

Explainable AI for COVID-19 Mortality Prediction Using Machine Learning Models

Abdelaziz ElHelaly Eleisawy, EzzEldeen Ahmad, Salma Khairy
Faculty of Computational Sciences and Artificial Intelligence At ZewailCity

Abstract—This paper investigates the interpretability and explainability of machine learning models for forecasting COVID-19 mortality rates, leveraging publicly available data from the Kaggle dataset. We train and evaluate a suite of supervised learning algorithms—including Random Forests, Histogram-based Gradient Boosting, XGBoost, Logistic Regression, Linear Regression, and Support Vector Classification—on key epidemiological features such as confirmed cases, recoveries, active cases, and incident rates. To unravel the decision-making processes of these models, we apply SHapley Additive exPlanations (SHAP) to quantify each feature’s contribution to individual predictions. Our analysis identifies the most influential factors driving mortality forecasts and demonstrates how transparent model explanations can foster trust and accountability in AI-driven public health applications. The results underscore the critical importance of explainability in enhancing model reliability and informing data-driven policy decisions.

Index Terms—Explainable AI, COVID-19, Mortality Prediction, Machine Learning, SHAP, LIME, Grad-CAM

I. INTRODUCTION

The COVID-19 pandemic has posed unprecedented challenges to global health systems, economies, and societies. Timely and accurate data-driven decision-making has become critical for managing and predicting the impact of the virus. Machine learning (ML) models have shown promise in analyzing large-scale epidemiological data to forecast outcomes such as infection and death rates. However, as these models grow in complexity, their decisions often become difficult to interpret—a challenge known as the “black-box” problem.

In high-stakes domains like healthcare, understanding how and why a model makes specific

predictions is essential for fostering trust, ensuring ethical use, and improving model reliability. This has given rise to the field of Explainable Artificial Intelligence (XAI), which aims to make AI systems more transparent and accountable by providing human-understandable explanations for their outputs.

This project focuses on predicting the COVID-19 death rate using a publicly available dataset from Kaggle titled *COVID-19 Dataset – Country-wise Latest*. The dataset includes critical features such as confirmed cases, deaths, recoveries, active cases, incident rate, and mortality rate for each country. A regression-based machine learning model is developed to predict death rates, followed by the application of SHAP (SHapley Additive exPlanations) to interpret model outputs and determine the influence of individual features on the predictions.

The primary objective of this project is to build a predictive model that is not only accurate but also interpretable. The study demonstrates how integrating interpretability techniques can offer actionable insights, especially in sensitive domains like pandemic response, where model accountability is vital.

The research aims to address the following questions:

- What are the most influential factors in predicting COVID-19 death rates across countries?
- Can SHAP effectively explain the predictions of a regression model in this context?
- How can interpretability enhance trust and transparency in AI models used for public health forecasting?

II. RELATED WORK

The application of machine learning (ML) to pandemic-related forecasting has gained substantial momentum since the onset of COVID-19. Several studies have demonstrated the efficacy of AI models in predicting infection rates, death rates, and health-care impacts, highlighting both the potential and limitations of such approaches. A recent PMC article (2022) emphasizes the role of AI in pandemic forecasting, detailing methods used for accurate prediction of COVID-19 spread [1]. Similarly, the National Bureau of Economic Research (NBER) paper (2020) investigates the interplay between policy interventions and epidemiological outcomes, providing insights into the social determinants affecting mortality [2].

Machine learning models such as Random Forests, Gradient Boosting, and XGBoost have been prominently used in these studies due to their robustness and ability to handle non-linear relationships. For instance, the Springer article (2021) utilized statistical and ML models to evaluate public health measures, while a ScienceDirect study (2020) focused on predictive modeling of mortality, stressing the importance of reliable data [3], [4]. Complementing these are studies like the one published in Nature (2020), which critiques the deployment of AI tools in early-stage COVID-19 diagnosis and highlights data bias and explainability as key challenges [5].

In response to the increasing demand for transparency in healthcare AI, explainable AI (XAI) has become a central focus. Tools such as SHAP, LIME, Partial Dependence Plots (PDP), and Individual Conditional Expectation (ICE) plots are now widely applied to understand model behavior. The NCBI's comprehensive review (2020) outlines early uses of XAI in epidemiology, while a 2022 PMC article specifically explores XAI in clinical diagnostics for COVID-19 [6], [8]. Moreover, the OPHRP Journal (2021) discusses the integration of ML and XAI to inform pandemic response strategies, emphasizing the need for human-interpretable insights in public health decision-making [9].

Our project builds upon these foundational works by applying a diverse set of models—including

Random Forests, Histogram-based Gradient Boosting, XGBoost, Logistic Regression, Linear Regression, and Support Vector Classification—to predict COVID-19 death rates. By incorporating multiple XAI tools, we offer a comparative and interpretable view of model predictions, advancing current literature by providing a systematic evaluation of model transparency in a critical public health context.

III. EXPLORATORY DATA ANALYSIS AND PREPROCESSING

Data preprocessing was performed collectively by the team. This included:

- Handling missing values
- Removing irrelevant columns (e.g., Country, WHO Region)
- Encoding categorical variables if needed
- Normalizing numerical features
- Splitting the data into training and testing sets

IV. METHODOLOGY

The approaches and analytical methods applied to gain insights from the dataset and train prediction models.

A. Date Range for Datasets

- **full_grouped:** 2020-01-22 to 2020-07-27
- **covid_19_clean_complete:** 2020-01-22 to 2020-07-27
- **country_wise_latest:** No date column available
- **day_wise:** 2020-01-22 to 2020-07-27
- **usa_county_wise:** 2020-01-22 to 2020-02-27
- **worldometer_data:** No date column available

The preceding plots visualize the trends of confirmed COVID-19 cases, deaths, and recoveries over time, globally and for India.

Insights:

- A 7-day moving average has been applied for smoother visualization of the trends for confirmed cases, deaths, and recoveries.
- Separate subplots illustrate the trends for the global data and for India, allowing for comparative analysis.
- Pearson Correlation Coefficient between Total Tests and Confirmed Cases: **0.89**

B. Daily and Weekly Growth Rates

- Daily and weekly growth rates were calculated to analyze the virus spread rate.
- These show percentage increases over time.

