

# Identification of Healthy Genomes Based On mutations Using PSI Blast

Tahmeena Fatima<sup>1</sup>, Singaraju Jyothi<sup>2</sup>, Dadala Mary Mamatha

<sup>1</sup>Research Scholar, Dept. of CSE, Sri Padamavati Mahila Visvavidyalayam, Tirupati  
tahmi.fatima18@gmail.com

<sup>2</sup>Professor, Department of Computer Science, Sri Padamavati Mahila Visvavidyalayam, Tirupati  
jyothi.spmvv@gmail.com

<sup>3</sup>Professor, Department of Bioscience & Sericulture, Sri Padamavati Mahila Visvavidyalayam, Tirupati  
dmmfulbrightucdavis@gmail.com

## Article Info

Volume 83

Page Number: 5803 - 5808

Publication Issue:

May - June 2020

## Abstract

The basis of unclear personal data is defended to predict the damages like stigmatisation and refinement are done by several methods from so many years. The rising challenge for people to use the anonymised information of genome is identified their potential using different efficient strategies. A genomic data is act as a general schema for genomic repositories. This dataset are literally gatherings of samples where every sample contains two parts there are region data and metadata. Each genomic set of data is connected with a data scheme in which former five attributes are fixed in order to represent the region coordinates and the sample identifier. The fixed region attributes consist of the chromosome which the region belongs to left and right ends within the chromosome and the value is denoted the DNA strand that contains the region. In recent times the area of computational genomics has been expanded to provide the essential component of the large international effort in genomics. For the analysing data science will permits the removal of applied visions from huge scales data. The techniques are used for identification of data is more difficult. However, genome can be healthy or unhealthy which is derived through the machine learning algorithm.

**Keywords:** Genome Database, Computational Genomics, Mutation.

## Article History

Article Received: 19 November 2019

Revised: 27 January 2020

Accepted: 24 February 2020

Publication: 17 May 2020

## I. INTRODUCTION

The simple form of a genome DNA is a complete set of organisms which includes its genes. DNA sequencing is an outline of a genomic data that refers to a person's genome which is derived from sequences of raw data that varies from individual or whole or a part of DNA.

By using NGS technology the genome [1] must sequenced minimum of 10 entities to cover the high entities. The genome was first sequenced by Watson to cover the 7.4× on 454 GS and Roche included 3.3 Million Single Nucleotide Polymorphisms (SNPs), around 82% of information is listed in Biotechnology Information SNP database (dbSNP) at national centre. Strangely the personal genome that follows the NGS technology [11] of SNP reports the same result with genome of 3 to 4 million and dbSNP is over lapped with 80 to 90%. This arrangement is so vigorous, in fact, it

may consider SNPs of 3 million with dbSNP of 80 to 90% concordance (depends on the culture of the gene) to be the 'gold standard' for detection of whole-genome sequencing (WGS) [3] SNP. Another suggestion for this genome arrangement is that contains the 0.5 million individual SNPs, where studies of WGS expand the growth of public database submission. In 2007 Watson genome is submitted completes skyrocketed dbSNP. As of 2010 February, over 100 million dbSNP are received from human, contains 23.7 matchless sequence variants are validated more than half. The backdrop of human genetics is speedily moving, fuelled by the advent of hugely parallel sequencing methodologies. New tools from Helicos Biosciences (Heliscope), Roche (454), Life Technologies (SOLiD) and Illumina (Genome Analyzer) produce millions of short sequence reads per run, creates the possible ways to sequence the entire human genomes in a matter of weeks. These 'next-

generation sequencing' (NGS) methods are previously been active to sequence the unconstitutional genomes of numerous individuals. Determined efforts like the Personal Genomes Project and the 1000 Genomes Project hopes to add thousands and more. The first five cancer genomes to be distributed revealed thousands of novel somatic mutations and involved new genes in tumour [10] enlargement and evolution. Our awareness on the genetic variants may find the disease vulnerability, response at treatment and other phenotype will repeatedly progress the study of DNA sequence is improved in human genome. The gene data will be kept in supercomputers that are highly secured and the IndiGen cards will be made anonymous. Thus, person's genetic information loss is not negotiated. then

A genome is a person's whole set of deoxyribonucleic acids or DNA, as well as genes contains base pairs of 3 billion. By sequencing the genome, researchers can discover the functions of genes and identify mutations responsible for any rare genetic diseases.

