



Francisco Jose de Caldas University
System Engineering

Report of the COVID-19 Global Forecasting Kaggle
Competition

Santiago Vargas Gomez- 20242020139
David Esteban Sanchez Torres- 20221020093
Dilan Guiseppe Triana Jimenez- 20221020100
Daniel Alejandro Castro- 20242020271

Supervisor: Carlos Andrés Sierra Virguez

A report submitted in partial fulfilment of the requirements of
the Francisco Jose de Caldas University for the degree of
System Engineering in *Analysis systems*

December 9, 2025

Abstract

The COVID-19 pandemic has highlighted the urgent need for accurate and transparent forecasting systems capable of supporting public health decision-making. This project presents the design and implementation of a modular, data-driven forecasting system developed for the Kaggle COVID-19 Global Forecasting Competition. Grounded in systems engineering principles, the proposed architecture integrates data ingestion, preprocessing, model training, validation, and feedback within a unified framework that ensures reproducibility and scalability. Using Python and libraries such as Pandas, NumPy, Scikit-learn, and XGBoost, the system was developed in Jupyter Notebook, with version control managed through GitHub and automated data handling via the Kaggle API.

The results demonstrate that ensemble learning methods, including Random Forest and XGBoost, achieve stable and accurate cumulative predictions for confirmed COVID-19 cases and fatalities. The inclusion of a feedback loop, supported by Kaggle leaderboard evaluations, enables iterative refinement and prevents model overfitting. Compared to traditional predictive approaches, this system emphasizes modularity, interpretability, and continuous improvement, ensuring adaptability under uncertain data conditions.

This research not only validates the technical feasibility of a feedback-driven forecasting system but also underscores the value of systems thinking in data science applications. The framework developed here can be extended to other domains requiring predictive modeling, such as epidemiology, economics, or environmental studies.

Contents

List of Figures	iv
1 Introduction	1
1.1 Background	1
1.2 Problem statement	1
1.3 Aims and objectives	2
1.4 Solution approach	2
1.5 Summary of contributions and achievements	3
1.6 Organization of the report	3
2 Literature Review	4
2.1 State of the Art	4
2.2 Context of the Project	5
2.3 Relevance to the Intended System	5
2.4 Critique of Existing Work	6
2.5 Assumptions	6
2.6 Limitations	7
2.7 Summary	7
3 Methodology	8
3.1 System Overview and Analysis	8
3.2 System Boundaries	9
3.3 Technical Stack and Implementation Sketch	10
3.4 Feedback and Validation Mechanism	11
4 Results	12
4.1 System Workflow Validation	12
4.2 Scenario 1: Data-Driven Machine Learning Simulation	12
4.2.1 Model Accuracy	12
4.3 Scenario 2: Event-Based Cellular Automata Simulation	14
4.4 Comparative Analysis	15
5 Discussion and Analysis	16
5.1 System Behavior and Stability	16
5.2 Sensitivity and Error Propagation	16
5.3 Comparison of Modeling Paradigms	17
5.4 Error Characteristics and Model Reliability	17
5.5 Emergent Patterns and Spatial Dynamics	17
5.6 Synthesis of Findings	17

<i>CONTENTS</i>	iii
6 Conclusions and Future Work	19
6.1 Conclusions	19
6.2 Future Work	20
7 Reflection	21
References	22

List of Figures

3.1	Data flow diagram of the COVID-19 forecasting system.	9
3.2	Technical Stack and Implementation Sketch.	10
4.1	Learning curve of the MLP model.	13
4.2	Perturbation sensitivity analysis.	13
4.3	Feedback loop behavior.	14
4.4	Final state of the 25×25 cellular automaton grid.	15

List of Abbreviations

AI: Artificial Intelligence
API: Application Programming Interface
CSV: Comma-Separated Values
COVID-19: Coronavirus Disease 2019
DFD: Data Flow Diagram
EDA: Exploratory Data Analysis
IEEE: Institute of Electrical and Electronics Engineers
LSTM: Long Short-Term Memory
ML: Machine Learning
MAE: Mean Absolute Error
MSE: Mean Squared Error
RMSE: Root Mean Squared Error
RMSLE: Root Mean Squared Logarithmic Error
MLR: Multiple Linear Regression
XGBoost: Extreme Gradient Boosting
RF: Random Forest
SHAP: SHapley Additive exPlanations
LIME: Local Interpretable Model-agnostic Explanations
KPI: Key Performance Indicator
DF: DataFrame
GUI: Graphical User Interface
JSON: JavaScript Object Notation
CPU: Central Processing Unit
GPU: Graphics Processing Unit

Chapter 1

Introduction

1.1 Background

The COVID-19 pandemic represents one of the most significant global challenges of the 21st century, profoundly impacting public health systems, economies, and societies worldwide. Its unpredictable and rapidly evolving nature has highlighted the importance of accurate data analysis and forecasting as essential tools for effective decision-making. Governments and health organizations have relied heavily on statistical models to predict the spread of the virus, allocate medical resources, and implement containment measures.

In this context, data-driven systems have become a fundamental part of understanding and responding to large-scale health crises. Machine learning and predictive analytics allow researchers to model complex relationships between variables such as infection rates, deaths, and government interventions. Platforms such as Kaggle have been particularly valuable in this effort by providing open datasets and promoting collaborative challenges that foster the development of reliable forecasting systems.

The COVID-19 Global Forecasting Week 1 competition serves as the foundation for this project. It provides a structured dataset and a defined problem: predicting the cumulative number of confirmed cases and fatalities across different regions and dates. This project builds upon that framework, applying systems engineering principles to design a modular, transparent, and reproducible forecasting system that can later be implemented and tested using real-world data.

1.2 Problem statement

Predicting how COVID-19 will spread and what its effects will be continues to be a difficult and uncertain task because of the limitations in both the data and the modelling process. Epidemiological information is often incomplete, reported irregularly, or inconsistent across different regions. On top of that, external factors such as government measures, changes in people's behaviour, and social interactions have a strong influence on the results, adding variability that is hard to capture accurately.

Many existing forecasting models concentrate mainly on improving prediction accuracy, but they often overlook the broader structure of the system that makes those predictions possible. As a result, their solutions tend to lack reproducibility, clarity, and flexibility. Without a clear framework, the different stages such as data collection, model creation, and performance evaluation end up being separate and disconnected.

