

SARSCoV2Analysis

Biology3579

2025-03-18

R-set-up

```
# Restoring project's environment  
renv::restore()
```

```
## - The library is already synchronized with the lockfile.
```

```
# Source library loading function  
source(here::here("functions", "libraries.R"))  
# Load necessary libraries for analysis  
load_libraries()
```

Introduction

The evolution of emerging pathogens, such as SARS-CoV-2, plays a critical role in shaping the trajectory of pandemics. Viral variants arise through genetic mutations, some of which confer increased transmissibility, immune escape, or enhanced replication fitness. These evolutionary changes impact the spread, dominance, and public health burden of different variants. The Alpha, Delta, and Omicron variants, for instance, demonstrated distinct epidemiological characteristics, leading to waves of infections with varying severity and control challenges.

Understanding the fitness advantage of viral variants—defined by their ability to outcompete predecessors in a given population—is essential for anticipating shifts in transmission dynamics. The basic and effective reproduction numbers, which quantify how efficiently a variant spreads, provide key epidemiological insights into variant dominance and the impact of interventions. By leveraging genomic surveillance data, it is possible to estimate these metrics and assess the selective advantage of emerging variants.

This analysis aims to:

Quantify the differences in fitness advantage and reproduction number among major SARS-CoV-2 variants during the COVID-19 pandemic in the UK. Estimate variant-specific growth rates using genomic sequencing data. Evaluate how changes in these metrics influence the effectiveness of epidemic control strategies.

Question 1: SARS-CoV-2 Major Lineages and Trends

This section analyses the major SARS-CoV-2 lineages circulating in England over time using data from the Sanger Institute's COG-UK programme. The dataset contains weekly counts of virus samples per lineage, providing insight into the evolution and spread of different variants.

The analysis consists of the following steps: 1. **Loading the dataset** – Importing the raw data from an online repository. 2. **Classifying major lineages** – Identifying key variants and grouping all others under a general category. 3. **Visualising lineage trends** – Generating stacked area plots to illustrate changes in lineage prevalence over time.

1.1 Load Sanger Data

The dataset is imported from a GitHub repository. The file `Genomes_per_week_in_England.csv` contains weekly counts of SARS-CoV-2 genomic sequences assigned to specific lineages.

```
# Load the raw dataset from an online repository (ensures reproducibility)
sanger_raw <- read.csv("https://raw.githubusercontent.com/Biology3579/SARSCoV2Assignment/main/data/Genomes_per_week_in_England.csv")

# Save a local copy of the dataset for future reference
write_csv(sanger_raw, here("data", "sanger_raw.csv"))
```

1.2 Process Sanger Data

Before conducting the analysis, the raw data must be processed to ensure that it is structured appropriately for statistical analysis and meaningful interpretation. The key processing steps include:

- Identification and classification of major SARS-CoV-2 lineages: Variants of concern and interest exhibit different transmission rates, immune evasion capabilities, and epidemiological impact. Classifying major lineages enables a clearer understanding of their relative dominance and replacement dynamics over time. This classification focuses on key variants such as Alpha (B.1.1.7), Delta (B.1.617.2), Omicron subvariants (BA.1, BA.2, BA.2.75, BA.4, BA.5, BA.5.3/BQ.1), and XBB, which have played significant roles in shaping the pandemic. By isolating these lineages, their individual growth trends and competitive advantages can be assessed.

All other other lineages - low-prevalence lineages with a much less significant epidemiological impact - are grouped as 'Other'. This approach helps avoid statistical noise and enhances the clarity of lineage trends.

- Formatting the dataset for further statistical analysis and visualisation: The raw dataset requires restructuring to facilitate robust statistical modeling and visualisation.

