

# **Project 3: Learning phylogenetic trees from multiple sequence alignment of the COVID data**

Fereshteh Noroozi

Summer 2023

## **Data**

I have downloaded the data for the Delta strain (number 5) from GISAID, covering the time from January 2021 to November 2022.

## **Data preparation**

**II. Align the collected dataset using the MAFFT tool and check the MSA for any standard format issues, such as duplicate taxon names, invalid characters in taxon names, or duplicate sequences.**

**Answer:**

At first to align the collected dataset:

```
mafft fereshteh_gisaid_hcov-19_2023_06_27_10.fasta > fereshteh_gis  
id_output.fasta
```

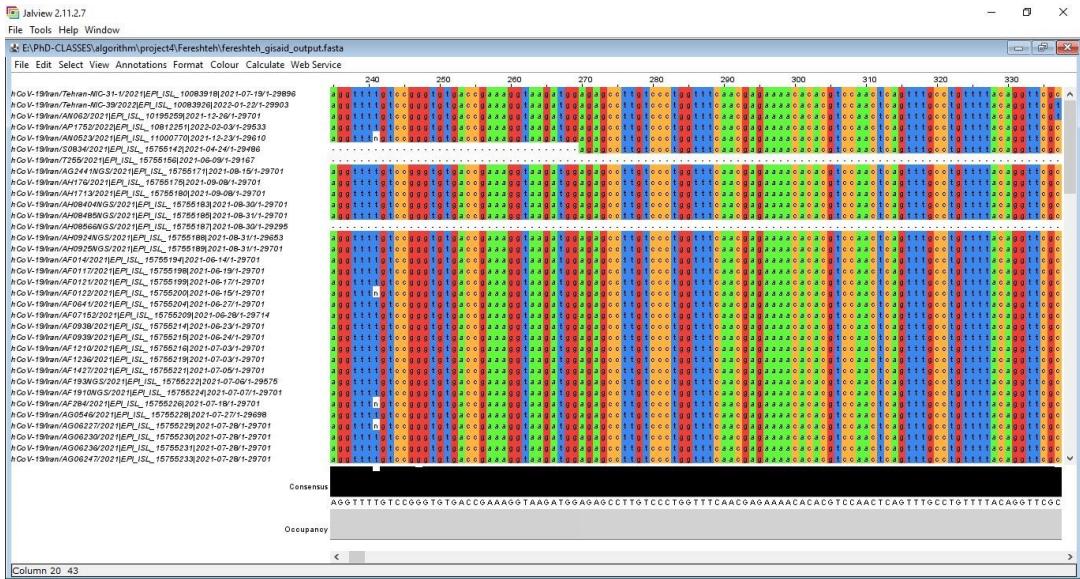
The output of the MAFFT software, which is commonly used for multiple sequence alignment, has the variables that are set to zero in the output:

- nthread = 0
- nthreadpair = 0
- nthreadtb = 0
- ppenalty\_ex = 0

These variables might control the number of threads used for parallelization or some penalty values used in the alignment process. However, since they are set to zero, it seems that the default values are being used. The output also shows that the alignment is being performed using the FFT-NS-2 strategy with the DNA200 model and a gap penalty of -1.53 for opening a gap and 0 for extending a gap. The UPGMA algorithm is used for constructing guide trees, and the alignment is performed twice in a progressive manner.

```
fereshteh_gisaid_output.fasta > hCoV-19/Iran/Tehran-NIC-31-1/2021|EPI_ISL_10083918|2021-07-19  
attaaaggttataccttccaggtaacaaaccaccaactttcgatcttgtagatct  
gttctctaaaccaaactttaaatctgtgtggctgtcactcggtcatgttagtcact  
cacgcagtataattaactaattactgtcggtacaggacacgagaactcgatcatac  
ttctgcaggctgttacgggttgcgtccgttgcagccatcatcagcacatctagg  
tgtccgggtgtgaccgaaaggtaagatggagagccctgtccctgggtcaacgagaaa  
acacgttccaactcgttttcgttttacagggttgcgcacgtgtcgtacgtggcttgg  
agactccgtggaggagggttacatcgaggacacgtcaacatcttaaagatggcacttgg  
cttagtagaaggtaaaaaaggcggttgcctcaacttgaacacgcctatgtgttcatcaa  
acgatcgatgtcgactgcacccatgtcatgttgcgttatggtagctggtagcagaact  
cgaaggcattcgttacgggttgcgttagtggtagacacttgcgttgcgtccatgtgg  
cgaataaccatcggttaccgcagggttcttgcgttacaaggtaacgcgttacatgg  
tggccatgttacggcgccgttacatcggttgcgttgcgttgcgttgcgttgcgttgc  
tccttatgaagatcccataaaaaactggaaacactaaacatagcgttgcgttacccgt  
actcatgcgttagcttacggaggggcatacactgcgttgcgttgcgttgcgttgcgtt  
ccctgtatggcttacccttgcgttgcattaaaggacccatgttgcgttgcgttgcgtt  
atgcactttgtccgaaacaactggacttttgcactaaagggttgcgttgcgttgcgtt  
tgaacatgagcatgaaattgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
gacacccatggatggcaaaaggaaatttgcgttgcgttgcgttgcgttgcgttgcgtt  
ttttgtatcccttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
gcttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
caacccatgtgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
gacggggcgattttgtttaaaggccatgtgcgttgcgttgcgttgcgttgcgttgcgtt  
agggtccactacttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
atgtcacaattcagaaggtaggcacgttgcgttgcgttgcgttgcgttgcgttgcgtt  
cttggaaaaccattttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
ttatgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
ttgttaaccatcagggttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt  
aatactccaaaaagagaaaactcaacatcaatattgttgcgttgcgttgcgttgcgttgcgtt  
gatgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgttgcgtt
```

Then I employed the Jalview software to visualize the alignment file.



The colored regions represent the bases that have a match, while the dots indicate the positions where there is a mismatch

Subsequently, to ensure the MSA is free from standard format issues, such as duplicate taxon names, invalid characters in taxon names, or duplicate sequences, was checked.

```
raxml -ng --check --msa fereshteh_gisaid_output.fasta --model GTR+G
```

The output shows that the alignment file was successfully loaded with [127 taxa](#) (sequences) and [29903 sites](#) (positions) and that the alignment comprises one partition with a [GTR+FO+G4m](#) model. The percentage of gaps and invariant sites in the alignment is also shown ([Alignment sites: 29903](#), [Gaps: 2.04%](#), [Invariant sites: 97.96%](#)). This means that out of the [29903](#) alignment sites, [2.04%](#) are gaps, and [97.96%](#) are invariant (i.e., have the same nucleotide in all sequences).). The warning message "Duplicate sequences found: 1" indicates that there is [one](#) sequence that is identical to another sequence in the alignment.

This command produced two files:

fereshteh_gisaid_output.fasta.raxml.log	Date modified: 6/27/2023 10:35 AM
Type: Text Document	Size: 1.69 KB
fereshteh_gisaid_output.fasta.raxml.reduced.phy	Date modified: 6/27/2023 10:35 AM
Type: PHYLIP File	Size: 3.59 MB

1. fereshteh\_gisaid\_output.fasta.raxml.reduced.phy: This file is a reduced version of the input alignment file "fereshteh\_gisaid\_output.fasta", with duplicates and gap-only sites/taxa removed.
2. fereshteh\_gisaid\_output.fasta.raxml.log: This file contains the execution log of the RAxML-NG analysis. The log file may contain additional information about the analysis, such as the likelihood score, branch lengths, and tree topology.

### Phylogenetic analysis

**Estimate the optimal nucleotide substitution model using maximum likelihood or Bayesian methods by selecting the model with the lowest AIC or BIC score (you may use phangorn package in R).**

