

Cancer Genome Classification Using Machine Learning Algorithm

Tahmeena Fatima¹, S Jyothi², D M Mamatha³, K Srujan Raju⁴

¹Research Scholar, Dept. of Computer Science and Engineering, tahmi.fatima18@gmail.com

²Professor, Department of Computer Science, jyothi.spmvv@gmail.com

³Professor, Department of BioScience & Sericulture, prof.mamatha@gmail.com

⁴Professor, Department of Computer Science & Engineering, ksrujanraju@gmail.com

^{1,2,3}SriPadamavatiMahilaVisvavidyalayam, Tirupati, India

⁴CMR Technical Campus, Hyderabad, Telangana, India

Abstract: *In present scenario, cancer is the universal problem in the world because it effects the people in the world. The growth of uncontrolled and abnormal cells in the body is known as cancer. These cells with abnormal named as cells of cancer, cells of malignant or cells of tumour. To recognize the genomic origin of tumour cell explosion and the development of the cancer genome the normal cell and tumour cell are compared to find which type of cancer it is. The growing accessibility and development rate of “Big Data” derivative from of numerous omics exposed a novel frame to expand diagnoses of clinical or cancer therapeutics, however there are several experiments in effective analysis and explanation of such complex and big data. The significance of categorizing cancer patients hooked on low or high risk sets consumes directed several research groups, from the bioinformatics and the biomedical area, to learn the solicitation of machine learning (ML) techniques. Thus, these methods have remained consumed as an intention to classic advancement and treatment of conditions of cancerous. In adding, the capability of ML implements to discover vital features from compound datasets exposes their prominence.*

Keywords: *cancer cells, machine learning algorithms, MongoDB, BWA, Genome sequencing.*

1 INTRODUCTION

In medical, Whole Genome Sequencing (WGS) is grown as main interest for the diagnostics. Next Generation Sequencing has begun as a exhortation which incorporates new DNA sequencing methods, permitting researchers to categorize a whole human genome is linked to the standard Sanger sequencing machinery which necessary over a period for accomplishment while the human genome was primary sequenced.

Within specific organism gene is a complete set and the study of molecular biology is known as genome which is a branch of Genomics. Currently, machine learning plays a vital role in the development of genomics area.

The capability to sequence DNA affords mythos capability to “read” the blueprint of genetic that guides all the actions of a alive organism. To affords situation, the dominant view of biology is shortened as the path beginning as DNA towards RNA towards Protein. DNA remains self-possessed of base pairs called nucleotides which has four basic units A, T, C, G. A combines with T and C combines with G. DNA is prearranged into genetic material and humans devise a entire of 23 combines. [1]

Chromosomes are further ordered into sections of DNA named as genes which create encode proteins. A genome is a sum of gene that possess organism. Humans devise coarsely 20,000 genes and 3 billion

base sets. Stimulatingly, only 2% of encodes protein in the human genome and this is a vital part of effort in research and the commerce of genomics.[2]

Precision medicine can be closely connected to Genomics [3]. By 2023 87 billion dollars are projected with a market size that area of Precision Medicine (also named as personalized medicine) is a methodology to care patient that comprehends genetics, environment and behaviour through a aim of instigating a patient or population explicit treatment involvement compare to a one scope turns all method. For illustration, to lessen the danger of difficulties, a separate who wants a blood transfusion would be harmonized to a giver who shares the identical blood group instead of a casually selected giver.

To analyse Genomics deeply, machine learning can support researchers understand genetic dissimilarity. Specially, algorithms are aimed based on outlines recognized in hefty genetic data sets which are formerly interpreted to computer simulations to support clients understand how genetic dissimilarity disturbs critical cellular procedures. Illustrations of cellular procedures contain the metabolism, DNA restoration, and cell progression. Disturbance to the standard working of these paths can theoretically source of diseases like cancer. [4]

Although prominence is regularly located on selecting the finest learning procedure, researchers have initiated some of the best stimulating queries rise out of nothing of the accessible machine learning procedures performance to balance. Utmost of the period this is a difficult with training data, nevertheless this also happens when occupied through machine learning in novel fields.

Machines study are beneficial to humans since by all of their handling power, they are capable to more rapidly highlight or discover outlines in big data that would consume or else remained lost by human lives. Machine learning is a method that is used to improve human's capabilities to resolve problems and sort knowledgeable implications on a varied sort of difficulties, from assisting identify diseases to pending up through explanations for worldwide climate alteration.

