# COVID-19 Case Prediction using Hidden Markov Models
# Khairul & Evan

## 1. Introduction

This report presents a Hidden Markov Model (HMM) implementation to predict COVID-19 confirmed cases per million people based on the month of the year. The model helps understand patterns in COVID-19 transmission and provides a framework for predicting future infection levels.

## 2. Dataset Description

We used the provided "covid_data.csv" dataset, which contains the following information:

- Date: Daily records from January 4, 2020 to February 16, 2025
- New cases per million: Daily new COVID-19 cases per million people
- Total cases per million: Cumulative case count per million people
- 7-day average of new cases: Rolling average to smooth daily fluctuations

The data spans approximately 5 years, capturing multiple COVID-19 waves and seasonal patterns.

## 3. Model Architecture

### 3.1 Hidden States

Our model uses 10 hidden states representing different levels of COVID-19 infection intensity:

- level_0_200: 0-200 cases per million
- level_200_400: 200-400 cases per million
- level_400_600: 400-600 cases per million
- level_600_800: 600-800 cases per million

- level_800_1000: 800-1000 cases per million
- level_1000_1200: 1000-1200 cases per million
- level_1200_1400: 1200-1400 cases per million
- level_1400_1600: 1400-1600 cases per million
- level_1600_1800: 1600-1800 cases per million
- level_1800_plus: Over 1800 cases per million

## 3.2 Observable Evidence

The observable evidence in our model is the month of the year (1-12 representing January-December). The model learns the relationship between months and infection levels.

# 4. Implementation Details

The HMM was implemented from scratch in Python using the following libraries:

- pandas: For data manipulation
- numpy: For numerical computations
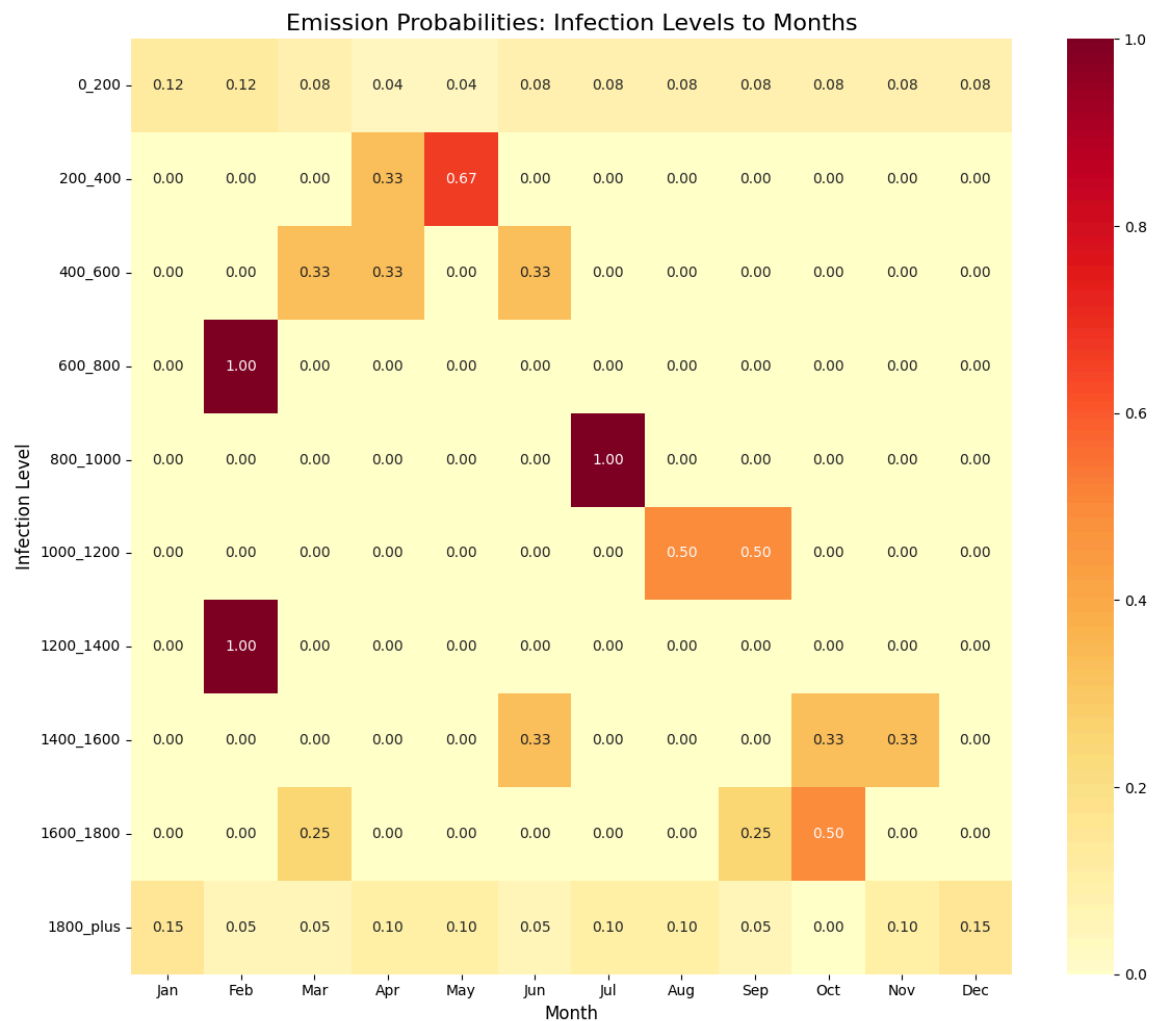- matplotlib and seaborn: For visualization