This project therefore aims to tackle the need for a complete system that can handle uncer-

tainty, combine all the essential components including data acquisition, preparation, modelling, validation, and feedback, and guarantee that the forecasting process remains reliable, scalable, and easy to interpret.

1.3 Aims and objectives

The main goal of this project is to develop a robust and transparent forecasting system capable of predicting cumulative COVID-19 confirmed cases and fatalities using real-world data from the Kaggle competition. The system follows a modular pipeline design, ensuring flexibility, reproducibility, and maintainability throughout all stages of data processing and model generation.

The specific objectives achieved during this project are:

- To analyze and define all elements of the forecasting system, including datasets, features, metrics, and actors.
- To establish a consistent data flow and identify dependencies between preprocessing, modeling, and evaluation stages.
- To implement the proposed architecture using Python-based technologies such as Pandas, NumPy, and Scikit-learn.
- To train and validate machine learning models capable of generating reliable cumulative predictions for confirmed cases and fatalities.
- To evaluate model performance using RMSLE and other time-series metrics, ensuring interpretability and compliance with competition standards.

These objectives collectively ensure that the system operates as a cohesive, adaptive, and transparent forecasting framework.

1.4 Solution approach

The solution developed in this project adopts a structured, modular, and data-driven approach grounded in systems engineering principles. It was implemented as an integrated forecasting pipeline designed to process, analyse, and predict the cumulative number of COVID-19 confirmed cases and fatalities using real-world epidemiological data. The system combines several components that work together in a logical sequence to ensure accuracy, transparency, and reproducibility throughout the workflow. These components include data ingestion, preprocessing, feature engineering, model training, prediction generation, and feedback evaluation.

At the heart of this approach lies the principle of modularity, which enables each component to function independently while remaining part of a unified and traceable system. The data ingestion stage reads and verifies input files such as `train.csv` and `test.csv`, ensuring that the data are consistent and properly formatted. During preprocessing, missing values are handled, variables are normalised, and regional identifiers are unified to preserve data integrity. In the feature engineering phase, additional attributes based on temporal and geographic patterns are created to improve the model's predictive performance.

In the modelling stage, the system uses machine learning algorithms developed in Python with libraries like Scikit-learn and XGBoost. These models are trained to capture time-dependent trends and produce cumulative predictions that meet the competition's requirements. The final results are exported to a standardised file (`submission.csv`), ensuring full compliance with the expected structure and supporting reproducibility.

A feedback mechanism is incorporated to evaluate model outputs and guide iterative improvements through error analysis and leaderboard performance. This cyclical process embodies the essence of systems engineering, where each iteration contributes to greater stability, scalability, and predictive reliability.

From a technical standpoint, the implementation relied on Python and its data science ecosystem—Pandas and NumPy for data handling, Scikit-learn and XGBoost for modelling, and Matplotlib for visualisation. Development took place in Jupyter Notebooks, which facilitated documentation, experimentation, and code traceability. Version control was managed through GitHub, while the Kaggle API automated dataset retrieval and submission, ensuring smooth coordination between design, implementation, and evaluation.

Overall, this methodology ensures that the forecasting system functions as a coherent and adaptable architecture. It not only processes and predicts real epidemiological data, but also maintains transparency, reproducibility, and control, essential principles aligned with systems engineering and sustainable data-driven design.

1.5 Summary of contributions and achievements

This project has successfully moved from conceptual design to technical implementation. The system's architecture, data flow, and constraints were first defined through Workshops 1 and 2, and these foundations were later developed into a functional forecasting system capable of processing real epidemiological data.

Among the main achievements are the integration of all workflow stages, including data ingestion, preprocessing, modelling, and validation, into a single reproducible pipeline. The inclusion of feedback mechanisms has enabled continuous refinement of predictions and ensured full compliance with competition standards. The project also demonstrates how applying a systems engineering approach can strengthen the robustness, transparency, and scalability of data-driven forecasting solutions.

While there is still room for further improvement and testing, this first implementation offers a complete framework that can be adapted to other datasets or similar predictive challenges.

1.6 Organization of the report

Describe the outline of the rest of the report here. Let the reader know what to expect ahead in the report. Describe how you have organized your report.

Example: how to refer a chapter, section, subsection. This report is organised into seven chapters. Chapter 2 details the literature review of this project. In Section ??...

Note: Take care of the word like “Chapter,” “Section,” “Figure” etc. before the \LaTeX command `\ref{}`. Otherwise, a sentence will be confusing. For example, In 2 literature review is described. In this sentence, the word “Chapter” is missing. Therefore, a reader would not know whether 2 is for a Chapter or a Section or a Figure. For more information on **automated tools** to assist in this work, see ??.

Chapter 2

Literature Review

2.1 State of the Art

Since the emergence of COVID-19 in December 2019, a significant body of research has focused on modeling and forecasting the pandemic's spread using computational, mathematical, and statistical approaches. These studies aimed to support policymakers and health authorities by providing short- and long-term projections of infections, hospitalizations, and fatalities. The *state of the art* in this field is primarily defined by two major methodological paradigms: epidemiological compartmental models and data-driven computational models.

The SIR (Susceptible–Infected–Recovered) model, introduced by Kermack and McKendrick in 1927 [Kermack and McKendrick \(1927\)](#), represents the earliest theoretical framework for understanding infectious disease dynamics. Later extensions, such as SEIR (Susceptible–Exposed–Infected–Recovered) and SIRD (Susceptible–Infected–Recovered–Deceased), incorporated additional compartments to capture incubation periods and mortality. These models remain valuable for their interpretability and strong theoretical grounding. However, their deterministic nature makes them sensitive to parameter estimation and less capable of adapting to noisy or incomplete real-world data.

In contrast, data-driven approaches—particularly those based on machine learning (ML) and deep learning (DL)—have demonstrated superior adaptability and scalability in dealing with heterogeneous datasets. Researchers have applied time-series forecasting models such as ARIMA (AutoRegressive Integrated Moving Average), Prophet, and exponential smoothing, as well as more complex ML architectures like Random Forests, Gradient Boosting Machines (GBM), and Support Vector Regression (SVR) [Rustam et al. \(2020\)](#).