A genome sequence does not end itself. The main task is to understand the genome [4] and what it contains and what are its function and how to use them. In previously it addressed as the mixture of experiment and computer analysis to control the regions of gene and locate them. This stated as the first method in genome sequencing. The second state is to understand the functions of the genome to certain range, simply last 30 years the molecular biology is stated in different ways. The only variation from the previous attention is directed at pathway expression of one gene is interconnected with another one. Now a days it becomes more common and relates to the whole genome expression.

## II. PROPOSED METHOD

### Finding the genes in a sequencing

After obtained the DNA sequence [12] of distinct cloned portion or a whole chromosome then applied numerous techniques to work on finding the genes present in the genome sequencing. These techniques are differentiated to simply examine the sequence,

by frequently using eye or computer, look at the genes which are associated with feature sequence, and techniques that locate the genes by tentative analysis of DNA sequence. Bioinformatics approach is the one of the technique of computer.

### Examination of sequence by locateGene

Examine the sequence is applied to locate the gene due to gene is not having series of nucleotides instead it has unique features. These structures define it a gene sequence or not, and the description is not influenced then it is non-coding DNA. At this point we did not understand the features and nature of the gene and examining of gene is not a complete evidence of gene location but it is quite potential tool and it is the initial technique applied to analyse the sequence of a new gene.

The genome sequence of human transcript [6] contains DNA sequence is copied from NCBI. The elements of mobile in human genome are recognised by Repeat Masker, elements of transposable sequence are recognised by Repbase Update. From NCBI genbank we can copy the information of pathology and tissue samples of transcript.

### About fasta file format

DNA sequence data frequently are contained in a file format called "fasta" format. Fasta format is simply a single line prefixed by the greater than symbol that contains annotations and another line that contains the sequence.

>information about the sequence

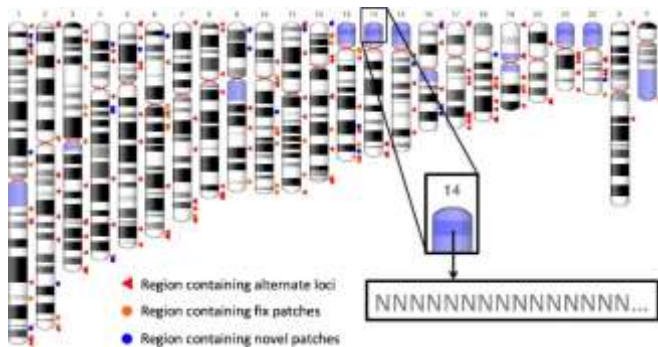
```
ATGTTTCGCATCACCAACATTGAGTTTCTTC  
CCGAATACCGACAAAAGGAGTCCAGGGAA
```

The file can contain one or many DNA sequences. There are lots of other formats, but fasta is the most common.

### About Grch38

The genome reference consortium (GRC) is providing the best possible reference assembly for human to generate multiple representations for regions that are too complex to be represented by a single path represented as Grch38 [16] as illustrated in figure 1. The GRC is stays silent due to its

mission of improving the human reference assembly genome, errors to be corrected and ensure by adding sequencing that provides the representation of human genome that meet the basic and clinical [2] needs of research. While the evaluation has new models and contains the sequence of a reference assembly of human and their process.



**Fig 1.** Shows the Grch38 data arrangement  
**Principle gathering**

- **ChromosomesCollected** for hg38 that contains the ranges of chr1–chr22, chrX, chrY and Mitochondrial (chrM).
- **Unlocalized** is a known belongs of a specific chromosome for orientation and are identified by random suffix.
- **Unplaced** is a unknown origin of sequences that are identified by chrU prefix.

## Mutation

Mutation specify how it makes the mutation [17] child by applying the random changes for a small individual. In mutation it provides the diversity to give space for searching broader. The mutation can specify the problems by taking the random add on distributed Gaussian mean with 0 at each entry of a vector parent which is unconstrained function. The Scale, Shrink are determined as the distribution of standard deviation.

The generation of standard deviation is determined by parameter of scale. The parent vector is coordinated at same intervals of standard deviation with vector  $v$  and its initial range is set as 2 by 1 and it derivation is given as  $Scale*(v(2)-v(1))$ .

If vector  $v$  initial range is set by two rows and the column as number of variables, at the coordinate  $i$  the standard deviation of parent vector is derived as  $Scale*(v(i,2) - v(i,1))$ .