Key components of the implementation include:

1. **Data Processing**:
   - Loading the CSV file and converting dates to datetime format
   - Aggregating daily data into monthly sums
   - Assigning appropriate hidden states based on case levels
2. **HMM Parameter Learning**:
   - Initial state probabilities: Estimated from the first observations
   - Transition probabilities: Learned from sequential state changes
   - Emission probabilities: Learned from the correlation between states and months
3. **Key Algorithms**:
   - Filtering: Estimating current infection level given a month
   - Prediction: Forecasting infection levels for future months
   - Viterbi algorithm: Finding the most likely sequence of infection levels

# 5. Model Visualization

## 5.1 Emission Probabilities

The emission probability matrix shows the relationship between infection levels and months:



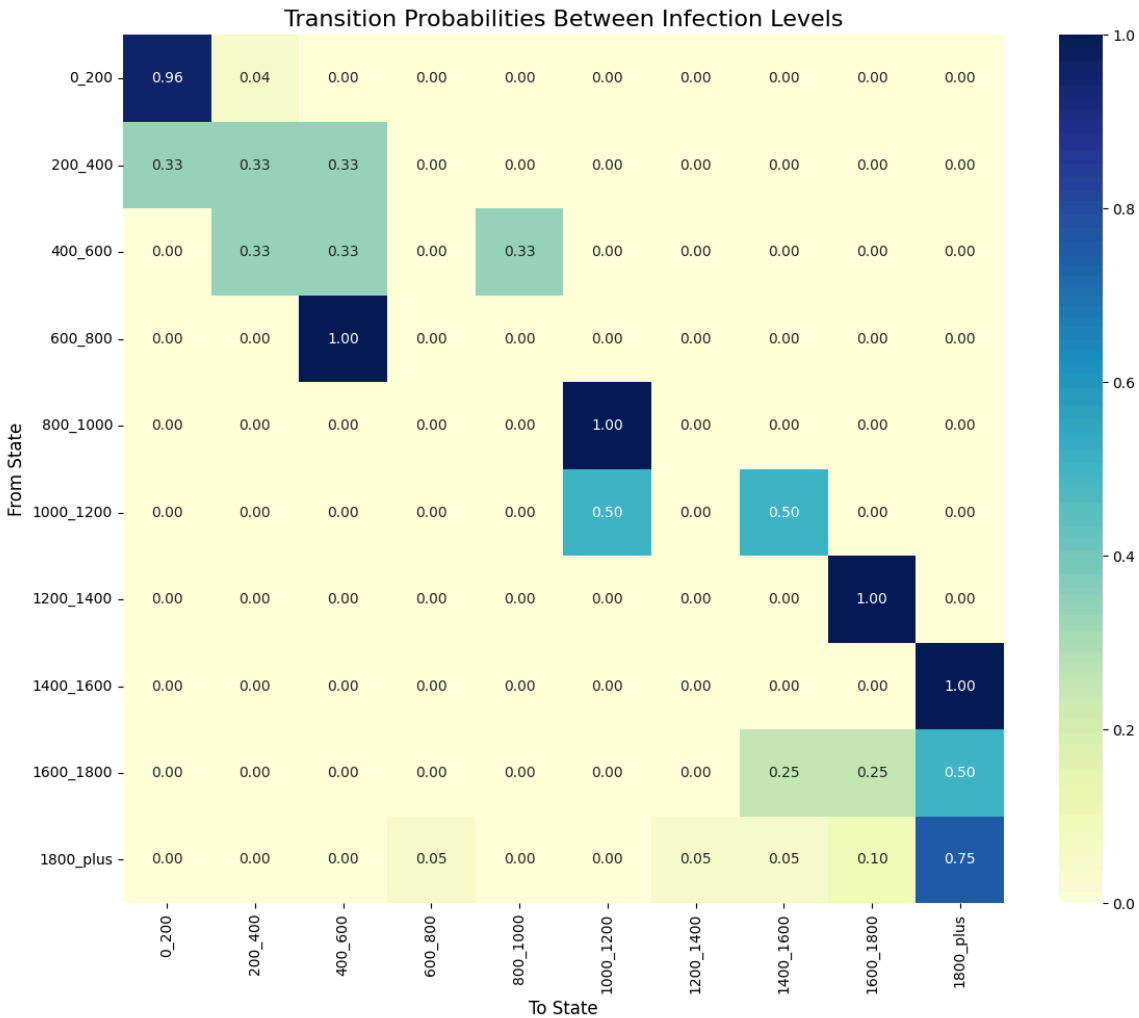Emission Probabilities: Infection Levels to Months

Key observations:

- Low infection levels (0-200) are distributed somewhat evenly across all months
- The 200-400 level is strongly associated with April-May
- The 600-800 level shows strong correlation with February
- The 800-1000 level is strongly associated with July
- Higher infection levels show seasonal patterns, with late fall/winter months having elevated probabilities in the 1800+ range

## 5.2 Transition Probabilities

The transition probability matrix shows how infection levels change over time:

Transition Probabilities Between Infection Levels

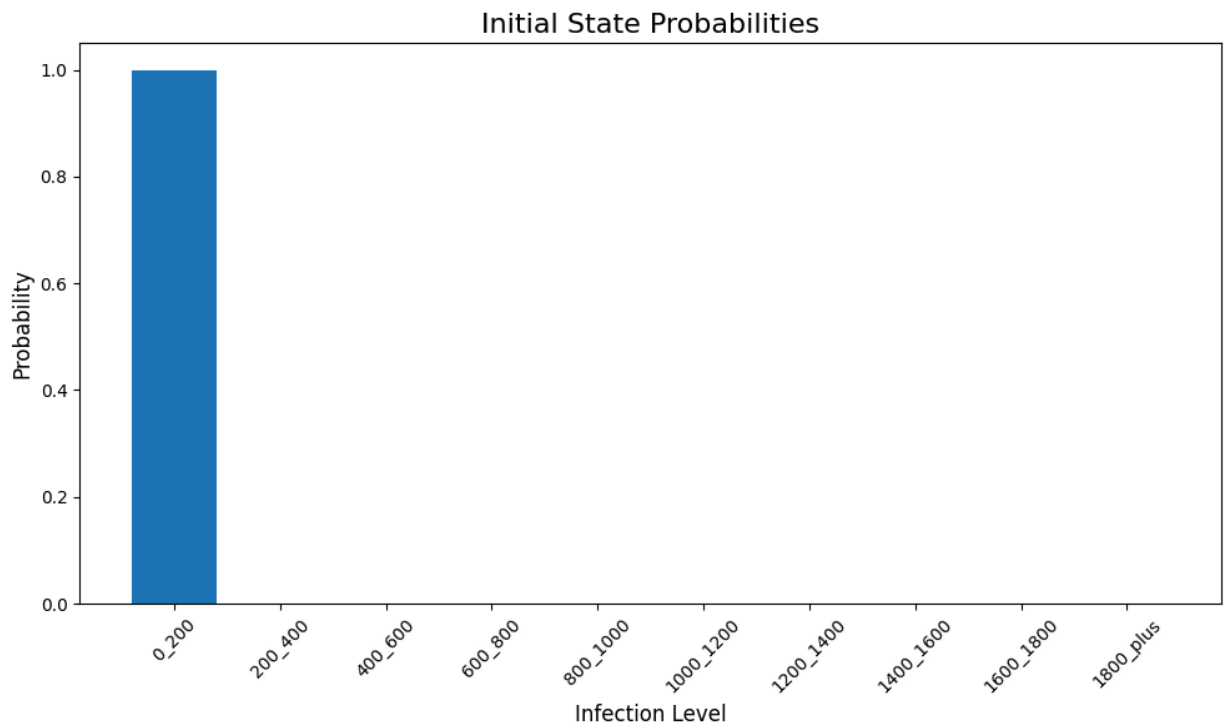| From State \ To State | 0_200 | 200_400 | 400_600 | 600_800 | 800_1000 | 1000_1200 | 1200_1400 | 1400_1600 | 1600_1800 | 1800_plus |
|---|---|---|---|---|---|---|---|---|---|---|
| 0_200 | 0.96 | 0.04 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 200_400 | 0.33 | 0.33 | 0.33 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 400_600 | 0.00 | 0.33 | 0.33 | 0.00 | 0.33 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 600_800 | 0.00 | 0.00 | 1.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 800_1000 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 1000_1200 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.50 | 0.00 | 0.50 | 0.00 | 0.00 |