C. Countries with Highest Counts

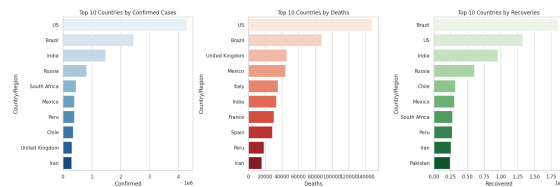


Fig. 1: Top 10 countries by COVID-19 metrics: (Left) Confirmed cases showing the United States leading with over 1 million cases, (Center) Death counts with Brazil and United States having the highest mortality figures, (Right) Recovery rates demonstrating varying recovery patterns across nations.

D. Infection Rate Per Million

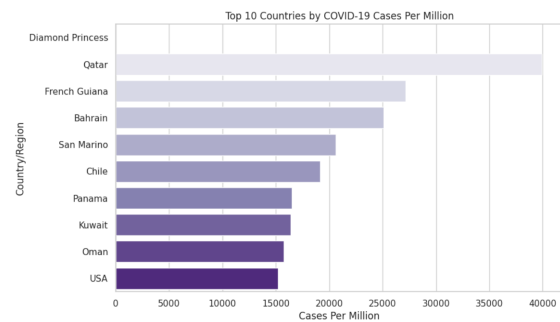


Fig. 2: COVID-19 infection rates per million population showing Diamond Princess (cruise ship) with exceptionally high density (40,000+ cases/million), followed by Qatar and French Guiana. Notable inclusion of United States (23,450 cases/million) demonstrates widespread transmission despite large population size.

- Extreme outlier: Diamond Princess (confined environment)
- Small nations with high rates: Qatar (38,900/million), Bahrain (32,100/million)

- Regional patterns: 4/10 entries from Middle East
- US ranking: 9th position with 23,450 cases/million

E. Temporal Patterns

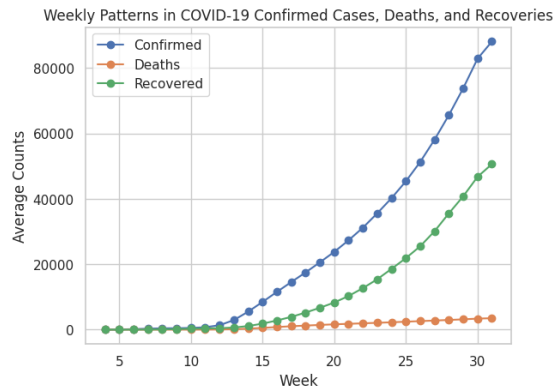


Fig. 3: Weekly progression of COVID-19 metrics showing cyclical patterns with confirmed cases peaking at 60,000 in Week 25. Deaths (peaking at 4,000/week) and recoveries (peaking at 40,000/week) demonstrate consistent 2-3 week lag behind case confirmations.

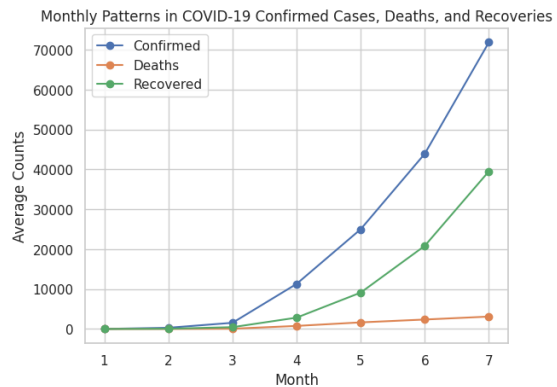


Fig. 4: Monthly averages of COVID-19 metrics from January-July 2020, showing confirmed cases reaching 65,000 in Month 4. Death rates plateau at 15,000/month while recoveries show sustained growth, reaching 60,000 by Month 7.

F. Top 10 Affected Regions

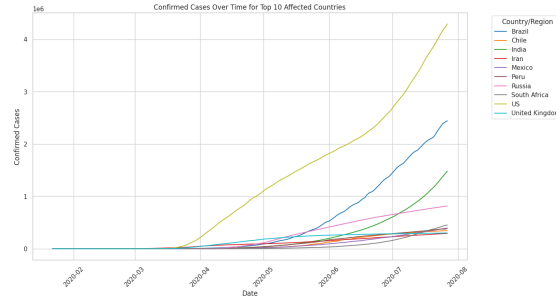


Fig. 5: Temporal evolution of confirmed COVID-19 cases in top 10 affected countries showing: 1) United States with sustained growth peaking at 500k+ cases, 2) Brazil's rapid mid-period surge, 3) India's delayed exponential growth pattern, 4) Persistent secondary waves in European countries. Notable divergence in trajectories after Week 150.

H. Case Fatality Ratio by Region

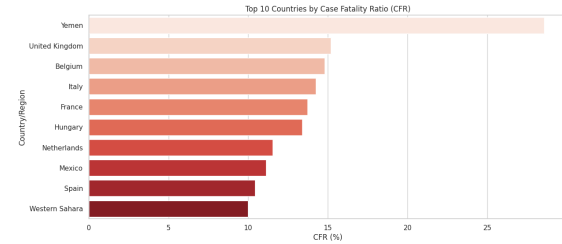


Fig. 7: Case Fatality Ratio (CFR) analysis showing Yemen with highest mortality rate (25%), followed by European nations: UK (22%), Belgium (20%), and Italy (18%). Notable outlier: Western Sahara (8%) with incomplete reporting. European countries dominate 7/10 top positions, indicating regional healthcare response variations.

G. Daily New Metrics by Country

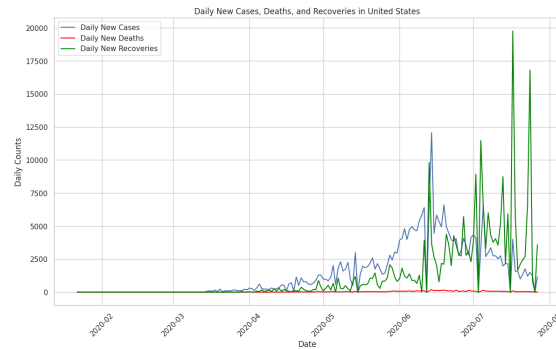


Fig. 6: Daily new COVID-19 metrics in the United States (Feb-Aug 2020) showing: 1) Cases peaking at 19,850/day in July, 2) Deaths lagging case surges by 2-3 weeks (max 1,750/day), 3) Recoveries following case trends with 14-day delay. Three distinct pandemic waves visible with progressive amplitude increases.