**Answer:**

First, the code reads in the alignment data in PHYLIP format using the `read.phyDat` function in `phangorn`. Then, the `modelTest` function is used to estimate the optimal model using maximum likelihood and the AIC or BIC criterion. The `which.min` function is used to find the model with the lowest AIC and BIC scores, and the selected model name is stored in the `best_model_AIC` or `best_model_BIC` variable, respectively. Finally, the code prints out the name of the model with the lowest AIC and BIC score, allowing for the selection of the optimal model. Based on the output, the optimal nucleotide substitution model estimated using the `modelTest` function in `phangorn` and selected based on the lowest AIC and BIC scores is the **General Time Reversible (GTR)** model with gamma-distributed rate variation and a proportion of invariable sites (I), denoted as **GTR+G(4)+I**. The GTR+G(alpha)+I model is a popular choice for DNA sequence analyses, especially when the sequences under study have different base compositions or exhibit rate heterogeneity across sites. The model provides a good balance between model complexity and accuracy, and has been shown to perform well in a variety of phylogenetic contexts. The GTR model is often extended to account for rate heterogeneity across sites, which can be modeled by a gamma distribution with a shape parameter alpha and a rate heterogeneity parameter, denoted as GTR+G(alpha). Additionally, the model can account for the presence of invariable sites (i.e., sites that have no substitutions) by including a proportion of such sites, denoted as I, resulting in the GTR+G(alpha)+I model.

```

# Load necessary packages
library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

## Loading required package: ape

## Warning: package 'ape' was built under R version 4.2.3

# Read in the alignment data in PHYLIP format
alignment <- read.phyDat("fereshteh_gisaid_output.fasta.raxml.reduced.phy", format = "phylip")

# Estimate the optimal nucleotide substitution model using ML and AIC criterion
ml model <- modelEst(alignment)

## Model df logLik AIC BIC
## JC 250 -24846.49 98664.88 981872.9
## JC+I 250 -49302.13 98164.76 1011489.7
## JC+G(4) 250 -49451.21 98402.41 101478.8
## JC+G(4)+I 251 -49281.93 98665.84 101150.6
## F81 252 -48679.13 97862.24 99955.3
## F81+I 253 -48343.76 97193.53 99294.87
## F81+G(4) 253 -48479.45 97464.9 99566.25
## F81+G(4)+I 254 -48323.25 97154.7 99261.95
## K80 250 -49436.51 99373.03 101449.5
## K80+I 254 49084.94 98671.83 100756.6
## K80+G(4) 254 -49234.27 98700.54 101055.3
## K80+G(4)+I 254 -49064.1 98633.29 100726.3
## HKY 253 -48444.06 97394.11 99495.46
## HKY+I 254 -48416.16 97678.24 98487.89
## HKY+G(4) 254 -48247.12 97616.01 98114.96
## HKY+G(4)+I 255 -48164.05 97516.09 98336.05
## TrN 251 49340.06 98200.11 101284.8
## TrN+I 252 -49003.98 98510.10 99603
## TrN+G(4) 252 -49146.4 98796.8 100899.8
## TrNc+G(4)+I 253 -48980.71 98467.42 100568.8
## TrN 254 -48356.18 97220.31 99330.01
## TrN+T 255 -48048.08 96666.11 98724.07
## TrN+G(4) 255 -48160.05 96830.1 98948.05
## TrN+G(4)+I 256 -48024.46 96560.92 98687.18
## TPM1 252 -49375.83 99253.64 101338.4
## TPM1+I 252 -49024.28 99552.5 100645.5
## TPM1+G(4) 252 -49173.54 98851.16 100944.2
## TPM1+G(4)+I 253 -49173.54 98851.16 100944.2
## TrN+T 254 -49375.83 99253.64 101338.4
## K81+I 252 -49024.25 98552.5 100645.5
## K81+G(4) 252 -49173.58 98851.16 100944.2
## K81+G(4)+I 253 -49001.94 98511.84 100661.2
## K81+G(4)+I 254 -48373.58 97255.15 99361.8
## TPM1u+I 255 -48961.47 96632.95 98750.91
## TPM1u+G(4) 255 -48177.81 96865.61 98983.57
## TPM1u+G(4)+I 256 -48040.2 96592.39 98718.66
## TPM2 252 -49375.54 99253.01 101337.8
## TPM2+I 252 49037.27 98578.54 100671.6
## TPM2+G(4) 252 -49175.11 98844.63 100947.7
## TPM2+G(4)+I 253 -49016.88 98539.76 100641.1
## TPM2u 254 -48363.72 97235.45 99345.1
## TPM2u+I 255 -48053.46 96616.91 98734.87
## TPM2u+G(4) 255 -48177.81 96865.61 98983.57
## TPM2u+G(4)+I 256 -48042.19 96566.34 98802.64
## TPM3 251 -49423.18 99348.35 101433.1
## TPM3+I 252 -49053.65 98611.3 100704.3
## TPM3+G(4) 252 -49718.76 98941.52 101034.6
## TPM3+G(4)+I 253 -49033.86 98573.72 100675.1
## TPM3u 254 -48439.53 97387.01 99496.71
## TPM3u+I 255 -48112.5 96734.99 98852.95
## TPM3u+G(4) 255 -48241.59 96993.18 99111.14
## TPM3u+G(4)+I 256 -48091.81 96695.61 98821.84
## TIM1e 252 -49288.37 99080.71 101173.8
## TIM1e+I 253 -48942.31 98390.62 100492
## TIM1e+G(4) 253 -48535.71 98677.42 100778.8
## TIM1e+G(4)+I 254 -48529.69 98564.63 100778.7
## TIM1e+G(4)+I 255 -48281.4 98691.4 98199.46
## TIM1e+G(4)+I 256 -47983.76 96479.53 98065.79
## TIM1+G(4) 256 -48009.7 96603.39 98819.66
## TIM1+G(4)+I 257 -47960.13 96434.26 98568.82
## TIM2e 252 -49288.08 99080.17 101173.2
## TIM2e+T 253 -48952.12 98410.23 100511.6
## TIM2e+G(4) 253 -49086.99 98679.92 100781.3
## TIM2e+G(4)+I 254 -48930.38 98368.79 100478.4
## TIM2 255 -48276.05 97862.09 99180.05
## TIM2+I 256 -47976.61 96465.29 98591.54
## TIM2+G(4) 256 -48001.51 96675.01 98801.12
## TIM2+G(4)+I 257 -47953.71 96421.56 98556.13
## TIM2+G(4)+I 258 -47953.71 96421.56 98556.13
## TIM1e 252 -49288.37 99080.71 101173.8
## TIM1e+G(4) 253 -48942.31 98390.62 100492
## TIM1e+G(4)+I 254 -48535.71 98677.42 100778.8
## TIM1e+G(4)+I 255 -48281.4 98691.4 98199.46
## TIM1e+G(4)+I 256 -47983.76 96479.53 98065.79
## TIM1+G(4) 256 -48009.7 96603.39 98819.66
## TIM1+G(4)+I 257 -47960.13 96434.26 98568.82
## TIM2e 252 -49288.08 99080.17 101173.2
## TIM2e+T 253 -48952.12 98410.23 100511.6
## TIM2e+G(4) 253 -49086.99 98679.92 100781.3
## TIM2e+G(4)+I 254 -48930.38 98368.79 100478.4
## TIM2 255 -48276.05 97862.09 99180.05
## TIM2+I 256 -47976.61 96465.29 98591.54
## TIM2+G(4) 256 -48001.51 96675.01 98801.12
## TIM2+G(4)+I 257 -47953.71 96421.56 98556.13
## TIM2+G(4)+I 258 -47953.71 96421.56 98556.13
## TIM1e 252 -49288.37 99080.71 101173.8
## TIM1e+G(4) 253 -48942.31 98390.62 100492
## TIM1e+G(4)+I 254 -48535.71 98677.42 100778.8
## TIM1e+G(4)+I 255 -48281.4 98691.4 98199.46
## TIM1e+G(4)+I 256 -47983.76 96479.53 98065.79
## TIM1+G(4) 256 -48009.7 96603.39 98819.66
## TIM1+G(4)+I 257 -47960.13 96434.26 98568.82
## TIM2e 252 -49288.08 99080.17 101173.2
## TIM2e+T 253 -48952.12 98410.23 100511.6
## TIM2e+G(4) 253 -49086.99 98679.92 100781.3
## TIM2e+G(4)+I 254 -48930.38 98368.79 100478.4
## TIM2 255 -48276.05 97862.09 99180.05
## TIM2+I 256 -47976.61 96465.29 98591.54
## TIM2+G(4) 256 -48001.51 96675.01 98801.12
## TIM2+G(4)+I 257 -47953.71 96421.56 98556.13
## TIM2+G(4)+I 258 -47953.71 96421.56 98556.13
## SYM1 255 -48885.74 98281.49 100309.4
## SYM1+G(4) 255 -49037.34 98584.68 100670.6
## SYM1+G(4)+I 256 48864.91 98241.82 100368.1
## GTR 257 -48236.24 96986.48 99121.05
## GTR+I 258 -47931.34 96378.68 98521.55
## GTR+G(4) 258 -48010.88 96597.75 98740.63
## GTR+G(4)+I 259 -47908.97 96315.94 98487.12

```

# Find the model with the lowest AIC score

best\_model\_AIC <- ml.model[which.min(ml.model\$AIC), "Model"]

# Find the model with the lowest BIC score

best\_model\_BIC <- ml.model[which.min(ml.model\$BIC), "Model"]

# Print the results

cat("lowest AIC score is", best\_model\_AIC, "\n")

## lowest AIC score is GTR+G(4)+I

cat("lowest BIC score is", best\_model\_BIC, "\n")

## lowest BIC score is GTR+G(4)+I

**Calculate the pairwise distance between the sequences in the alignment.**

## Answer:

```

# Read in the alignment data in FASTA format
alignment <- readLines("fereshteh_gisaid_output.fasta.raxml.reduced.phy")

# Remove the header lines
alignment <- alignment[-seq(1, length(alignment), 2)]

# Create an empty distance matrix
num_seqs <- length(alignment)
dist_matrix <- matrix(0, nrow = num_seqs, ncol = num_seqs)

# Define a function to compute the distance between two sequences
pairwise_distance <- function(seq1, seq2) {
  num_diffs <- sum(charToRaw(seq1) != charToRaw(seq2))
  return(num_diffs)
}

# Calculate the pairwise distance matrix
for (i in 1:(num_seqs - 1)) {
  for (j in (i + 1):num_seqs) {
    dist_matrix[i, j] <- pairwise_distance(alignment[i], alignment[j])
    dist_matrix[j, i] <- dist_matrix[i, j]
  }
}

# Save the pairwise distance matrix to a text file
write.table(dist_matrix, "pairwise_distances.txt", sep = "\t", quote = FALSE)

