# Machine Learning in Bioinformatics: COVID-19 Mutation Prediction

[1]AMALA JAMES, [2]AMBILY JACOB

[1]MSc Scholar, [2]Assistant Professor

[1,2]Department of Computer Science,

[1,2]St Joseph's College (Autonomous), Irijnalakuda, Thrissur, Kerala

**Abstract:** Antibody therapeutics and vaccines are among our pis aller to finish the raging COVID-19 pandemic. They, however, are susceptible to over 5000 mutations on the spike (S) protein uncovered by a Mutation Tracker supported over 200 000 genome isolates. Recent variants of the virus within the UK, South Africa, and Brazil seem to spread more easily, which have the potential to steer to more hospitalizations and deaths. RNA sequence analysis of emerging SARS-CoV-2 infection is effective for tracking viral evolution and developing novel diagnostic tools. Furthermore, SARS-CoV-2 sequence analysis can provide insight into potential antigenic drift events that cause strain speciation and changing clinical outcomes. The tactic is definitely adaptable to research potential mutations of the virus, ensuring the simplest possible vaccines are quickly identified. This will give humans a big advantage over evolving mutations, with the model accomplishing vaccine design cycles that when took months or years during a matter of seconds or minutes. This text is an effort to predict mutations of novel coronavirus by applying machine learning techniques over bioinformatics.

**Keywords- COVID-19, Machine Learning, Sequential Pattern Mining, Corpus, Genome.**

## I. INTRODUCTION

The Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) virus, also mentioned as COVID-19, was first reported in Wuhan, China in December 2019. Later it had been declared as a plague by the earth Health Organization (WHO). There are 179,686,071 confirmed cases of COVID-19, including 3,899,172 deaths, reported to WHO. Genome of SARS-CoV-2 has been sequenced by different groups around the world which revealed multiple strains of the virus and showed that its genome is 79% almost just like the SARS-CoV-1 and 50% to the MERS-CoV (Middle East Respiratory Syndrome Coronavirus), respectively. The evolution process of virus isn't completely mentioned because it changes slowly compared to other viruses, thus, giving less mutations to review. On the typical, the coronavirus accumulates about two changes per month in its genome. Most of the changes within the COVID-19 genome structure won't affect how the virus behaves, but a few of may influence the disease's transmissibility or severity. it is vital to understand the pattern within the virus mutates also as its

mutation rate. The mutation rate of any virus could also be a critical parameter to understand viral evolution. Genomic sequence and mutation analysis are crucial and accurate information on the mutation rate can play a crucial role within the assessment of possible drugs/vaccination strategies. During this regard, we propose an algorithm which can be used to analyze genome sequences for variations and also to review the mutation rate.

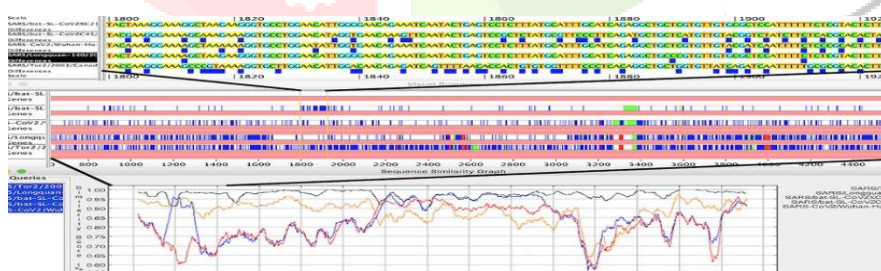The objective of this paper includes:

- To find interesting patterns in COVID-19 genome sequences, we apply sequential pattern mining.
- To evaluate if subsequent nucleotide bases(s) are often predicted in COVID-19 genome sequences.
- To analyze mutations in genome sequence, we propose an algorithm to hunt out the mutations that takes place in genome sequences also because the mutation rate.

## II.    METHODOLOGY

### i.    Genome assessment gadgets by VBRC:

The Viral Bioinformatics Research Centre (VBRC) is an accomplished resource collected expressly for virologists to work with the close to examination of viral genomes. Viral Genome Organizer envisions and dissects the relationship of characteristics inside various all out viral genomes. The mechanical assembly allows the customer to exchange protein or DNA game plans and can show START/STOP codons for 6-traces similarly as open arrangement edges and other customer described outcomes. Genome Annotation Transfer Utility is a gadget for explaining genomes using information from a reference genome. It obliges instinctive clarification, normally explaining characteristics that are essentially equivalent to the reference disease anyway leaving others for a human decision. The VBRC was created for dsDNA contaminations anyway have been adapted to Covid. SARS-CoV-2 and solidly related contaminations have been added to the informational index. VBRC mechanical assemblies will help with imagining contrasts between Covid groupings at different levels of objective (see Figure 1).