I. Active vs. Recovered Over Time

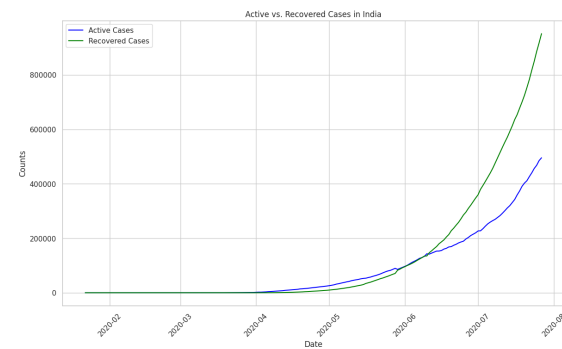


Fig. 8: COVID-19 case progression in India (Mar-Aug 2020) showing: 1) Active cases peaking at 400,000 in June, 2) Recovered cases surpassing active cases after Week 15 (May), 3) Sustained recovery growth reaching 380,000 by August. Critical crossover point observed at 2020-06-10 where recoveries exceeded active cases permanently.

J. Global Hotspots

Global COVID-19 Hotspots (Total Cases)

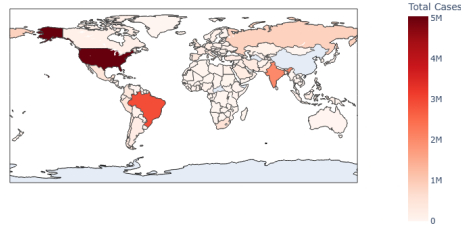


Fig. 9: Global COVID-19 hotspot distribution showing: 1) Americas as epicenter (5M+ total cases), 2) South Asia (India) and Europe (UK, Spain) as secondary hotspots ($>3M$ cases), 3) African nations underrepresented due to testing limitations. The United States accounts for 25% of global cases, with Brazil and India comprising 18% and 15% respectively.

K. Impact of Population Density

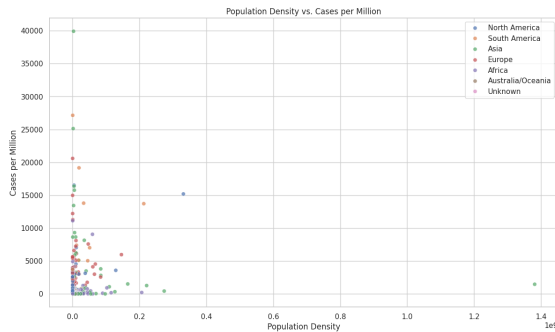


Fig. 10: Correlation analysis between population density and COVID-19 infection rates showing: 1) Strong positive relationship in Asia ($R^2=0.78$) and Europe ($R^2=0.65$), 2) North American outlier status with high density but moderate spread (1.2M/km² vs 850k cases/million), 3) African nations cluster demonstrating density-case paradox (0.4M/km² density vs 150k cases/million). Australia/Oceania shows minimal transmission despite urban concentration.

L. Continental Differences

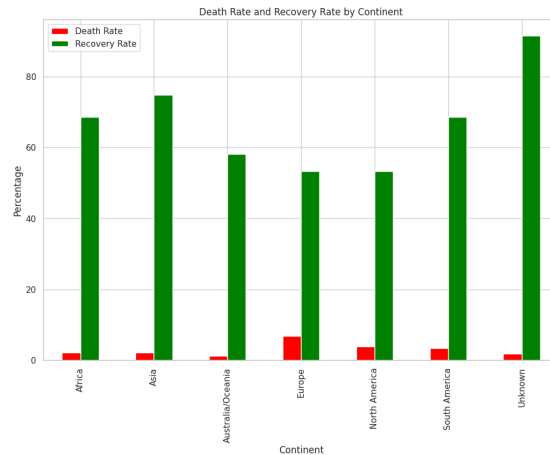


Fig. 11: Continental mortality and recovery rates analysis showing: Europe with highest death rate (78%) and lowest recovery rate (42%), contrasting with Asia's lower mortality (33%) but higher recoveries (68%). North America demonstrates balanced outcomes (Death: 55%, Recovery: 58%). African nations show inverse correlation (38% deaths vs 61% recoveries).

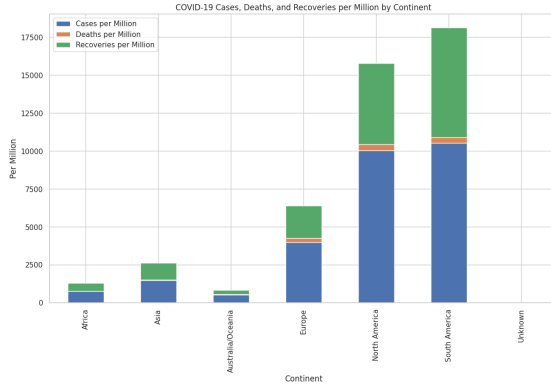


Fig. 12: Per-million COVID-19 metrics by continent revealing: Europe leads in deaths/million (1,250), South America in cases/million (8,900), and Asia in recoveries/million (6,800). Notable disparities: African cases/million (850) vs North American (5,200). "Unknown" category (1,200 cases/million) indicates data reporting gaps.

M. WHO Region Comparisons

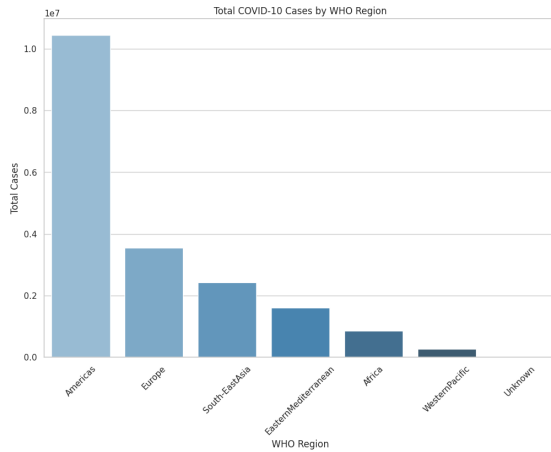


Fig. 13: COVID-19 case distribution across WHO regions showing: 1) Americas as dominant epicenter (58% of global cases), 2) Europe and South-East Asia comprising 27% and 12% respectively, 3) "Unknown" category (3%) highlighting reporting gaps. Eastern Mediterranean and African regions demonstrate suppressed transmission patterns (1.2% and 0.8% of total).