A more detailed breakdown of these processing steps is provided in the `cleaning_and_curating.R` script, from which the processing functions are sourced.

```
# Source the data processing functions
source(here("functions", "cleaning_and_curating.R"))

# Pipe to process data concisely
sanger_analysis_data <- sanger_raw %>%
  clean_sanger_data(variants = c("B.1.1.7", "B.1.617.2", "BA.1", "BA.2", "BA.2.75", "BA.4", "BA.5", "BA.5.3", "XBB.1.5", "XBB.1.6", "XBB.1.9", "XBB.1.10", "XBB.1.11", "XBB.1.12", "XBB.1.13", "XBB.1.14", "XBB.1.15", "XBB.1.16", "XBB.1.17", "XBB.1.18", "XBB.1.19", "XBB.1.20", "XBB.1.21", "XBB.1.22", "XBB.1.23", "XBB.1.24", "XBB.1.25", "XBB.1.26", "XBB.1.27", "XBB.1.28", "XBB.1.29", "XBB.1.30", "XBB.1.31", "XBB.1.32", "XBB.1.33", "XBB.1.34", "XBB.1.35", "XBB.1.36", "XBB.1.37", "XBB.1.38", "XBB.1.39", "XBB.1.40", "XBB.1.41", "XBB.1.42", "XBB.1.43", "XBB.1.44", "XBB.1.45", "XBB.1.46", "XBB.1.47", "XBB.1.48", "XBB.1.49", "XBB.1.50", "XBB.1.51", "XBB.1.52", "XBB.1.53", "XBB.1.54", "XBB.1.55", "XBB.1.56", "XBB.1.57", "XBB.1.58", "XBB.1.59", "XBB.1.60", "XBB.1.61", "XBB.1.62", "XBB.1.63", "XBB.1.64", "XBB.1.65", "XBB.1.66", "XBB.1.67", "XBB.1.68", "XBB.1.69", "XBB.1.70", "XBB.1.71", "XBB.1.72", "XBB.1.73", "XBB.1.74", "XBB.1.75", "XBB.1.76", "XBB.1.77", "XBB.1.78", "XBB.1.79", "XBB.1.80", "XBB.1.81", "XBB.1.82", "XBB.1.83", "XBB.1.84", "XBB.1.85", "XBB.1.86", "XBB.1.87", "XBB.1.88", "XBB.1.89", "XBB.1.90", "XBB.1.91", "XBB.1.92", "XBB.1.93", "XBB.1.94", "XBB.1.95", "XBB.1.96", "XBB.1.97", "XBB.1.98", "XBB.1.99", "XBB.1.100", "XBB.1.101", "XBB.1.102", "XBB.1.103", "XBB.1.104", "XBB.1.105", "XBB.1.106", "XBB.1.107", "XBB.1.108", "XBB.1.109", "XBB.1.110", "XBB.1.111", "XBB.1.112", "XBB.1.113", "XBB.1.114", "XBB.1.115", "XBB.1.116", "XBB.1.117", "XBB.1.118", "XBB.1.119", "XBB.1.120", "XBB.1.121", "XBB.1.122", "XBB.1.123", "XBB.1.124", "XBB.1.125", "XBB.1.126", "XBB.1.127", "XBB.1.128", "XBB.1.129", "XBB.1.130", "XBB.1.131", "XBB.1.132", "XBB.1.133", "XBB.1.134", "XBB.1.135", "XBB.1.136", "XBB.1.137", "XBB.1.138", "XBB.1.139", "XBB.1.140", "XBB.1.141", "XBB.1.142", "XBB.1.143", "XBB.1.144", "XBB.1.145", "XBB.1.146", "XBB.1.147", "XBB.1.148", "XBB.1.149", "XBB.1.150", "XBB.1.151", "XBB.1.152", "XBB.1.153", "XBB.1.154", "XBB.1.155", "XBB.1.156", "XBB.1.157", "XBB.1.158", "XBB.1.159", "XBB.1.160", "XBB.1.161", "XBB.1.162", "XBB.1.163", "XBB.1.164", "XBB.1.165", "XBB.1.166", "XBB.1.167", "XBB.1.168", "XBB.1.169", "XBB.1.170", "XBB.1.171", "XBB.1.172", "XBB.1.173", "XBB.1.174", "XBB.1.175", "XBB.1.176", "XBB.1.177", "XBB.1.178", "XBB.1.179", "XBB.1.180", "XBB.1.181", "XBB.1.182", "XBB.1.183", "XBB.1.184", "XBB.1.185", "XBB.1.186", "XBB.1.187", "XBB.1.188", "XBB.1.189", "XBB.1.190", "XBB.1.191", "XBB.1.192", "XBB.1.193", "XBB.1.194", "XBB.1.195", "XBB.1.196", "XBB.1.197", "XBB.1.198", "XBB.1.199", "XBB.1.200", "XBB.1.201", "XBB.1.202", "XBB.1.203", "XBB.1.204", "XBB.1.205", "XBB.1.206", "XBB.1.207", "XBB.1.208", "XBB.1.209", "XBB.1.210", "XBB.1.211", "XBB.1.212", "XBB.1.213", "XBB.1.214", "XBB.1.215", "XBB.1.216", "XBB.1.217", "XBB.1.218", "XBB.1.219", "XBB.1.220", "XBB.1.221", "XBB.1.222", "XBB.1.223", "XBB.1.224", "XBB.1.225", "XBB.1.226", "XBB.1.227", "XBB.1.228", "XBB.1.229", "XBB.1.230", "XBB.1.231", "XBB.1.232", "XBB.1.233", "XBB.1.234", "XBB.1.235", "XBB.1.236", "XBB.1.237", "XBB.1.238", "XBB.1.239", "XBB.1.240", "XBB.1.241", "XBB.1.242", "XBB.1.243", "XBB.1.244", "XBB.1.245", "XBB.1.246", "XBB.1.247", "XBB.1.248", "XBB.1.249", "XBB.1.250", "XBB.1.251", "XBB.1.252", "XBB.1.253", "XBB.1.254", "XBB.1.255", "XBB.1.256", "XBB.1.257", "XBB.1.258", "XBB.1.259", "XBB.1.260", "XBB.1.261", "XBB.1.262", "XBB.1.263", "XBB.1.264", "XBB.1.265", "XBB.1.266", "XBB.1.267", "XBB.1.268", "XBB.1.269", "XBB.1.270", "XBB.1.271", "XBB.1.272", "XBB.1.273", "XBB.1.274", "XBB.1.275", "XBB.1.276", "XBB.1.277", "XBB.1.278", "XBB.1.279", "XBB.1.280", "XBB.1.281", "XBB.1.282", "XBB.1.283", "XBB.1.284", "XBB.1.285", "XBB.1.286", "XBB.1.287", "XBB.1.288", "XBB.1.289", "XBB.1.290", "XBB.1.291", "XBB.1.292", "XBB.1.293", "XBB.1.294", "XBB.1.295", "XBB.1.296", "XBB.1.297", "XBB.1.298", "XBB.1.299", "XBB.1.300", "XBB.1.301", "XBB.1.302", "XBB.1.303", "XBB.1.304", "XBB.1.305", "XBB.1.306", "XBB.1.307", "XBB.1.308", "XBB.1.309", "XBB.1.310", "XBB.1.311", "XBB.1.312", "XBB.1.313", "XBB.1.314", "XBB.1.315", "XBB.1.316", "XBB.1.317", "XBB.1.318", "XBB.1.319", "XBB.1.320", "XBB.1.321", "XBB.1.322", "XBB.1.323", "XBB.1.324", "XBB.1.325", "XBB.1.326", "XBB.1.327", "XBB.1.328", "XBB.1.329", "XBB.1.330", "XBB.1.331", "XBB.1.332", "XBB.1.333", "XBB.1.334", "XBB.1.335", "XBB.1.336", "XBB.1.337", "XBB.1.338", "XBB.1.339", "XBB.1.340", "XBB.1.341", "XBB.1.342", "XBB.1.343", "XBB.1.344", "XBB.1.345", "XBB.1.346", "XBB.1.347", "XBB.1.348", "XBB.1.349", "XBB.1.350", "XBB.1.351", "XBB.1.352", "XBB.1.353", "XBB.1.354", "XBB.1.355", "XBB.1.356", "XBB.1.357", "XBB.1.358", "XBB.1.359", "XBB.1.360", "XBB.1.361", "XBB.1.362", "XBB.1.363", "XBB.1.364", "XBB.1.365", "XBB.1.366", "XBB.1.367", "XBB.1.368", "XBB.1.369", "XBB.1.370", "XBB.1.371", "XBB.1.372", "XBB.1.373", "XBB.1.374", "XBB.1.375", "XBB.1.376", "XBB.1.377", "XBB.1.378", "XBB.1.379", "XBB.1.380", "XBB.1.381", "XBB.1.382", "XBB.1.383", "XBB.1.384", "XBB.1.385", "XBB.1.386", "XBB.1.387", "XBB.1.388", "XBB.1.389", "XBB.1.390", "XBB.1.391", "XBB.1.392", "XBB.1.393", "XBB.1.394", "XBB.1.395", "XBB.1.396", "XBB.1.397", "XBB.1.398", "XBB.1.399", "XBB.1.400", "XBB.1.401", "XBB.1.402", "XBB.1.403", "XBB.1.404", "XBB.1.405", "XBB.1.406", "XBB.1.407", "XBB.1.408", "XBB.1.409", "XBB.1.410", "XBB.1.411", "XBB.1.412", "XBB.1.413", "XBB.1.414", "XBB.1.415", "XBB.1.416", "XBB.1.417", "XBB.1.418", "XBB.1.419", "XBB.1.420", "XBB.1.421", "XBB.1.422", "XBB.1.423", "XBB.1.424", "XBB.1.425", "XBB.1.426", "XBB.1.427", "XBB.1.428", "XBB.1.429", "XBB.1.430", "XBB.1.431", "XBB.1.432", "XBB.1.433", "XBB.1.434", "XBB.1.435", "XBB.1.436", "XBB.1.437", "XBB.1.438", "XBB.1.439", "XBB.1.440", "XBB.1.441", "XBB.1.442", "XBB.1.443", "XBB.1.444", "XBB.1.445", "XBB.1.446", "XBB.1.447", "XBB.1.448", "XBB.1.449", "XBB.1.450", "XBB.1.451", "XBB.1.452", "XBB.1.453", "XBB.1.454", "XBB.1.455", "XBB.1.456", "XBB.1.457", "XBB.1.458", "XBB.1.459", "XBB.1.460", "XBB.1.461", "XBB.1.462", "XBB.1.463", "XBB.1.464", "XBB.1.465", "XBB.1.466", "XBB.1.467", "XBB.1.468", "XBB.1.469", "XBB.1.470", "XBB.1.471", "XBB.1.472", "XBB.1.473", "XBB.1.474", "XBB.1.475", "XBB.1.476", "XBB.1.477", "XBB.1.478", "XBB.1.479", "XBB.1.480", "XBB.1.481", "XBB.1.482", "XBB.1.483", "XBB.1.484", "XBB.1.485", "XBB.1.486", "XBB.1.487", "XBB.1.488", "XBB.1.489", "XBB.1.490", "XBB.1.491", "XBB.1.492", "XBB.1.493", "XBB.1.494", "XBB.1.495", "XBB.1.496", "XBB.1.497", "XBB.1.498", "XBB.1.499", "XBB.1.500", "XBB.1.501", "XBB.1.502", "XBB.1.503", "XBB.1.504", "XBB.1.505", "XBB.1.506", "XBB.1.507", "XBB.1.508", "XBB.1.509", "XBB.1.510", "XBB.1.511", "XBB.1.512", "XBB.1.513", "XBB.1.514", "XBB.1.515", "XBB.1.516", "XBB.1.517", "XBB.1.518", "XBB.1.519", "XBB.1.520", "XBB.1.521", "XBB.1.522", "XBB.1.523", "XBB.1.524", "XBB.1.525", "XBB.1.526", "XBB.1.527", "XBB.1.528", "XBB.1.529", "XBB.1.530", "XBB.1.531", "XBB.1.532", "XBB.1.533", "XBB.1.534", "XBB.1.535", "XBB.1.536", "XBB.1.537", "XBB.1.538", "XBB.1.539", "XBB.1.540", "XBB.1.541", "XBB.1.542", "XBB.1.543", "XBB.1.544", "XBB.1.545", "XBB.1.546", "XBB.1.547", "XBB.1.548", "XBB.1.549", "XBB.1.550", "XBB.1.551", "XBB.1.552", "XBB.1.553", "XBB.1.554", "XBB.1.555", "XBB.1.556", "XBB.1.557", "XBB.1.558", "XBB.1.559", "XBB.1.560", "XBB.1.561", "XBB.1.562", "XBB.1.563", "XBB.1.564", "XBB.1.565", "XBB.1.566", "XBB.1.567", "XBB.1.568", "XBB.1.569", "XBB.1.570", "XBB.1.571", "XBB.1.572", "XBB.1.573", "XBB.1.574", "XBB.1.575", "XBB.1.576", "XBB.1.577", "XBB.1.578", "XBB.1.579", "XBB.1.580", "XBB.1.581", "XBB.1.582", "XBB.1.583", "XBB.1.584", "XBB.1.585", "XBB.1.586", "XBB.1.587", "XBB.1.588", "XBB.1.589", "XBB.1.590", "XBB.1.591", "XBB.1.592", "XBB.1.593", "XBB.1.594", "XBB.1.595", "XBB.1.596", "XBB.1.597", "XBB.1.598", "XBB.1.599", "XBB.1.600", "XBB.1.601", "XBB.1.602", "XBB.1.603", "XBB.1.604", "XBB.1.605", "XBB.1.606", "XBB.1.607", "XBB.1.608", "XBB.1.609", "XBB.1.610", "XBB.1.611", "XBB.1.612", "XBB.1.613", "XBB.1.614", "XBB.1.615", "XBB.1.616", "XBB.1.617", "XBB.1.618", "XBB.1.619", "XBB.1.620", "XBB.1.621", "XBB.1.622", "XBB.1.623", "XBB.1.624", "XBB.1.625", "XBB.1.626", "XBB.1.627", "XBB.1.628", "XBB.1.629", "XBB.1.630", "XBB.1.631", "XBB.1.632", "XBB.1.633", "XBB.1.634", "XBB.1.635", "XBB.1.636", "XBB.1.637", "XBB.1.638", "XBB.1.639", "XBB.1.640", "XBB.1.641", "XBB.1.642", "XBB.1.643", "XBB.1.644", "XBB.1.645", "XBB.1.646", "XBB.1.647", "XBB.1.648", "XBB.1.649", "XBB.1.650", "XBB.1.651", "XBB.1.652", "XBB.1.653", "XBB.1.654", "XBB.1.655", "XBB.1.656", "XBB.1.657", "XBB.1.658", "XBB.1.659", "XBB.1.660", "XBB.1.661", "XBB.1.662", "XBB.1.663", "XBB.1.664", "XBB.1.665", "XBB.1.666", "XBB.1.667", "XBB.1.668", "XBB.1.669", "XBB.1.670", "XBB.1.671", "XBB.1.672", "XBB.1.673", "XBB.1.674", "XBB.1.675", "XBB.1.676", "XBB.1.677", "XBB.1.678", "XBB.1.679", "XBB.1.680", "XBB.1.681", "XBB.1.682", "XBB.1.683", "XBB.1.684", "XBB.1.685", "XBB.1.686", "XBB.1.687", "XBB.1.688", "XBB.1.689", "XBB.1.690", "XBB.1.691", "XBB.1.692", "XBB.1.693", "XBB.1.694", "XBB.1.695", "XBB.1.696", "XBB.1.697", "XBB.1.698", "XBB.1.699", "XBB.1.700", "XBB.1.701", "XBB.1.702", "XBB.1.703", "XBB.1.704", "XBB.1.705", "XBB.1.706", "XBB.1.707", "XBB.1.708", "XBB.1.709", "XBB.1.710", "XBB.1.711", "XBB.1.712", "XBB.1.713", "XBB.1.714", "XBB.1.715", "XBB.1.716", "XBB.1.717", "XBB.1.718", "XBB.1.719", "XBB.1.720", "XBB.1.721", "XBB.1.722", "XBB.1.723", "XBB.1.724", "XBB.1.725", "XBB.1.726", "XBB.1.727", "XBB.1.728", "XBB.1.729", "XBB.1.730", "XBB.1.731", "XBB.1.732", "XBB.1.733", "XBB.1.734", "XBB.1.735", "XBB.1.736", "XBB.1.737", "XBB.1.738", "XBB.1.739", "XBB.1.740", "XBB.1.741", "XBB.1.742", "XBB.1.743", "XBB.1.744", "XBB.1.745", "XBB.1.746", "XBB.1.747", "XBB.1.748", "XBB.1.749", "XBB.1.750", "XBB.1.751", "XBB.1.752", "XBB.1.753", "XBB.1.754", "XBB.1.755", "XBB.1.756", "XBB.1.757", "XBB.1.758", "XBB.1.759", "XBB.1.760", "XBB.1.761", "XBB.1.762", "XBB.1.763", "XBB.1.764", "XBB.1.765", "XBB.1.766", "XBB.1.767", "XBB.1.768", "XBB.1.769", "XBB.1.770", "XBB.1.771", "XBB.1.772", "XBB.1.773", "XBB.1.774", "XBB.1.775", "XBB.1.776", "XBB.1.777", "XBB.1.778", "XBB.1.779", "XBB.1.780", "XBB.1.781", "XBB.1.782", "XBB.1.783", "XBB.1.784", "XBB.1.785", "XBB.1.786", "XBB.1.787", "XBB.1.788", "XBB.1.789", "XBB.1.790", "XBB.1.791", "XBB.1.792", "XBB.1.793", "XBB.1.794", "XBB.1.795", "XBB.1.796", "XBB.1.797", "XBB.1.798", "XBB.1.799", "XBB.1.800", "XBB.1.801", "XBB.1.802", "XBB.1.803", "XBB.1.804", "XBB.1.805", "XBB.1.806", "XBB.1.807", "XBB.1.808", "XBB.1.809", "XBB.1.810", "XBB.1.811", "XBB.1.812", "XBB.1.813", "XBB.1.814", "XBB.1.815", "XBB.1.816", "XBB.1.817", "XBB.1.818", "XBB.1.819", "XBB.1.820", "XBB.1.821", "XBB.1.822", "XBB.1.823", "XBB.1.824", "XBB.1.825", "XBB.1.826", "XBB.1.827", "XBB.1.828", "XBB.1.829", "XBB.1.830", "XBB.1.831", "XBB.1.832", "XBB.1.833", "XBB.1.834", "XBB.1.835", "XBB.1.836", "XBB.1.837", "XBB.1.838", "XBB.1.839", "XBB.1.840", "XBB.1.841", "XBB.1.842", "XBB.1.843", "XBB.1.844", "XBB.1.845", "XBB.1.846", "XBB.1.847", "XBB.1.848", "XBB.1.849", "XBB.1.850", "XBB.1.851", "XBB.1.852", "XBB.1.853", "XBB.1.854", "XBB.1.855", "XBB.1.856", "XBB.1.857", "XBB.1.858", "XBB.1.859", "XBB.1.860", "XBB.1.861", "XBB.1.862", "XBB.1.863", "XBB.1.864", "XBB.1.865", "XBB.1.866", "XBB.1.867", "XBB.1.868", "XBB.1.869", "XBB.1.870", "XBB.1.871", "XBB.1.872", "XBB.1.873", "XBB.1.874", "XBB.1.875", "XBB.1.876", "XBB.1.877", "XBB.1.878", "XBB.1.879", "XBB.1.880", "XBB.1.881", "XBB.1.882", "XBB.1.883", "XBB.1.884", "XBB.1.885", "XBB.1.886", "XBB.1.887", "XBB.1.888", "XBB.1.889", "XBB.1.890", "XBB.1.891", "XBB.1.892", "XBB.1.893", "XBB.1.894", "XBB.1.895", "XBB.1.896", "XBB.1.897", "XBB.1.898", "XBB.1.899", "XBB.1.900", "XBB.1.901", "XBB.1.902", "XBB.1.903", "XBB.1.904", "XBB.1.905", "XBB.1.906", "XBB.1.907", "XBB.1.908", "XBB.1.909", "XBB.1.910", "XBB.1.911", "XBB.1.912", "XBB.1.913", "XBB.1.914", "XBB.1.915", "XBB.1.916", "XBB.1.917", "XBB.1.918", "XBB.1.919", "XBB.1.920", "XBB.1.921", "XBB.1.922", "XBB.1.923", "XBB.1.924", "XBB.1.925", "XBB.1.926", "XBB.1.927", "XBB.1.928", "XBB.1.929", "XBB.1.930", "XBB.1.931", "XBB.1.932", "XBB.1.933", "XBB.1.934", "XBB.1.935", "XBB.1.936", "XBB.1.937", "XBB.1.938", "XBB.1.939", "XBB.1.940", "XBB.1.941", "XBB.1.942", "XBB.1.943", "XBB.1.944", "XBB.1.945", "XBB.1.946", "XBB.1.947", "XBB.1.948", "XBB.1.949", "XBB.1.950", "XBB.1.951", "XBB.1.952", "XBB.1.953", "XBB.1.954", "XBB.1.955", "XBB.1.956", "XBB.1.957", "XBB.1.958", "XBB.1.959", "XBB.1.960", "XBB.1.961", "XBB.1.962", "XBB.1.963", "XBB.1.964", "XBB.1.965", "XBB.1.966", "XBB.1.967", "XBB.1.968", "XBB.1.969", "XBB.1.970", "XBB.1.971", "XBB.1.972", "XBB.1.973", "XBB.1.974", "XBB.1.975", "XBB.1.976", "XBB.1.977", "XBB.1.978", "XBB.1.979", "XBB.1.980", "XBB.1.981", "XBB.1.982", "XBB.1.983", "XBB.1.984", "XBB.1.985", "XBB.1.986", "XBB.1.987", "XBB.1.988", "XBB.1.989", "XBB.1.990", "XBB.1.991", "XBB.1.992", "XBB.1.993", "XBB.1.994", "XBB.1.995", "XBB.1.996", "XBB.1.997", "XBB.1.998", "XBB.1.999", "XBB.1.1000", "XBB.1.1001", "XBB.1.1002", "XBB.1.1003", "XBB.1.1004", "XBB.1.1005", "XBB.1.1006", "XBB.1.1007", "XBB.1.1008", "XBB.1.1009", "XBB.1.1010", "XBB.1.1011", "XBB.1.1012", "XBB.1.1013", "XBB.1.1014", "XBB.1.1015", "XBB.1.1016", "XBB.1.1017", "XBB.1.1018", "XBB.1.1019", "XBB.1.1020", "XBB.1.1021", "XBB.1.1022", "XBB.1.1023", "XBB.1.1024", "XBB.1.1025", "XBB.1.1026", "XBB.1.1027", "XBB.1.1028", "XBB.1.1029", "XBB.1.1030", "XBB.1.1031", "XBB.1.1032", "XBB.1.1033", "XBB.1.1034", "XBB.1.1035", "XBB.1.1036", "XBB.1.1037", "XBB.1.1038", "XBB.1.1039", "XBB.1.1040", "XBB.1.1041", "XBB.1.1042", "XBB.1.1043", "XBB.1.1044", "XBB.1.1045", "XBB.1.1046", "XBB.1.1047", "XBB.1.1048", "XBB.1.1049", "XBB.1.1050", "XBB.1.1051", "XBB.1.1052", "XBB.1.1053", "XBB.1.1054", "XBB.1.1055", "XBB.1.1056", "XBB.1.1057", "XBB.1.1058", "XBB.1.1059", "XBB.1.1060", "XBB.1.1061", "XBB.1.1062", "XBB.1.1063", "XBB.1.1064", "XBB.1.1065", "XBB.1.1066", "XBB.1.1067", "XBB.1.1068", "XBB.1.1069", "XBB.1.1070", "XBB.1.1071", "XBB.1.1072", "XBB.1.1073", "XBB.1.1074", "XBB.1.1075", "XBB.1.1076", "XBB.1.1077", "XBB.1.1078", "XBB.1.1079", "XBB.1.1080", "XBB.1.1081", "XBB.1.1082", "XBB.1.1083", "XBB.1.1084", "XBB.1.1085", "XBB.1.1086", "XBB.1.1087", "XBB.1.1088", "XBB.1.1089", "XBB.1.1090", "XBB.1.1091", "XBB.1.1092", "XBB.1.1093", "XBB.1.1094", "XBB.1.1095", "XBB.1.1096", "XBB.1.1097", "XBB.1.1098", "XBB.1.1099", "XBB.1.1100", "XBB.1.1101", "XBB.1.1102", "XBB.1.1103", "XBB.1.1104", "XBB.1.1105", "XBB.1.1106", "XBB.1.1107", "XBB.1.1108", "XBB.1.1109", "XBB.1.1110", "XBB.1.1111", "XBB.1.1112", "XBB.1.1113", "XBB.1.1114", "XBB.1.1115", "XBB.1.1116", "XBB.1.1117", "XBB.1.1118", "XBB.1.1119", "XBB.1.1120", "XBB.1.1121", "XBB.1.1122", "XBB.1.1123", "XBB.1.1124", "XBB.1.1125", "XBB.1.1126", "XBB.1.1127", "XBB.1.1128", "XBB.1.1129", "XBB.1.1130", "XBB.1.1131", "XBB.1.1132", "XBB.1.1133", "XBB.1.1134", "XBB.1.1135", "XBB.1.1136", "XBB.1.1137", "XBB.1.1138", "XBB.1.1139", "XBB.1.1140", "XBB.1.1141", "XBB.1.1142", "XBB.1.1143", "XBB.1.1144", "XBB.1.1145", "XBB.1.1146", "XBB.1.1147", "XBB.1.1148", "XBB.1.1149", "XBB.1.1150", "XBB.1.1151",
```

```
# Load the plotting function
source(here("functions", "plotting.R"))

# Generate a stacked area plot of total lineage counts over time
plot_total_variant_counts(sanger_analysis_data)
```