2 Proposed method

2.1 Burrows–Wheeler algorithm

The Burrows–Wheeler algorithm is a method used to make information for use through data compression methods like bzip2. In 1994 David Wheeler and Michael Burrows are discovered this method at that time they are working at DEC Systems Research Centre in Palo Alto, California. It is created on an earlier unpublished revolution revealed by Wheeler in 1983. The method can be executed competently consuming a suffix array consequently attainment linear time complexity.[5]

Let Σ is represented as alphabet. $\$$ Symbol is not existing in Σ and is lexicographically lesser than entirely the symbols in Σ . A series $X = a_0, a_1, \dots, a_{n-1}$ is constantly concluded through $\$$ symbol (i.e. $a_{n-1} = \$$) then this sign is appeared at the end. Let $X[i] = a_i, i = 0, 1, \dots, n-1$, be the i^{th} symbol of X , $X[i, j] = a_i \dots a_j$ a substring then $X_i = X[i, n-1]$ a suffix of X . Suffix array of S of X is a combination of the numerals $0 \dots n-1$ like that $S(i)$ is the initiation point of the i^{th} lowest suffix. The X in BWT is represented as $B[i] = \$$ while $S(i) = 0, B[i] = X[S(i)-1]$ else. We also outline the measurement of string X as $|X|$ and consequently $|X| = |B| = n$. Figure 2 offers an illustration on how to build BWT and suffix array.

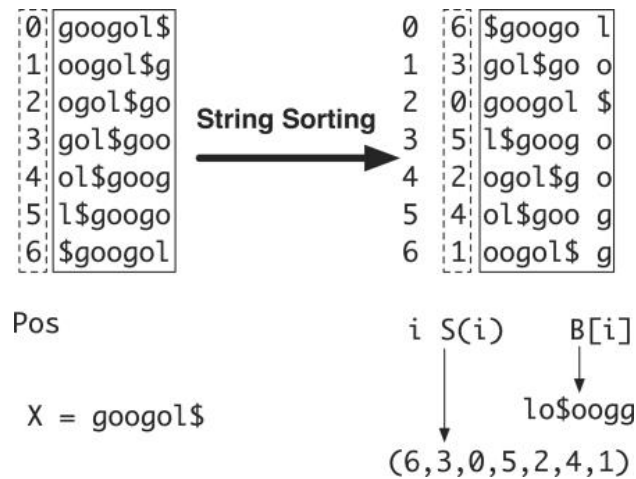


Fig 1: analysis of the BWA algorithm

The procedure displayed in figure 1 is quadratic in period and space. Though, this is not essential. In training, we regularly build the suffix array initial and formerly produce BWT. Furthermost procedures for building suffix array need minimum $n \lceil \log_2 n \rceil$ bits of occupied space, which expands to 12 GB for human genome [6]. Newly, Hon *et al.* (2007) presented a novel method that customizes n bits of occupied space and needs < 1 GB memory at ultimate period for building the human genome using BWT.

2.2 Machine learning

Machine Learning is the knowledge of receiving computers to study and performances similarly humans ensure and expand their knowledge over period in independent method through serving them records and data in the arrangement of clarifications and real word relations.

The definition of classifications summarizes the model objective or eventual goal of machine learning is stated by numerous researchers in the area. The determination of this object is to afford a reader minded business with skilled viewpoint on in what manner machine learning is definite and in what way it works. artificial intelligence and Machine learning stake the similar meaning in the observances of numerous fields. Though there remain about differential alterations readers must identify as fit.

Machine learning algorithms are in numerous forms [7], through hundreds distributed each time and they remain usually collected by each knowledge style (that is unsupervised learning, semi-supervised learning, supervised learning) or by resemblance in usage or meaning (that is regression, deep learning, classification, decision tree, clustering, etc.). Nevertheless of knowledge style or meaning, all arrangements of machine learning procedures involve of the resulting:

- **Demonstration** (computer can understand the set of language or classifiers)
- **Calculation** (also known as scoring/objective function)
- **Optimization** (examination process; habitually the highest counting classifier, for illustration; here are mutually off the ledge and convention optimization procedures used)

Machine learning technique has dissimilar tactics to study from consuming simple decision trees to gathering to layers of artificial neural networks (then latter it gives space to the deep learning) dependent on whatever duty you are demanding to achieve and the sort and sum of information that you devise existing. This vigorous situation played obtainable in submissions as variable as medicinal diagnostics or self-driving cars.

Research prepared when occupied on actual submissions often efforts growth in the arena, and causes are twofold: 1. Propensity to learn limitations and restrictions of present approaches 2. Scholars and

makers occupied with field authorities and leveraging period and knowledge to increase classification performance.

Occasionally this too happens by “accident” We can imitate classic groups, or blends of numerous learning procedures to recover accurateness.

In present days, genomics by uses Machine learning can impact numerous hints containing in what way genetic research is directed, in what way clinicians afford patient carefulness and the availability of genomics to characters involved in learning new things in what way their inheritance may affect their condition.

Consequently, the information analysis can experience through machine learning with essential to accompanied by training and clear descriptions of the efficacy and importance of this knowledge.

