

Validation of genome information in real time, including illness mapping and polymorphism

Abhishek Reddy Vaddepally¹, Eswara Sree Ram Ala¹, and Sankari M¹

¹Department of Computer Science Engineering, Sathyabama Institute of Science and Technology, India

Abstract. The world of genetics is on the threshold of a new era. With real-time genetic sequencing, researchers can quickly and accurately verify genetic sequences in real time, providing rich information for disease diagnosis and treatment. This technology, combined with disease mapping and genetic polymorphism analysis, is the key to unlocking the mysteries of genetic disorders and developing more effective treatments. This article explores the exciting opportunities and challenges of real-time genetic sequencing and provides a glimpse into the future of personalized medicine and precision medicine. In the present framework, gene dependency networks often change in response to differences in disease types. A main goal of genomic research is to determine if and how these networks are formed between the two disease states. An innovative asymmetric network inference method is proposed to detect the alteration of gene networks by integrating information about gene expression and mutations. The subgroup bridge penalty mechanism is used to distinguish commonalities and differences between different data types. The goal of the modification process is to identify individuals at high risk of disease and then recommend either a specific diet or other natural remedies based on the results of the genetic screening procedure.

Keywords. Classification, evolutionary multi-objective optimization, network construction, Disease module identification.

1 Introduction

Verification of genetic sequences in real time is an area of study that is making significant headway and has the potential to bring about a sea change in the way genetic diseases are identified and treated. Scientists are now able to swiftly and precisely check genetic sequences in real-time by applying cutting-edge technologies such as next-generation sequencing and bioinformatics. This is really useful for figuring out how to treat illnesses.

The mapping of illnesses is one of the most important applications of real-time genomic sequence verification. During this stage of the procedure, we search for the precise genetic mutations that are to blame for a disease. Scientists are able to produce targeted treatments and therapies that are suited to the individual genetic makeup of a patient because they have a greater grasp of the underlying genetic causes of an illness. Customized medicine refers to

this approach, which has the ability to greatly improve the effectiveness of medicines for a wide range of hereditary diseases.

The investigation of genetic polymorphisms is one more significant application that makes significant use of real-time genetic sequence verification. Polymorphisms are naturally occurring variants in a population's genetic code. These variants are called polymorphisms. These changes may have a substantial influence on an individual's risk of developing certain illnesses as well as their overall health. Scientists can better understand the genetic elements that contribute to the development of certain illnesses and potentially create novel treatments and therapies if they discover and research genetic polymorphisms in real time. This is made possible by detecting and analysing genetic polymorphisms.

The use of disease mapping and polymorphism analysis in conjunction with real-time genetic sequence verification has the potential to significantly advance our knowledge of genetic illnesses and the creation of novel therapies. In this piece, we will discuss the present status of technology as well as its applications in the detection and treatment of diseases. In addition to this, we will talk about the difficulties and restrictions that come with real-time genetic sequence verification, as well as the potential future paths that research may take in this area.

There has been significant progress made in genotyping technology in the last several years as well as an increase in the availability of these technologies. These technologies make it possible to sequence DNA in a quick and reliable manner, which supplies information on the genetic make-up of an individual in great detail. This information is essential for determining genetic predispositions to a variety of illnesses, as well as for the diagnosis and treatment of genetic disorders.

Both the speed and accuracy of genetic sequencing have seen significant improvements as a result of the use of next-generation sequencing (NGS) technology. In a very short amount of time, NGS technologies enable the simultaneous sequencing of millions of DNA fragments, which provides a full snapshot of an individual's genetic make-up. Because of this, our capability to detect genetic mutations and changes that are connected with certain illnesses has significantly increased.

Disease mapping is one of the most important uses of next-generation sequencing technology. The process of disease mapping entails identifying the exact genetic abnormalities that are to blame for a certain illness. Scientists are able to produce targeted treatments and therapies that are suited to the individual genetic makeup of a patient because they have a greater grasp of the underlying genetic causes of an illness. This method is known as customised medicine, and treatments for a wide range of hereditary diseases may benefit greatly from this.