**Figure 1:**

**Bioinformatics devices accelerate SARS-CoV-2 assessment. A couple devices are depicted in Table 1.**

**Table 1:**

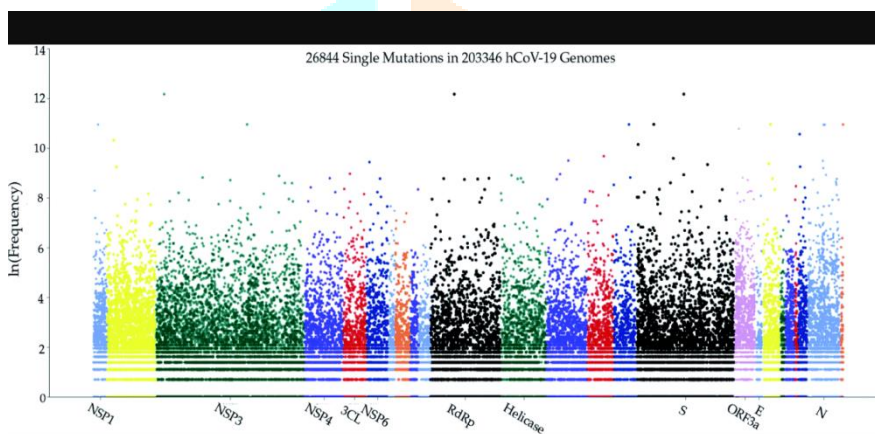| Tool | Advancing SARS-CoV-2 investigation by |
|---|---|
| poreCov | reducing drawn-out bioinformatics bottlenecks in getting ready sequencing |
| VADR | validation and remark of SARS-CoV-2 progressions |
| COPASI | modelling the components of the pandemic and effect of interventions |
| COVIDSIM | analysing effects of contact decline measures and guide political dynamic |
| CoV-GLUE | tracking changes gathering in the SARS-CoV-2 genome |
| VBRC genome assessment gadgets | visualizing contrasts between Covid courses of action at different levels of objective |
| VIRULIGN | fast, codon-right unique progression plan and remark of contamination genomes |
| Rfam COVID-19 | annotating coordinated RNAs in Covid progressions and expecting discretionary plans |

## ii.    Covidex: Arrangements free subtyping utilizing ML:

Viral subtypes address bunches among secludes from the worldwide populace of a characterized animal groups. Subtypification is significant for concentrates on infection the study of disease transmission, advancement and pathogenesis. Most subtype grouping strategies require the arrangement of the info information against a bunch of pre-characterized subtype reference successions. These strategies can be computationally costly, especially for long successions, for example, SARS-CoV-2 ( ≈≈30 kb per genome). To handle this issue, AI apparatuses might be utilized for infection subtyping. Covidex was created as an open-source arrangement free AI subtyping device. It's anything but a sparkly application that permits quick and exact (out-of-pack blunder rate < 1.5 %) order of viral genomes in pre-characterized bunches. For SARS-CoV-2, the default transferred model depends on Nextstrain and GISAID information. Covidex depends on a quick execution of arbitrary backwoods prepared over a k-mer information base. Via preparing the order calculations over k-mer recurrence vectors, Covidex significantly lessens computational and time necessities and can characterize many SARS-CoV-2 genomes like a flash. Consequently, with regards to the flow worldwide pandemic where the quantity of accessible SARS-CoV-2 genomes is developing dramatically, SARS-CoV-2 exploration can profit with this particular device intended to decrease the time required in information examination essentially.

### iii. Changes on the spike protein:

Infection changes are presented by normal determination, replication component, cell climate, polymerase loyalty, quality altering, arbitrary hereditary float, late epidemiological highlights, have resistant reactions, and so forth Remarkably, seeing how transformations have changed the SARS-CoV-2 construction, work, infectivity, action, and harmfulness is critical for thinking of life-saving systems in infection control, regulation, avoidance, and drug. Genome sequencing, SNP calling, and phenotyping give an effective way to parse changes from countless viral examples. In this work, we recovered more than 200 000 complete SARS-CoV-2 genome arrangements from the GISAID information base and made a constant intelligent SARS-CoV-2 Mutation Tracker to report more than 26 000 extraordinary single transformations alongside their change recurrence on SARS-CoV-2 as of January 20, 2021. Figure 2 is a screen capture of online Mutation Tracker. It depicts the appropriation of transformations on the total coding area of SARS-CoV-2. The y-pivot shows the regular log recurrence for every change at a particular position.