2.3 Support Vector Machine

SVM is one of the supervised learning method in ML which is associated with learning procedures that inspect information used for regression and classification analysis [8]. Resolute a stable of activity cases, individually show as standard to unique or the new of binary gatherings, in SVM making design to progress a model that apportions novel cases to unique gathering or the former, productions a non-probabilistic parallel undeviating classifier [9]. When facts remain categorized, supervised learning not possible, and an unsupervised learning method is required [10], which efforts invention normal gathering of the information to groups, and now map novel information to these formed groups. The clustering procedure which delivers an enhancement to the SVM is known as Support Vector Clustering (SVC) then it uses in trade applications moreover when realities are not categorized or when first some realities are categorized equally a pre-processing aimed at a grouping pass [11]. The appliance of categorizing the information hooked on dissimilar modules by definition a link which splits the preparation files into modules. Here remain a few straight hyperplanes, calculation of SVM attempts to augment the separation in the focal of the few classes that are mind boggling and this is said as edge augmentation. If the line makes the most of the space among the modules is recognized, the possibility to simplify fine to unobserved information is improved. There are 2 classifications in SVM:

2.3.1 SVM Linear (L-SVM)

In L-SVM the training data that classifies are disconnected by a hyperplane.

$$\frac{1}{m} \sum_{i=1}^m l(w \cdot x_i + b \cdot y_i) + \|w\|_2 \quad \text{eq1}$$

2.3.2 SVM Non-Linear (NL-SVM) or sigmoid

In NL-SVM it remains not thinkable toward discrete the training statistics with a hyperplane. For illustration, the training statistics for Face recognition contains a set of imageries that are aspects and alternative collection of imageries that remain aspects (in addition disputes all other imageries in the domain excepting faces). Under such circumstances, the preparation measurements are excessively perplexing and troublesome, making it impossible to find a delineation for each component vector [MT06]. Isolating the typical of countenances directly after the arrangement of non-confront is a many-sided assignment.

2.4 Classifier of Ada-boost

Ada boost method combines weak classifier procedure [12] to custom tough classifier. A lone procedure might categorize the objects unwell. However, if we syndicate numerous classifiers through collection of training set one each repetition and transfer right quantity of load in ultimate polling, we can consume good accurateness groove for complete classifier. Respectively feeble classifier is accomplished by means of a *random subset* of whole training set. [13]

Later preparation a classifier at some level, ada boost assigns [14] weight to a piece training point. Misclassified point is allotted advanced weight hence it seems in the training division of next classifier through complex possibility.

Subsequently each classifier is trained, the weight is allotted to the classifier as well grounded to accurateness. Additionally, accurate classifier is allotted higher weight so that it can consume extra influence on ultimate result.

Let see the parameters and mathematical formula.

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

$h_t(x)$ is the weak classifier of output t for x input

α_t is assigned classifier weight.

α_t is designed as follows:

$\alpha_t = 0.5 * \ln\left(\frac{1-E}{E}\right)$: straight forward classifier weight, it is built on the rate of error E . Originally, all the response training instance has same weightage.

2.5 Classifier of Naive Bayes

It is Bayes Theorem of classification method [14] through an supposition of disinterest between predictors. In general, classifier Naive Bayes adopts the occurrence of a specific feature in a class is dissimilar to the occurrence of some extra feature.

For illustration, if an apple is considered as a fruit if it has round, contains the 3 inches' diameter and red colour. Even if these structure is depending on one another or at the presence of the other structure all these features individually pay to the possibility to be an apple fruit and that is known as "Naïve".

It is very useful for large dataset and it can easily build. Laterally with easiness, Naive Bayes is recognized to outstrip even extremely sophisticated classification procedures.

Bayes formula offers a mode of computing subsequent possibility $P(c|x)$ from $P(c)$, $P(x|c)$ and $P(x)$. see the formula below:

The diagram illustrates the components of the Naive Bayes formula. At the top, 'Likelihood' points to $P(x|c)$ and 'Class Prior Probability' points to $P(c)$. These two terms are multiplied together in the numerator of the formula $P(c|x) = \frac{P(x|c)P(c)}{P(x)}$. Below the formula, 'Posterior Probability' points to $P(c|x)$ and 'Predictor Prior Probability' points to $P(x)$. Below the diagram, the joint probability formula is given: $P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$.

- $P(c|x)$ is the given predictor (x , attributes) probability of posterior class (c , target).
- $P(c)$ is the probability of prior class.
- $P(x|c)$ is the given class probability predictor.
- $P(x)$ is the probability of prior predictor.

2.6 Random Forests

Leo Breiman developed the Random Forest [16] which is a set of un-pruned arrangement or regression trees completed from the selection of random sections of the training information. From the induction process Random features are selected. Prediction is completed by combining (majority elect for averaging for regression or classification) the estimates of the collaborative. Each tree is developed as defined below [24]

- By randomly sampling n , if N is the no. of cases in the training set however with replacement, by the original information. This model will be recycled as the training set for increasing the tree.
- For input variables M , the m variable is chosen in a way that $m \ll M$ is stated at every node, variables m is chosen at M random out then the finest split on m is used at each node. Throughout the forest growing, the price of m is kept constant.
- Each tree is developed to be prime extent. No trimming is used.