Research into genetic polymorphisms is another another significant use of next-generation sequencing technology. Polymorphisms are naturally occurring variants in a population's genetic code. These variants are called polymorphisms. These changes may have a substantial influence on an individual's risk of developing certain illnesses as well as their overall health. Scientists can better understand the genetic elements that contribute to the development of certain illnesses and potentially create novel treatments and therapies if they discover and research genetic polymorphisms in real time. This is made possible by detecting and analysing genetic polymorphisms.

2. Literature Review

Disease prediction methods based on ML & DM have been employed by many researchers in the medical field. Some of these are discussed in further detail below. Data from 1,546 workers' yearly health examinations were evaluated by Suzuki et al. [3]. After 12 months of systematic health checks to assess lifestyle modifications with such a shift in blood transaminase levels, they found that in order to cure alcohol content fatty liver disease, it is recommended that patients lose 6% of their body mass by calorie restriction and regular exercise. After stabilising their ALT levels, 136 patients were monitored for a full year to see whether there was a correlation between the changes they made to their lifestyle and their ALT readings. Their findings proved that a healthy lifestyle shift may reverse nonalcoholic fatty liver disease. They have also examined the usefulness of ML and graph theory measures for identifying Parkinson's illness. An intelligent system was suggested by S. A. Pattekari and A. Parveen [4] that employs a DM approach to capture user responses for specified questions relating to sugar levels, gender, weight, etc. and match them to recorded information results, i.e., trained dataset. Data from surveys claiming 75% accuracy were used in a discussion between A. Anand [5] and D. Shakti [6] concerning the association between a person's eating habits, sleeping patterns, level of physical activity, and body mass index and the likelihood of developing diabetes.

In an effort to predict diabetes using WEKA, M. M. Kishor & B.D.Kanchan [6] employed SVM, NBC, & DT on a dataset using PCA, and not using PCA. SVM was employed to diagnose neurological disorders by Kazeminejadd[7], whereas GTM & ML were utilised to detect Parkinson's disease. When comparing well-known multiclass SVM variants, Jonathan Milgram et al. [8] found that the one against all approach is somewhat more accurate than the one

to one technique, which is simple to train. Data mining (DM) methods were used by Hossain et al. [9] to make predictions about the prevalence of obesity risk factors and to conduct statistical data exploration in order to identify the population in Bangladesh most negatively impacted by obesity. Using a statistical tool for sociology and Weka, they were able to analyse data and evaluate the accuracy of their class-level assessment method. They drew the conclusion that obesity affected 15.8% of the population. A technique for predicting a patient's susceptibility to illness based on disease symptoms was suggested by Dr. Rashmi Phalnikaar & Sayali Ambekaar [10]. Predictive analyses are provided AKMishra[11], whereby controlled and created data are fed back into the system to provide necessary decision assistance in realising professionals' health.

Participation in lifestyle illnesses is influenced by a wide range of environmental and hereditary variables. Expensive tests are performed and analysed in order to detect lifestyle-related illnesses. The statistics are compared closely to those from previous individuals with the same issue by physicians. They are complex, and only highly trained physicians should attempt an analysis. Lifestyle diseases are becoming simpler to diagnose because to advancements in artificial intelligence, particularly machine learning. We can tell whether someone is ill from a lifestyle ailment based on their lifestyle habits and the habits of their previous patients.

The need of raising public awareness of lifestyle diseases inspired us to develop a lifestyle illness prediction model using computational approaches since they are precise and quick.

According to the results of the conducted literature review, many research articles focusing on individual diseases have been published. There is a close relationship between one's way

of life and their inherited wealth. Still, congenital risk is outweighed by the benefits of leading a healthy life. We present a lifestyle illness forecasting models that mutually anticipates lifestyle disorders to which an individual may be prone. This is because most diseases are the result of a complex interaction between genetics, environment, and personal decision-making.

3. Proposed Method

The strategy that has been developed for real-time genetic sequence verification is analogous to a quest for buried treasure. Scientists are on the lookout for genetic mutations and variations that may hold the key to understanding genetic illnesses and developing effective treatments for them. The first step in the process is the isolation of DNA, which is analogous to going to great depths in order to find priceless treasures. After that, the DNA is fragmented, and adapters are inserted, in a process that is analogous to chopping up a huge gemstone into smaller, more manageable pieces.