**Figure 2:**



26844 Single Mutations in 203346 hCoV-19 Genomes

### iv. Analyzing COVID-19 Genome with SPM and progression assumption methods:

SPM have been expected to mine models in coordinated sequential data. SPM contains recognizing huge models in a lot of discrete progressions. It involves two essential parts:

1. Corpus headway: COVID-19 genome progressions are changed into a corpus of discrete groupings, where each whole genome game plan is changed over into a course of action of nucleotides.
2. Learning using SPM and Sequence Prediction Techniques: SPM computations are applied on the corpus to discover frequently happening nucleotides, progressive associations among nucleotides, and to expect the accompanying nucleotides base(s) of a plan.

**Corpus advancement:**

**Definition 1 (Nucleotide base set):** Let NB = {A, C, G, T} be the arrangement of all unmistakable nucleotide bases. The documentation |NB| indicates the set cardinality. Thus, |NB| = 4 as there are 4 unmistakable nucleotide

**Definition 2 (COVID-19 genome grouping):** CGS = ⟨NB1, NB2,..., NBn⟩ , to such an extent that NBi⊆NBi⊆ NB (1 ≤ I ≤ n).

**Utilizing SPM and grouping forecast procedures:**

**SPM:**

Subsequent to setting up the corpus, different SPM methods can be applied to discover designs that show up in genome groupings. The most well-known measure to assess designs in design mining is the help measure. SPM utilizing the help measure is known as the assignment of continuous SPM.

**Definition 3 (Support):** The help of a genome (sub-) succession $S\alpha$ in a corpus CGSC is the absolute number of arrangements that contain $S\alpha$. It is signified as $sup(S\alpha)$ and characterized as: $sup(S\alpha)=|\{S|S\alpha\sqsubseteq S \wedge S\in CGSC\}|sup(S\alpha)=|\{S|S\alpha\sqsubseteq S \wedge S\in CGSC\}|$.

**Definition 4 (Frequent SPM in a genome grouping corpus)**: May there be a genome grouping corpus CGSC and a client characterized least help limit minsup, with the end goal that minsup > 0. The assignment of incessant SPM in a CGSC is to identify all regular genome aftereffects. A genome aftereffect S is successive if $sup(S) \geq$ minsup.

**Sequence forecast methods:**

Another learning task acted in this examination is to assemble grouping expectation models utilizing the COVID-19 genome successions to check whether the plan of nucleotide bases is unsurprising. The applied models incorporate CPT+, CPT, DG, AKOM, Mark1, TDAG and LZ78.

**Regular nucleotide sets**:

The Apriori calculation for FIM was first applied on the corpus to discover the oftentimes happening nucleotide base sets. Apriori takes a corpus and a minsup limit as info and yields the continuous nucleotide base sets. A post-preparing step was then performed to keep just continuous itemsets containing solitary nucleotides or having a different of three nucleotides (the length of a codon).
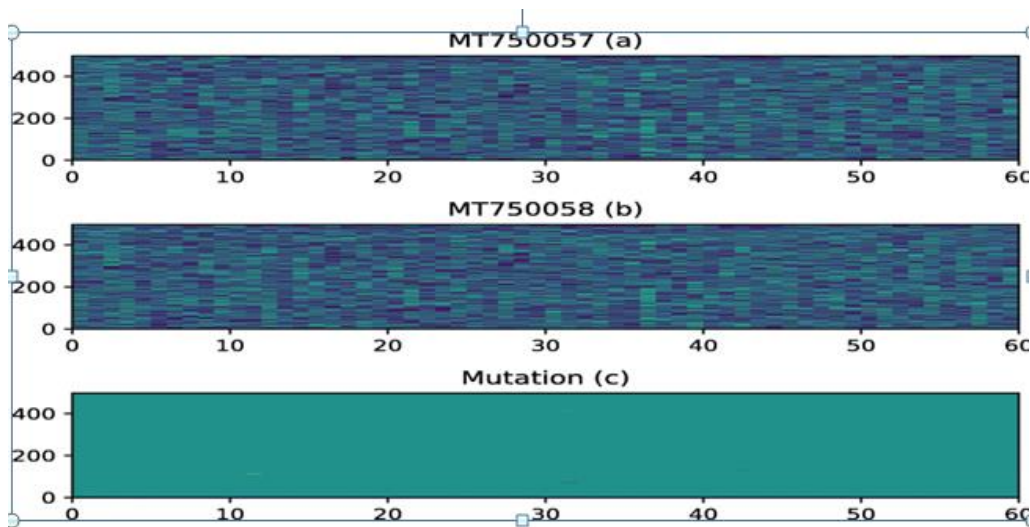
**Incessant successive examples:**