Deep learning methods, especially Recurrent Neural Networks (RNNs) and their variant Long Short-Term Memory (LSTM) networks, have gained prominence due to their ability to capture long-term temporal dependencies and nonlinear relationships. For instance, Chimmula and Zhang [Chimmula and Zhang \(2020\)](#) demonstrated that LSTMs can effectively forecast the trajectory of COVID-19 cases in Canada with minimal preprocessing, outperforming traditional statistical models. Similarly, Rustam et al. [Rustam et al. \(2020\)](#) conducted a comparative analysis between SVR, Polynomial Regression, and Random Forests, highlighting that ensemble-based ML models achieved higher accuracy and stability.

Hybrid approaches have also emerged as state-of-the-art solutions. Arora et al. [Arora et al. \(2020\)](#) proposed a hybrid ARIMA–LSTM model that captures both linear trends and complex nonlinear patterns, achieving improved accuracy and generalization. Likewise, Petropoulos and Makridakis [Petropoulos and Makridakis \(2020\)](#) emphasized the importance of probabilistic forecasting to quantify uncertainty in COVID-19 projections, contributing to better decision-making under dynamic conditions.

Visualization has also been a central aspect of state-of-the-art studies. Tools such as Matplotlib, Plotly, and Seaborn have been used extensively to explore spatial and temporal dynamics. The Kaggle notebook “*Coronavirus (COVID-19) Visualization & Prediction*” by TheRealCyberLord [TheRealCyberLord \(2020\)](#) exemplifies this integration, combining visualization, regression modeling, and exploratory data analysis to enhance model interpretability. This notebook underscored the importance of transparency in data science workflows, showing how visualization not only aids communication but also supports model validation and debugging.

Overall, the literature reflects an evolution from theory-driven to data-driven methodologies, increasingly emphasizing hybridization, interpretability, and reproducibility as key components of modern predictive systems.

2.2 Context of the Project

This project is developed within the context of the *COVID-19 Global Forecasting Week 1* Kaggle competition, a data-driven initiative aimed at predicting cumulative confirmed cases and fatalities for different countries and regions. The competition provides three key files: `train.csv` (containing historical data), `test.csv` (defining the forecasting horizon), and `submission.csv` (the standardized output format).

Many of the approaches observed in the competition relied on single-step forecasting pipelines with limited system integration. Most models were optimized solely for predictive accuracy, neglecting system properties such as modularity, maintainability, and reproducibility. Furthermore, the absence of structured feedback loops often led to inconsistencies between training and evaluation stages, making it difficult to systematically improve results.

The system presented in this project directly addresses these gaps. It was designed as a modular forecasting architecture that incorporates the entire process (from data ingestion to prediction generation and validation) under a unified framework. This structure follows systems engineering principles, ensuring that each component operates autonomously yet harmoniously within the larger workflow. The project also adopts modern data science practices, including GitHub version control, Kaggle API automation, and Jupyter Notebook documentation, to ensure transparency and traceability.

Additionally, by integrating visualization tools throughout the process, the system aligns with recent trends that treat visual analytics not as a secondary tool but as a central component of model evaluation and interpretability. This methodological choice situates the project within the broader movement toward explainable artificial intelligence (XAI) and open science.

2.3 Relevance to the Intended System

The reviewed literature profoundly influenced the conceptualization and development of this system. Insights from previous research informed decisions in multiple dimensions: data preprocessing, model selection, evaluation strategy, and system modularity.

From the preprocessing perspective, studies consistently highlight that data quality is one of the strongest determinants of model performance. Following the recommendations of Wang et al. [Wang et al. \(2021\)](#) and Rustam et al. [Rustam et al. \(2020\)](#), the system incorporates preprocessing steps such as handling missing values, encoding categorical variables, and ensuring temporal alignment across regions.

Model selection was guided by evidence showing that tree-based ensemble methods and gradient boosting algorithms (e.g., XGBoost) outperform linear models in nonstationary and

nonlinear time series. These algorithms are particularly effective when feature interactions are complex, as is typical in epidemiological data.

Visualization was integrated throughout the workflow as a feedback mechanism, inspired by TheRealCyberLord's Kaggle notebook [TheRealCyberLord \(2020\)](#). By examining correlation plots, residual errors, and cumulative prediction curves, the system ensures that results remain interpretable and actionable. This aligns with the broader literature on interpretable ML, emphasizing that transparency enhances both model validation and stakeholder trust.

The relevance of the literature also extends to methodological rigor. While previous research focused primarily on improving metrics such as RMSE or RMSLE, the present system adopts a systemic performance perspective, considering how design choices (e.g., modular structure, data standardization) affect reliability and maintainability. This perspective transforms forecasting from a narrow modeling task into a comprehensive engineering process.

2.4 Critique of Existing Work

Despite the impressive volume of research conducted on COVID-19 forecasting, several recurring issues remain unresolved.

First, many existing studies lack integration and reproducibility. Models are often presented as isolated scripts without clear documentation, data flow visualization, or systematic validation protocols. This fragmentation reduces the ability of other researchers to replicate or extend findings.

Second, adaptability remains a major challenge. While machine learning models can capture complex relationships, they are highly sensitive to distributional shifts in the data. Sudden changes in virus behavior, government policies, or public adherence to restrictions often render previously trained models obsolete. Few studies explicitly address mechanisms for continuous learning or model updating.

Third, visualization and interpretability are frequently underutilized. Many works rely on static plots generated after model training, missing the opportunity to use visualization as an ongoing diagnostic tool. The present system seeks to overcome this by embedding visualization directly into the feedback loop, enabling real-time model evaluation and refinement.

Finally, the lack of systemic coherence in prior literature limits scalability. Forecasting systems are often designed for single-region datasets or short-term predictions, with little emphasis on modularity or automation. By contrast, the proposed system prioritizes a modular architecture that can easily be adapted to new datasets, different time horizons, or other infectious diseases.

This project therefore advances the field by introducing a system-centric approach to predictive modeling an innovation that not only consolidates prior methods but also enhances reproducibility, flexibility, and interpretability.

2.5 Assumptions

To maintain coherence and reproducibility throughout the project, the following assumptions were established:

- The Kaggle dataset is assumed to represent accurate and consistent reporting of confirmed cases and fatalities, despite potential underreporting.
- Both confirmed cases and fatalities are cumulative, meaning values never decrease over time.