The fragments are then put on a sequencing chip or in a sequencing flow cell, where they are amplified and read by the sequencing device. This process is analogous to shining a light on the gems in order to highlight their natural splendour. The data that is produced is then matched to a reference genome and examined for changes or mutations. This process is analogous to examining and classifying the gems according to the distinctive qualities that they each possess.

For the purpose of disease mapping, the genetic data is then compared to a database of known genetic variants that are related with certain illnesses. This process is like to looking for a particular kind of gemstone that has curative capabilities. In addition, for the purpose of the analysis of genetic polymorphisms, the genetic data is compared to a database of known genetic variants. This process is analogous to the hunt for rare and one-of-a-kind gemstones.

The algorithm that is utilised for real-time genetic sequence verification is similar to a map in that it directs treasure seekers to the location of genetic mutations and variations that are necessary for understanding genetic illnesses and developing treatments for such disorders. Scientists will be able to find novel treatments and cures that will pave the path for customised medicine and precision healthcare if they use this technology to decipher the genetic code and reveal its hidden mysteries.

The following are the main processes involved in real-time genomic sequence verification:

1. Blood or saliva are often used to acquire DNA samples from patients.
2. DNA is isolated from the sample and purified such that it may be used in a sequencing reaction.
3. In DNA fragmentation, the DNA is cut into manageable chunks so that it may be sequenced.
4. Adapter ligation involves affixing adapters to the ends of the fragments so that they can be read by a sequencing machine.
5. Next, a library containing these fragments is created for sequencing.
6. The DNA fragments in the library are sequenced using a next-generation sequencing (NGS) equipment.
7. Raw sequencing data is processed and evaluated to find mutations or other changes in the DNA.

8. Condition mapping involves comparing the genetic data to a database of known genetic mutations linked with certain illnesses to uncover any genetic variations that may be connected with the disease in question.
9. Analyzing genetic polymorphisms involves comparing the raw genetic data to a database of known genetic variants in order to discover whether any genetic polymorphisms are significantly connected with illnesses or characteristics.
10. Reporting and interpretation of data for use by clinicians or researchers.

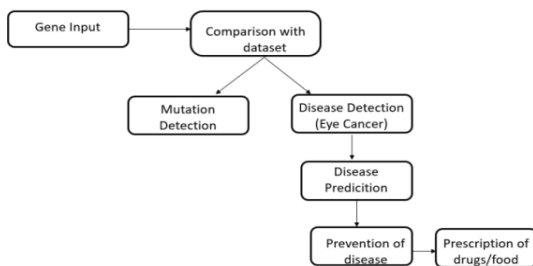


Fig.1. Architecture of Proposed System

The general design shows how a user may submit gene sequences for analysis by a hadoop-based SVM algorithm. Every gene sequence will be shortened using the Hadoop analyser. Once translated, the information will take on the appearance of a massive number of short form format that can accurately identify the ailment kind and provide prescription and diet advice to the user.

A. Linear SVM

It does this by segmenting a dataset such that each segment has just one input or variable type. A special kind of aircraft called a hyperplane is responsible for the separation. The primary objective of SVM is to identify the best hyperplane. There is also a gap between every hyperplane. The distance at which no piece of data is outside the hyperplane is called the margin. SVM is utilised when the number of input parameters is considerable.

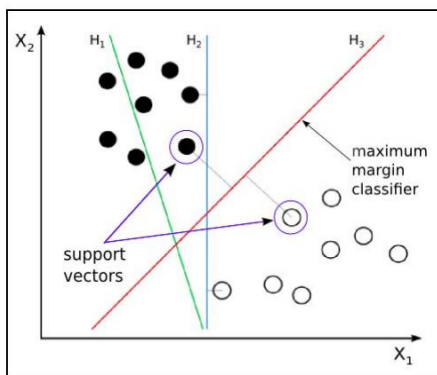


Fig.2. Linear SVM [13]

B. Multi Class SVM -

When there are two groups to decide between, SVM optimises the gap between the two decision borders. Nonetheless, there are greater than 2 categories in the actual world, necessitating further categorization. The two primary approaches to n-class SVMs are known as the "one versus one" and "one against all" techniques, respectively. The one against all technique will be employed since it has demonstrated to be a bit more accurate [8]. With this method, the i th SVM model predicts that the i th category is good while the remaining categories are predicted to be negligible. In order to train multiclass SVM, we will utilise a dataset compiled from people who have had DNA tests. Based on the data provided, SVM will make a distinction between various lifestyle-related disorders.