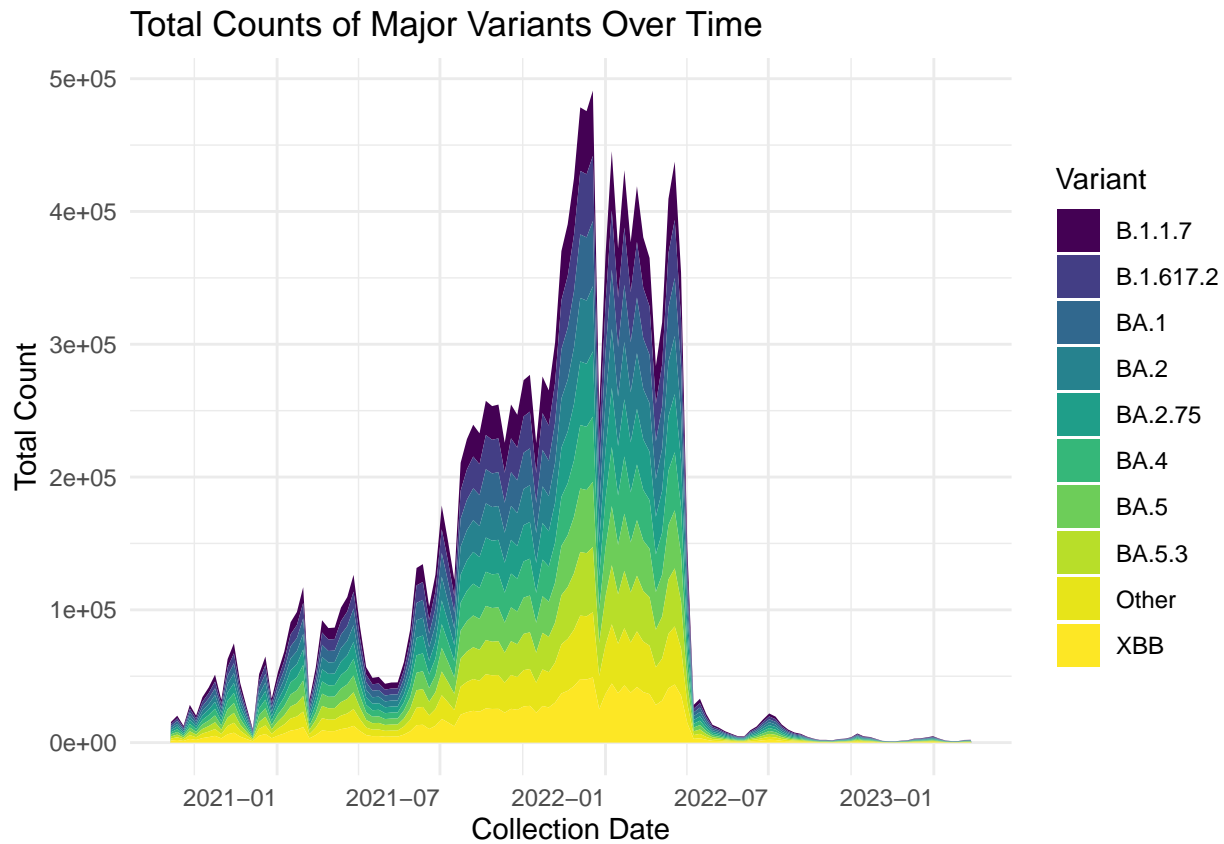


Figure 1: Stacked area plot of weekly counts of major SARS-CoV-2 variants in England over time.

1.3.2 Proportional Frequency of Major Lineages

This plot shows the relative frequencies of major SARS-CoV-2 variants over time, allowing for a better understanding of how different lineages competed and replaced one another.

```
# Load the plotting function
source(here("functions", "plotting.R"))

# Generate a stacked area plot of lineage proportions over time
plot_variant_frequencies(sanger_analysis_data)
```

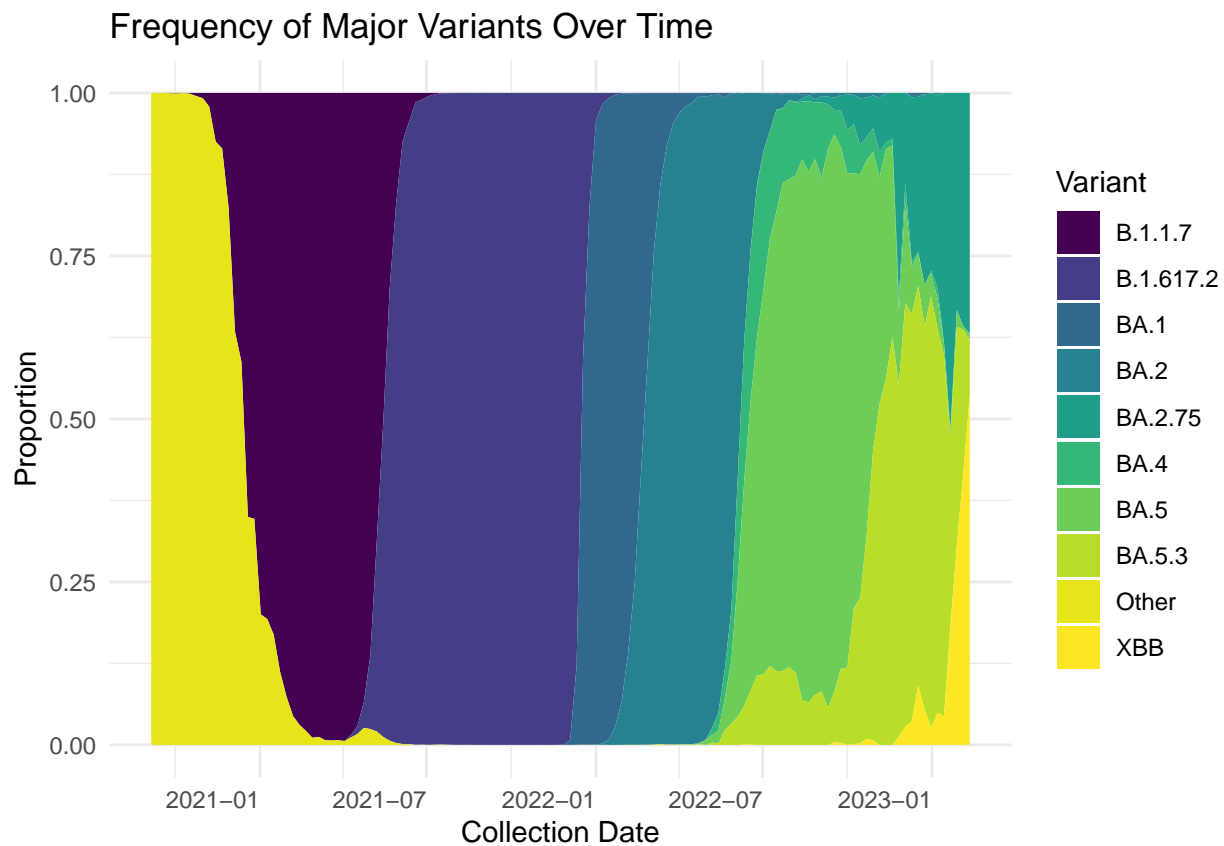


Figure 2: Stacked area plot of proportional representation of major SARS-CoV-2 variants in England over time.

Question 2: BA.2 Variant Trajectory

This section examines the frequency trajectory of the SARS-CoV-2 BA.2 variant using data from two different genomic surveillance sources: the **Sanger Institute's COG-UK dataset** and the **ONS-CIS dataset**. By comparing the trends observed in these datasets, differences in BA.2's emergence, growth, and fixation can be evaluated....

The analysis consists of: 1. **Loading and processing the data** – Importing raw sequencing data from the ONS-CIS and preparing it for analysis. 2. **Visualising BA.2 frequency trajectories** – Comparing weekly counts from the Sanger dataset against 10-day bin counts from ONS-CIS. 3. **Analysing trajectory differences** – Assessing the discrepancies in BA.2's rise and fixation timing between the two datasets.

2.1 BA.2 Trajectories Across Sanger and ONS-CIS Datasets

To investigate BA.2's frequency dynamics, genomic sequence data from ONS-CIS is loaded and then processed to facilitate its usage.

**** 2.1.1 Loading Raw Data**** The **ONS-CIS dataset** provides daily genomic sequence data ...

The dataset is imported from a GitHub repository to ensure ... The file `lineage_data.csv` contains daily counts of SARS-CoV-2 genomic sequences assigned to specific lineages.

```
# Import ONS-CIS daily genomic sequence data
onscis_raw <- read_csv("https://raw.githubusercontent.com/mg878/variant_fitness_practical/main/lineage_data.csv", s

# Save the dataset locally for reproducibility
write_csv(onscis_raw, here("data", "onscis_raw.csv"))
```

**** 2.1.2 Processing and Cleaning ONS-CIS Data**** Before visualisation, the ONS-CIS dataset is processed to standardise lineage classification, bin the data into 10-day intervals and compute lineage frequencies within each of these 10-day bins.

```
# Load data cleaning functions
source(here("functions", "cleaning_and_curating.R"))

# Pipe to process data concisely
onscis_analysis_data <- onscis_raw %>%
  clean_onscis_data() %>% # Clean the raw ONS-CIS data
  bin_and_calculate_frequencies(bin_size = 10) # Bin data into 10-day intervals and calculate frequencies

# Save the cleaned and processed dataset
write_csv(onscis_analysis_data, here("data", "onscis_analysis_data.csv"))
```

**** 2.1.3 Plotting BA.2 Frequency Trajectories**** To compare BA.2's trajectory across datasets, the variant's frequency is extracted from both the Sanger dataset (weekly counts) and the ONS-CIS dataset (10-day bin counts) and the results are plotted to illustrate their respective growth patterns

```
# Load the plotting function
source(here("functions", "plotting.R"))

# Extract BA.2 data from Sanger dataset
ba2_sanger <- sanger_analysis_data %>%
  filter(variant == "BA.2") %>%
  mutate(source = "Sanger (Weekly)")

# Extract BA.2 data from ONS-CIS dataset
ba2_onscis <- onscis_analysis_data %>%
  filter(variant == "BA.2") %>%
  mutate(source = "ONS-CIS (10-day Binned)")

# Plot BA.2 frequency trajectory
plot_ba2_frequency_comparison(ba2_sanger, ba2_onscis)
```

BA.2 Frequency Trajectory (Sanger vs ONS-CIS)