```

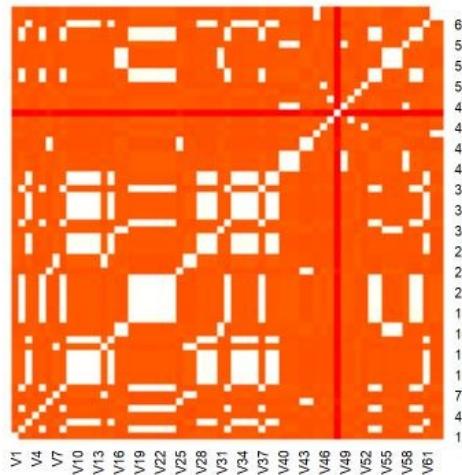
```
# Save the pairwise distance matrix to a text file  
write.table(dist_matrix, "pairwise_distances.txt", sep = "\t", quote = FALSE, col.names = NA)
```

This R code calculates the pairwise distance between the sequences in the alignment. It does so by reading in the alignment data in FASTA format, removing the header lines, and then creating an empty distance matrix. It defines a function `pairwise_distance` to compute the distance between two sequences, and then calculates the pairwise distance matrix by looping through all pairs of sequences and applying the `pairwise_distance` function. Finally, it saves the resulting pairwise distance matrix to a text file.

A heatmap is a good choice if we want to visualize the overall structure of the distance matrix and identify patterns or clusters. In a heatmap, the matrix is represented as a grid of colored squares, with darker colors indicating higher distances. We can also add clustering to the rows and/or columns of the heatmap to group similar points together.

```
# Read in the pairwise distance matrix
dist_matrix <- as.matrix(read.table("pairwise_distances.txt"))

# Create a heatmap of the distance matrix
heatmap(dist_matrix, Rowv = NA, Colv = NA, col = rev(heat.colors(256)),
        scale = "none", margins = c(10, 10))
```



## Generate a first tree topology based on pairwise distances.

**Answer:**

```
# Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

# Read in the pairwise distance matrix
dist_matrix <- as.matrix(read.table("pairwise_distances.txt", header = FALSE, fill = TRUE))

# Perform the neighbor-joining algorithm
tree <- nj(dist_matrix)

## Warning in storage.mode(m) <- "numeric": NAs introduced by coercion

# Save the tree topology to a file
write.tree(tree, "pairwise_distances.nwk")
```

This code performs a neighbor-joining algorithm on the matrix using the `nj()` function from the `ape` package, and saves the resulting tree topology to a Newick file "pairwise\_distances.nwk" using the `write.tree()` function.

Neighbour-joining (NJ) is an algorithm for building phylogenetic trees based on pairwise distances between taxa. The basic idea behind the NJ algorithm is to iteratively join pairs of taxa that are the "closest" in terms of their pairwise distances, until a complete tree is constructed.

```
pairwise_distances.nwk
1 ((v36:10.13043229,v37:10.86956771):1.585071087,(v29:13.19999371,v34:12.80000629):1.414928913,(((v11:3.604402184,v12:7.395597816):0.8161770276,v
```

Then to plot my tree:

```
# Load the ape package
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

# Read in the pairwise distance matrix
dist_matrix <- as.matrix(read.table("pairwise_distances.txt"))

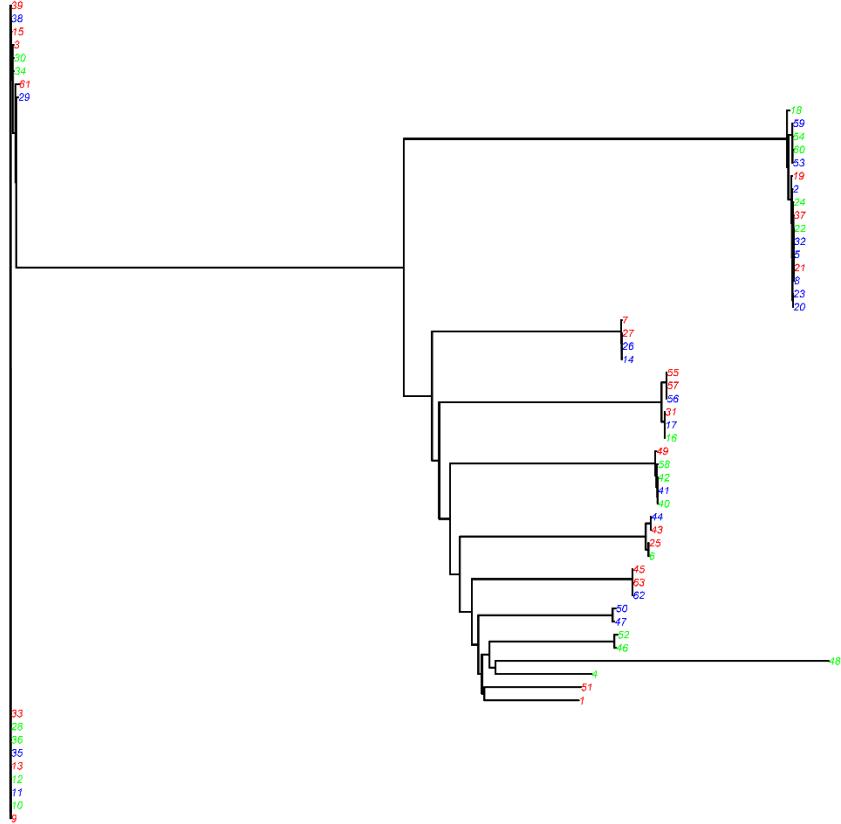
# Convert the distance matrix to a distance object
dist_obj <- dist(dist_matrix)

# Generate a tree topology using the neighbor-joining algorithm
tree <- nj(dist_obj)

# Define a vector of colors for each group
group_colors <- c("red", "blue", "green")

# Assign a color to each tip label based on the group variable
group <- c("Group1", "Group2", "Group1", "Group3", "Group2", "Group3")
tip_colors <- group_colors[as.factor(group)]

# Plot the tree with larger font size, longer branches, and colored tip labels
plot(tree, cex = 0.3, edge.length = 5, tip.color = tip_colors)
```



## Optimize the phylogenetic tree topology.

**Answer:**

The bellow code optimizes the tree topology using maximum likelihood estimation (`pml()`) from the `phangorn` package. The model used is General Time Reversible (GTR) with gamma-distributed rate heterogeneity and an estimated proportion of invariable sites. The `optInv`, `optGamma`, and `optParams` options are set to `TRUE` to optimize the proportion of invariable sites, the shape parameter of the gamma distribution, and the base frequencies, respectively. The `trace` option is set to 0 to suppress the output of optimization details.

```

# Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the nucleotide alignment
alignment <- read.phyDat("fereshteh_gisaid_output.fasta.raxml.reduced.phy", format = "phylip")

# Convert the alignment to a DNABin object
alignment_bin <- as.DNABin(alignment)

# Calculate pairwise distances between sequences
distances <- dist.dna(alignment_bin, model = "raw")

# Generate a tree topology based on pairwise distances
tree <- nj(distances)

# Optimize the phylogenetic tree topology using ML estimation
fit <- pml(tree, data = alignment, model = "GTR", optInv = TRUE, optGamma = TRUE,
            optParams = TRUE, trace = 0)

# Save the optimized tree topology to a file
write.tree(fit$tree, "optimized_tree.nwk")

