 IS3
 T TACACTC CGTTGGT AATATGG AATGTGT TTGAAAT
 GCTACCC G G
 IS3 LS12 LS13/RS
 LS14
 G GTTTTTT AATGTGT CTCTXXX



Before Mutation

After Mutation in

Logical Sequence

Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number	Logical Sequence (LS)	Repeated Sequence (RS)	Identification Number
LS _{1,1}	RS	0.182464	LS _{1,1}	RS	0.182464
LS _{1,2}		0.442375	LS _{1,2}		0.442375
LS _{1,3}		0.672457	LS _{1,3}		0.672
LS _{1,4}		0.671137	LS _{1,4}		0.671
LS _{1,5}	RS	0.182464	LS _{1,5}	RS	0.182464
LS _{1,6}		0.475453	LS _{1,6}		0.475453
LS _{1,7}		0.627014	LS _{1,7}		0.627014
LS _{1,8}		0.150465	LS _{1,8}		0.150465
LS _{1,9}	RS	0.182464	LS _{1,9}	RS	0.182464
LS _{1,10}		0.239480	LS _{1,10}		0.361546
LS _{1,11}		0.731645	LS _{1,11}		0.731645
LS _{1,12}		0.411033	LS _{1,12}		0.411033
LS _{1,13}	RS	0.182464	LS _{1,13}	RS	0.182464

[**Result:** Change in polypeptide sequence might change the shape or function of the protein, depending on where in the sequence occurs.]

6. CONCLUSION

As an attempt to automate the genetic finger printing for precise identification of individuals from their DNA sample, the Neural-fuzzy Pattern Recognition System implemented using the concept of fuzzy resonance theory mapping discussed in the above work classifies the sequences to identify a unique number from the given sample that actually includes nucleotide basis of adenine, guanine, cytosine, thymine and uracil which are represented by fuzzy values respectively. Any type of mutation for instance, gene mutations in the Japanese HNPCC (Hereditary Non Polyposis Colorectal Cancer) which triggers HNPCC tumor that could not be detected even by PCR-SSCP (Polymerase chain reaction-Single strand conformation polymorphism) can be easily detected by subjecting the sample to gene sequencing process and analyzed using the proposed system.

Further development can be extended by training patterns in DNA protein represented by suitable fuzzy equivalent to classify and predict the protein structure in the protein folding problem and also above technique can be used in the areas where feature extraction is to be done in genetic engineering with suitable modification.

7. REFERENCES

1. Richard O. Duda, Peter E.Hart, David G. Stork, "Pattern classification"- Second Edition", John Wiley and sons, 2006.
2. John Hertz, Anders Krogh, and Richard G. Palmer. "Introduction to the Theory of Neural Computation". Addison Wesley, Redwood City, A, 2008.
3. Stephen J. Hanson, Jack D. Cowan, and C. Lee Giles, "Advances in Neural Information Processing Systems", volume 5, Morgan Kaufmann San Mateo CA, 2009.

4. "Advances in Neural Networks issn-2006", Third international symposium on neural networks, Springer Berlin Heidelberg, New York publications.
5. Robert Schalkoff, "Pattern Recognition: Statistical, Structural and Neural Approaches, 2007, John Wiley and sons.
6. Carpenter, G.A. and S. Grossberg, "A Massively Parallel Architecture for a self-organizing Neural Pattern Recognition Machine", Computer Vision, Graphics and Image Processing, 37, PP. 54-115.
7. Carpenter, G.A. and S. Grossberg, and J.H. Reynolds (2010), "ARTMAP: Supervised Real Time Learning and Classification of Non-stationary Data by a Self-organizing Neural Network". Vol. 4, pp. 565-588.
8. Phipps Arabie, Lawrence J. Hubert, and Geert De Soete, editors, "Clustering and Classification". World Scientific, River Edge, NJ.
9. Stephen, Krawetz, David D. Womble, "Introduction to Bioinformatics A Theoretical and Practical Approach", Human Press Inc.,
10. David W. Mount, David W. Mount, "Bio informatics Sequence and Genome analysis"- Second Edition, Cold Spring Harbor Laboratory Press, New York.
11. Norah Rudin, Keith Inman, "An Introduction forensic DNA Analysis", 2002-CRC Press.
12. Donald R. Tsveter. "The Pattern Recognition Basis of Artificial Intelligence". IEEE Press, New York, page 117, Computational Intelligence and Bio inspired Systems, 8th international work conference on artificial neural networks, iwann-2005proceedings.
13. Julie A. Ayala-Gross, "DNA Analysis: The best method for Human Identifications", National University, San Diego – 2001.
14. Joe Nickell and John F. Fischer, "Crime Science Methods of Forensic Detection", 1999. University Press of Kentucky.
15. John O. Savino, Brent E Turvey, "Rape Investigation Hand book", 2005, Elsevier Inc.,
16. David E. Newton, "DNA Evidence and Forensic science"- 2008 facts on file, Inc. <http://www.factsonfile.com>.
17. Jorg T. Epplen Thomas Lubjuhn, Birkhauser, "DNA Profiling and DNA Finger Printing", Verlag Publication.

18. Simon Eastaeal, Neil Mc Lead, Ken, Harwood , “DNA Profiling Principles, Pitfalls and Potential”, Academic Publishers, Inc.,
19. Des Higgins, willie Taylor, “Bioinformatics Sequence, Structure and data banks”, Oxford University Press, 2000.
20. “Bioinformatics for geneticists”, Michael R. Barnes , Second Edition, John Wiley & Sons Ltd.,
21. Andreas D. Buxevanis, “Bioinformatics-A practical Guide to the Analysis of genes and proteins”, second edition, A John wiley & sons, Inc., Publication.
22. Charles L. Valon, “New developments in Mutation Research”, Nova science publishers Inc New York, 2007.
23. “Oxidative Damage to Nucleic Acids”, Springer science press, New York.
24. Richard G. H. Cotton, Edward Edkins, Sue Forrest “Mutation detection”, IRL Press at Oxford University Press.
25. Graham R. Taylor “Laboratory methods for the detection of mutations and polymorphisms in DNA”, CRC Press, 2007 - Science.
26. S.N Sivanandam, “Introduction to neural networks and MATLAB-6.0” ,Tata McGraw-Hill publishing company,2006