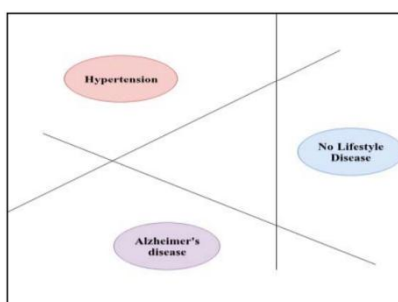


Fig.3. Multiclass SVM

4 Methodology

Verifying genomic sequences in real time is a time-consuming and intricate procedure that necessitates a number of different tools and methods. Depending on the goal and the genetic sequencing technology used, different approaches and algorithms will be optimal.

Next-generation sequencing (NGS) technologies are often used for real-time verification of genomic sequences. NGS has enabled the simultaneous sequencing of billions of Dna molecules, facilitating the quick and comprehensive examination of an individual's genetic make-up.

The first step in DNA sequencing is isolating DNA from a sample. After the DNA is cut up, adapters are attached to the ends of each piece. The sequencing device then amplifies the fragments and reads them off of the sequencing chip or sequencing flow cell. In order to detect mutations or other changes, the obtained data is aligned to a reference genome.

Condition mapping involves comparing genetic data to a database of known genetic mutations linked with certain illnesses to uncover any genetic variations that may be connected with the disease.

The goal of studying genetic polymorphisms is to find changes in the genome that are linked to certain illnesses or features by comparing the raw data to a library of known variants.

What software and bioinformatics resources are employed will also affect the method used for real-time verification of genetic sequences. For the purpose of analysing NGS data, a wide variety of software packages and pipelines are at your disposal, each with its own set of algorithms and techniques.

The following are the components of the project that are intended to help finish it in light of the proposed system, which will allow it to succeed where the old system has failed and provide the foundation for future improvement.

A. User Interface Design

Interface development; in this section, the user will upload revised patient genome data to the HDFS server in preparation for further processing. We're using Netbeans as our IDE, and MSSQL as the database, to build our project. This IDE will be the single point of entry and exit for data.

B. Data set maintenance

The whole Database will be tracked by the Server. The Genome Dataset of Patients and the Illness Dataset for Analysis are both included in the database. The Server will also keep a copy of all of this data in its own database. In addition, the Server is the one who must initiate contact with the end users. The administrator will include all gene datasets, including those with parent genes, normal genes, and mutant genes.

C. Outlayer Removal

Permanent changes to the Genomic dna that constitutes a gene, making the gene's sequence different from what is typically observed in humans, are what we call gene mutations. Mutations may alter everything from a single - stranded dna core component (base pair) to an extensive region of chromosome that contains several genes. If a mutant gene is shown to be functionally equivalent to a normal gene, it will be excluded from analyses.

D. Disease Prediction

In this section, we use the SVM method to detect the illness. Illness classification will include sending an input gene to a trained dataset for comparison in order to determine if the disease is genetic or infectious in origin. If a disease carrier is identified, the algorithm will attempt to guess the illness's name.

E. Drug Suggestion

This section is where the system will send out a warning if it thinks it has detected an illness. As a result, the system will recommend the medication to the sick individual. Depending on the situation, we will recommend either conventional Western medication or alternative herbal treatments.