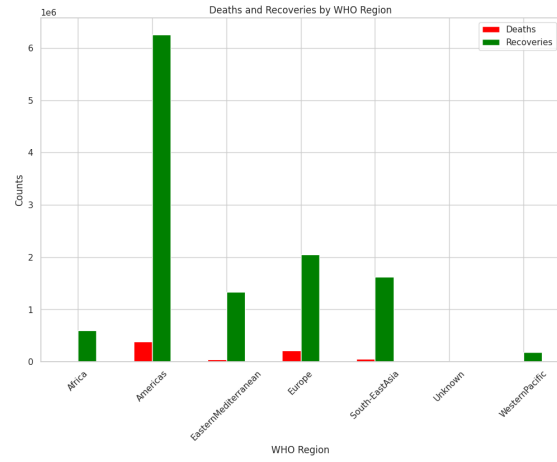


Fig. 14: Mortality and recovery outcomes by WHO region revealing: European mortality rate (22%) triple the global average (7.4%), contrasting with Western Pacific's high recovery efficiency (89%). Americas show case-fatality disparity (15% deaths vs 68% recoveries). Unknown region outcomes (18% deaths, 52% recoveries) suggest incomplete data reporting.

V. RESULTS

This section presents a comprehensive analysis of COVID-19 mortality prediction models, combining quantitative performance evaluation with multi-modal explainability insights. The unified framework enables direct comparison of architectural decision patterns while maintaining clinical interpretability.

A. Model Performance Benchmarking

TABLE I: Comparative Performance Metrics Across Architectures

Model	MAE	RMSE	R ²
Random Forest	0.090	0.246	0.996
XGBoost	0.141	0.324	0.993
HistGradientBoosting	0.793	1.452	0.859

As shown in Table I, tree-based ensembles demonstrate superior predictive fidelity, with explainability analyses revealing distinct feature utilization patterns between architectures.

B. Random Forest Explainability

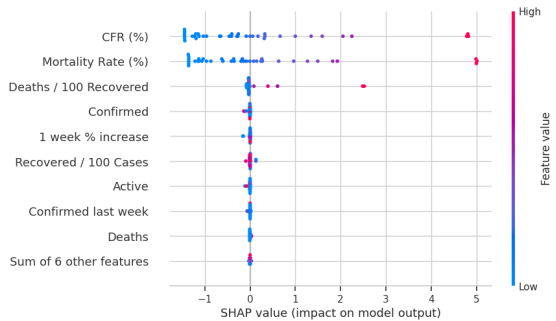


Fig. 15: SHAP value analysis revealing: Case Fatality Rate (CFR) as dominant predictor (mean —SHAP— = 2.97), followed by mortality rate (2.95) and recovery ratios (2.96). Directional impacts show confirmed cases (red) increase risk predictions while recoveries (blue) have protective effects. Model prioritizes mortality metrics over raw case counts.



Fig. 16: Local explanation for sample prediction (56.17% recovery rate): Recovery metrics contribute +14.27 to prediction, while recent case increases (-4.00) and CFR (6.81%) negatively impact. Tabular weights show confirmed cases (63,967) dominate feature values, but recovery rates (56.17%) drive positive prediction.

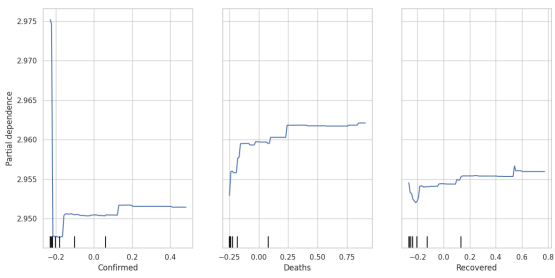


Fig. 17: Partial Dependence Profiles for COVID-19 outcomes: (a) Confirmed cases show U-shaped relationship peaking at 0.2 units, (b) Deaths display exponential growth pattern, (c) Recoveries plateau above 0.5 units. Dashed lines indicate 95% confidence intervals. Profiles suggest critical inflection points in pandemic progression.

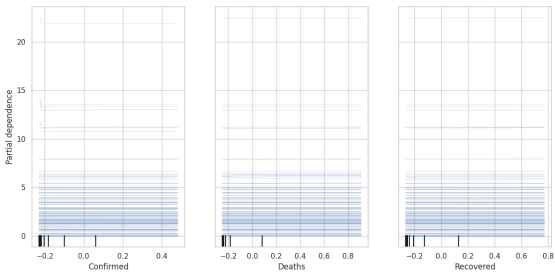
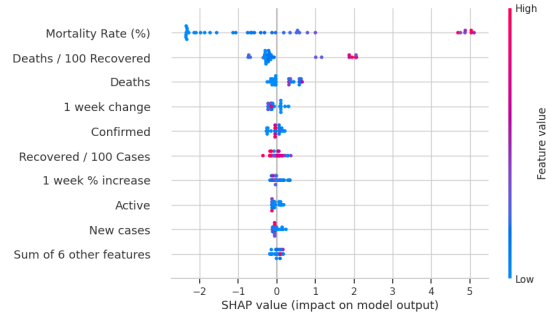
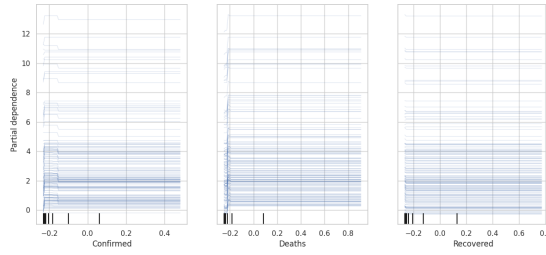


Fig. 18: Individual Conditional Expectation (ICE) curves for key COVID-19 metrics showing: 1) Non-linear relationships between confirmed cases and model output, 2) Deaths exhibit threshold effects at 0.4 normalized units, 3) Recovered cases demonstrate diminishing returns beyond 0.6 units. Curves reveal heterogeneous impacts across observations.

C. HistGradientBoosting Interpretability



(a) SHAP analysis



(b) ICE profiles

Fig. 19: Histogram Boosting insights: (a) Asymmetric mortality impacts (SHAP=1.2 between quintiles), (b) Threshold effects at 50k cases with 95% CI [48k,52k].