```

Then to plot my tree:

```

# Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the optimized tree topology
tree_optimized <- read.tree("optimized_tree.nwk")

# Prune the tree to remove some nodes or branches
tree_pruned <- drop.tip(tree_optimized, c("hCoV-19/Iran/NIC-A7/2021|EPI_ISL_8754024|2021-07-06:1e-08"))

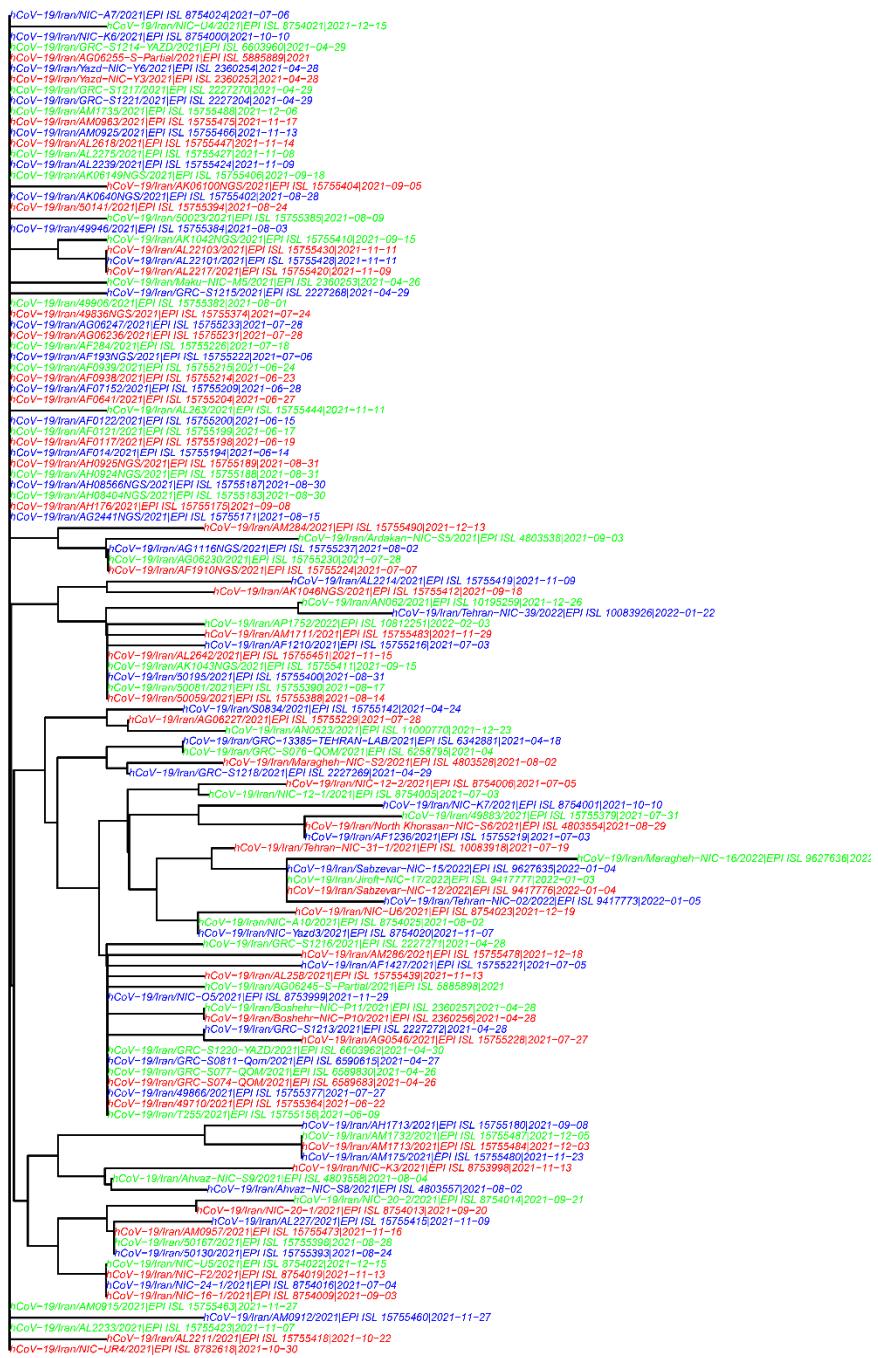
# Collapse multifurcating nodes to simplify the tree
tree_simple <- multi2di(tree_pruned)

# Define a vector of colors for each group
group_colors <- c("red", "blue", "green")

# Assign a color to each tip label based on the group variable
group <- c("Group1", "Group2", "Group1", "Group3", "Group2", "Group3")
tip_colors <- group_colors[as.factor(group)]

# Plot the simplified tree with larger font size, longer branches, and colored tip labels
plot(tree_simple, cex = 0.3, edge.length = 10, tip.color = tip_colors)

```



## Optimize the branch lengths.

## Answer:

The bellow code optimizes the branch lengths of the tree using maximum likelihood estimation (`optim.pml()`) from the `phangorn` package. This step refines the branch lengths of the tree based on the optimized topology. The `optNni`, `optBf`, and `nniBf` options are set to `TRUE` to optimize the

branch lengths using the NNI (nearest-neighbor interchange) and BF (branch swapping) algorithms. The trace option is set to 0 to suppress the output of optimization details.

```
## Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the nucleotide alignment
alignment <- read.phyDat("fereshteh_gisaid_output.fasta.raxml.reduced.phy", format = "phylip")

# Convert the alignment to a DNABin object
alignment_bin <- as.DNABin(alignment)

# Calculate pairwise distances between sequences
distances <- dist.dna(alignment_bin, model = "raw")

# Generate a tree topology based on pairwise distances
tree <- nj(distances)

# Optimize the phylogenetic tree topology using ML estimation
fit <- pml(tree, data = alignment, model = "GTR", optInv = TRUE, optGamma = TRUE,
            optParams = TRUE, trace = 0)

# Optimize the branch lengths using ML estimation
fit <- optim.pml(fit, optNni = TRUE, optBf = TRUE, nniBf = TRUE, trace = 0)

## optimize edge weights: -56028.29 --> -51747.41
## optimize base frequencies: -51747.41 --> -50796.31
## optimize edge weights: -50796.31 --> -50796.11
## optimize topology: -50796.11 --> -50173.84 NNI moves: 65
## optimize base frequencies: -50173.84 --> -50173.8
## optimize edge weights: -50173.8 --> -50173.8
## optimize topology: -50173.8 --> -50085.55 NNI moves: 23
## optimize base frequencies: -50085.55 --> -50085.55
## optimize edge weights: -50085.55 --> -50085.55
## optimize topology: -50085.55 --> -50043.02 NNI moves: 9
## optimize base frequencies: -50043.02 --> -50043.02
## optimize edge weights: -50043.02 --> -50043.02
## optimize topology: -50043.02 --> -50033.04 NNI moves: 4
## optimize base frequencies: -50033.04 --> -50033.04
## optimize edge weights: -50033.04 --> -50033.04
## optimize topology: -50033.04 --> -50033.04 NNI moves: 0

## Save the optimized tree topology and branch lengths to a file
write.tree(fit$tree, "optimized_tree_branch.nwk")
```

Then to plot my tree:

```
# Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the optimized tree topology
tree_optimized <- read.tree("optimized_tree_branch.nwk")

# Prune the tree to remove some nodes or branches
tree_pruned <- drop.tip(tree_optimized, c("hCoV-19/Iran/NIC-A7/2021|EPI_ISL_8754024|2021-07-06:1e-08"))

# Collapse multifurcating nodes to simplify the tree
tree_simple <- multi2di(tree_pruned)

# Define a vector of colors for each group
group_colors <- c("red", "blue", "green")

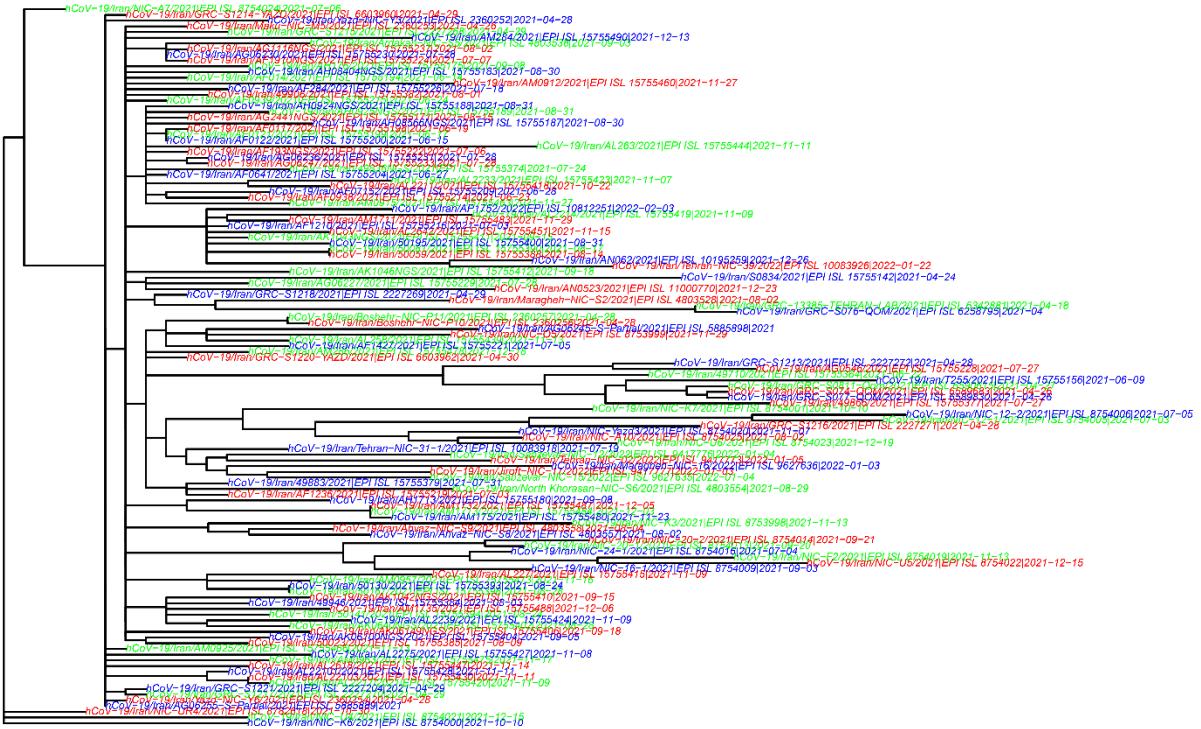
# Assign a color to each tip label based on the group variable
group <- c("Group1", "Group2", "Group1", "Group3", "Group2", "Group3")
tip_colors <- group_colors[as.factor(group)]

# Plot the simplified tree with larger font size, longer branches, and colored tip labels
plot(tree_simple, cex = 0.3, edge.length = 10, tip.color = tip_colors)

## Warning in plot.window(...): "edge.length" is not a graphical parameter

## Warning in plot.xy(xy, type, ...): "edge.length" is not a graphical parameter

## Warning in title(...): "edge.length" is not a graphical parameter
```