The Shrink standard deviation generates to control the parameters of Shrink. If the vector initial range is set with 2 by 1 at  $K^{th}$  generation of standard deviation, parent vector is same at all coordinates generates  $\sigma_k$ , and the recursive formula is derived as

$$\sigma_k = \sigma_{k-1} \left( 1 - Shrink \frac{1}{kGenerations} \right)$$

The vector with two rows and columns as number of variables, coordinate  $i$  as standard deviation of parent vector at  $K^{th}$  generation and the recursive formula of  $\sigma_{i,k}$  is derived as

$$\sigma_{i,k} = \sigma_{i,k-1} \left( 1 - Shrink \frac{1}{kGenerations} \right)$$

If 1 is set as Shrink, then the shrink algorithm will coordinate each value of standard deviation until it reaches the last generation or become 0. The shrink causes the grow of standard deviation at negative value. The scale, shrink has the default value as 1.

## Algorithm of Machine Learning

To improve the computer experience the machine learning (ML) algorithms are adapted. These machine learning is the union of artificial intelligence that are concerned to develop the algorithms and techniques. It has a so many application areas like search engines, natural language process, stock market analysis, bioinformatics, medical diagnosis etc. The amount required to analyse the biological data is exploded through many machine learning algorithms that are implemented for data explosion. Hence, ML becomes an important tool for bioinformatics [15][18].

## III.EXPERIMENTAL ANALYSIS

Primarily, install the required packages called Bio python in python programming language because this a tool is used for computational molecular biology. Python is an object oriented, interpreted, flexible language that is becoming increasingly popular for scientific computing because it is easy to learn and has a very clear syntax and can easily extended with modules. For analysing the normal and abnormal genome initially taken the GRCh38 DNA sequencing which is a normal clinical data [8] used to compare the data. For comparison different types of genome sequencing is taken. Then



analysed the gene sequencing by counting the A,T,C,G's in a genome and evaluate the GC content in percentage to know whether the given genome is human genome or not because the human consists 30% to 60% [14] range across 100 kbp as illustrated in table 1. The GC ratio [5] in a gene is a coding region that evaluates the length of sequence that is highly proportional to G,C content. Mainly the GC content is calculated because it has strong interaction than A, T which has ability to form a three hydrogen bonds whereas A, T has two hydrogen bonds due to this melting point is high for G, C content along with a longer piece of DNA. DNA has GC content high then it hard to perform the amplification of PCR, and also hard to design a primer and specifies is as long enough. The DNA polymerases optimal temperature is maintained below the melting point.

Secondly, mutation is applied on the genome sequencing because DNA sequence may make some error. For study, a DNA base of a gene is compared with each other called as mutation or variation in a gene. The general code has integral redundancies this does not make much effect on protein that made a gene. In other case the error might in the base of codon that specifies the amino acid in a protein. Because amino acid is not a crucial part of the protein. If protein has the crucial part of amino acid, then it may be defective and not works as well and this mutation is done for different genome by taking the reference clinical genome Grch38 and found the genome sequence mutation as shown in figure 2.

**Table 1.** GC% of a genome classification

NAME	GC%
<b>Normal Clinical Data</b>	47.010464062922175
<b>Genome 1</b>	35.304752066115704
<b>Genome 2</b>	39.76670201484623
<b>Genome 3</b>	40.70054945054945
<b>Genome 4</b>	60.61759516007058



**Fig 2. Mutation count**



**Fig 3.** Sequence alignment to find healthy and unhealthy genome

Finally, to conclude the normal and abnormal genome sequencing by applying the blast tool which uses the machine learning algorithm in python to alignment the sequence to know the genome is healthy or unhealthy. Various approaches are applied to predict taken the homo sapiens whole genome sequencing called Grch38 which is a clinical data for finding the alignment of sequencing to implement accurate statistical model of sequence and makes the efficient search of a large sequencing database by Blast tool [13]. The sequencing search method is based on the extension of the statistics

unmapped local alignment [7] for High Scoring Segment Pairs (HSP) and is highly efficient at searching as it uses heuristics to reduce the search space. And uses PSI Blast tool for finding the position specific score matrix to form a multiple alignment. This captures the preserved pattern in alignment and stores the matrix score of each position in the alignment highly conserved positions receive high score zero. The profile is used to compare the clinical data and genome sequencing data to search and detect the sequence matching using the position specific matrix. The evaluated sequencing from the search threshold score is specified and refined for another round of search. The process is continued until the state has no new sequences are detected this makes PSI more capable of detecting the distant sequence alignment. This evolution of whole genome sequencing alignment is as illustrated in figure 3. This PSI Blast user terms are in two page format and by seeing we can analyse the genome is healthy or not.