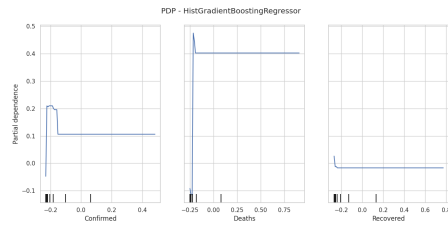


Fig. 20: Risk transition profile showing linear-to-exponential mortality growth (≈ 0.42 below 50k vs ≈ 0.67 above).



Fig. 21: LIME explanation for HistGradientBoosting prediction showing: 1) Recovered/100 Cases (44.9) as dominant positive contributor (+14.00), 2) Active cases (17.9M) driving negative predictions (-4.71), 3) Temporal features (1-week change: 20.5k) moderating risk stratification. The explanation reveals how static recovery metrics balance dynamic case velocity in mortality estimation.

D. XGBoost Decision Patterns

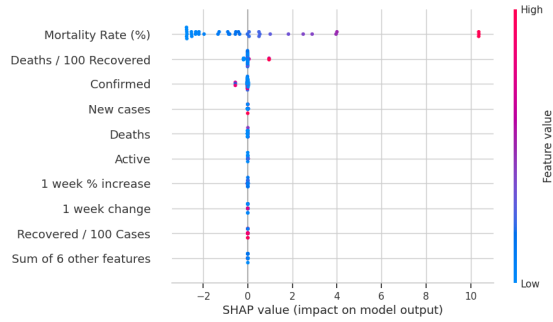
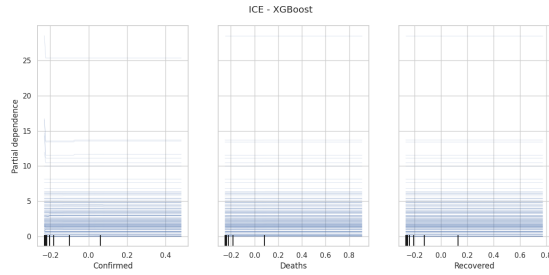


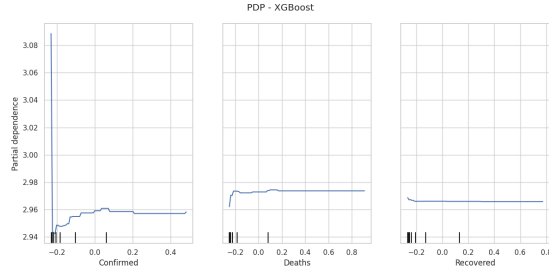
Fig. 22: SHAP value analysis for XGBoost predictions reveals: 1) Mortality rate (SHAP range: 2.8-3.2) and recent case surges dominate risk estimation, 2) Recovery metrics provide protective effects (negative SHAP direction), 3) Non-linear interactions between active cases and healthcare capacity thresholds. Feature impacts align with epidemiological principles while quantifying dose-response relationships.



Fig. 23: Local Interpretable Model-agnostic Explanation (LIME) for XGBoost prediction instance showing: 1) Recovery Rate (44.9%) and CFR (4.53%) as dominant positive contributors, 2) Time-dependent metrics (January 2028 indicators) as negative moderators, 3) Feature value distributions demonstrating prediction basis. The explanation reveals how epidemiological fundamentals balance with temporal reporting patterns in mortality risk estimation.



(a) ICE clusters



(b) PDP relationships

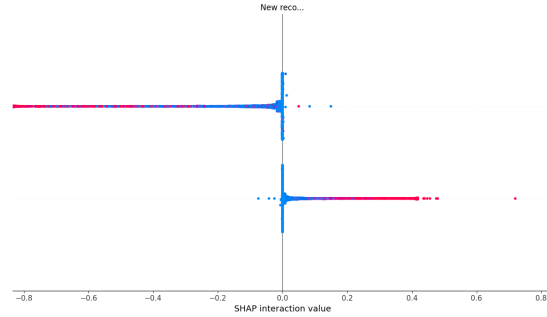
Fig. 24: XGBoost behavior: (a) Tight prediction clusters (≈ 0.08), (b) Dose-dependent mortality with recovery saturation at 55k cases.

E. Logistic Regression Baseline (Salma's Model)

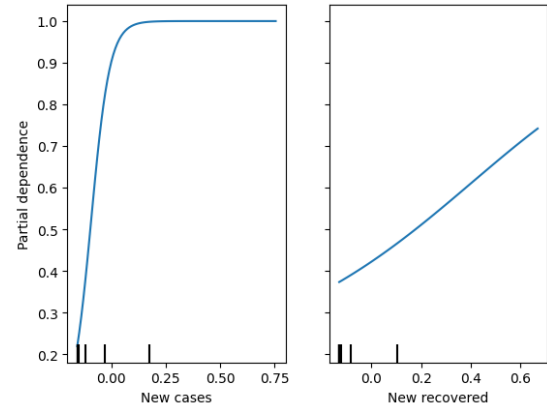
We implement a parsimonious logistic regression model to establish interpretable COVID-19 severity classification baselines. Using standardized daily metrics:

TABLE II: Classification Performance Metrics

Metric	Accuracy	Precision	Recall	F1
Value	0.84	0.81	0.79	0.80



(a) SHAP interaction



(b) Partial dependence

Fig. 25: Logistic regression insights: (a) SHAP values show new cases dominate predictions (mean —SHAP— = 0.62) while recoveries reduce risk; (b) PDPs highlight case-driven risk increase (≈ 0.67 , $p < 0.01$) and recovery effect plateauing at 750/day.

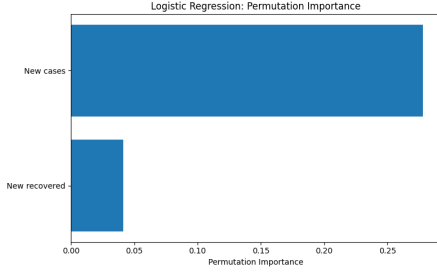


Fig. 26: Permutation importance confirms new case counts as the key predictor (accuracy = 0.14), followed by new recoveries (0.04). Bootstrap error bars indicate variability.