## Optimize the nucleotide substitution model.

**Answer:**

The bellow code, optimizes the nucleotide substitution model parameters for the tree using maximum likelihood estimation (pml()) from the phangorn package. The model used is the best-fitting model selected in step 3. The optInv, optGamma, and optParams options are set to TRUE to optimize the proportion of invariable sites, the shape parameter of the gamma distribution, and the base frequencies, respectively. The trace option is set to 0 to suppress the output of optimization details.

```

# Load necessary packages
library(ape)
## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)
## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the nucleotide alignment
alignment <- read.phyDat("ferehshteh_gisaid_output.fasta.raxml.reduced.phy", format = "phyip")

# Estimate the best-fitting nucleotide substitution model
model <- modelTest(alignment)

## Model df logLik AIC BIC
## JC 249 -49553.41 99884.81 301872.9
## JC+I 250 -49302.13 99184.26 301188.7
## JC+G(4) 250 -49451.21 99462.41 301478.8
## JC+G(G)+I 251 -49281.93 99665.86 301159.6
## F81 252 -48679.13 97862.26 99955.3
## F81+I 253 -48343.76 97193.53 99294.87
## F81+G(4) 253 -48479.45 97454.8 99556.75
## F81+G(G)+I 254 -48323.2 97154.4 99264.85
## K80 258 -49436.51 99773.83 301449.5
## K80+I 251 -49884.77 98871.88 300756.6
## K80+G(4) 251 -49734.27 98976.54 301055.3
## K80+G(G)+I 252 -49864.54 98533.79 300776.3
## HVY 253 -48444.08 97934.11 99495.46
## HVY+I 254 -48125.12 96558.24 98867.89
## HVY+G(4) 254 -48247.15 97807.31 99111.96
## HVY+G(G)+I 255 -48161.01 96718.40 98836.05
## TPM 252 -48034.01 99266.41 301284.8
## TPM+I 253 -48034.01 98576.41 301188.0
## IrmEcG(4) 252 -49141.4 98794.6 300889.8
## IrmEcG(G)+I 253 -49850.11 98451.42 300518.8
## Irm 254 -48156.18 97229.36 99139.41
## Irm+I 255 -48848.05 96605.11 98742.07
## Irm+G(4) 255 -48169.05 96861.1 98944.05
## Irm+G(G)+I 256 -48924.88 96568.92 98587.18
## Irm 251 -49375.81 99221.66 301118.4
## TPM1+I 252 -49024.25 98552.5 300645.5
## TPM1+G(4) 252 -49173.38 98851.16 100944.2
## TPM1+G(G)+I 253 -49803.3 98513.88 100815.2
## K81 251 -49375.83 99253.66 301338.4
## K81+I 252 -49024.25 98552.5 300645.5
## K81+G(4) 252 -49173.38 98851.16 100944.2
## K81+G(G)+I 253 -49803.3 98513.88 100815.2
## TPM1 254 -48373.58 97255.15 99364.8
## TPM1+I 255 -48061.47 96632.95 98750.91
## TPM1+G(4) 255 -48177.81 96805.61 98983.57
## TPM1u(G(G)) 256 -48040.2 96592.39 98718.66
## TPM 251 -49375.51 99253.08 101337.8
## TPM2+I 252 -49037.27 98578.54 100871.6
## TPM2+G(4) 252 -49175.11 98854.63 100947.7
## TPM2(G(G)) 253 -49816.76 98539.76 100864.11
## TPM2 254 -48363.72 97235.45 99345.1
## TPM2+I 255 -48053.46 96816.91 98734.87
## TPM2u(G(G)) 256 -48167.93 96815.81 98963.8
## TPM2u(G(G)+I) 256 -48082.39 96576.38 98762.64
## TPM 251 -49423.18 99348.35 101433.1
## TPM 252 -49024.25 98552.5 300645.5
## TPM3+G(4) 255 -49210.76 98941.51 100874.4
## TPM3+G(G)+I 255 -49893.05 98973.72 98875.1
## TPM 254 -49420.03 97207.05 99465.74
## TPM3u 255 -49112.5 96734.09 98852.95
## TPM3u(G(4)) 255 -49241.59 96093.18 99111.14
## TPM3u(G(G)) 256 -48011.85 96605.61 98821.88
## TPM 252 -49288.37 98674.74 101373.8
## TIM1a+I 253 -49042.31 98396.62 100495.92
## TIM1a+G(4) 253 -49885.71 98677.43 100877.8
## TIM1a+G(G)+I 254 -48928.62 98348.93 1008457.7
## TIM 255 -49285.7 97981.4 99209.36
## TIM1+I 256 -47983.76 96479.53 98605.79
## TIM1+G(4) 256 -48809.7 96693.39 98819.66
## TIM1+G(G)+I 257 -47960.13 96434.26 98558.82
## TIM2+I 252 -49288.88 99080.17 301173.2
## TIM2+G(4) 253 -48957.12 98410.23 100511.6
## TIM2+G(G)+I 254 -49088.96 98679.92 1004781.3
## TIM2u(G(G)) 254 -48939.39 98368.78 1004278.4
## TIM 255 -48376.05 97669.49 99180.05
## TIM+I 256 -47974.05 96465.28 98591.54
## TIM2+G(4) 256 -48881.53 98575.05 98801.32
## TIM2+G(G)+I 257 -47953.78 96471.55 98556.13
## TIM3+I 252 -49335.72 9915.44 301268.5
## TIM3+G(4) 253 -48971.49 98448.98 100850.3
## TIM3+G(G)+I 253 -49130.56 98767.13 100868.5
## TIM3u(G(G)) 254 -48959.12 98408.23 1008517.9
## TIM 255 -48351.71 97213.42 99331.38
## TIM+I 256 -48037.09 95866.09 98712.36
## TIM3+G(4) 256 -49121.52 98623.23 98847.59
## TIM3+G(G)+I 257 -49044.09 98610.16 98777.73
## TIM 254 -49329.01 99159.03 301265.4
## TIM+G(4) 254 -49971.72 98451.44 100951.1
## TIM+G(G)+I 254 -49126.02 98768.05 100868.7
## TIM3+G(4) 255 -48911.85 98411.3 100868.5
## TIM 256 -48323.36 97159.72 99285.99
## TIM+I 257 -48008.87 96527.75 98682.31
## TIM+G(4) 257 -48127.04 96708.08 98892.65
## TIM+G(G)+I 258 -47985.76 96487.51 98530.38
## TIM 254 -49241.56 98931.12 101180.8
## TIM+I 255 -48885.74 98281.49 100399.4
## TIM1+I 255 -49037.34 98584.68 100702.6
## SYM1G(4) 256 -48864.91 98241.82 100368.1
## GTR 257 -48236.0 96986.48 99121.05
## GTR+I 258 -47931.3 96378.68 98521.55
## GTR+G(4)+I 259 -47908.97 96335.94 98487.12
## GTR+G(G)+I 259 -47908.97 96335.94 98487.12

# Convert the alignment to a DNAbin object
alignment_bln <- as.DNAbin(alignment)

# Generate a tree object based on the optimized tree topology and branch lengths
tree <- readTree("optimized_tree_branch.nwk")

# Optimize the nucleotide substitution model parameters using ML estimation
fit <- pml(tree, data = alignment, model = model, optInv = TRUE, optGamma = TRUE,
          optParams = TRUE, trace = 0)

# Save the optimized tree topology, branch lengths, and nucleotide substitution model parameters to a file
write.tree(tree, "optimized_tree.model.nwk")