Then, at that point, SPM calculations were applied to discover covered up consecutive connections between nucleotides. The CM-SPAM calculation was executed, which needs to set a minsup edge. CM-SPAM was designed to discover just coterminous successive examples since designs that skip nucleotides would be difficult to decipher, and designs that are not a numerous of three nucleotides (the size of a codon) were sifted through.

## v.　Coronavirus Genome transformation investigation

The SARS-CoV-2 infection is continually advancing by changing nucleotides, to adjust to new conditions. The matplotlib library was utilized to create the plots for two genome arrangements and the changed succession . In the event that there is any transformation in the two genome arrangements, the changed succession will have brilliant spots. Assuming there is no transformation, the will be dim with no splendid spots (Figure 3).

**Figure 3:**



## III.   CONCLUSION

This paper proposed two ways to deal with research and break down COVID-19 genome successions. In the principal approach, design mining procedures are utilized to discover regular nucleotide bases in the arrangements, their incessant examples and the successive connection between such examples. Besides, different succession expectation models were assessed on genome groupings, where AKOM (All-K-Order-Markov) performed better compared to other cutting edge calculations. In the subsequent methodology, a calculation was proposed to dissect changes in COVID-19 genome successions. The calculation finds the location(s) in COVID-19 strains where the nucleotide bases are changed to compute the transformation rate. The methodologies introduced in this paper are not restricted to the SARS-CoV-2 infection. They could be utilized for the investigation of other human infections as well.

## IV.   REFERENCES

[1] World Health Organization, http://www.euro.who.int/en/healthtopics/health-emergencies/coronavirus-covid-19/news/news/2020/3/who-announces-covid-19-outbreak-a-pandemic, last seen June 20, 2021.

[2] A. Alkhateeb and L. Rueda. "Zseq: an approach for preprocessing next-generation sequencing data." Journal of Computational Biology 24.8 (2017): 746--755.

[3] Li, Bo, and Colin N. Dewey. "RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome." BMC bioinformatics 12.1 (2011): 323.

[4] Iulian Derecichei and Govindaraja Atikukke, https://dl.acm.org/doi/10.1145/3388440.3415991, last seen June 27 2021.

[5] M. Saqib Nawaz, Philippe Fournier-Viger, Abbas Shojaee and Hamido Fujita, https://link.springer.com/article/10.1007/s10489-021-02193-w , last seen June 27 2021.

[6] Mousavizadeha l, ghasemi s (2020) genotype and phenotype of covid-19: their roles in pathogenesis. j microb immuno infect. https://doi.org/10.1016/j.jmii.2020.03.022

[7] Zihayat m, davoudi h, an a (2017) mining significant high utility gene regulation sequential patterns. bmc syst biol 11:109.

[8] Kawade dr, oza ks (2013) exploration of dna sequences using pattern mining. j biomed informa 2:144–148.

[9] Singer j, et al. (2020) cov-glue: a web application for tracking sars-cov-2 genomic variation. preprints: 2020060225

[10] franziska hufsky, kevin lamkiewicz, alexandre almeida, abdel aouacheria, cecilia arighi, alex bateman, jan baumbach, niko beerenwinkel, christian brandt, marco cacciabue, sara chuguransky, oliver drechsel, robert d finn, adrian fritz, stephan fuchs, georges hattab, anne-christin hauschild, dominik heider, marie hoffmann, martin hölzer, stefan hoops, lars kaderali, ioanna kalvari, max von kleist, renó kmiecinski, denise kühnert, gorka lasso, pieter libin, markus list, hannah f löchel, maria j martin, roman martin, julian matschinske, alice c mchardy, pedro mendes, jaina mistry, vincent navratil, eric p nawrocki, áine niamh o'toole, nancy ontiveros-palacios, anton i petrov, guillermo rangel-pineros, nicole redaschi, susanne reimering, knut reinert, alejandro reyes, lorna richardson, david l robertson, sepideh sadegh, joshua b singer, kristof theys, chris upton, marius welzel, lowri williams, manja marz, *Briefings in Bioinformatics*, Volume 22, Issue 2, March 2021, Pages 642–663, https://doi.org/10.1093/bib/bbaa232