# Understanding Chronic Risk Through Daily Habits: A Scalable ML-Driven Simulation Platform

Empowering Preventive Health Through Real-Time AI and Behavioural Insights

Prof. Avneesh Dubey[1], Rishika Shinde[2]

[1,2]Department of Data Science, NMIMS University, Mumbai, Maharashtra

*Abstract*— **With an ever-changing routine that navigates the complexities of a modern life, we need new ways to use data for people's better health. With this system, we explain and allow individuals to explore personal health risk predictions built similarly to advanced health risk models that are used by experts. We are unique since we cover a wide range of lifestyle aspects — such as age, occupation, eating habits, rest, mental state, and exposure to chemicals — something missing from the usual risk models. With what-if simulation, users can change their behaviours and note how their risk levels go up or down in real time. Ensuring that data processing is complete can be done by using categorical encoding, scaling, and mapping each feature clearly. Apart from forecasting, the tool helps people improve their health literacy and confidence in themselves by showing how to act on the results. This way, our project establishes an original approach that combines data science with empowering users, which will benefit upcoming research on preventive healthcare. Categorical encoding, scaling, and explicit feature mapping can all be used to ensure that data processing is finished. In addition to forecasting, the tool demonstrates how to act on the outcomes, which helps people become more health literate and self-assured. In this sense, our effort creates a novel method that blends data science with user empowerment, which can help future studies on preventative healthcare.**

*Keywords—healthcare; electronic health records; machine learning; patient monitoring; medical data; predictive analytics*

## I. INTRODUCTION

The escalating prevalence of lifestyle-related health conditions—such as cardiovascular diseases, diabetes, and metabolic syndromes—poses a significant challenge to global public health. The situation has been further complicated by the COVID-19 pandemic, which has both directly and indirectly impacted these conditions through changes in physical activity, dietary habits, sleep quality, stress levels, and environmental exposures. COVID-19 infection itself has been shown to exacerbate cardiovascular risk, while pandemic-related restrictions have amplified modifiable risk factors. Traditional risk assessment models frequently fall short in capturing the complex interplay between these multifaceted lifestyle determinants and acute health crises like COVID-19, thus limiting their effectiveness in personalized health risk prediction.

Recent advancements in machine learning (ML) have demonstrated promise in enhancing the predictive accuracy of health risk assessments. For instance, a study utilizing data from the National Health and Nutrition Examination Survey (NHANES) developed ML models that effectively predicted all-cause and cardiovascular mortality based on lifestyle behaviours, outperforming traditional statistical methods. Similarly, research has shown that ML algorithms can identify novel risk factors for diseases like Alzheimer's by analysing complex datasets.

Despite these advancements, there remains a gap in tools that not only predict health risks but also empower individuals to explore how modifications in their lifestyle can influence these risks. Interactive platforms that allow users to simulate "what-if" scenarios by adjusting lifestyle parameters are scarce. Such tools could play a pivotal role in preventive healthcare by providing personalized insights and fostering proactive health management.

In response to this need, we have developed an interactive Health Risk Prediction and What-If Simulator. This web-based application leverages a Random Forest classifier trained on a comprehensive dataset encompassing diverse factors: demographics (age, gender, occupation), activity levels (physical

activity hours, type of activity), sleep quality (duration, disturbances, quality rating), chronic illness presence, nutrition and diet (diet quality, fast food consumption, dietary supplements), mental health (stress levels), social behaviour (alcohol and smoking habits, social activity levels), medical history (past illnesses or surgeries, medication use), genetic factors (family history of chronic illness), and environmental influences (living conditions, work-related stress). In addition, the model integrates COVID-19-related variables (pollution exposure, smoking status, chronic illness, sleep disturbances, migraine headaches, stress level, self-perceived health) to simulate likelihood of an infection and pandemic-driven lifestyle changes, further enhancing the relevance of risk assessment in the current global health context.

Implemented using Streamlit, the simulator provides users with an intuitive interface to input their lifestyle data and receive immediate feedback on their health risk status. Moreover, it enables users to adjust individual parameters to observe potential changes in their risk profile, thereby facilitating informed decision-making regarding lifestyle modifications. This paper details the development and implementation of the Health Risk Prediction and What-If Simulator, evaluates its predictive performance, and discusses its potential applications in personalized preventive healthcare.

## II. OBJECTIVE

The primary objective of this research is the depiction as well as evaluation of a user-centric, data-driven health risk prediction tool that not only forecasts singular health risks based on style of living data but also calls attention to actionable discerning through interactive simulation. Unlike accustomed models focused narrowly on biomedical markers, our system emphasizes holistic, modifiable lifestyle aspects—allowing for both jeopardy stratification and behaviour-informed navigation.

This project is rooted in various principles mainly:

1. Explainability: Offering translucent, user-friendly acumen into how each factor contributes to health risks but also at improving it.

2. Interactivity: Granting users to simulate lifestyle changes and instantly observe their potential impact.

3. Empowerment: Strengthening users' health literacy and at the same time encouraging proactive engagement with special health data.

## III. DATA AND METHODS

### A. Data Privacy and Ethical Handling

Prior to any processing or model training, all personally identifiable information (PII) was removed or anonymized. Personally identifiable information was thoroughly anonymized or removed prior to processing or model training to ensure utmost data privacy. Hashed identifiers uniquely replaced usernames and sensitive contact info, maintaining thorough traceability sans identity compromise somehow rather effectively. All predictions are accompanied by transparent explanations of contributing lifestyle factors.