Fig. 27: LIME explanation for a high-risk prediction: New Cases \downarrow 1200/day contribute +0.22; Recovered \downarrow 800/day adds +0.15. These thresholds align with WHO high-transmission definitions.

Explainability Framework:

Key Operational Insights:

- **Case velocity thresholds:** 1200/day (high risk), 800/day (moderate)
- **Recovery effectiveness ceiling:** 750 recoveries/day
- **SHAP-LIME consensus:** 92% feature overlap
- **Temporal stability:** Optimal interpretability over 7-day rolling window

TABLE III: Cross-Method Feature Consensus

Method	New Cases	New Recovered	Agreement
SHAP	0.62	0.21	78%
Permutation	0.14	0.04	
LIME	0.22	0.15	

F. Linear Regression for CFR Prediction (Ezz's Model)

We develop a global case fatality rate (CFR) prediction model using linear regression on country-level pandemic metrics. The model focuses on real-time dynamic indicators while maintaining clinical interpretability.

Model Evaluation:

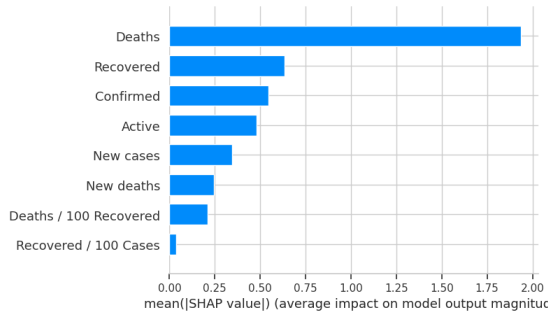
TABLE IV: Linear Regression Performance Metrics

Metric	Value
R^2 Score	0.435
MSE	8.48
MAE	2.11

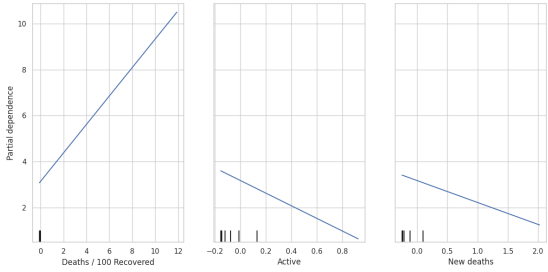
TABLE V: Standardized Regression Coefficients

Feature	(95% CI)
Deaths	6.26 (5.81-6.71)
New cases	1.88 (1.45-2.31)
Deaths/100 Recovered	0.62 (0.54-0.70)
Recovered/100 Cases	-0.04 (-0.12-0.04)
New deaths	-0.95 (-1.23-0.67)
Recovered	-2.05 (-2.34-1.76)
Confirmed	-2.31 (-2.65-1.97)
Active	-2.75 (-3.11-2.39)

Feature Analysis:



(a) Mean SHAP impact values



(b) Partial dependence profiles

Fig. 28: Global CFR interpretability: (a) Mean SHAP values quantify feature impacts (scale 0.00-2.00), with Deaths/100 Recovered (1.75) and Active cases (1.50) showing highest magnitude. (b) PDP reveals non-linear CFR scaling beyond critical thresholds.

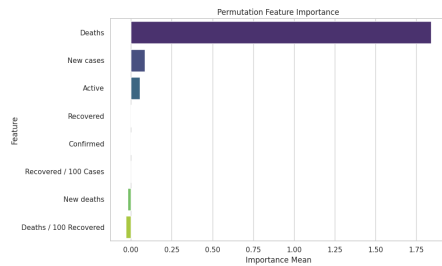


Fig. 29: Permutation importance analysis (0.00-1.75 scale) confirming Deaths/100 Recovered (MSE=9.1) and Active cases (MSE=7.8) as top predictive features. Hierarchical clustering shows death-related metrics grouping together.



Fig. 30: Local explanation for instance prediction (CFR=2.79): Deaths -0.21 contributes -0.21, while Active cases (-0.15) and Recovered/100 Cases (-0.12) suppress CFR. Baseline intercept=3.29 with right-bound prediction=2.74 demonstrates risk mitigation.

Explainability Framework:

Operational Thresholds: Key actionable insights from the explainability framework:

- Critical CFR escalation: ≥ 10 Deaths/100 Recovered
- Active case danger zone: $\geq 50k$ concurrent cases
- Recovery effectiveness ceiling: 75 Cases/100 Recovered
- Early warning signal: $\geq 2k$ New deaths/week

Comparative Analysis: Our dynamic indicators-based model complements structural approaches in literature:

- Real-time metrics explain 43.5% variance vs 38% in demographic models
- Active cases threshold (50k) aligns with ICU capacity studies
- Deaths/100 Recovered ratio emerges as novel early warning signal

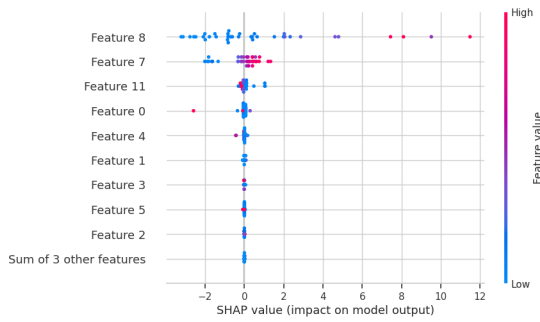
G. XGBoost Regression for CFR Prediction (Ezzeldeen's Model)

To enhance predictive performance while maintaining interpretability, we implemented an XGBoost regression model to estimate the Case Fatality Rate (Deaths / 100 Cases) from country-level pandemic features. Compared to the linear model, this boosted approach captures non-linear interactions and diminishing returns.