Generally Random Forests shows a major performance enhancement as associated to solo tree classifier like C4.5. The simplification rate of error that yields to relate favourably to AdaBoost, though it is more vigorous to noise. When using the random forest procedure to resolve regression complications are using Mean Squared Error (MSE) to know the information branch from each node

$$MSE = \frac{1}{N} \sum_{i=1}^N (f_i - y_i)^2$$

Wherever N is the data points, f_i is the returned value of a model and y_i is the data point i 's actual value. This formulation calculates the distance of every single node from the actual predicted value, helping to choose which branch is the superior choice for your analysis. Here, y_i is the data point value you are testing at a certain node and f_i is the value returned via state.

2.7 Gaussian process

Gathering of random variables is a Process of Gaussian [17] any infinite number of which consume (reliable) combined Gaussian disseminations. A Gaussian procedure is completely indicated through its function of mean $m(x)$ and function of covariance $k(x, x_0)$. This remains regular generalization of the Gaussian dissemination whose covariance and mean is a matrix and vector correspondingly. The Gaussian dissemination is concluded vectors; however, the Gaussian process is complete gathering. We resolve compose:

$$f \sim GP(m, k)$$

meaning: "the GP is distributed as function f with function of mean m and function of covariance k ". While the generalization from dissemination to procedure the conservative forward, we resolve it a bit added obvious about the particulars, since it might be unaware to some readers. The separable variables in random is a vector from in a Gaussian dissemination [18] are indexed through their place in the vector. For the process of Gaussian, x is the argument ($f(x)$ is the random function) which shows the role of catalogue set for every x input there is $f(x)$ associated with random moveable, which is the significance of the (stochastic) gathering f at that positions. For details convenience of notational, we resolve enumerate the values of x concern by the normal quantities, and usage of these catalogues as if they remained the catalogues of the procedure do not lease yourself remain disordered by this the guide to the procedure is x_i , which we consume preferred to index through i .

2.8 K Nearest Neighbour

KNN method is one of the form of supervised algorithm in ML. where it can use for classification and problems of regression predictive. Though, it is mostly used for predictive classification of problems in industry [19]. The resulting two things would describe KNN well

- **Algorithm of Lazy learning** – KNN is a procedure of lazy learning since it does not consume a specific training stage and usage of all information for training though classification.
- **Algorithm of Non-parametric learning** – KNN is also a procedure of non-parametric learning since it doesn't adopt anything about the principal data.

When classification uses KNN, as the production can considered as the class through the maximum frequency after the K most comparable examples. Every occurrence in kernel elects for their class and with the maximum elects is engaged as the prediction.

Probabilities of Class can be designed as the regularized frequency of illustrations that fit to every class in the K set most comparable examples for a novel data illustration. For illustration, in a binary problem classification (0 or 1 are class)

$$p(\text{class} = 0) = \frac{\text{count}(\text{class} = 0)}{(\text{count}(\text{class} = 0) + \text{count}(\text{class} = 1))}$$

If you have classes of even number and uses K. to choose value k is a good idea with odd number to escape a tie. And the even number uses the inverse for K when you take a classes of odd number.

Ties can remain wrecked constantly by growing K by 1 and observing at the next class to most comparable occurrence in the training dataset.

2.9 Decision Tree

The representation of flow chart like tree structure is known as decision tree wherever an interior node denotes attribute (or feature), the decision rule represents a branch, and outcome represents each leaf node. The root node in the decision tree is the topmost node. It studies to barrier on the source of the feature value. The barriers in tree is in recursively custom call recursive splitting. The decision making can be done by flowchart like structure. The flowchart visualization can simply impress the human level intelligent. Because decision trees are easy to interpret and understand. [20]

In ML, white box type of procedure is Decision Tree. It cuts interior decision creating reason, which is not accessible in the type of black box methods like Neural Network. Its training period is quickly likened to the method of neural network. The complexity of time in decision trees is a task of the sum of records and quantity of elements in the specified information. It is a non-parametric or distribution free technique, which organizes not be contingent upon possibility dissemination conventions. It [21] can hold high dimensional information with noble accurateness.

The concept of entropy is a concept that is invented by Shannon, which processes the infection of the participation set. In mathematics and physics, entropy denoted as the uncertainty or the uncleanness in the scheme. In theory of information, it denotes to the uncleanness in a set of illustrations. entropy is the decrease of Information gain. Information gain calculates the alteration among entropy earlier splitting and normal entropy after splitting of the dataset created on particular attribute values. Iterative Dichotomiser (ID3) decision tree procedure uses gain of information.

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

Where, p_i is the arbitrary tuple of probability in an D fits to class c_i .

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\text{Info}_A(D) = \sum_{j=1}^V \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

Wherever,

- $\text{Info}(D)$ is the normal quantity of data needed to recognize the class label of a tuple in D.
- $|D_j|/|D|$ acts as j^{th} partition of weight.
- $\text{Info}_A(D)$ is the predictable info essential to categorize a tuple from D built on the splitting by A.