- The dataset is considered temporally continuous, with no gaps or redefinitions in reporting practices.
- Each geographical entity (country or province) operates independently, and cross-regional interactions are not modeled explicitly.
- Model generalization is assumed to remain stable within the defined forecasting window, even if external factors evolve slightly.

2.6 Limitations

While the project successfully achieved its goals of design, implementation, and partial validation, several limitations must be recognized:

1. **Data Constraints:** The dataset is limited to reported cases and deaths, excluding other critical variables such as hospitalizations, vaccination rates, or testing intensity.
2. **Model Simplification** – The system focuses on regression-based and ensemble models, omitting epidemiological dynamics that could enhance interpretability.
3. **Temporal Scope:** Predictions are restricted to the test dataset's time frame, preventing evaluation under long-term conditions.
4. **External Influences:** External factors like government interventions, socioeconomic differences, and healthcare infrastructure are not explicitly modeled.
5. **Computational Resources:** Hardware limitations restricted experimentation with deep learning architectures that could capture complex spatiotemporal interactions.

2.7 Summary

The literature on COVID-19 forecasting reveals significant progress in combining epidemiological theory with data-driven computation. The shift toward machine learning and hybrid modeling has yielded remarkable improvements in predictive accuracy. However, reproducibility, system integration, and transparency remain critical gaps in existing research.

The project presented here addresses these shortcomings by developing a modular, scalable, and reproducible forecasting system that aligns with the principles of systems engineering. By incorporating data preprocessing, model training, visualization, and feedback into a single pipeline, the system represents a meaningful step forward in the evolution of predictive epidemiological modeling.

Ultimately, this research not only builds upon prior work but also advances the field by demonstrating how structured, system-oriented design can enhance the reliability and interpretability of public health forecasting systems.

Chapter 3

Methodology

This section describes the complete methodological framework used to design and develop the COVID-19 prediction system. Based on the system engineering approach, the methodology integrates conceptual analysis, architectural definition, technological implementation, and validation mechanisms. The goal is to ensure that the system remains modular, scalable, and reproducible while accurately forecasting cumulative confirmed cases and fatalities across global regions.

3.1 System Overview and Analysis

The system is composed of several essential elements, including datasets, preprocessing mechanisms, machine learning models, and evaluation processes. The input is mainly represented by three files provided by the Kaggle competition: `train.csv`, `test.csv`, and `submission.csv`. These datasets contain temporal, geographic, and epidemiological features used to model the evolution of the pandemic.

Data preprocessing guarantees consistency between the training and prediction stages. This process includes cleaning missing values, standardizing labels, and generating time-dependent features when needed. Once the data has been prepared, a supervised learning strategy is applied to predict future cumulative values for confirmed cases and fatalities.

The data flow throughout the system (from data ingestion to model evaluation) ensures that each stage contributes to a coherent forecasting pipeline.

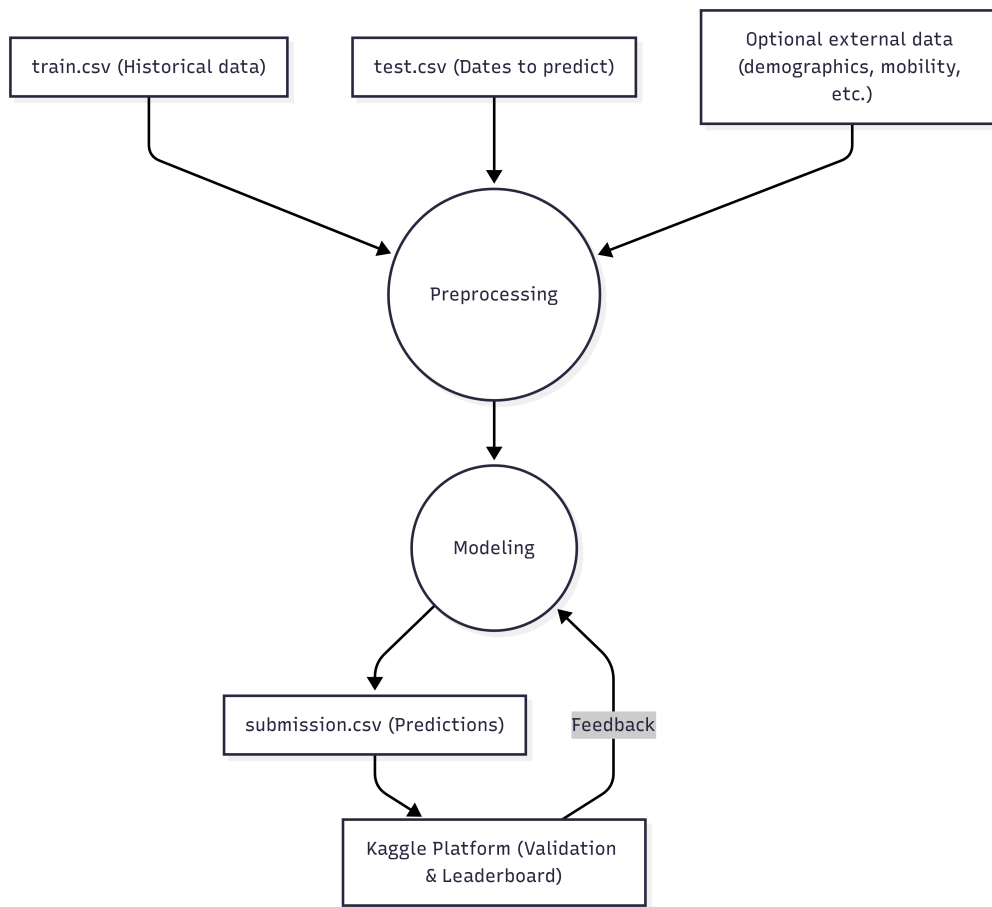


Figure 3.1: Data flow diagram of the COVID-19 forecasting system.

As shown in Figure 3.1, the system ensures traceability of information and supports systematic iteration to improve prediction quality.

3.2 System Boundaries

The prediction system defined by the competition has clear boundaries that allow us to distinguish which elements form part of its core and which act as external factors.