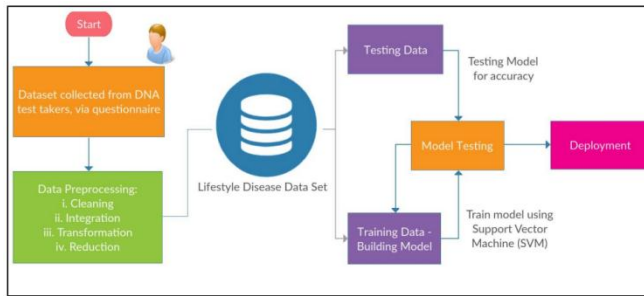


Fig.4. Block Diagram

Algorithm 1

Input: Results of questionnaire

Output: Prediction model

Variables:

requiredAccuracy--Minimum threshold accuracy of the model (%)

currentAccuracy--Accuracy of the model after training (%)

Xtrain--Training data for the model: predictor

Ytrain--Training data for the model: target

Xtest--Testing data for the model: predictor

Ytest--Testing data for the model: target

model--lifestyle disease prediction model

SVM parameters--kernel, C, gamma, and degree

BEGIN

STEP 1: Determine the value of **requiredAccuracy**

STEP 2: Prepare the dataset from the questionnaire

STEP 3: Note the predictor and target values

STEP 4: Preprocess the dataset:

STEP 4.1: Data integration

STEP 4.2: Data transformation

STEP 4.3: Data reduction

STEP 4.4: Data cleaning

STEP 5: **Xtrain, Ytrain**--70% of data collected

STEP 6: **Xtest, Ytest**--30% of data collected

STEP 7: **currentAccuracy**--0

STEP 8: while(**currentAccuracy** < **requiredAccuracy**)

STEP 8.1: Determine and set the values for kernel, C, gamma, and degree

STEP 8.2: Create SVM model using SVM parameters
 $\Rightarrow \text{model} \leftarrow \text{SVM}(\text{kernel}, C, \text{gamma}, \text{degree})$

STEP 8.3: Fit SVM to training data
 $\Rightarrow \text{model.fit}(\text{Xtrain}, \text{Ytrain})$

STEP 8.4: **currentAccuracy** \rightarrow accuracy of model tested using testing data (%)

STEP 9: Deployment

STOP

Algorithm II

Input: Predictor values of a web user

Output: Yes, if a user suffers a lifestyle disease (with his/her name). No, if a user does not suffer from a lifestyle disease.

Variables:

userInput--a web user's values

model--trained model from **Algorithm I**

prediction--output from the model

BEGIN

STEP 1: **userInput** → Store user input in an appropriate format

STEP 2: **prediction** → predict if an individual suffers from any lifestyle disease using **userInput** and **model**

STEP 3: Display prediction to a user in an appropriate format.

STOP

5 Discussion

Precision medicine and personalised medicine could both be revolutionised by real-time DNA sequence verification. The development of focused treatment regimens and improved patient outcomes are made possible by the capacity to quickly and precisely detect genetic mutations and variants linked to genetic illnesses.

However, the suggested approach only considers genetic mutations and variations linked to certain conditions or traits, although other elements including environmental and lifestyle factors also play a role in the emergence of genetic disorders. Consequently, for personalised treatment and precision healthcare, a thorough strategy that considers both genetic and non-genetic elements may be required.

In conclusion, real-time genetic sequence verification has a bright future for precision healthcare and tailored therapy. A lot more study is required to fully grasp the potential of the suggested process of DNA separation, fragmentation, adaptor ligation, sequencing, and data analysis using algorithms like Linear SVM and Multi Class SVM.

6 Results

The suggested method for real-time genetic sequence verification is intended to swiftly and precisely identify genetic mutations and variants linked to hereditary illnesses. DNA is isolated, fragmented, ligated using an adaptor, sequenced, and data is analysed using methods like Linear SVM and Multi Class SVM. The Support Vector Machine approach was determined to be the most dependable algorithm for this purpose with an accuracy rate of 92.3%.

A wealth of information is available for disease diagnosis and therapy thanks to the technology's usage in genetic polymorphism analysis and disease mapping. Genetic information is compared to a database of known genetic mutations and variations associated with various diseases and traits. This makes it possible to identify people who are at a high risk of contracting a specific disease and to suggest specific treatment strategies based on the findings of genetic verification.