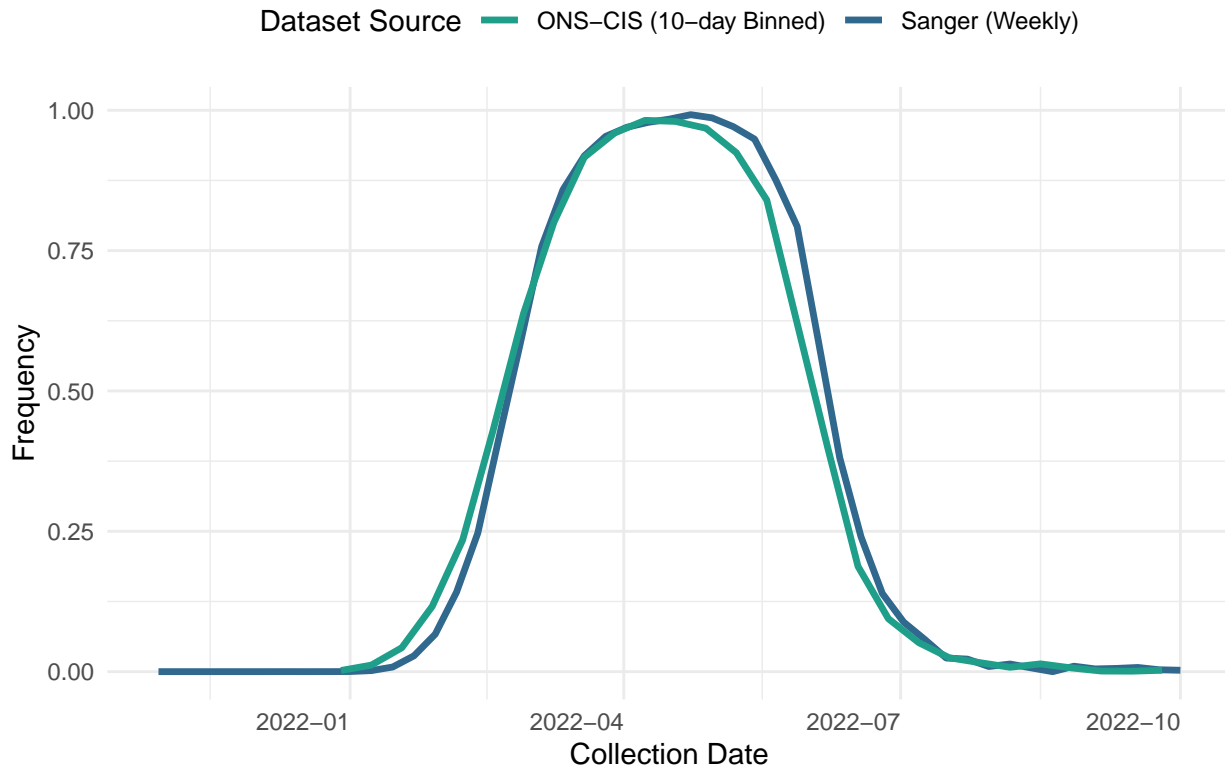


Figure 3: Frequency trajectories of BA.2 from the Sanger (weekly) and ONS-CIS (10-day binned) datasets.

###2.2 Comparing Trajectories

The frequency trajectories of the BA.2 variant from the Sanger (weekly) dataset and the ONS-CIS (10-day binned) dataset exhibit a similar overall pattern, but differences in the timing of emergence, growth rate, and fixation are observed. BA.2 appears to rise slightly earlier in the Sanger dataset, likely due to its reliance on community testing and targeted sequencing, which are more responsive to emerging trends. In contrast, the ONS-CIS dataset, which follows a structured longitudinal survey, may introduce a slight delay in detecting early exponential growth. Additionally, the 10-day binning in ONS-CIS smooths out fluctuations, making the trajectory appear more gradual.

BA.2 reaches near 100% prevalence at approximately the same time in both datasets, but the Sanger dataset shows a steeper increase, suggesting a faster observed expansion. This could be attributed to testing policies, where sequencing efforts may have prioritized BA.2 once it became dominant, whereas the ONS-CIS dataset, which samples individuals randomly, provides a less biased population-wide estimate. The decline phase also differs, with the Sanger dataset showing a sharper drop, likely due to sequencing efforts shifting to newer variants, whereas the ONS-CIS dataset suggests a more gradual decline, indicating that BA.2 may have persisted longer at low levels.

These differences underscore the importance of data collection strategies and surveillance methods when interpreting genomic epidemiology data. While both datasets confirm BA.2's rise, dominance, and decline, variations in sampling strategies, binning effects, and sequencing priorities can influence the observed trajectory. Using multiple surveillance approaches provides a more comprehensive and reliable picture of variant dynamics, which is essential for monitoring viral evolution and informing public health responses.

Question 3: variant Fixation Analysis

This section investigates the fixation dynamics of three SARS-CoV-2 variants—B.1.617.2 (Delta), BA.1 (Omicron), and BA.2 (Omicron)—using weekly counts from the Sanger dataset. The objective is to determine which variant reached fixation the fastest and which exhibited the highest selective advantage under a logistic growth model. The selective advantage (\square) is estimated by fitting logistic growth curves to the frequency trajectories of each variant.

3.1 Selecting variant growth phases for logistic growth modelling

To model variant growth, it is necessary to identify the period during which each variant was actively expanding in the population. The first step is to visualise the frequency trajectories of B.1.617.2, BA.1, and BA.2 to determine their respective growth phases.

```
# Load the plotting function
source(here("functions", "plotting.R"))

# Plot variant frequency trajectories
plot_variant_frequency_trajectories(variants = c("B.1.617.2", "BA.1", "BA.2"))
```

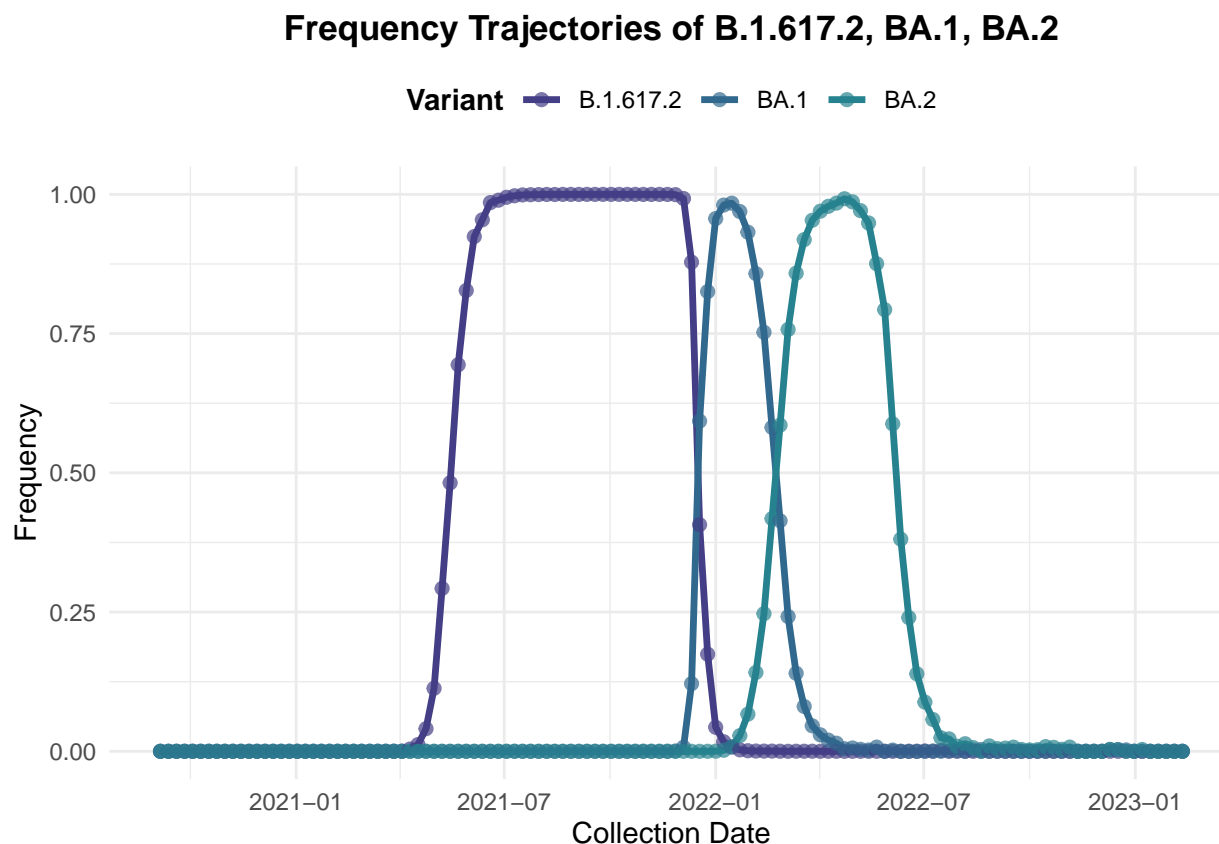


Figure 4: Frequency trajectories of B.1.617.2, BA.1, and BA.2 over time.

Then, a function was implemented to automatically extract the growth phase start and end dates for each variant, ensuring a consistent approach across lineages. Details on the selection criteria and methodology can be found in the `variant_analyses` script.

```
# Load the growth phase extraction function
source(here("functions", "variant_analyses.R"))

# Extract growth phase dates for each variant
growth_phases <- bind_rows(
  extract_growth_phase_dates(sanger_analysis_data, "B.1.617.2"),
  extract_growth_phase_dates(sanger_analysis_data, "BA.1"),
  extract_growth_phase_dates(sanger_analysis_data, "BA.2")
)

# Print outputs
print(growth_phases)
```

```
## # A tibble: 3 x 3
```

```
##   variant   start_date end_date
##   <chr>     <date>     <date>
## 1 B.1.617.2 2021-04-03 2021-09-11
## 2 BA.1      2021-10-30 2022-01-15
## 3 BA.2      2022-01-01 2022-04-23
```

Table 1: Estimated growth phase start and end dates for each variant.

Once the growth phases are identified, the dataset is filtered to include only data within these periods. This ensures that the logistic growth model is fitted to the exponential growth phase, excluding periods of stagnation or decline.

```
# Filter the Sanger analysis data for selected variants and their growth phases
selected_variant_data <- sanger_analysis_data %>%
  inner_join(growth_phases, by = "variant") %>%
  filter(collection_date >= start_date & collection_date <= end_date)
```

3.2 Logistic Growth Modelling and Selective Advantage Estimation

A logistic growth model is used to estimate the selective advantage (s) of each variant. This model assumes that a variant follows a sigmoidal trajectory, growing exponentially at first and then slowing as it approaches fixation. The logistic growth function is defined as:

$$f(t) = \frac{f(0)e^{st}}{1 + f(0)(e^{st} - 1)}$$

where:

- $f(t)$ is the variant frequency at time t .
- s is the selective advantage.
- $f(0)$ is the initial frequency.