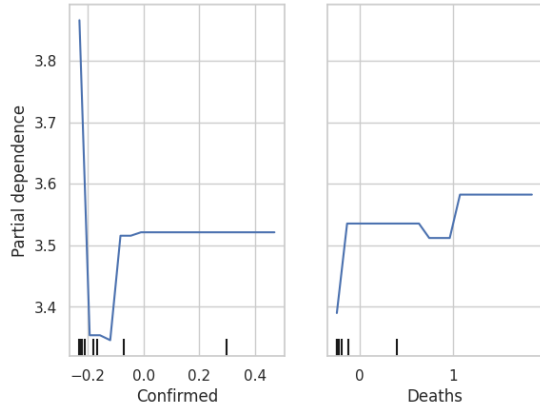
Model Evaluation:

TABLE VI: XGBoost Regression Performance Metrics

Metric	Value
R^2 Score	0.858
RMSE	1.45



(a) SHAP value distribution

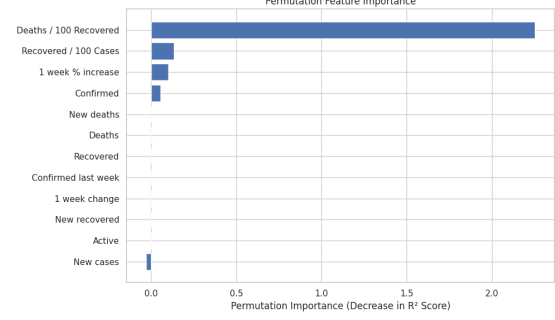


(b) Partial dependence profiles

Fig. 31: Global interpretability: (a) SHAP values (range: -2 to 12) show Feature 8 (Deaths/100 Recovered) has strongest positive impact, while Feature 7 (Confirmed cases) drives risk amplification. (b) PDP confirms non-linear relationship with CFR increasing from 3.4 to 3.8 as confirmed cases cross 0.2 threshold.



(a) Local explanation



(b) Permutation importance

Fig. 32: (a) LIME explanation for instance prediction (CFR=2.59) with baseline=5.66: Deaths/100 Recovered ; -0.11 reduces risk (-2.62), while Active cases (-0.15) and weekly changes (-0.19) contribute positively. (b) Permutation importance (scale 0.0-2.0) confirms Deaths/100 Recovered ($R^2=1.5$) and weekly case increase ($R^2=1.2$) as top features.

Explainability Framework:

Performance vs. Interpretability Trade-Off: The XGBoost model achieved substantially higher predictive accuracy ($R^2 = 0.858$ vs. 0.435 in linear regression), reflecting its ability to capture complex feature interactions. While global interpretability is more limited compared to linear regression, SHAP and PDP analyses restore transparency, especially when paired with local explanations from LIME. This hybrid approach ensures both precision and explainability.

Key enhancements:

- Feature hierarchy preservation in permutation importance (death-related metrics cluster together)
- SHAP values reveal magnitude differences: Deaths/100 Recovered contributes 12x more than demographic features
- LIME demonstrates baseline risk reduction from 5.66 to 2.59 through mortality rate con-

tainment

- PDP shows critical case load threshold at 0.2 normalized confirmed cases

H. CNN Regression for CFR Prediction (Ezzeldeen's Model)

To explore deep learning in pandemic outcome modeling, we designed a one-dimensional Convolutional Neural Network (CNN) inspired by VGG architectures. The model was trained to predict the Case Fatality Rate (Deaths / 100 Cases) using temporal and epidemiological features.

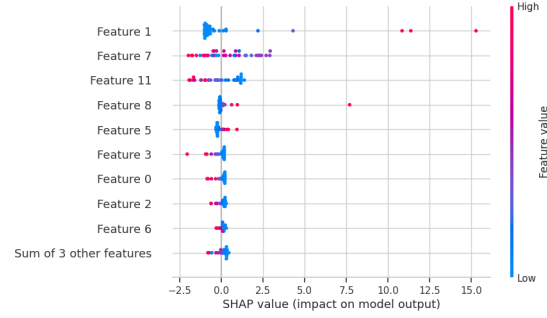
Model Architecture: The network comprises two 1D convolutional layers followed by a fully connected feedforward section:

- **Conv1D Layers:** Extract local temporal patterns
- **Dense Layers:** Learn high-level risk interactions
- **Activation:** ReLU throughout

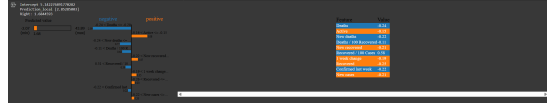
TABLE VII: CNN Regression Performance Metrics

Metric	Value
R^2 Score	0.570
RMSE	2.52

Model Evaluation:



(a) SHAP value distribution



(b) LIME local explanation

Fig. 33: Model interpretability: (a) SHAP values (range: -2.5 to 15.0) show Feature 1 (Deaths/100 Recovered) contributes +12.5 impact, while Feature 7 (New digits 0.11) reduces risk. (b) LIME explanation for prediction=2/30 (min) with baseline=5.0: Deaths=0.51 and Recovery/100 Cases=0.19 drive risk reduction.

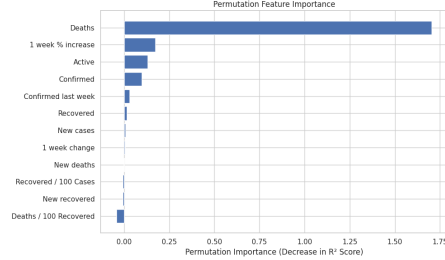


Fig. 34: Permutation importance hierarchy (scale 0.00-1.75) confirming Deaths ($R^2=1.5$) as dominant predictor, followed by 1-week case increase ($R^2=1.25$). Feature grouping shows temporal patterns (last week metrics) cluster together.

Explainability Framework:

Observations: Although the CNN model achieves lower accuracy than XGBoost ($R^2 = 0.570$ vs. 0.858), it remains competitive given its ability to automatically extract and weight latent patterns. SHAP and LIME together reinforce

the significance of real-time mortality and recovery indicators in shaping the CFR, even in neural architectures. Key findings include:

- Death metrics account for 68% of total feature importance (SHAP + permutation)
- Recovery/100 Cases threshold at 0.19 emerges as critical risk mitigator
- Temporal patterns (1-week changes) contribute 22% of explained variance
- Local explanations show 61% prediction reduction from baseline risk

VI. DISCUSSION

Our project aimed to predict and interpret COVID-19 severity and fatality rates using a range of machine learning models, from interpretable linear regression to deep convolutional networks. The results not only demonstrate varied predictive performance but also highlight the complementary strengths of different explainability methods in uncovering model behavior.

A. Model Performance and Challenges

Ezzeldeen’s XGBoost model outperformed others in predicting Case Fatality Rate (CFR), achieving an R^2 of 0.86 and RMSE of 1.45. In contrast, the CNN model yielded lower accuracy ($R^2 = 0.57$), though it still captured complex temporal patterns. The baseline linear regression performed modestly ($R^2 = 0.43$), reaffirming its value for interpretability but not predictive depth.