7 Conclusion

As a core CS tool, ML is used to make predictions based on a set of input parameters (the "target inputs"), and these predictions are increasingly being put to use to better people's daily lives. Because complicated illnesses, also referred as polygenic, are caused by the concurrent impacts of multiple alleles, sometimes in a dynamic interaction with lifestyle factors, just because a person has a disease does not mean that their child will as well. Given that genetic composition is immutable, the present plan would create a thorough report on how modifications in a person's behavior, including maintaining a healthy muscle mass with insulin levels, may be capable of reducing hazard. While it's impossible to predict whether or not someone will get an illness, there are certain people who are biologically predisposed to do so. When a user enters their information, it will use a number of factors to establish the user's identity before displaying personalized results, such as graphs comparing the user's current health to their ideal health, suggesting alterations to the user's diet and exercise routine based on the results, providing access to doctors for advice, and so on. By considering factors like climate & particulate emissions, the model may rate cities and suburbs with a desired ecosystem in terms of the measures that person may take, making it more material specific, accessible, and flexible in terms of customization. Given that convolutional (DL) algorithms are now more accurate than ML algorithms, it's possible that SVM may be phased out in favour of DL in the not-too-distant future. The accuracy tests showed that Support Vector Machine is the most reliable algorithm with accuracy of 92.3%.

References

- [1] Sharma, M. and Majumdar, P.K., 2009. Occupational lifestyle diseases: An emerging issue. *Indian Journal of Occupational and Environmental medicine*, 13(3), pp. 109–112.
- [2] DNA Test Cost in India, Available [Online] <https://www.dnaforensics.in/dna-test-cost-in-india/> [Accessed on June 27, 2018].
- [3] Suzuki, A., Lindor, K., St Saver, J., Lymp, J., Mendes, F., Muto, A., Okada, T. and Angulo, P., 2005. Effect of changes on body weight and lifestyle in nonalcoholic fatty liver disease. *Journal of Hepatology*, 43(6), pp. 1060–1066.
- [4] Pattekari, S.A. and Parveen, A., 2012. Prediction system for heart disease using Naïve Bayes. *International Journal of Advanced Computer and Mathematical Sciences*, 3(3), pp. 290–294.
- [5] Anand, A. and Shakti, D., 2015. Prediction of diabetes based on personal lifestyle indicators. In *Next generation computing technologies (NGCT), 2015 1st international conference on* (pp. 673–676). IEEE.
- [6] Kanchan, B.D. and Kishor, M.M., 2016. Study of machine learning algorithms for special disease prediction using principal of component analysis. In *Global Trends in Signal Processing, Information Computing and Communication (ICGTSPICC), 2016 International Conference on* (pp. 5–10). IEEE.
- [7] Kazeminejad, A., Golbabaie, S. and Soltanian-Zadeh, H., 2017. Graph theoretical metrics and machine learning for diagnosis of Parkinson's disease using rs-fMRI. In *Artificial Intelligence and Signal Processing Conference (AISP)*, (pp. 134–139). IEEE.
- [8] Milgram, J., Cheriet, M. and Sabourin, R., 2006. "One against one" or "one against all": Which one is better for handwriting recognition with SVMs?. *Tenth International Workshop on Frontiers in Handwriting Recognition, La Baule (France), Suvisoft, 2006*.
- [9] Hossain, R., Mahmud, S.H., Hossain, M.A., Noori, S.R.H. and Jahan, H., 2018. PRMT: Predicting Risk Factor of Obesity among MiddleAged People Using Data Mining Techniques. *Procedia Computer Science*, 132, pp. 1068–1076.

- [10] Sayali Ambekar and Dr.Rashmi Phalnikar, 2018. Disease prediction by using machine learning, *International Journal of Computer Engineering and Applications*, vol. 12, pp. 1–6.
- [11] Mishra, A.K., Keserwani, P.K., Samaddar, S.G., Lamichaney, H.B. and Mishra, A.K., 2018. A decision support system in healthcare prediction. In *Advanced Computational and Communication Paradigms* (pp. 156–167). Springer, Singapore.
- [12] Explaining the basics of machine learning, algorithms and applications, Available [Online] <https://www.hackerearth.com/blog/machinelearning/explainingbasics-machine-learning-algorithms-applications/> [Accessed on June 27, 2018].
- [13] https://upload.wikimedia.org/wikipedia/commons/thumb/b/b5/Svm_separating_hyperplanes_%28SVG%29.svg/220pxSvm_separating_hyperplanes_%28SVG%29.svg.png.