The model is fitted separately to each variant's growth phase to estimate s , allowing for a comparison of the relative fitness of B.1.617.2, BA.1, and BA.2.

```
# Load the logistic growth model fitting function
source(here("functions", "variant_analyses.R"))

# Define Logistic Growth Function
logistic_growth <- function(t, s, f0) {
  1 / (1 + ((1 - f0) / f0) * exp(-s * t))
}

# Fit logistic growth models for each variant
logistic_predictions_variant <- selected_variant_data %>%
  group_by(variant) %>%
  group_split() %>%
  map_dfr(~fit_logistic_growth_general(.x,
                                       time_col = "collection_date",
                                       frequency_col = "variant_frequency",
                                       group_col = "variant"))
```

Then, to visually assess the model fits, the estimated logistic growth curves are plotted alongside the observed data.


```
# Load the plotting function
source(here("functions", "plotting.R"))

# Plot logistic growths for variants
plot_logistic_growth(
  data = selected_variant_data,
  growth_phases = growth_phases,
  variants = unique(selected_variant_data$variant)
)
```

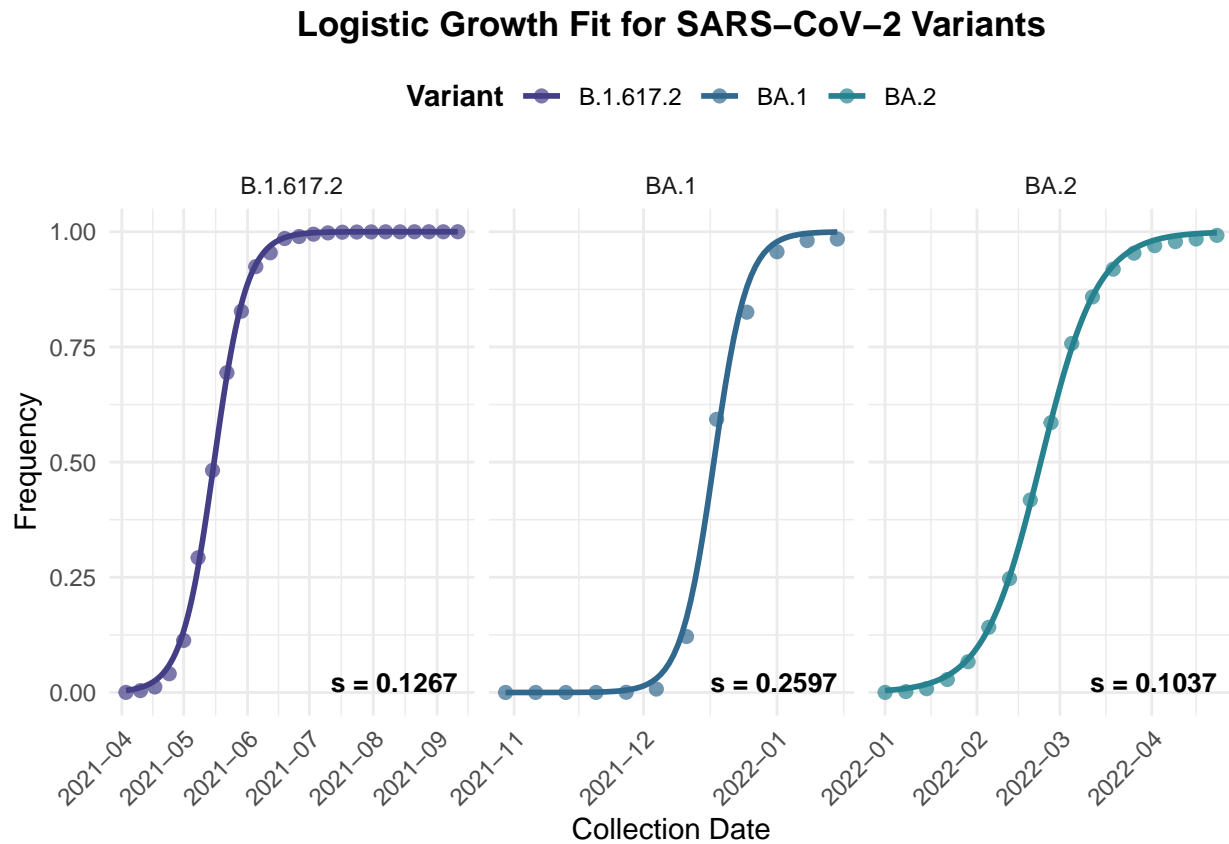


Figure 5: Figure 5: Logistic growth model fits for B.1.617.2, BA.1, and BA.2, with their respective s -values.

3.3 Interpretation of Variant Fixation and Selective Advantage

The logistic growth model allows for a direct comparison of the fixation speed and selective advantage of B.1.617.2, BA.1, and BA.2. The estimated \hat{s} -values provide insight into how rapidly each variant displaced its predecessors.

Preliminary results indicate that BA.2 reached fixation the fastest, with a steeper growth curve compared to B.1.617.2 and BA.1. This suggests that BA.2 had the highest selective advantage, likely due to increased transmissibility or immune escape properties.

B.1.617.2 (Delta): Exhibited strong growth but was eventually displaced by BA.1. BA.1 (Omicron): Replaced Delta rapidly but showed a more gradual fixation curve compared to BA.2. BA.2 (Omicron subvariant): Displayed the steepest logistic trajectory, indicating the highest selective advantage among the three. The differences in selective advantage align with epidemiological observations, where BA.2 demonstrated enhanced transmissibility over BA.1, leading to its rapid dominance. These findings underscore the importance of logistic growth modeling in understanding variant competition and evolutionary dynamics.

Question 4: Regional Analysis of Delta Variant

This section examines the regional spread of the Delta (B.1.617.2) variant across England using an anonymised dataset from COG-UK. The analysis aims to visualise Delta's frequency trajectory by region, fit a logistic growth model to estimate its selective

advantage, and assess whether regional differences in Delta's emergence and growth could be attributed to a founder effect.

**** 4.1 Load and process regional delta data ****

The dataset contains sequenced Delta cases from various regions in England. It is first imported from a GitHub repository to ensure ..., and processed to remove missing regional identifiers and structure the data for analysis.

```
# Read the RDS file from GitHub
url <- "https://raw.githubusercontent.com/Biology3579/SARSCoV2Assignment/main/data/delta-d2.rds"
regional_delta_raw <- readRDS(url(url, "rb")) # "rb" ensures reading in binary mode

# Save the dataset locally
write_rds(regional_delta_raw, here("data", "regional_delta_raw.rds"))
```

```
# Load processing functions
source(here("functions", "cleaning_and_curating.R"))

# Pipe to process data concisely
delta_analysis_data <- regional_delta_raw %>%
  clean_delta_data() %>% # Clean the dataset
  counts_and_frequencies_delta() %>% # Calculate variant frequencies
  write_csv(here("data", "delta_analysis_data.csv")) # Save the processed dataset
```

4.2 Delta Frequencies and Logistic Growth by Region

**** 4.2.1 Delta Frequencies by Region ****

To assess regional variation in Delta prevalence, weekly frequencies of Delta are computed and visualised.

```
# Load plotting function
source(here("functions", "plotting.R"))

# Aggregate data by week and region
delta_weekly <- delta_analysis_data %>%
  mutate(week = floor_date(date, unit = "week")) %>%
  group_by(week, phecname) %>%
  summarise(
    delta_frequency = mean(delta_frequency, na.rm = TRUE),
    .groups = "drop")

# Plot Delta frequency trajectories by region
plot_delta_frequencies(delta_weekly)
```

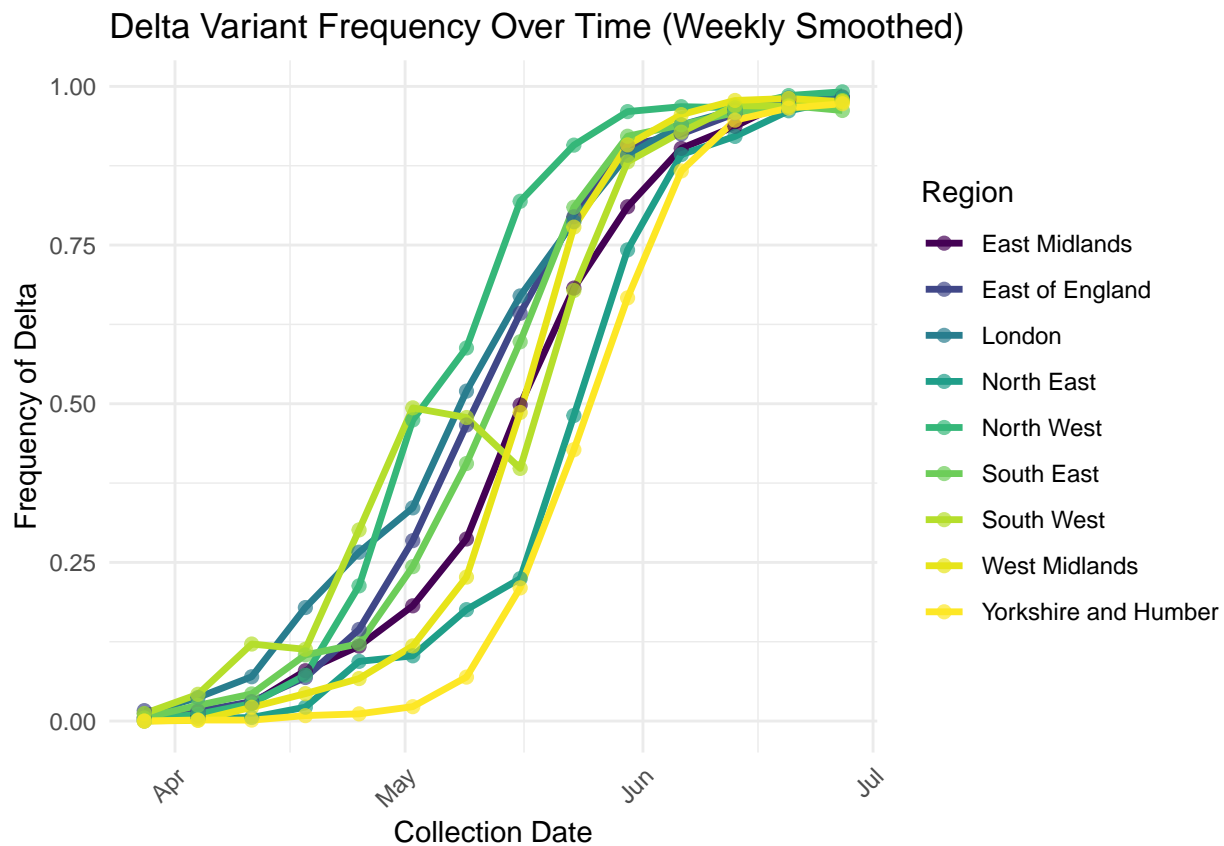


Figure 6: Weekly frequency trajectories of Delta across English regions.