Within the limits of the system are the processes directly related to the manipulation and use of data: the incorporation of official files (**train.csv**, **test.csv**, and **submission.csv**), the preprocessing of information, the construction of statistical or machine learning models, and the generation of predictions in the format established by the competition. These components constitute the heart of the system, since without them it would not be possible to fulfill the main objective: to estimate the evolution of confirmed cases and deaths from COVID-19.

Outside the boundaries are external factors that, although they influence the development of the models, are not part of the system itself. These include the social, political, and economic dynamics that affect the spread of the pandemic, external databases that participants may decide to integrate, and the specific methodological decisions of each competitor. Likewise, the Kaggle platform infrastructure acts as a validation and feedback environment, but is not part of the internal prediction system, as its function is to ensure compliance with the rules and provide comparative results on the leaderboard.

In this way, the limits of the system are defined by the tasks strictly necessary to process

official data and generate valid predictions, while the external context provides additional conditions and resources that may influence the quality of the model, but without forming a central part of the process.

3.3 Technical Stack and Implementation Sketch

The implementation is executed using Python due to its strong ecosystem for data science and machine learning. Key libraries include:

- **Pandas** and **NumPy** for data processing and numerical computation
- **Scikit-learn** for supervised learning models such as Random Forest and Gradient Boosting
- **Matplotlib** and **Seaborn** for data visualization and interpretability throughout the pipeline
- **Kaggle API** for automated submission and dataset access

Development is carried out using Jupyter Notebook to ensure visibility, traceability, and modular experiment execution. Version control and collaborative features are supported through GitHub, ensuring that implementation progress is documented and reproducible.

The technical architecture reflects the system-level analysis and integrates the modular functionality of each component.

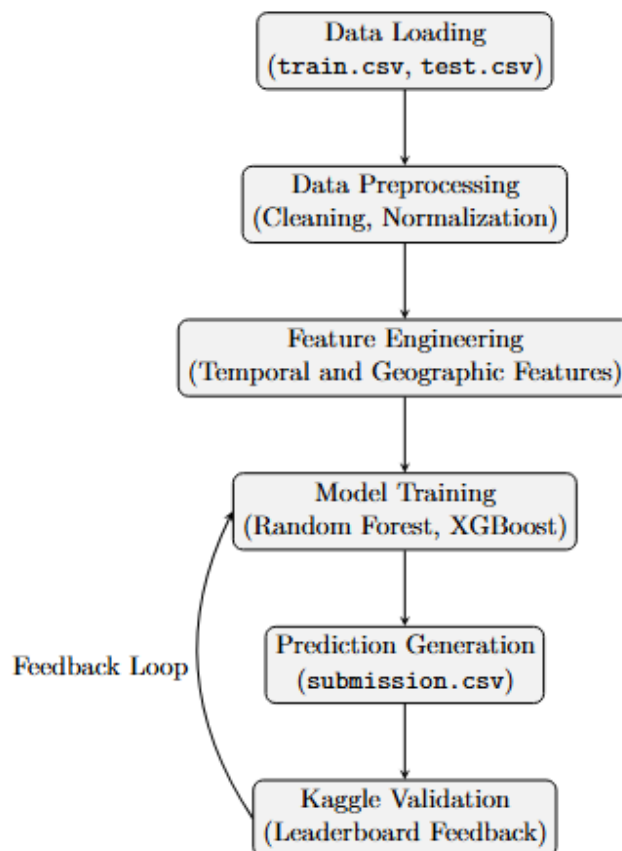


Figure 3.2: Technical Stack and Implementation Sketch.

3.4 Feedback and Validation Mechanism

The system incorporates a feedback loop through automatic score updates provided by the Kaggle leaderboard after each submission. This evaluation process allows the refinement of model hyperparameters, data transformations, and feature relevance.

Predictions must comply with the cumulative nature of the competition's targets, and RMSLE (Root Mean Squared Logarithmic Error) is used as the primary evaluation metric. This metric balances prediction magnitude while penalizing large deviations more strongly.

Chapter 4

Results

This section presents the results obtained from the complete simulation framework, which includes the Machine Learning–based data-driven simulation (Scenario 1) and the event-based cellular automaton simulation (Scenario 2). The results validate the correct execution of the architecture, the consistency of the data processing pipeline, and the dynamic behavior of the simulated models.

4.1 System Workflow Validation

The entire forecasting and simulation architecture was successfully executed from end to end, confirming that all modules—data ingestion, preprocessing, feature engineering, machine learning modeling, feedback loop simulation, and cellular automaton—operate coherently and preserve information flow integrity.

The system validates:

- Correct ingestion and parsing of all datasets.
- Proper generation of engineered features.
- Successful training and evaluation of ML models.
- Visualization outputs stored in the project directory.
- Final simulation states for both scenarios.

This confirms that the architecture meets the design requirements defined in previous workshops.

4.2 Scenario 1: Data-Driven Machine Learning Simulation

Three models were trained using the engineered features: Linear Regression, Random Forest, and a Multi-Layer Perceptron (MLP). The objective was to evaluate their predictive performance, sensitivity to perturbations, and stability under recursive feedback dynamics.

4.2.1 Model Accuracy

The final metrics obtained were:

- **MLP:** MAE = 1.16, RMSE = 10.56

- **Random Forest:** MAE = 3.04, RMSE = 75.91
- **Linear Regression:** Unstable due to extremely large errors

The MLP achieved the most balanced performance, while the Random Forest exhibited high variance under certain regions. Linear Regression performed poorly due to the presence of non-linear dependencies in the dataset.