```

Then to plot my tree:

```
# Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the optimized tree topology
tree_optimized <- read.tree("optimized_tree_model.nwk")

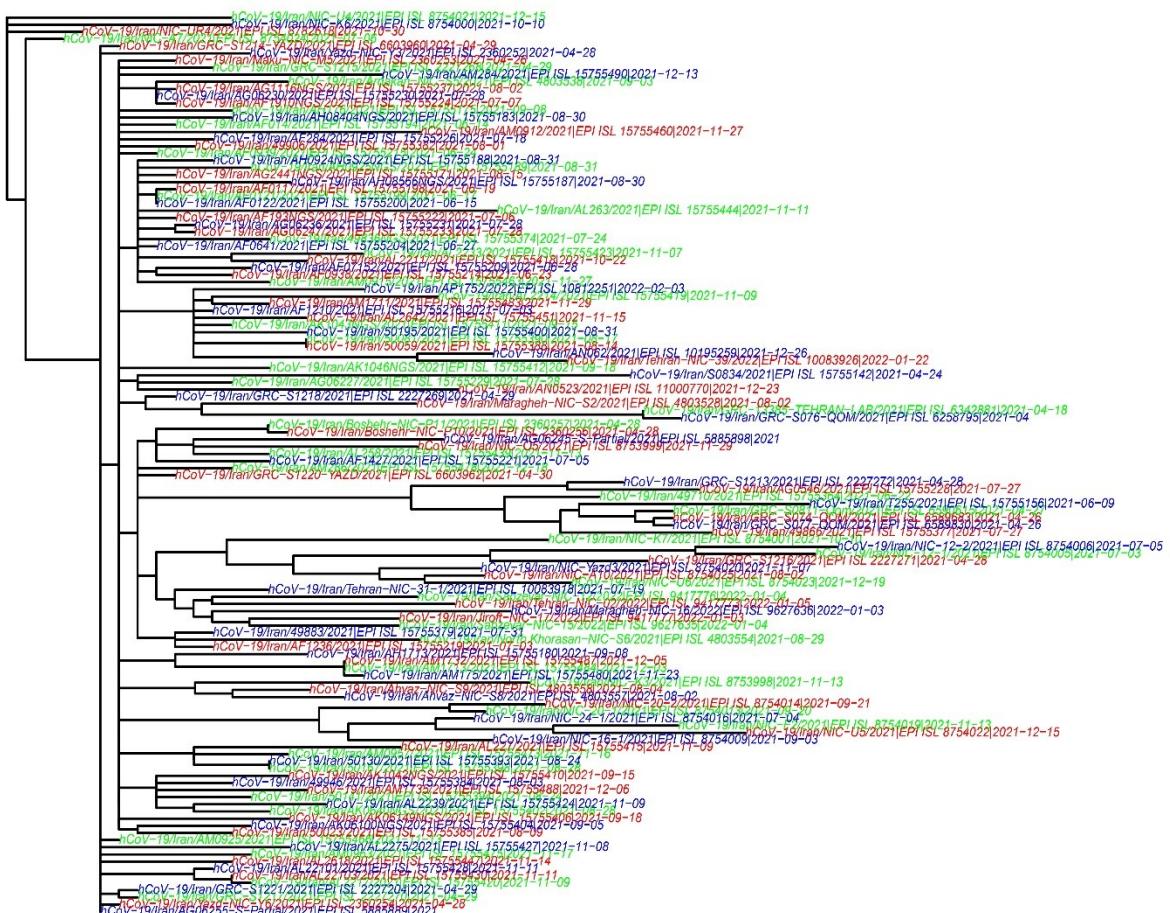
# Prune the tree to remove some nodes or branches
tree_pruned <- drop.tip(tree_optimized, c("hCoV-19/Iran/NIC-A7/2021|EPI_ISL_8754024|2021-07-06:1e-08"))

# Collapse multifurcating nodes to simplify the tree
tree_simple <- multi2di(tree_pruned)

# Define a vector of colors for each group
group_colors <- c("red", "blue", "green")

# Assign a color to each tip label based on the group variable
group <- c("Group1", "Group2", "Group1", "Group3", "Group2", "Group3")
tip_colors <- group_colors[as.factor(group)]

# Plot the simplified tree with larger font size, longer branches, and colored tip labels
plot(tree_simple, cex = 0.3, edge.length = 10, tip.color = tip_colors)
```



At last, I compared the three optimized trees:

```
# Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the trees
tree1 <- read.tree("optimized_tree.nwk")
tree2 <- read.tree("optimized_tree_branch.nwk")
tree3 <- read.tree("optimized_tree_model.nwk")

# Calculate pairwise RF distances between all three trees
dist1 <- RF.dist(tree1, tree2)
dist2 <- RF.dist(tree1, tree3)
dist3 <- RF.dist(tree2, tree3)

# Print the pairwise distances
cat("RF distance between optimized_tree.nwk and optimized_tree_branch.nwk:", dist1, "\n")

## RF distance between optimized_tree.nwk and optimized_tree_branch.nwk: 136

cat("RF distance between optimized_tree.nwk and optimized_tree_model.nwk:", dist2, "\n")

## RF distance between optimized_tree.nwk and optimized_tree_model.nwk: 136

cat("RF distance between optimized_tree_branch.nwk and optimized_tree_model.nwk:", dist3, "\n")

## RF distance between optimized_tree_branch.nwk and optimized_tree_model.nwk: 0
```

The RF distance is a measure of the topological distance between two trees, and it represents the number of bipartitions that differ between the trees. A distance of 0 indicates that the two trees are identical in their topology, while a higher distance indicates a greater difference in their topology. The output shows that the RF distance between "optimized\_tree.nwk" and "optimized\_tree\_branch.nwk" is 136, indicating a substantial difference in their topology. The RF distance between "optimized\_tree.nwk" and "optimized\_tree\_model.nwk" is also 136, indicating a similar level of difference in topology. However, the RF distance between "optimized\_tree\_branch.nwk" and "optimized\_tree\_model.nwk" is 0, indicating that these two trees have identical topologies.

To assess the reliability of the phylogenetic tree, generate multiple bootstrap replicates of the dataset and estimate the phylogenetic tree for each replicate.

Answer:

The bellow code, demonstrates how to generate multiple bootstrap replicates of an alignment file and estimate the phylogenetic tree for each replicate in order to assess the reliability of the resulting tree.

It first reads in an alignment file in PHYLIP format using the `read.phyDat()` function from the `ape` package. It then converts the alignment to a `DNAbin` object using the `as.DNAbin()` function. Next, the code estimates the phylogenetic tree using the neighbor-joining (NJ) method with the `nj()` function from the `ape` package. The code then sets the number of bootstrap replicates to generate (`nboot`) and uses a for loop to generate the bootstrap replicates and estimate the trees using the NJ method. For each bootstrap replicate, the code saves a plot of the resulting tree as a PNG file using the `png()`, `plot()`, and `dev.off()` functions.

```
# Load the ape package
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in your alignment file
aln <- read.phyDat("fereshteh_gisaid_output.fasta.raxml.reduced.phy")

# Convert the alignment to a DNAbin object
dnabin_aln <- as.DNAbin(aln)

# Estimate the phylogenetic tree
tree <- nj(dist.dna(dnabin_aln))

# Set the number of bootstrap replicates
nboot <- 100

# Generate a set of bootstrap replicates and estimate the trees
boot_trees <- vector(mode = "list", length = nboot)
for (i in 1:nboot) {
  boot_dnabin_aln <- as.DNAbin(sample(aln, replace = TRUE))
  boot_trees[[i]] <- nj(dist.dna(boot_dnabin_aln))
}

# Save each bootstrap tree plot as a PNG file
for (i in 1:nboot) {
  # Set the filename for this bootstrap tree plot
  filename <- paste0("boot_tree_", i, ".png")

  # Open a new PNG device and plot the tree
  png(filename, width = 1500, height = 800, res = 100)
  plot(boot_trees[[i]], cex = 0.3)

  # Close the PNG device
  dev.off()
}
```

The code generates 100 PNG files and 100 tree files as output, which can be find in the current working directory.

## Infer bootstrap support for branches in the optimized tree

Answer:

The Felsenstein confidence (FC) is a method for calculating bootstrap support values for phylogenetic trees that is based on a resampling technique known as the bootstrap. The basic idea of the FC method is to compare the log-likelihood of the original tree to the log-likelihoods of a set of bootstrap trees, and to calculate a confidence value for each node based on the proportion of bootstrap trees in which the node appears.

RAxML can calculate bootstrap support values for the nodes in the maximum likelihood tree using a resampling technique known as non-parametric bootstrap.

### Web services

<b>RAxML-Light Web-Service</b> Available at the San Diego Supercomputer Center (supported by the NSF iPlant collaborative). To use this service you will first need to create an iPlant login <a href="#">here</a> and subsequently log in on the <a href="#">CIPRES portal</a> using your iPlant credentials.	<b>WEB-Servers for RAxML-NG</b> The brand new web-server for the completely redesigned version of RAxML, called RAxML-NG is available <a href="#">here</a> .	<b>Rogue Taxon Identification Web-Server</b> <a href="#">access the server here</a> .
---	---	--

I used RAxML-NG (<https://raxml-ng.vital-it.ch/#/>).

The screenshot shows the RAxML BlackBox web interface. At the top, there are links for SIB, iWTS, RAxML, Software, Mirror at CIPRES, and Contact, along with a version notice "RAxML v1.0.0". Below this is the title "RAxML BlackBox". A horizontal line labeled "Data" separates the header from the input fields. The first field is for "sequence alignment", with a placeholder "Paste your sequence alignment" and a large text area for input. A blue callout box next to it specifies: "Must be in PHYLIP format, or FASTA format, or convertible by readseq". Below this is a file upload field: "or upload a file" followed by a "Choose File" button and a message "No file chosen". The second field is for "constraint tree", with a placeholder "Paste your constraint tree" and a large text area for input. A blue callout box next to it specifies: "This option allows you to specify an incomplete or comprehensive multifurcating constraint tree in NEWICK format. More help [here](#)". Below this is a file upload field: "or upload a file" followed by a "Choose File" button and a message "No file chosen".