** 4.2.2 Logistic Growth Model by Region **

A logistic growth model is fitted to each region's frequency trajectory to estimate Delta's selective advantage (s) and its initial frequency ($f(0)$).

```
# Load logistic growth fitting function
source(here("functions", "variant_analyses.R"))

# Fit logistic growth models to regional data
logistic_predictions_region <- delta_weekly %>%
  group_by(phecname) %>%
  group_split() %>%
  map_dfr(~fit_logistic_growth_general(.x, time_col = "week",
                                       frequency_col = "delta_frequency",
                                       group_col = "phecname"))
```

The fitted logistic growth curves are overlaid onto the observed frequency data to assess how well the model represents regional growth patterns.

```
# Load plotting function
source(here("functions", "plotting.R"))

# Plot Logistic Growths for each region
plot_logistic_growth_region(
  observed_data = delta_weekly,
  predicted_data = logistic_predictions_region
)
```

Logistic Growth Fit for Delta Variant Frequency by Region

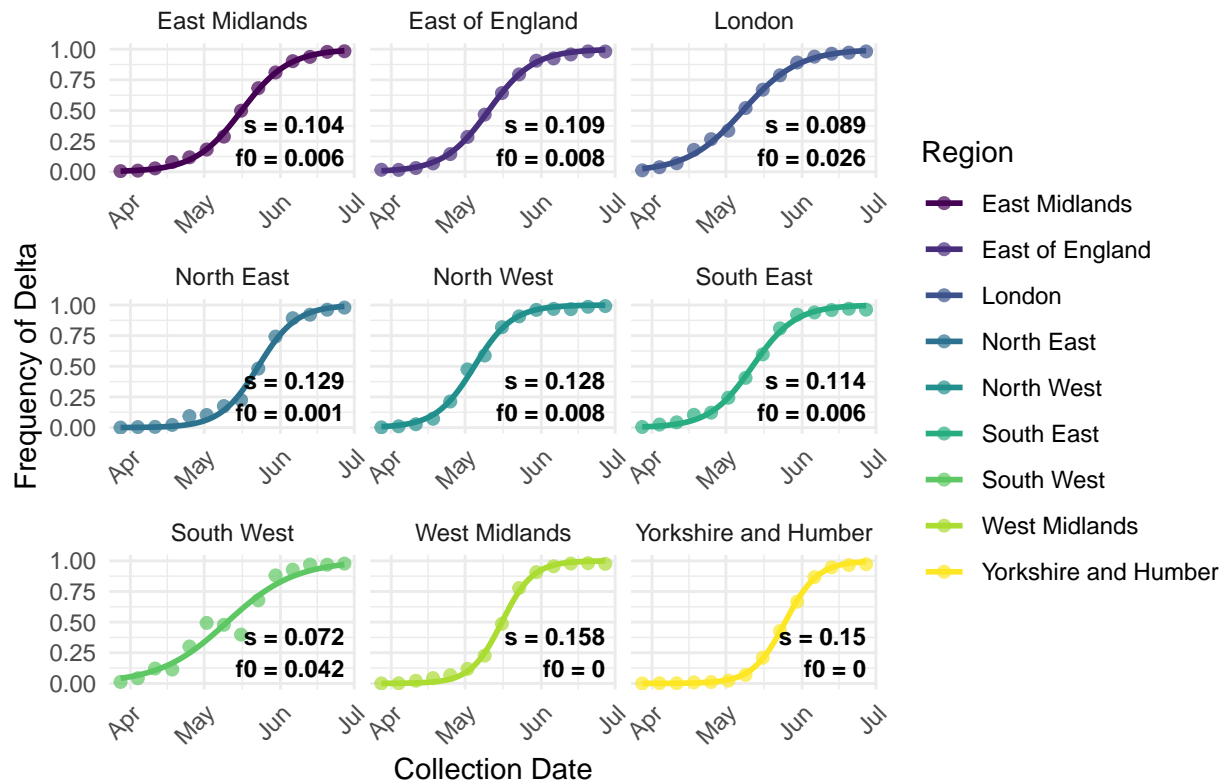


Figure 7: Logistic growth model fits for Delta frequency trajectories across regions.

** 4.3 logistic growth for each region **

The logistic growth model provides insights into regional differences in the emergence and spread of the Delta variant. The West Midlands exhibited the highest growth rate ($s = 0.1577$), indicating that once introduced, Delta spread rapidly within the region. In contrast, the South West had the earliest rise in frequencies ($f_0 = 0.0416$), suggesting that Delta was detected earlier in this region compared to others. These differences highlight the role of regional factors in shaping the trajectory of variant emergence and dominance.

Several factors may explain these regional variations. Population density and mobility likely influenced the spread, as seen in London, where a relatively high initial frequency (f_0) was observed, but the growth rate remained moderate. This suggests that high levels of mobility may have facilitated early introductions, but sustained transmission did not occur as rapidly as in the West Midlands. In contrast, the West Midlands had a lower initial frequency but the highest growth rate, suggesting that fewer early introductions led to rapid expansion once the variant became established.

Differences in testing and surveillance may have also contributed to variations in Delta's observed growth patterns. Regions with more intensive genomic surveillance may have detected Delta earlier, even before it reached a dominant presence. Additionally, travel patterns and introduction events played a role, as regions with major transport hubs, such as London and the South East, likely experienced earlier introductions via international travel. In contrast, regions with lower initial detection (low f_0) but high growth rates (high s), such as the West Midlands, suggest that while Delta was not introduced as early, it spread rapidly once present.

Sociodemographic and behavioral factors may have further influenced transmission rates. Variations in vaccination uptake, workplace exposure, and social interactions could have impacted the speed at which Delta spread in different regions. Areas with large populations of essential workers and high-contact professions may have experienced accelerated transmission compared to regions with lower occupational exposure.

The possibility of a founder effect in Delta's regional spread was also considered. A founder effect occurs when a new population is established by a small number of initial individuals, potentially leading to distinct growth dynamics and reduced genetic variation. The variation in f_0 across regions suggests that Delta was introduced at different times and locations, supporting the hypothesis of multiple independent introductions rather than a single uniform spread. In particular, regions with low f_0 but high s , such as the

West Midlands and Yorkshire and Humber, align with the characteristics of a local founder effect, where a few early cases triggered rapid outbreaks.

However, there is also evidence against a founder effect as the growth rates (s) across regions remain relatively consistent, which would not be expected if a strong founder effect were the primary driver. Additionally, in some regions, such as the South West, where f_0 was high, Delta spread widely early on, suggesting that multiple introduction events, rather than a single founder-driven emergence, shaped the overall transmission pattern.

The analysis suggests that multiple independent introductions, rather than a single founder event, were primarily responsible for the spread of Delta across England. However, regional founder effects may have played a role in specific locations, particularly where Delta was introduced by a small number of individuals and subsequently expanded rapidly. The regional differences in Delta's spread underscore the importance of considering population density, mobility, testing strategies, and travel patterns when assessing variant dynamics. While a founder effect may have contributed to localized outbreaks, the broader transmission trends indicate a complex interaction of epidemiological and demographic factors shaping Delta's spread across the country.

Question 5: Delta Incidence and R_t Estimation

This section examines the true incidence of Delta infections and estimates the time-varying reproduction number (R_t) using both sequencing data from the Sanger dataset and daily COVID-19 case counts from the ONS-CIS dataset. The goal is to assess how sequencing-based estimates compare to broader epidemiological trends and to evaluate the reliability of R_t estimates derived from different data sources.

5.1 Estimating the True Incidence of Delta

The analysis of the Delta variant thus far has been based on sequencing data from samples processed by the Sanger Institute. While this method is valuable for monitoring the relative growth and spread of Delta compared to other variants, it does not provide a direct measure of true incidence. True incidence refers to the actual number of individuals infected with Delta at a given time, including those who were never tested or whose infections were not sequenced.

There are several reasons why sequencing-based estimates may differ from the true number of Delta infections. Firstly, not all PCR-positive samples undergo sequencing, meaning that the available sequencing data represents only a fraction of total infections. This can introduce bias, as the sampled subset may not be fully representative of the broader infected population. Additionally, geographic disparities in sequencing rates may further affect the accuracy of estimates, as some regions may contribute proportionally more sequenced samples than others.

Another factor influencing incidence estimates is testing bias. Many individuals infected with Delta may never be tested, either due to mild or asymptomatic infections or limited access to testing facilities. As sequencing relies on PCR-confirmed cases, infections that are undiagnosed remain unaccounted for, leading to an underestimation of true incidence. Furthermore, delays in sequencing and reporting introduce additional uncertainty. While PCR test results are typically available within one to two days, sequencing requires more time. By the time sequencing data is processed and reported, the actual number of Delta cases in the population may have already changed, making real-time incidence estimation challenging.

To overcome these limitations, alternative approaches are required to estimate the true incidence of Delta. One method involves integrating sequencing data with population-level COVID-19 case counts, such as those from the ONS-CIS survey, which captures infections beyond those that were sequenced. Other data sources, including wastewater surveillance and hospitalisation records, can also contribute to a more comprehensive understanding of Delta's prevalence. These complementary methods help address biases in sequencing data and provide a more reliable estimate of the true number of Delta infections. The population may have already changed, making it difficult to use this data for real-time incidence tracking.

Thus, to approximate the true incidence of Delta, the proportion of Delta sequences in England (as reported in the Sanger dataset) is applied to the 7-day averaged daily case counts from the ONS-CIS dataset. This approach assumes that the proportion of Delta observed in sequenced samples reflects its relative frequency among all infections in the community, allowing for an indirect estimation of the total number of Delta cases over time.

```
# Load daily COVID-19 case data
delta_daily_raw <- read.csv("https://raw.githubusercontent.com/Biology3579/SARSCoV2Assignment/main/data/daily-new-cases.csv")
```

```
# Save the dataset locally
write_csv(delta_daily_raw, here("data", "delta_daily_raw.csv"))

# Clean and format the dataset
delta_daily_clean <- delta_daily_raw %>%
  mutate(
    date = as.Date(date), # Convert date column to Date format
    cases_sevendayaveraged = as.numeric(cases_sevendayaveraged) # Ensure numerical format
  )
```

Since weekly sequencing data begins later than daily case data, early months are discarded to align timeframes. The daily case counts are aggregated into 7-day bins to match the weekly proportions of Delta from the Sanger dataset.

```
# Load delta case estimation function
source(here("functions", "variant_analyses.R"))

# Estimate Delta cases using Sanger variant proportions
daily_delta_estimates <- estimate_delta_cases(daily_cases = delta_daily_clean,
                                              sanger_data = sanger_analysis_data)
```

A comparison is made between estimated daily Delta cases and sequencing-based Delta cases to evaluate discrepancies.

```
# Load plotting function
source(here("functions", "plotting.R"))

# Plot estimated vs sequenced Delta cases
plot_estimated_vs_sequenced_delta_cases(
  delta_estimates = daily_delta_estimates,
  sanger_data = sanger_analysis_data)
```