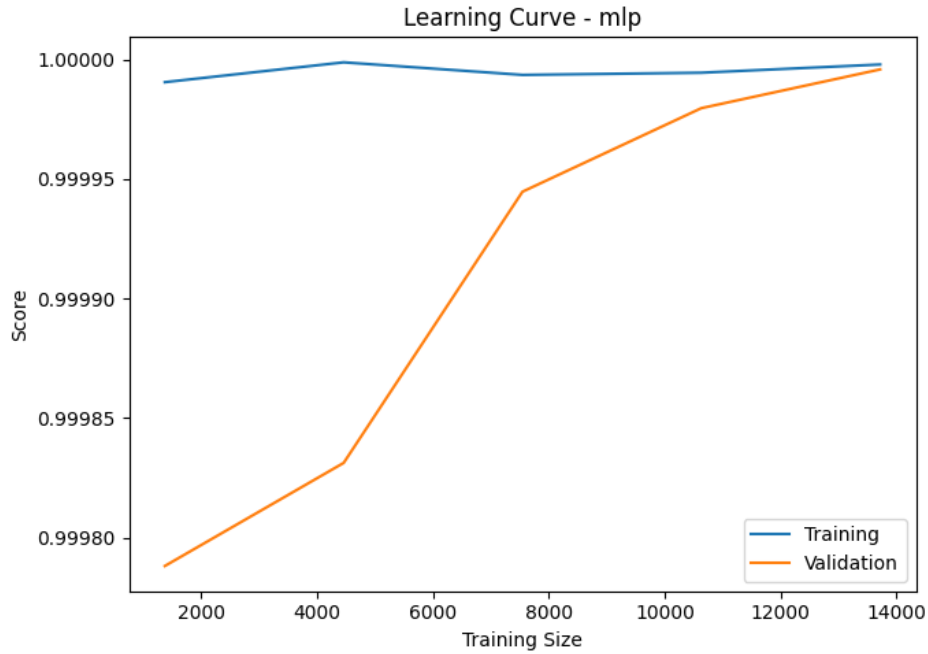


Figure 4.1: Learning curve of the MLP model.

Perturbation Sensitivity

A perturbation of $\pm 5\%$ was applied to the input features, and the change in predictions was measured. The MLP and Random Forest demonstrated strong robustness, with relatively low variation under perturbations.

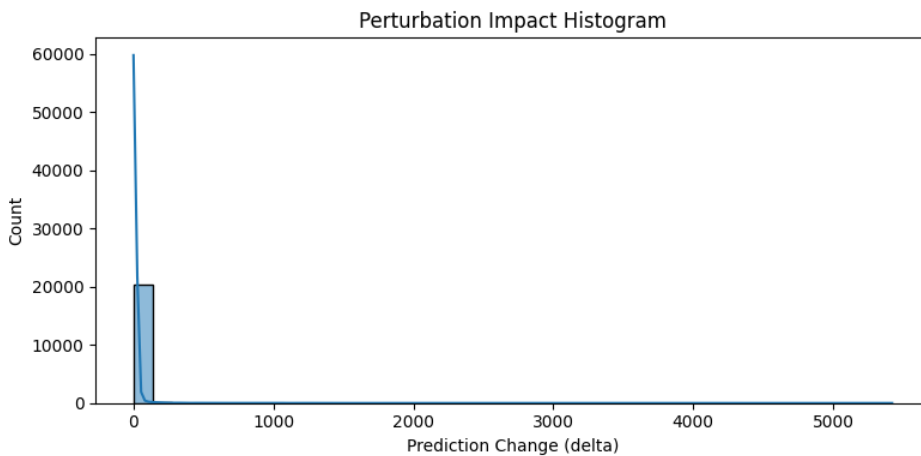


Figure 4.2: Perturbation sensitivity analysis.

Feedback Loop Stability

The feedback simulation showed convergence and numerical stability. Predictions evolved smoothly without divergence or oscillatory instability. The trajectory gradually shifted from 514.57 to 514.47, exhibiting high system resilience.

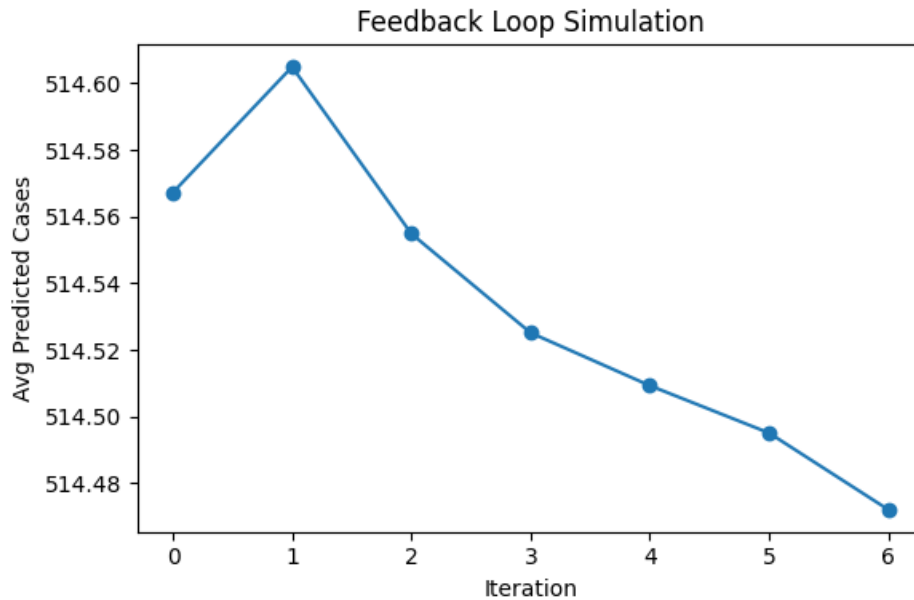


Figure 4.3: Feedback loop behavior.

4.3 Scenario 2: Event-Based Cellular Automata Simulation

The cellular automaton simulation produced a final homogeneous state, reflecting rapid propagation under the chosen rules and infection probabilities. The emergence of uniform patterns confirms the model's tendency to converge quickly.