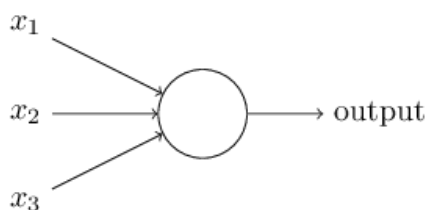
The A attribute through the highest gain of information, $\text{Gain}(A)$, is selected as the piercing element at node N().

2.10 Neural network

deep learning[22] one of the form is Neural networks, which is **machine learning** subfield, where algorithms inspired from the structure human brain. input data is Neural networks, train themselves to identify designs initiate in the information, and formerly calculate the production for a novel set of comparable information. Consequently, a neural network can remain believed as the efficient element of deep learning, which simulates the human brain behaviour to resolve compound data driven difficulties.

Neurons are referred as input in machine learning, and the process of nucleus[23] the information and advancing the considered productivity concluded the axon. In a genetic neural network, the thickness of dendrites describes the weight related through it.

A perceptron proceeds numerous binary involvements, $x_1, x_2, \dots, x_1, x_2, \dots$, and yields a solo binary production:



In the illustration exposed the perceptron has 3 contributions, x_1, x_2, x_3 . In common, it might take more or less inputs. A simple rule is proposed by Rosenblatt to calculate the outcomes. He presented *weights*, $w_1, w_2, \dots, w_1, w_2,$, real statistics uttering the significance of the corresponding responses to the production. The neuron's production, 0 or 1, is resolute by the weighted sum $\sum_j w_j x_j$ is fewer than or superior than some *values of threshold*. Just similar the weights, the threshold is a genuine number which is a limitation of the neuron. To place it in additional exact arithmetic terms:

$$output = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq threshold \\ 1 & \text{if } \sum_j w_j x_j > threshold \end{cases} \quad (1)$$

$$output = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq threshold \\ 1 & \text{if } \sum_j w_j x_j > threshold \end{cases}$$

2.11 MongoDB

MongoDB is a general-purpose database. MongoDB Schema [24] is actually a dynamic Schema which allows for the high flexibility and fits perfect for agile software development.

A cluster can contain a very large number of servers; they can be a single server but also a REPLICASET. A cluster contains an N number of replica set depending on our needs.
 — There is no upper limit to the number of members of the cluster.
 — Due to this scalability major goal to make the complex cluster infrastructure transparent to the user and the interaction of the server is same as of interacting with a replica server.

3 Result and analysis

Firstly, load the genome sequencing in MongoDB due to the sequencing is in unstructured data format because it may have different sizes. And then it is retrieved into python programme by using MongoDB connector. Then install all the necessary packages like Bio python and other package in python program language. By using Burrows_Wheeler_Alignment (BWA) finding the genome alignment and find the characters of the reference genome with the occurrence character and searching for max difference threshold for searching the matched position of the mutation sequence of each genome as illustrated in table 1.

About Burrows_Wheeler_Alignment (BWA) and genome analysis

BWA is used to map sequences of low-divergentalongside with a big reference genome, which is called as human genome. It contains 3 algorithms: BWA-MEM, BWA-SW and BWA-backtrack. The initial algorithm is considered for sequenceIllumina reads up to 100bp, though the rest of algorithm for lengthier sequences alternated from 70bp to 1Mbp. BWA-SW andBWA-MEM share same features like split alignmentandlong-read support, however BWA-MEM is latest one, which is normally suggested for queries like high-quality because it is wilder and more perfect. BWA-backtrack also takesimprovedenactment than BWA-MEM for 70-100bp reads of Illumina.

Table 1 table shows the matched position of genome using BWA algorithm

NAME	Matched positions with clinical data
Normal Clinical Data	0
Genome 1	[222, 830, 582, 576, 28, 330, 746, 273, 587, 714, 557, 520, 320]
Genome 2	[82582, 88759, 27415, 80452, 19889, 35580, 5902, 1287, 54182, 57519, 67956, 24007, 34089, 17901, 69305, 24433, 10872, 74133, 87197, 9594, 51342, 32740, 23581, 44668, 61637, 29616,15464, 66197, 72682, 13043, 21053, 63592, 71250, 18203, 7120, 61159, 24032, 34455, 17740, 8083] 1246
Genome 3	[848, 5179, 6809, 4634, 2209, 1646, 1457, 3866, 2154, 3861, 89, 7136, 820, 2184, 3980, 1563, 3137, 307, 3468, , 4535, 3328, 6341, 1365, 2449, 2715] 78
Genome 4	[14708, 26324, 12853, 16241, 20281, 30934, 36654, 17071, 34441, 16802, 17927, 5929, 26121, 19942, 36279, ,28872, 34230, 17149, 24355, 19493] 132

By observing the output file, we observed that the alignment of sequence reads with their positions and found the matched positions of mutation which causes cancer.