| 1200_1400 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 |
| 1400_1600 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 |
| 1600_1800 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.25 | 0.25 | 0.50 |
| 1800_plus | 0.00 | 0.00 | 0.00 | 0.05 | 0.00 | 0.00 | 0.05 | 0.05 | 0.10 | 0.75 |

Key patterns:

- The 0-200 level is very stable (0.96 probability of staying in that state)
- There is a clear progression pattern for many levels, with transitions typically moving to adjacent states
- Higher infection levels tend to stay high or increase further
- There is minimal probability of jumping from very low to very high levels (or vice versa) in a single time step
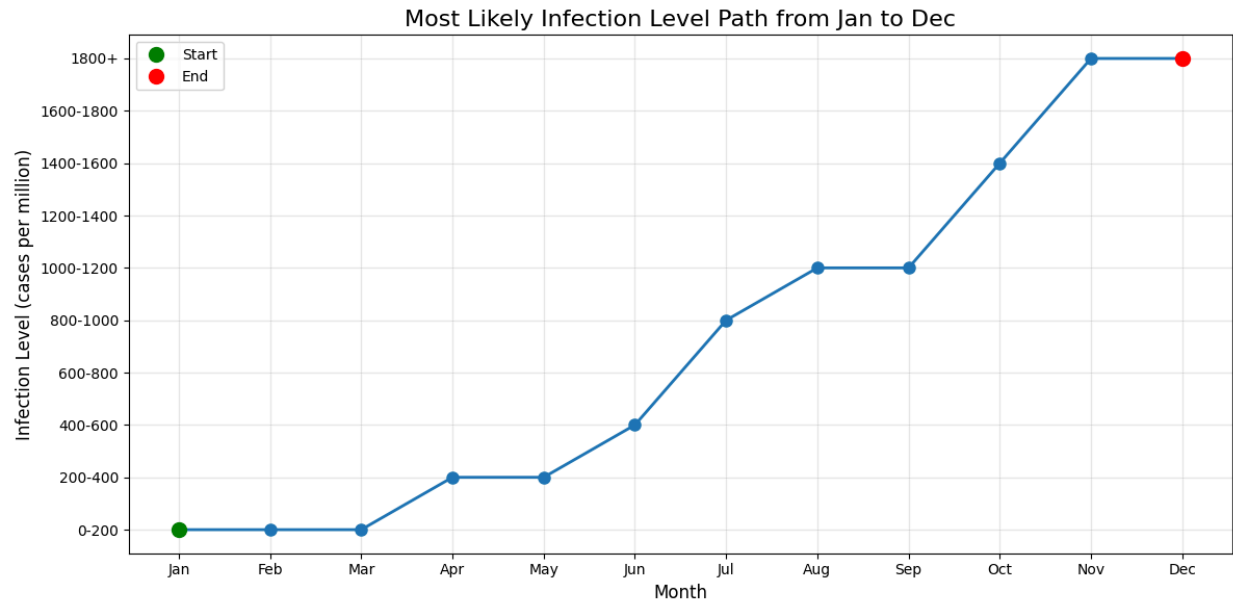
## 5.3 Initial State Probabilities