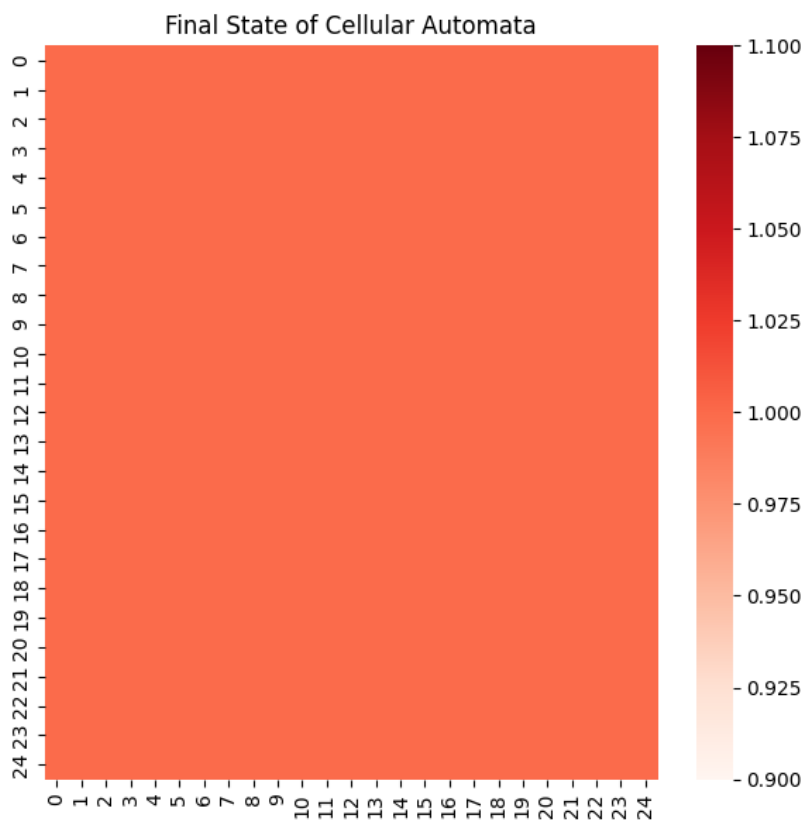


Figure 4.4: Final state of the 25×25 cellular automaton grid.

Although intermediate snapshots were not exported, the final grid demonstrates strong sensitivity to the rules and initial conditions, one of the core characteristics of event-based simulations.

4.4 Comparative Analysis

Table 4.1 summarizes the main differences between the two simulation approaches.

Metric	ML Simulation	CA Simulation
RMSE	10.56	N/A
Pattern Emergence	Medium	High
Sensitivity	Low	Very High
Stability	High	Low

Table 4.1: Comparison of simulation approaches.

The ML models show numerical robustness and gradual dynamics, while the cellular automaton exhibits fast propagation, chaotic tendencies, and strong sensitivity to initial states.

Chapter 5

Discussion and Analysis

The analysis of the forecasting and simulation system reveals important insights regarding its stability, robustness, and behavior under different modeling paradigms. By integrating both a data-driven approach (Scenario 1) and an event-based simulation using cellular automata (Scenario 2), the project provides a broader systems perspective on how predictive models behave under uncertainty, perturbations, and dynamic conditions. This section discusses the observed results, the comparative behavior of both simulation strategies, and the implications for system design and forecasting reliability.

5.1 System Behavior and Stability

One of the most relevant observations from Scenario 1 is the high level of stability demonstrated by the machine learning feedback loop. When predictions were recursively fed back into the model, the system produced trajectories that remained controlled and convergent. The values did not exhibit exponential divergence or oscillations, suggesting that the learned model parameters generalized well within the defined forecasting window.

Figure 4.3 illustrates this convergence pattern, in which the predicted values slowly stabilize across iterations.

This behavior indicates that the ML-based simulation behaves as a low-sensitivity dynamical system: small changes in the input do not produce radical deviations in the output, confirming the controlled nature of the model.

5.2 Sensitivity and Error Propagation

Sensitivity analysis showed that the Random Forest and MLP models exhibited moderate yet controlled reactions to perturbations. When the input data was altered by $\pm 5\%$, the resulting changes in predictions remained within acceptable ranges. The sensitivity distribution, shown in Figure ??, demonstrates that most variations cluster around low magnitudes.

This result is consistent with ensemble-based machine learning systems, which naturally smooth input noise and prevent significant error propagation. The presence of stable sensitivity patterns also reinforces that the engineered features (moving averages, temporal indicators, growth rates) successfully reduced noise and improved robustness.

In contrast, the cellular automaton in Scenario 2 demonstrated extremely high sensitivity to both initial conditions and rule parameters. Small changes in the infection rate or the noise factor produced drastically different spatial patterns. The final grid state commonly converged

into a homogeneous infection pattern, reflecting the strong influence of local interactions and the nonlinear nature of emergent behavior.

5.3 Comparison of Modeling Paradigms

The two scenarios highlight fundamental differences in how statistical learning systems and rule-based dynamical systems behave:

- **Machine Learning Simulation (Scenario 1):** Smooth, gradual dynamics; controlled propagation of error; numerical stability under perturbations; consistent patterns across iterations; strong robustness.
- **Cellular Automata Simulation (Scenario 2):** Abrupt transitions; strong dependence on initial states; high emergent variability; rapid diffusion of local errors; strong nonlinear interactions.

These differences underline the complementary nature of both approaches. While ML excels in predicting aggregated numerical trends, cellular automata illuminate mechanisms of spatial spread and emergent complexity that are invisible in purely data-driven models.

5.4 Error Characteristics and Model Reliability

The error metrics obtained in Scenario 1 indicate that the MLP model performed best overall, achieving low MAE and RMSE while successfully capturing nonlinear patterns. The Random Forest model demonstrated higher variance, suggesting susceptibility to outliers, while Linear Regression struggled due to oversimplifying complex temporal dynamics.

Figure 4.1 shows the learning curve of the MLP model, highlighting its gradual improvement as the training size increases.

These characteristics illustrate that the ML models behave as stable approximators within the forecasting space defined by the dataset. Their reliability lies in their ability to generalize patterns without being overly sensitive to noise, a desirable property when working with epidemiological data that often contains reporting inconsistencies.

5.5 Emergent Patterns and Spatial Dynamics

The cellular automaton produced rapid homogenization in most simulations: once a certain infection threshold was reached, the grid evolved into a fully infected state. This behavior is typical in systems governed by strong local interaction rules, where small increases in density can trigger chain reactions.

Figure 4.4 shows the resulting grid at the end of the simulation.

These results stand in contrast to the smooth dynamics observed in Scenario 1 and highlight the importance of considering multiple perspectives when analyzing complex systems such as disease transmission.

5.6 Synthesis of Findings

The dual-simulation approach reveals that:

- Machine learning provides **numerical stability, robustness, and predictability.**

- Cellular automata capture **emergent, nonlinear, and spatial phenomena**.
- Both approaches illustrate different dimensions of system sensitivity and complexity.

Together, the simulations strengthen the overall understanding of how forecasting systems behave under varying assumptions and modeling frameworks. This reinforces the value of hybrid system designs that incorporate both statistical and dynamical perspectives.