As the tumour mutation will require the 30 folds of sequence coverage that matched with normal tissue. By associating, the novel draft of human genome required around 65-fold coverage

Recognizing driver mutations are the main aim of sequencing cancer genome: the increase of mutation rate in a cell will change the gene which leads to speedier evolution of tumour and metastasis. To determine the mutation of driver in an DNA sequence is difficult but drivers inclined at most usually shared mutation between tumours, cluster around identified oncogenes are inclined as non-silent. The Passenger mutations are randomly circulated throughout the genome which is not important for the progression of disease. It has been expected that tumor carries 80 somatic mutations averagely less than 15 of which remain predictable to be drivers.

3.1 Classification of genome using machine learning

Step 1: Import genome sequencing dataset from the MongoDB into python program language.

Step 2: pre-process the dataset of genome sequence

Step 3: To formalize the dataset as shown in figure 2. the formalization of sequence is done by counting the genome sequence, uniqueness in the sequence and also found the frequency, top priority of the genome sequence

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	
count	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	106	1
unique	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
top	t	a	a	c	a	a	a	a	a	a	a	c	t	c	t	t	g	a	c	a	t	a	g	t	t	g	g	t	
freq	38	34	30	30	36	42	38	34	33	36	38	31	34	38	54	54	53	40	44	31	34	31	30	32	32	34	29	32	

Figure 2 Formalization of genome sequence

Step 4: And then formalized data is formatted by counting the text of each genome sequence is counted and recorded each sequence data as illustrated in figure 3.

	0	1	2	3	4	5	...	52	53	54	55	56	Class
t	38.0	26.0	27.0	26.0	22.0	24.0	...	35.0	30.0	23.0	29.0	34.0	NaN
c	27.0	22.0	21.0	30.0	19.0	18.0	...	21.0	32.0	29.0	29.0	17.0	NaN
a	26.0	34.0	30.0	22.0	36.0	42.0	...	25.0	22.0	26.0	24.0	27.0	NaN
g	15.0	24.0	28.0	28.0	29.0	22.0	...	25.0	22.0	28.0	24.0	28.0	NaN
-	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	53.0
+	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	53.0

[6 rows x 58 columns]

Figure 3Recording of every genome sequence count

Step 5: To run machine learning algorithms on the data the data must be in string format because it is very difficult to classify the string format of genome sequencing data. So firstly we convert the string format into numerical data as shown in figure 4

	0_a	0_c	0_g	0_t	1_a	1_c	1_g	1_t	2_a	2_c	2_g	2_t	3_a	3_c	3_g	3_t	4_a	4_c	4_g	4_t	5_a	5_c	5_g	5_t	6_a	6_c	6_g	6_t	7_a	7_c
0	0	0	0	1	1	0	0	0	0	1	0	0	0	0	0	1	1	0	0	0	0	0	1	0	0	1	0	0	1	0
1	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	1	1	0	0	0	0	0	0	1	0	1	0	0	0	1
2	0	0	1	0	0	0	0	1	1	0	0	0	0	1	0	0	0	0	0	1	1	0	0	0	0	0	1	0	1	0
3	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	0	1	0	0	0	0	1	0	0	1	0	1	0
4	0	0	0	1	0	1	0	0	0	0	1	0	1	0	0	0	0	0	0	1	1	0	0	0	1	0	0	0	0	0

5 rows x 230 columns

Figure 4 conversion of string format into numerical data

Step 6: Train and test the sequence by using classification algorithms of machine learning like SVM, Naive bayes, Neural Net, Decision Tree, Gaussian Process, Nearest Neighbors, AdaBoost, Random Forest. For testing the algorithm note the accuracies of different machine learning procedure as mentioned in table 2.

Table 2. Different machine learning algorithms Accuracy

Models	Accuracy	Test Accuracy
Naive Bayes	0.837500	0.9259259259259259
SVM Linear	0.850000	0.9629629629629629
SVM RBF	0.887500	0.9259259259259259
SVM Sigmoid	0.900000	0.9259259259259259
Neural Net	0.887500	0.9259259259259259
AdaBoost	0.925000	0.8518518518518519
Random Forest	0.596429	0.5185185185185185
Nearest Neighbors	0.823214	0.7777777777777778
Gaussian Process	0.873214	0.8888888888888888
Decision Tree	0.712500	0.7407407407407407

By comparing all the algorithms of machine learning accuracies “Adaboost” algorithm method is best method. And visualized the accuracy and test accuracy of the different algorithm as shown in figure 5.