By performing the sequence alignment of whole genome sequencing identified regions which has similar and consequence function relationship between the sequences on our different genome data. The clinical data is a healthy genome by using this data we compare aligned sequences of genome are represented within a matrix and by seeing gaps we conclude it as an unhealthy genome.

#### IV. CONCLUSION

Every genome contains unique mapping sequence are in small number of individuals then assemble these together to form a complete sequence. For analysis these sequence firstly find the genome sequence is homo sapiens by calculating the A, T, G, C count and GC percent due to Genome melting rate is limited and it is well recognized and not essentially follows the variance melting is only answerable for connection among GC under representation content in human genome. Then by using the blast tool the genome is compared with clinical data and formed the sequence alignment. By sequence alignment of whole genome sequencing identified regions which has similar and

Consequence function relationship between the sequences. Aligned sequences of genome are represented within a matrix and by seeing gaps we conclude it an unhealthy genome. In further by comparing sequence can be analysed to determine the disease or what is the problem in genome.

#### REFERENCE

- [1] Bloss, C. S. *et al.* A genome sequencing program for novel undiagnosed diseases. *Genetics in medicine: official journal of the American College of Medical Genetics* **17**, 995–1001, (2015).
- [2] Lee, H. *et al.* Clinical exome sequencing for genetic identification of rare Mendelian disorders. *Jama* **312**, 1880–1887, (2014).
- [3] Lam, H. Y. *et al.* Performance comparison of whole-genome sequencing platforms. *Nature biotechnology* **30**, 78–82, (2012).
- [4] Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nature genetics* **46**, 912–918. (2014).
- [5] Adams, R. L., and R. Eason. Increased GC content of DNA stabilizes methyl CpG dinucleotides. *Nucleic Acids Res.* **12**:5869–587. (1984)
- [6] Bernardi, G. : The human genome: organization and evolutionary history. *Annu. Rev. Genet.* **29**:445–476. (1995)
- [7] Altschul SF; Gish W (1996). *Local Alignment Statistics. Meth. Enz.* Methods in Enzymology. 266. pp. 460–480.
- [8] Wright C, Middleton A, Burton H *et al.*: Policy challenges of clinical genome sequencing. *BMJ*. Nov (2013).
- [9] The fantom project charts an atlas of gene activity over the human body. [www.science.ku.dk/english/press/news/2014/fantom/](http://www.science.ku.dk/english/press/news/2014/fantom/). Accessed: 2017-08-08.
- [10] Greenblatt MS, Bennett WP, Hollstein M, Harris CC. Mutations in the p53 tumor suppressor gene: clues to cancer etiology and molecular pathogenesis. *Cancer Res.* **54**(18):4855–78. (1994)
- [11] Treangen TJ, Salzberg SL. Repetitive DNA and next-generation sequencing: computational challenges and solutions. *Nat Rev Genet.* **13**:36–46. (2011)
- [12] Eid J, Fehr A, Gray J, *et al.* Real-time DNA sequencing from single polymerase molecule. *Science*. **323**:133–138. (2009)

- [13] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, and D.J Lipman. A basic local alignment search tool. *Journal of Molecular Biology*, 215: 403-410. (1990)
- [14] Masahiko Mizuno, Minoru Kanehisa, Distribution profiles of GC content around the translation initiation site in different species, *Federation of European Biochemical Societies*, (1994)
- [15] Cock,P.J., Antao,T., Chang,J.T., Chapman,B.A., Cox,C.J., Dalke,A., Friedberg,I., Hamelryck,T., Kauff,F., Wilczynski,B. et al.:Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25, 1422–1423.(2009)
- [16] Schneider,V.A., Graves-Lindsay,T., Howe,K., Bouk,N., Chen,H.C., Kitts,P.A., Murphy,T.D., Pruitt,K.D., Thibaud-Nissen,F., Albracht,D. et al.: Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. *Genome Res.*, 27, 849–86.(2017)
- [17] Amar D, Izraeli S, Shamir R. Utilizing somatic mutation data from numerous studies for cancer research: proof of concept and applications. *Oncogene*. 36:3375–83. (2017)
- [18] Bhargavi, P., Lohitha Lakshmi, K., Jyothi, S. (2020), Gene sequence analysis of breast cancer using genetic algorithm *Emerging Research in Data Engineering Systems and Computer Communications* pp 177-190 (AISC, volume 1054), Springer Series