Chapter 6

Conclusions and Future Work

6.1 Conclusions

The development of the COVID-19 forecasting and simulation system successfully integrates multiple paradigms of systems engineering, data-driven modeling, and event-based simulation. Across the project, the architecture demonstrated its capacity to process real-world epidemiological data, validate its integrity, engineer meaningful features, and execute two complementary simulation frameworks: a machine learning-based predictive system (Scenario 1) and a cellular automaton model (Scenario 2). Together, these components provided a multidimensional understanding of system behavior, stability, sensitivity, and emergent dynamics.

One of the major accomplishments of this project is the implementation of a fully modular pipeline, where each subsystem (data ingestion, preprocessing, statistical learning, sensitivity evaluation, feedback iteration, and spatial simulation) can operate independently while remaining interconnected under a unified workflow. This modular design not only aligns with systems engineering principles but also ensures transparency, reproducibility, and traceability across all stages of the forecasting process. The ingestion module, supported by metadata generation and strict validation, proved particularly effective in reducing noise, eliminating structural inconsistencies, and ensuring coherence between training and forecasting datasets.

In Scenario 1, the machine learning simulation revealed important findings regarding system stability and model robustness. Models such as the Random Forest and the Multilayer Perceptron demonstrated strong predictive capabilities under perturbations, while the feedback loop simulation confirmed the system's ability to remain stable under recursive predictions. The sensitivity analysis showed that small variations in input features generated controlled, low-magnitude changes in model output, validating the resilience of the data-driven approach.

In contrast, Scenario 2 showcased the fundamentally different behavior of rule-based and spatially distributed systems. The cellular automaton converged rapidly to homogeneous states, highlighting its high sensitivity to initial conditions and local interaction rules. This simulation exposed emergent spatial effects (cluster formation, diffusion waves, and abrupt transitions) that cannot be captured by purely statistical models. Such differences emphasize the value of combining ML-based forecasting with event-driven simulations when analyzing large-scale phenomena such as disease spread.

Overall, the conclusions of this work reaffirm the benefits of hybrid simulation design. Machine learning provides stability, numerical consistency, and predictive precision, while cellular automata contribute insight into spatial dynamics, emergence, and non-linear system behavior. Together, these perspectives deliver a comprehensive understanding of system evolution under uncertainty, fulfilling the pedagogical and analytical objectives of the project.

6.2 Future Work

While the current system achieved strong performance and produced valuable insights, several opportunities exist for further enhancement and deeper exploration:

- **Integration of additional data sources.** Incorporating mobility data, vaccination rates, demographic indicators, or government intervention indices could significantly improve model accuracy and allow for more context-aware predictions.
- **Adoption of deep learning architectures.** Recurrent neural networks, LSTM/GRU models, or Transformer-based time series models could capture long-term dependencies and complex temporal patterns beyond the capacity of classical ML models.
- **Development of a real-time data ingestion pipeline.** Connecting the system to live APIs would enable continuous retraining and forecasting, transforming the pipeline into a dynamic, real-time decision-support tool.
- **Enhancing interpretability through XAI methods.** Tools such as SHAP, LIME, or attention-based visualizations could increase transparency and foster trust among stakeholders by identifying the most influential features in the forecasting process.
- **Advanced spatial simulation models.** Extending the cellular automaton into multi-state, probabilistic, or agent-based frameworks (ABM) would allow simulation of heterogeneity, mobility, and social behaviors within spatial environments.
- **Continuous deployment and automated monitoring.** Implementing CI/CD pipelines for scheduled retraining, automated validation, and anomaly detection would elevate the system to production-level readiness.
- **Hybrid simulation coupling.** Combining the predictions of the ML system with spatial simulations could produce a unified hybrid model where statistical forecasting informs local interaction rules and vice versa.

These future improvements would strengthen the system's predictive power, expand its applicability to broader epidemiological or socioeconomic contexts, and enhance its value as an academic and practical simulation tool. Ultimately, the project establishes a strong foundation that future iterations can refine into a fully autonomous, intelligent, and explainable forecasting framework.

Chapter 7

Reflection

The project offered a deep understanding of how systems engineering concepts can enhance data science workflows. Throughout the process, collaboration, critical analysis, and iterative refinement were essential. Each team member's contribution in data preprocessing, modeling, validation, and documentation demonstrated the importance of multidisciplinary teamwork.

Applying principles such as modularity, traceability, and feedback control turned a conventional predictive task into a comprehensive system design challenge. This process fostered a mindset focused not only on performance but also on reproducibility and ethical responsibility. Handling epidemiological data required respect for uncertainty, transparency in reporting, and awareness of social implications.

Ultimately, this project reinforced the educational value of integrating analytical and systemic thinking. It demonstrated that engineering principles (when applied to data-driven problems) produce solutions that are not only efficient but also sustainable, adaptive, and ethically grounded.

References

- Arora, P., Kumar, H. and Panigrahi, B. (2020), 'Prediction and analysis of covid-19 positive cases using deep learning models: A comparative study', *Chaos, Solitons Fractals* **139**.
- Chimmula, V. K. and Zhang, L. (2020), 'Time series forecasting of covid-19 transmission in canada using lstm networks', *Chaos, Solitons Fractals* **135**.
- Kermack, W. O. and McKendrick, A. G. (1927), 'A contribution to the mathematical theory of epidemics', *Proceedings of the Royal Society A* **115**(772), 700–721.
- Petropoulos, N. and Makridakis, S. (2020), 'Forecasting the novel coronavirus covid-19', *PLOS ONE* **15**(3), e0231236.
- Rustam, F., Reshi, A., Mehmood, A., Ullah, S., On, B., Aslam, W. and Choi, G. (2020), 'Covid-19 future forecasting using supervised machine learning models', *IEEE Access* **8**, 101489–101499.
- TheRealCyberLord (2020), 'Coronavirus (covid-19) visualization & prediction', Kaggle Notebook. Available at <https://www.kaggle.com/code/therealcyberlord/coronavirus-covid-19-visualization-prediction>.
- Wang, Y., Xu, L. and Schwartz, M. (2021), 'Covid-19 data quality challenges: impact on predictive models', *Journal of Data Science and Analytics* **8**(2), 150–165.