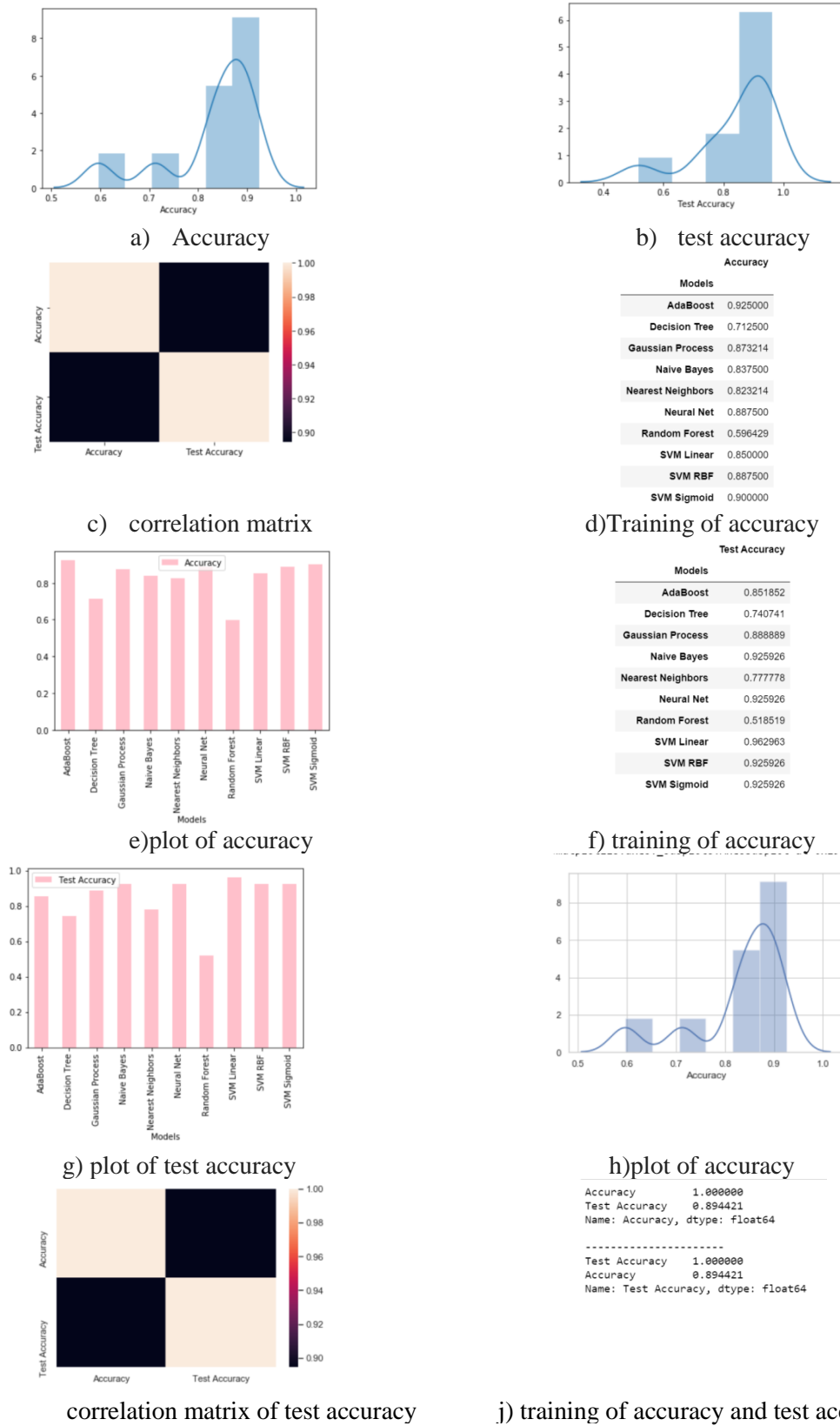


Figure 5 visualization of accuracy

By combining the accuracy visualization of figure 5 we visualized both the accuracy in one graph as shown in figure 6 by graph we concluded that Adaboost algorithm is best for the classification of genome.

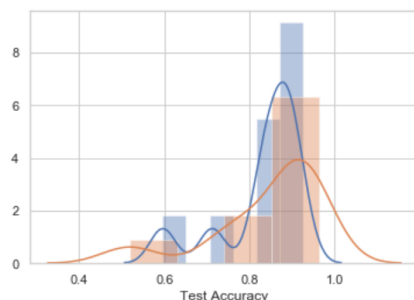


Figure 6 visualization of different machine learning algorithms accuracy

4. Conclusion

The sequencing technology like massive parallel sequencing to prospect the genome sequencing list of each DNA base in a genome, a genome map recognizes the innovations. A map of genome is detailed less and helps in navigating the genome sequence around the genome. By applying Burrows_Wheeler_Alignment (BWA) founded the genome alignment and found the characters of the reference genome with the occurrence character and searching for max difference threshold by analysing the output of sequence reads is aligned with their positions and found the matched positions of mutation which causes cancer. Because cancer genome would need 30-fold coverage of sequence in a genome and a corresponding normal flesh. By associating, the novel draft of human genome required around 65-fold coverage. By applying different algorithms of machine learning to the sequence dataset to formalize which algorithm is best to train and test the classification algorithms. Here Naive Bayes, Support Vector Machine(SVM), Ada Boost, Nearest Neighbours, Random Forest, Decision Tree etc., algorithms are applied and founded the accuracy for each classifier. By comparing accuracy of each classifier Ada Boost method has the high accuracy and visualization.