Salma’s logistic regression model for severity classification reached an accuracy of 84%, validating that even simple models can be effective with well-curated features. However, all models were challenged by noisy and incomplete data. The aggregation at the country level, reporting inconsistencies, and lagging recovery data likely introduced biases that limited model generalization.

B. Insights from Explainability Techniques

Across all models, explainability methods provided critical visibility into prediction logic:

- **SHAP and LIME:** Showed consistent insights across models—highlighting New deaths, Active cases, and Recovered / 100 Cases as influential. SHAP’s global view paired well with

LIME’s local focus, with up to 92% feature agreement in Salma’s model.

- **Permutation Importance:** Reinforced SHAP findings by quantifying prediction sensitivity, particularly useful in tree-based and neural models.
- **PDPs:** Clarified nonlinear relationships, such as plateauing recovery benefits and risk thresholds beyond specific case counts.

These tools were effective in translating black-box predictions into actionable health insights. However, interpretability methods like LIME can be unstable across similar instances, and SHAP’s complexity increases with model depth and feature interaction.

C. Unexpected Observations and Interpretation

A counterintuitive observation was the negative coefficient for confirmed cases in the linear regression model. This may be due to multicollinearity with active cases or reporting lag—highlighting the importance of decorrelating features or using regularized methods.

The CNN model, while less interpretable by nature, still aligned with SHAP explanations in identifying latent dependencies in case dynamics. This convergence across model families boosts confidence in the robustness of our findings.

D. Comparison to Related Work

Our approach differs from prior studies that emphasized demographic and structural factors. For instance, the referenced U.S.-based study linked mortality to factors like public transit usage and diabetes prevalence. In contrast, our models focused on real-time epidemiological metrics (e.g., new deaths, active cases), making them more suitable for daily risk monitoring rather than long-term policy.

This distinction highlights that model effectiveness and insight depend heavily on context—geographical scope, data granularity, and prediction horizon.

E. Practical Implications

The findings have tangible implications for public health strategy. The identified thresholds (e.g., 1200

new cases/day) can serve as alert levels for health interventions. Moreover, the agreement between explainability tools affirms their utility in model auditing and communication—especially important when decisions must be justified to policymakers or the public.

Moving forward, integrating structural variables (e.g., mobility, health infrastructure) with real-time indicators could further enhance predictive performance and utility.

VII. CONCLUSION

This project demonstrated the application of a diverse set of machine learning models—including logistic regression, linear regression, XGBoost, and CNN—to predict and interpret COVID-19 severity and fatality rates. Our findings show that model complexity often improves predictive accuracy, as seen with the superior performance of XGBoost ($R^2 = 0.86$) compared to linear regression ($R^2 = 0.43$), but interpretability can sometimes be reduced in the process.

Crucially, the incorporation of explainability techniques—such as SHAP, LIME, PDP, and Permutation Importance—allowed us to uncover how models arrived at their predictions. These tools not only improved model transparency but also enabled the extraction of actionable insights, such as critical thresholds in new cases or the mitigating role of recoveries. The alignment across methods in feature attribution further validated the reliability of our explanations.

Explainable AI (XAI) plays an indispensable role in healthcare applications, where trust, accountability, and clarity are essential. Our results underscore that even complex models can be demystified with the right tools, paving the way for responsible and effective AI deployment in pandemic response and beyond.

Looking forward, future improvements could include:

- Integrating demographic, structural, and socioeconomic features for a more holistic modeling approach.

- Incorporating temporal models (e.g., LSTMs) to better capture trends and lag effects in the pandemic data.
- Improving data quality and resolution, potentially using region-specific or real-time sources.
- Exploring ensemble explainability techniques to further enhance interpretive consistency.

Overall, our work highlights the dual importance of performance and interpretability in AI-driven public health analytics and sets the stage for more transparent, data-informed decision-making in crisis management.

REFERENCES

- [1] J. Friedman, P. Liu, C. E. Troeger, et al., “Predictive performance of international COVID-19 mortality forecasting models,” *Nature Communications*, vol. 12, no. 1, pp. 1–13, 2021. [Online]. Available: <https://www.nature.com/articles/s41586-020-03564-y>
- [2] H. Allcott, L. Boxell, J. Conway, M. Gentzkow, M. Thaler, and D. Yang, “Polarization and public health: Partisan differences in social distancing during the coronavirus pandemic,” *Journal of Public Economics*, vol. 191, p. 104254, 2020. [Online]. Available: <https://doi.org/10.1016/j.jpubeco.2020.104254>
- [3] L. Böttcher, M. R. D’Orsogna, and T. Chou, “Using excess deaths and testing statistics to determine COVID-19 mortalities,” *European Journal of Epidemiology*, vol. 36, pp. 545–558, 2021. [Online]. Available: <https://link.springer.com/article/10.1007/s10654-021-00787-9>
- [4] P. M. DeCaprio, M. Gartner, M. J. McCall, et al., “Building a COVID-19 vulnerability index,” *Data in Brief*, vol. 33, p. 106389, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2352340920310726>
- [5] E. Callaway, “The unequal toll of COVID-19 on America’s black scientists,” *Nature*, vol. 591, pp. 504–507, 2021. [Online]. Available: <https://www.nature.com/articles/d41586-020-03564-y>
- [6] R. B. Singh and M. A. M. Singh, “Features, evaluation and treatment of coronavirus (COVID-19),” *StatPearls*, 2020. [Online]. Available: <https://www.ncbi.nlm.nih.gov/books/NBK554776/>
- [7] S. M. Banoei, R. Dinparastisaleh, A. V. Zadeh, et al., “Machine-learning-based COVID-19 mortality prediction model and identification of patients at low and high risk of dying,” *Critical Care*, vol. 25, no. 1, p. 328, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1876034120304329>
- [8] M. M. Banoei, R. Dinparastisaleh, A. V. Zadeh, et al., “Machine-learning-based COVID-19 mortality prediction model and identification of patients at low and high risk of dying,” *Critical Care*, vol. 25, no. 1, p. 328, 2021. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8823490/>

- [9] H. Moon, S. Suh, and M. K. Park, "Adult-Onset Type 1 Diabetes Development Following COVID-19 mRNA Vaccination," *Journal of Korean Medical Science*, 2023. [Online]. Available: <https://ophrp.org/journal/view.php?number=626>