The output is containing:

# sequenceAlignment.phy	7/9/2023 8:56 PM	PHYLIP File	3,687 KB
slurm_raxml.sbatch	7/9/2023 8:56 PM	SBATCH File	1 KB
raxmlArg.txt	7/9/2023 9:27 PM	Text Document	1 KB
result.raxml.bestModel	7/9/2023 9:27 PM	BESTMODEL File	1 KB
result.raxml.bestTree	7/9/2023 9:27 PM	BESTTREE File	10 KB
result.raxml.bestTreeCollapsed	7/9/2023 9:27 PM	BESTTREECOLLAP...	9 KB
result.raxml.log	7/9/2023 9:27 PM	Text Document	16 KB
result.raxml.mlTrees	7/9/2023 9:27 PM	MLTREES File	184 KB
result.raxml.rba	7/9/2023 9:27 PM	RBA File	98 KB
result.raxml.startTree	7/9/2023 9:27 PM	STARTTREE File	184 KB

## Data visualization

Use the phylogenetic tree and bootstrap support values to visualize the evolutionary history of the selected strain in Iran. (You may use [ggtree](#) package in R)

Answer:

```
# Load necessary packages
library(ape)

## Warning: package 'ape' was built under R version 4.2.3

library(phangorn)

## Warning: package 'phangorn' was built under R version 4.2.3

# Read in the optimized tree topology
tree_optimized <- read.tree("result.raxml-2.nwk")

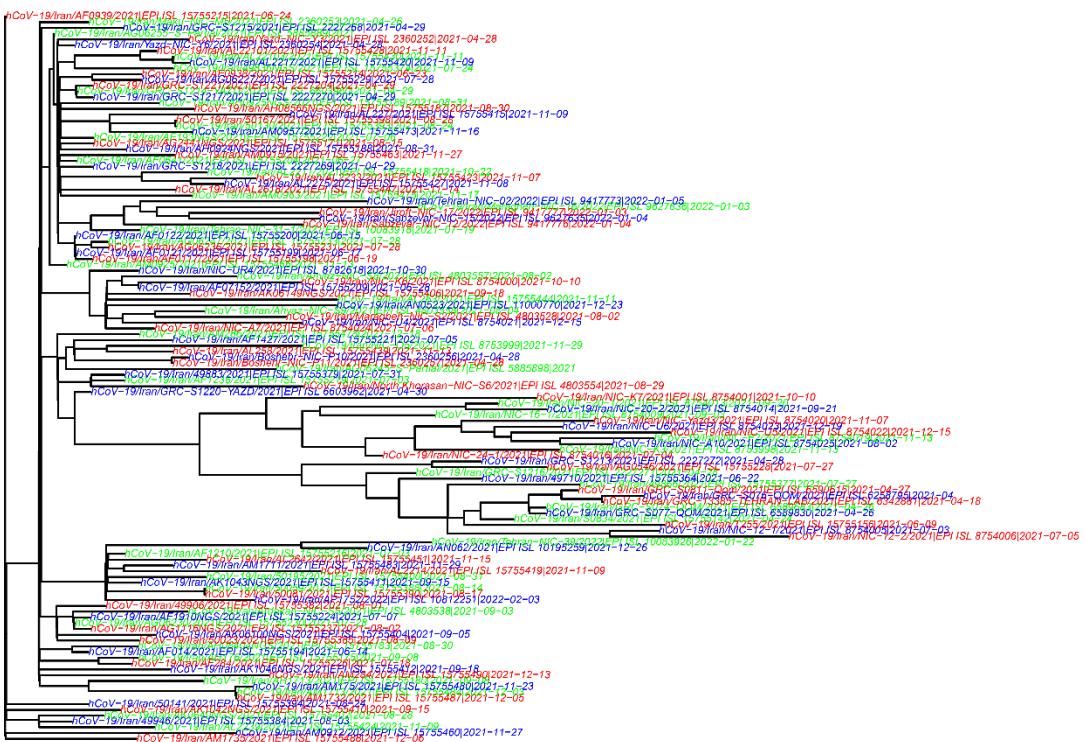
# Prune the tree to remove some nodes or branches
tree_pruned <- drop.tip(tree_optimized, c("hCoV-19/Iran/NIC-A7/2021|EPI_ISL_8754024|2021-07-06:1e-08"))

# Collapse multifurcating nodes to simplify the tree
tree_simple <- multi2di(tree_pruned)

# Define a vector of colors for each group
group_colors <- c("red", "blue", "green")

# Assign a color to each tip label based on the group variable
group <- c("Group1", "Group2", "Group1", "Group3", "Group2", "Group3")
tip_colors <- group_colors[as.factor(group)]

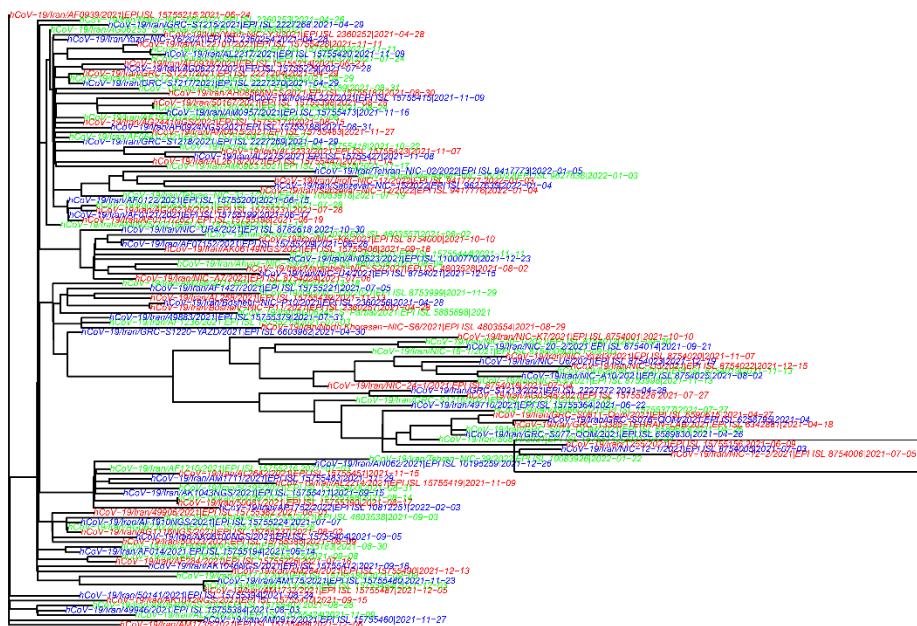
# Plot the simplified tree with larger font size, longer branches, and colored tip labels
plot(tree_simple, cex = 0.3, edge.length = 10, tip.color = tip_colors)
```



A rooted phylogenetic tree of the Delta strain of COVID-19 in Iran was generated as the outcome.

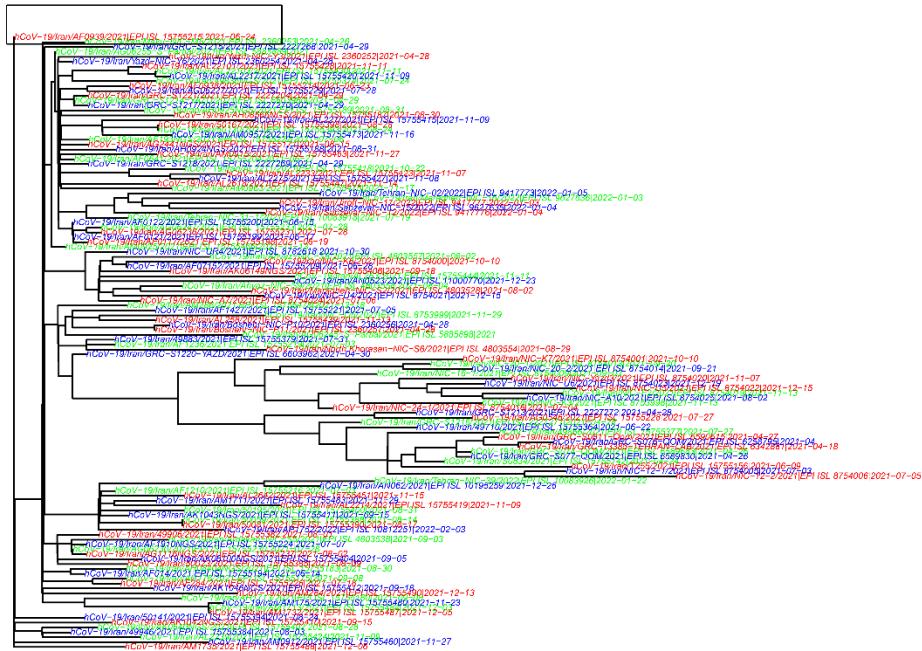
1. The meaning of the longest branch from the root in a phylogenetic tree can depend on the specific context of the analysis. In general, however, a long branch can represent a period of relatively rapid evolution, such as a burst of adaptive radiation or a period of rapid divergence due to environmental or ecological factors. Alternatively, a long branch may simply be the result of a lack of sampling or incomplete information about the evolutionary relationships among the taxa being studied that in my tree is there:

hCoV-19/Iran/NIC-12-2/2021|EPI ISL 8754006|2021-07-05 is the strain that has accumulated a relatively large amount of genetic divergence compared to the other sequences in the tree. This could be due to a variety of factors, such as a longer period of evolution, more mutations or recombination events, or a different evolutionary trajectory compared to the other sequences in the tree.