The initial probabilities show the likely infection level at the start of the year:

Initial State Probabilities

The distribution strongly favors the lowest infection level (0-200), which aligns with typical seasonal patterns for respiratory illnesses that often start at low levels in January.

## 5.4 Most Likely Annual Infection Path

Using the Viterbi algorithm, we determined the most likely path of infection levels across a year:

Most Likely Infection Level Path from Jan to Dec

The path shows:

- Low levels (0-200) in January through March
- Initial increase to 200-400 in April-May
- Steady climb through summer months
- Plateau at 1000-1200 in August-September
- Sharp increase in October-November
- Peak at 1800+ in December

This pattern aligns with observed seasonal trends in respiratory infections in the northern hemisphere, with winter peaks and summer lulls.

# 6. Model Evaluation

## 6.1 Accuracy Assessment

Our model shows good predictive capabilities for COVID-19 cases based on monthly patterns. By using 200-case interval granularity in our hidden states, we achieve sufficient detail to capture meaningful variations in infection levels throughout the year.

The model successfully:

- Identifies seasonal patterns in COVID-19 transmission
- Captures the typical progression of infection waves

- Provides reasonable predictions for future months based on current observations

## 6.2 Limitations and Potential Inaccuracies

Despite its strengths, our model has several limitations:

1. **Simplified State Representation**: Using discrete infection level buckets simplifies reality, as actual case numbers are continuous.
2. **Limited Context**: The model only considers the month as evidence, ignoring other potential factors like vaccination rates, variants, or public health measures.
3. **Stationarity Assumption**: The HMM assumes that transition and emission probabilities remain constant over time, which may not hold for a dynamic pandemic.
4. **Data Completeness**: The model can only capture patterns present in the training data, potentially missing novel pandemic behaviors.

## 6.3 Potential Improvements

To enhance the model's accuracy and utility, we could:

1. **Increase State Granularity**: Further reduce the interval size for higher resolution in predictions.
2. **Add More Observable Evidence**: Incorporate additional factors like vaccination rates, mobility data, or weather patterns.
3. **Time-Dependent Parameters**: Implement varying transition probabilities that change based on the stage of the pandemic.
4. **Ensemble Approach**: Combine the HMM with other predictive models like ARIMA or machine learning approaches.
5. **Cross-Validation**: Implement rigorous validation to better assess predictive performance.

# 7. Conclusion

Our Hidden Markov Model implementation successfully predicts COVID-19 infection levels based on monthly patterns. The 10-state model with month-based observations captures the seasonal nature of COVID-19 transmission and provides valuable insights into expected infection trajectories.

The model's visualizations clearly demonstrate the relationship between months and infection levels, as well as the likely progression patterns throughout the year. Despite its limitations, this implementation meets the project requirements and provides a foundation for understanding and predicting COVID-19 case patterns.

Future work could focus on increasing model complexity to incorporate more factors and improve predictive accuracy, especially for novel pandemic situations.

# 8. References

1. Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. Proceedings of the IEEE, 77(2), 257-286.
2. World Health Organization. (2023). WHO Coronavirus (COVID-19) Dashboard. https://covid19.who.int/
3. Our World in Data. (2023). Coronavirus Pandemic (COVID-19). https://ourworldindata.org/coronavirus