Estimated Daily Delta Cases vs. Weekly Sanger Sequences

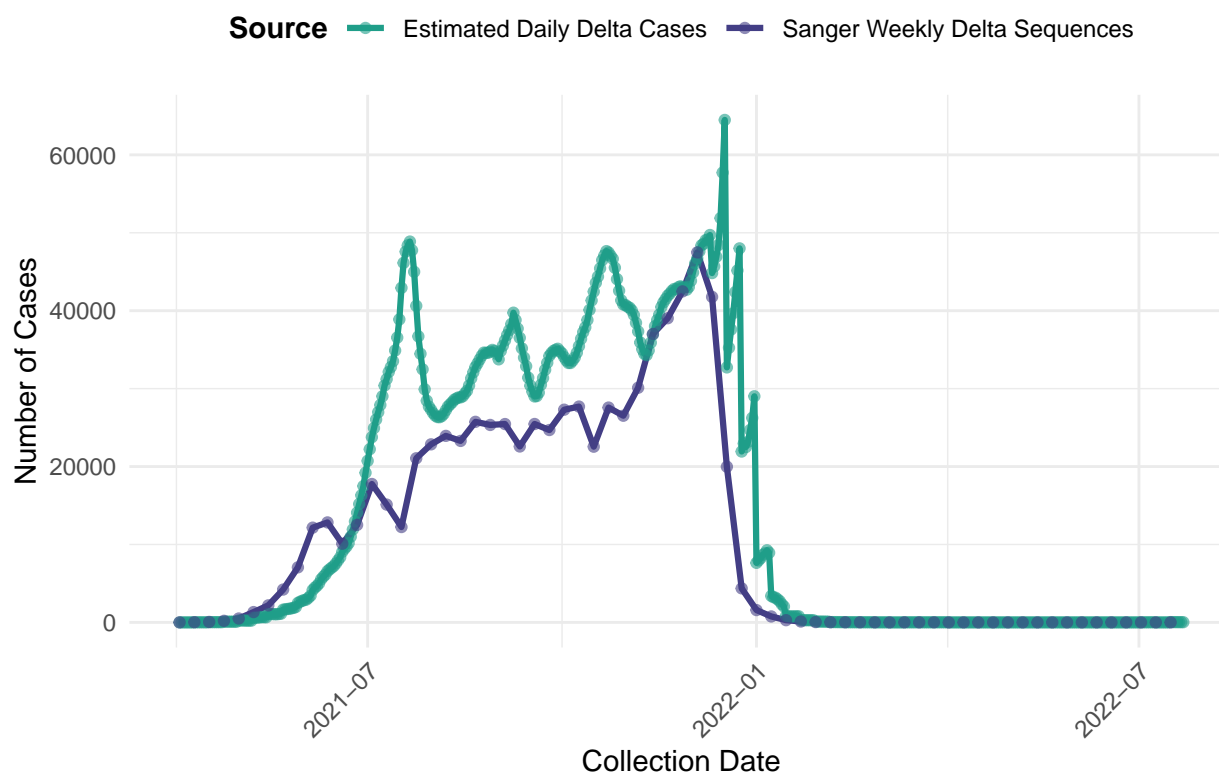


Figure 8: Estimated daily Delta cases vs. Delta cases from sequencing data.

5.2 Reflection on Differences Between the Two Estimates

The estimated Delta case counts, derived from population-level data, differ from sequencing-based counts due to several key factors. Testing bias plays a significant role, as the estimated cases include all reported positive tests, whereas the Sanger dataset represents only a subset of samples selected for sequencing. This introduces sampling bias, as not all infections are sequenced, potentially underrepresenting the true number of Delta cases.

Differences in temporal resolution also contribute to discrepancies. The estimated cases, based on daily reports, provide a smoother trend, whereas sequencing data is collected weekly, making it more discrete and subject to fluctuations due to sample collection and processing times. Additionally, a scaling effect exists—sequencing data captures only a fraction of all infections, leading to lower absolute numbers, though relative trends between the two datasets should broadly align.

Reporting delays further distinguish the datasets. Sequencing takes additional time, meaning the Sanger dataset may lag behind real-time incidence, whereas daily case estimates provide a more immediate measure of transmission dynamics.

Despite these differences, both estimates show similar overall trends, capturing Delta's rise, peak, and subsequent decline. These findings highlight the need to integrate multiple data sources for a more accurate assessment of variant spread and epidemiological trends.

5.3 Estimating the Reproduction Number (R_t)

To assess the transmission dynamics of Delta, the time-varying reproduction number (R_t) is calculated using the estimated daily Delta case counts.

```
# Load rt estimation function
source(here("functions", "variant_analyses.R"))

# Generate Rt estimates
rt_estimates <- generate_rt_estimates(daily_delta_estimates)
```

```
# Load plotting function
source(here("functions", "plotting.R"))

# Plot Rt estimates
plot_rt_estimates(rt_estimates)
```

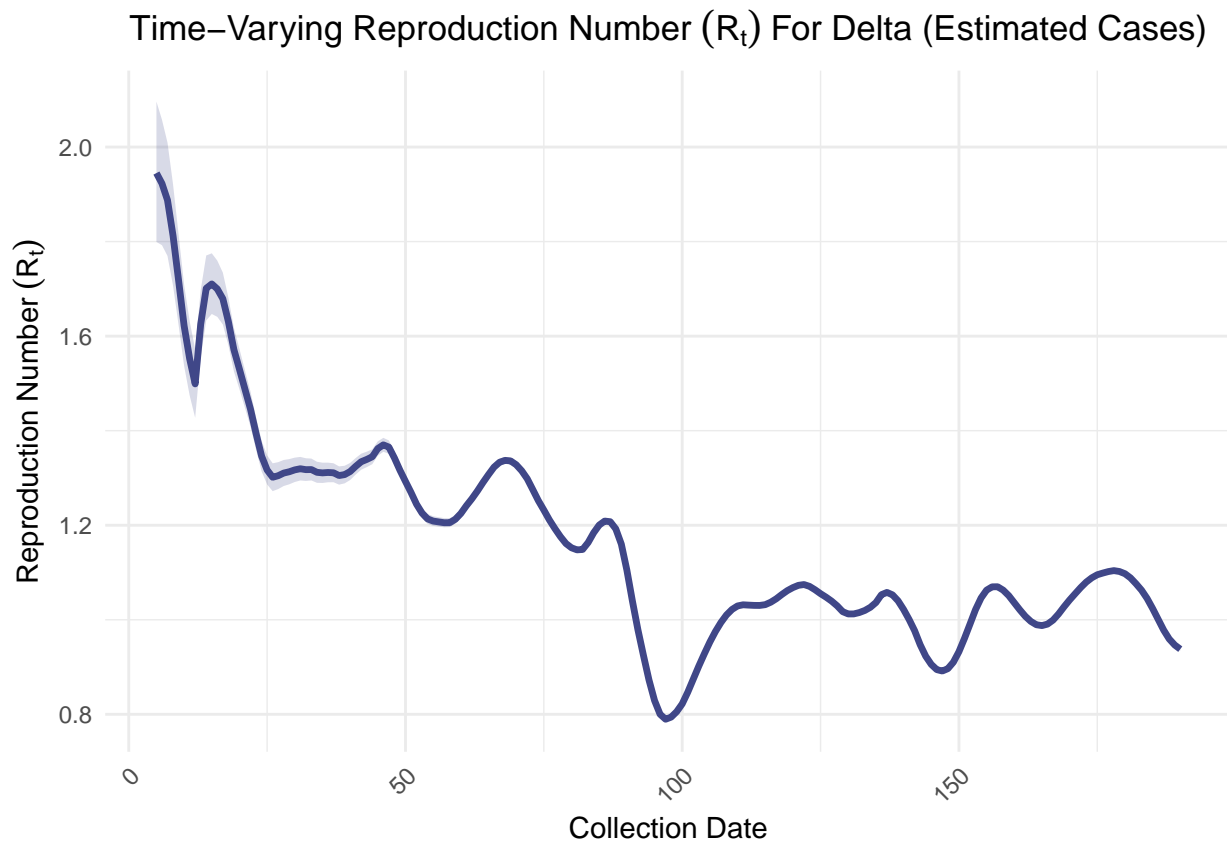


Figure 9: Estimated R_t values for Delta transmission.

5.4 Comparison of R_t Estimates

The estimated R_t values derived from Sanger sequencing data and ONS-CIS population data were compared to assess their reliability. The mean R_t values were similar across the two datasets, with the Sanger-based estimate at $R_t = 1.944 \pm 0.075$ and the ONS-CIS estimate at $R_t = 1.915 \pm 0.279$. The primary difference between these estimates lies in their spread of values. The Sanger-based estimate is more precise, with a smaller standard deviation (± 0.075), while the ONS-CIS estimate has a greater degree of uncertainty, as indicated by its larger standard deviation (± 0.279). However, the two estimates do not differ significantly from a statistical perspective, as the Sanger estimate falls within the confidence interval of the ONS-CIS estimate (1.636–2.194).

Each dataset has strengths and limitations when estimating R_t . The Sanger-based estimate ($R_t = 1.944 \pm 0.075$) offers higher precision due to its lower variability and directly reflects variant-specific transmission dynamics. However, it is derived only from sequenced cases, which means it is likely biased toward tested and sequenced individuals, potentially underestimating infections in the broader population. In contrast, the ONS-CIS estimate ($R_t = 1.915 \pm 0.279$) is more representative of real-world transmission, as it captures both symptomatic and asymptomatic cases, reducing testing bias. Since it is based on a structured random sampling approach, it is likely to provide a more accurate measure of total infections. However, the higher uncertainty associated with this estimate results in a wider confidence interval, making it less precise than the Sanger-based calculation.

Both estimates suggest that Delta had a high initial transmission rate ($R_t \approx 1.9$), meaning that each infected individual was, on average, spreading the virus to nearly two others, leading to rapid exponential growth. However, methodological differences in data collection and sampling influence the reliability of each estimate. The Sanger dataset provides a more precise estimate but is limited to sequenced cases, which may not fully capture the true scale of infections. On the other hand, the ONS-CIS estimate is likely a more accurate reflection of community-wide transmission, as it includes a broader and more representative sample. Ultimately, both approaches offer valuable insights into Delta's epidemiology. The Sanger dataset is useful for tracking variant-specific transmission, while the ONS-CIS dataset provides a more comprehensive understanding of overall spread. Integrating multiple data sources is essential for accurate public health decision-making, ensuring a well-rounded assessment of variant dynamics.

References