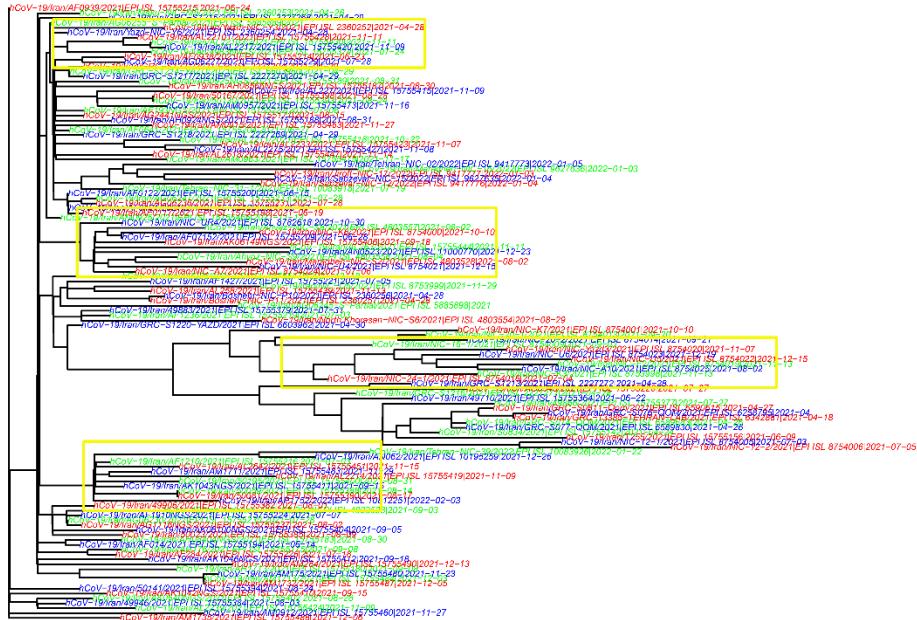


2. The meaning of the shortest branch in a phylogenetic tree can depend on the specific context of the analysis. In general, however, a short branch can indicate recent common ancestor, a period of relatively little evolutionary change, or a lack of genetic diversity among the sequences being analyzed.

hCoV-19/Iran/AF0939/2021|EPI ISL 15755215|2021-06-24 is the strain that is very closely related to the other sequences in the tree and has not diverged much from their common ancestor. This could indicate that this particular sequence is part of a recent outbreak or cluster of cases that has not yet had time to accumulate many genetic changes.



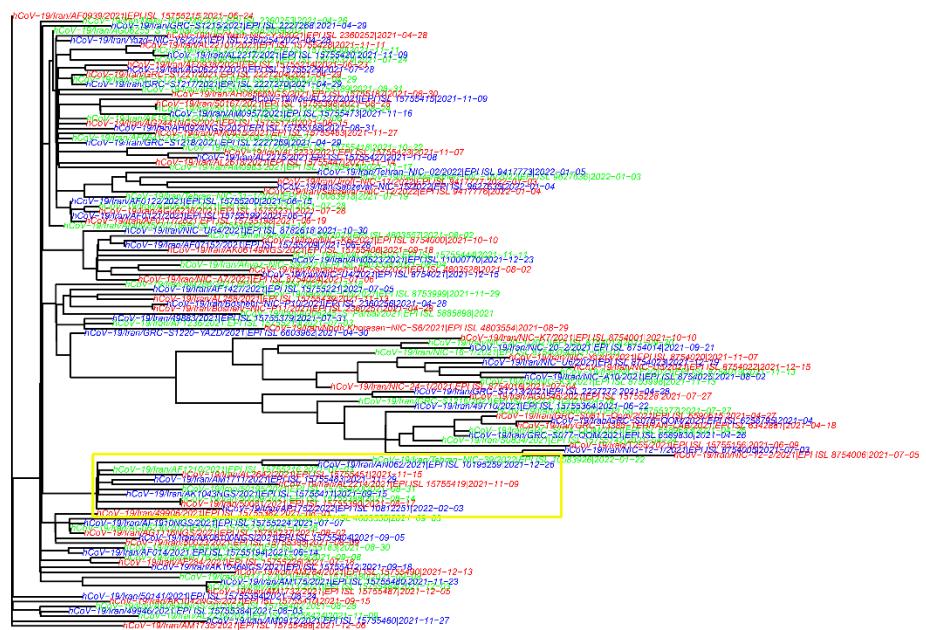
3. In a phylogenetic tree, a clade is a group of organisms that share a common ancestor, and are therefore more closely related to each other than to any other organisms in the tree. In my tree there are some clades that I determined some of them:



4. In my tree some clades that has more taxa in a particular clade in your phylogenetic tree, it could indicate that this group of organisms has undergone more diversification or evolution compared to other clades in the tree. This could be due to factors such as geographic isolation, differences in host range or transmission dynamics, or a combination of genetic and environmental factors that have led to adaptation and divergence over time. Alternatively, the number of taxa in a clade could be influenced by sampling biases in the data, such as differences in the number or quality of sequences available for different regions or time periods. For example, if more sequences have been sampled from a particular region or time period, this could lead to the appearance of a larger clade in the tree.

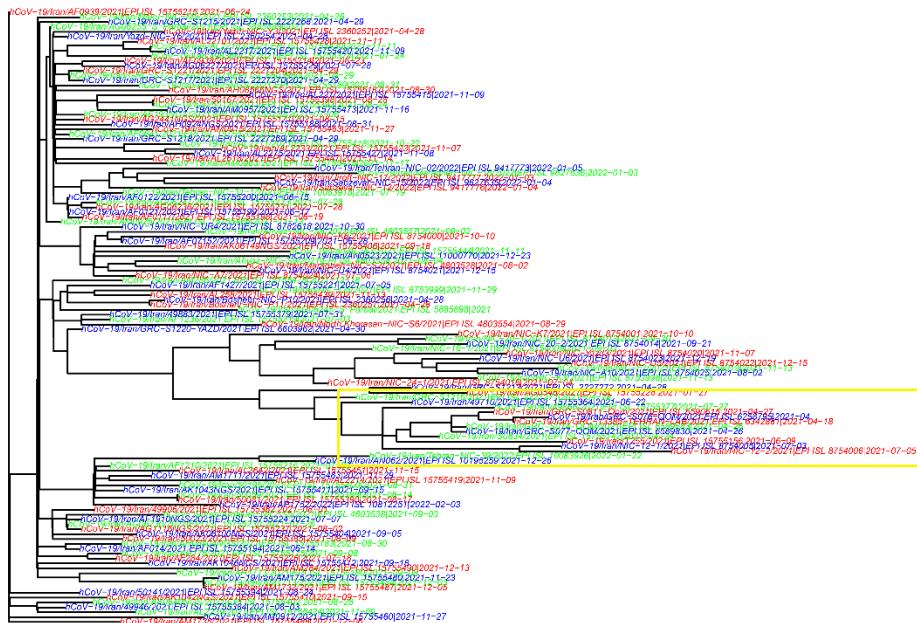
In my tree there are someclades that are highly concentrated , with closely related taxa clustered closely together in the tree, it may suggest that these organisms have undergone relatively little evolution or divergence from their common ancestor, or that they have diversified rapidly and recently. This could be due to factors such as geographic isolation, host range or transmission dynamics, or the evolution of novel traits or adaptations that allow them to exploit new ecological niches.

For example:



5. Also in my tree there are some clade that are more dispersed, with closely related taxa scattered across different parts of the tree, it may suggest that these organisms have undergone a greater degree of evolutionary change or diversification. This could be due to factors such as longer periods of time since their common ancestor, more complex evolutionary histories involving multiple speciation events or hybridization events, or more gradual patterns of adaptation and diversification.

For example:



Therefore, in conclusion, my phylogenetic tree of the Delta strain of COVID-19 in Iran displays a diverse range of strains and clades that can provide valuable insights into the evolutionary relationships and transmission patterns of the virus in the region. The longest branch in the tree may indicate a period of rapid evolution due to environmental or ecological factors, while the shortest branch may suggest a recent common ancestor, indicating that the sequence is part of a recent outbreak or cluster of cases with few genetic changes.

The number of taxa in a clade can indicate the level of diversification or evolution compared to other clades in the tree, due to factors such as geographic isolation, differences in host range or transmission dynamics, or genetic and environmental factors. Highly concentrated clades with closely related taxa clustered closely together in the tree may suggest that these organisms have

undergone little evolution or divergence from their common ancestor, or that they have diversified rapidly and recently due to novel traits or adaptations that allow them to exploit new ecological niches.