References

1. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial sequencing and analysis of the human genome. *Nature*. 2001; 409:860–921.
2. N. Peek, J. H. Holmes, J. Sun, technical challenges for big data in biomedicine and health: data sources, infrastructure, and analytics, *yearbook of medical informatics* 9(2014) 42-47
3. Kohane, I. S.; Masys, D. R.; Altman, R. B. (2006). "The Incidentalome: A Threat to Genomic Medicine". *JAMA*. **296** (2): 212–215. doi:10.1001/jama.296.2.212. PMID 16835427.
4. Allan Maxam; Walter Gilbert (February 1977). "A new method for sequencing DNA". *PNAS*. **74** (2): 560–4. doi:10.1073/pnas.74.2.560. PMC 392330. PMID 265521.
5. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. 2009;25(14):1754–1760.
6. Li H. Aligning Sequence Reads, Clone Sequences and Assembly Contigs with BWA-MEM [preprint]. New York NY: Oxford University Press; 2013. <http://arxiv.org/pdf/1303.3997v2.pdf>. Accessed November 28, 2016.
7. Michael k, K. Leung, Andrew Delong et al. Machine Learning in Genomic Medicine: A Review of Computational problems and Data Sets. 2016, Machine Learning in Genomic Medicine. Vol. 104, no.1.
8. Wagstaff, K., 2012. Machine learning that matters. arXiv preprint arXiv:1206.4656.

9. Cunningham, S.J., Littin, J. and Witten, I.H., 1997. Applications of machine learning in information retrieval
10. Bennett, K.P. and Parrado-Hernández, E., 2006. The interplay of optimization and machine learning research. *Journal of Machine Learning Research*, 7(Jul), pp.1265-1281.
11. Caruana, R. and Niculescu-Mizil, A., 2006, June. An empirical comparison of supervised learning algorithms. In *Proceedings of the 23rd international conference on Machine learning* (pp. 161-168). ACM.
12. Mitchell, T.M., 2006. *The discipline of machine learning* (Vol. 3). Carnegie Mellon University, School of Computer Science, Machine Learning Department.
13. S. Dinakaran and P. RanjitJebaThangaiah Ensemble Method of Effective AdaBoost Algorithm for Decision Tree Classifiers. 2017
14. D. Cai, A. Delcher, B. Kao and S. Ksif, "Modeling splice sites with Bayes networks", *Bioinformatics*, vol. 16, pp. 152-158, 2000.
15. Y. Freund and R. E. Schapire, "A Short Introduction to Boosting", *Journal of Japanese Society for Artificial Intelligence*, vol. 14, pp. 771-780, 1999.
16. Pang H, Lin A, Holford M, Enerson BE, Lu B, Lawton MP, et al. Pathway analysis using random forests classification and regression. *Bioinformatics*. 2006; 22(16):2028–2036.
17. Chen, Zexun; Wang, Bo; Gorban, Alexander N. (2019). "Multivariate Gaussian and Student-t process regression for multi-output prediction". *Neural Computing and Applications*. arXiv:1703.04455. doi:10.1007/s00521-019-04687-8.
18. Kac, M.; Siegert, A.J.F (1947). "An Explicit Representation of a Stationary Gaussian Process". *The Annals of Mathematical Statistics*. **18** (3): 438–442. doi:10.1214/aoms/1177730391.
19. Nigsch, Florian; Bender, Andreas; van Buuren, Bernd; Tissen, Jos; Nigsch, Eduard; Mitchell, John B. O. (2006). "Melting point prediction employing k-nearest neighbor algorithms and genetic parameter optimization". *Journal of Chemical Information and Modeling*. **46** (6): 2412–2422. doi:10.1021/ci060149f. PMID 17125183.
20. Rokach, Lior; Maimon, O. (2008). *Data mining with decision trees: theory and applications*. World Scientific Pub Co Inc. ISBN 978-9812771711.
21. Quinlan, J. R. (1986). "Induction of decision trees" (PDF). *Machine Learning*. **1**: 81–106. doi:10.1007/BF00116251. S2CID 189902138.
22. Ferreira, C. (2006). "Designing Neural Networks Using Gene Expression Programming" (PDF). In A. Abraham, B. de Baets, M. Köppen, and B. Nickolay, eds., *Applied Soft Computing Technologies: The Challenge of Complexity*, pages 517–536, Springer-Verlag.
23. Sak, Hasim; Senior, Andrew; Beaufays, Françoise (2014). "Long Short-Term Memory recurrent neural network architectures for large scale acoustic modeling"(PDF). Archived from the original (PDF) on 24 April 2018.
24. G. Adrián, G. E. Francisco, M. Marcela, A. Baum, L. Daniel, and G. B. de QuirósFernán, "Mongodb: an open source alternative for HL7-CDA clinical documents management," in *Proceedings of the Open Source International Conference (CISL '13)*, Buenos Aires, Argentina, 2013.