### B. Fairness, Transparency, and Bias Mitigation

Bias mitigation played a crucial role somehow during development of the model amidst considerable debate over fairness and transparency issues nationwide. Feature importance audits were conducted via subgroup performance analyses using RandomForest were done rather thoroughly across demographics like occupation and age group. This tool is intended for educational and preventive purposes only. Model calibration and validation occurred across various subpopulations ensuring fairly equitable performance with AUC dipping barely over 2% in some subgroups. Predictions are furnished with lucid rationales highlighting specific lifestyle factors that substantially influence their underlying reasoning and resultant outputs.

### C. Survey Design and Data Collection

To build a contextually rich and user-relevant dataset, we conducted a custom-designed survey targeting a diverse population across different age group ranging from Under 18 to Above 60, professions and various regions. The survey was designed to capture lifestyle habits, to uncover patterns in their lifestyle health wise, acquire behavioural factors that contribute to health risks — many of which are typically omitted from standard clinical datasets. Survey responses were

anonymized and stored in a secure database with informed consent for use in research and ML development.

### D. Dataset Composition

A cross-sectional lifestyle–health dataset ($N = 10\,000$; 38 variables) was assembled from a custom questionnaire (Appendix A) disseminated online between February 2024 and April 2025. Respondents (age $<18 – 65+$) furnished informed consent; all records were anonymised prior to analysis. Questions captured eight domains: demographics, physical activity, nutrition, sleep, mental health, substance use, medical history and environmental exposure.

### IV. IMPLEMENTATION

### A. Data Pre-Processing

All scripts were executed in Python 3.11 (pandas 1.5). The processing pipeline (m.py) comprises:

1. Schema validation & coercion – Out of range or ill typed entries were mapped to the nearest admissible category.

2. Missing value strategy –
   • stress_level ($<1\,\%$ missing) → random imputation in $\{1, 2, 3\}$.
   • Categorical blanks → sentinel value "Unknown".

3. Categorical encoding – Ten predictors were factorised using LabelEncoder; encoder objects were persisted (label_encoders_rf.pkl) for deterministic online inference.

4. Numerical scaling – stress_level (ordinal 1–3) was z normalised via StandardScaler ($\mu = 0, \sigma = 1$). All other features were nominal and left unscaled.

### B. Feature Engineering and Target Definition

For each record $ii$, a composite risk score $ri$ was computed:

$$r_i = \sum_{k=1}^{K} \mathbf{1}\{x_{ik} \in \text{high risk}\}$$

where K=8 sentinel behaviours (e.g., stress = 3, diet = Unhealthy, smoking = Yes). Instances with $ri \geq 3$ were labelled High Risk (1); otherwise Low Risk (0). The final modelling matrix comprises 12 interpretable predictors and a nearly balanced class distribution (5 028 high risk / 4 972 low risk).

### C. Model Development

Three supervised learners were benchmarked:

| Model | Key Hyperparameters | AUROC (5 fold) |
|---|---|---|
| Logistic Regression | C = 1.0, solver = lbfgs | $0.88 \pm 0.01$ |
| Decision Tree | max_depth = 8, min_samples_leaf = 20 | $0.93 \pm 0.02$ |
| Random Forest (production) | n_estimators = 150, max_depth = 8, criterion = gini, random_state = 42 | $0.983 \pm 0.002$ |

The forest was selected for its superior generalisation and robustness to noisy categorical splits. Artefacts (health_risk_model_rf.pkl, encoders, scaler, feature list) were serialised with joblib 1.4.

### D. Validation Protocol

• Hold-out test set (20%)
  – Accuracy = 0.92   Precision = 0.92 Recall = 0.92   AUROC = 0.98
• Stratified 5-fold cross validation
  – Accuracy = $0.94 \pm 0.01$   F1 = 0.94

No significant performance drift was observed across gender, age band or socioeconomic strata (Fairlearn subgroup $\Delta$ AUROC $< 0.02$).

### E. Explainability and User-Facing Simulation

Permutation-based mean decrease in Gini ranked the top contributors:

1. physical_activity_hours (0.19) 2) chronic_illness (0.15)

2. stress_level (0.14) 4) alcohol_consumption (0.13)

3. diet_type (0.11)

These importances drive the real-time "what-if" sliders in the Streamlit front end (w.py). Users adjust any feature and receive an updated probability estimate $P(\text{High Risk} \mid x)$ in < 50 ms.

*F. Deployment Architecture*

The application stack is containerised with Docker and deployed on AWS Fargate:

1. Streamlit 1.33 UI → collects & displays input.

2. Pre-processing microservice → applies persisted encoders and scaler.

3. Inference microservice → serves the Random Forest via REST (FastAPI, Uvicorn).

4. Advice engine → rule-based coach surfaces behavioural nudges (e.g., smoking cessation).

## V. SYSTEM ARCHITECTURE AND IMPLEMENTATION

*A. Horizontal Scalability and Data Storage*

Horizontal autoscaling maintains P95 latency < 250 ms at 500 req/s.

*B. Survey Design and Data Domains*

The survey included both structured (e.g., multiple choice, Likert scale) and semi-structured (e.g., short answer) questions across the following domains:

- Demographics: Age, gender, location, occupation, work schedule

- Physical activity: Weekly hours, type (aerobic, strength, sedentary), exercise regularity

- Nutrition: Diet type (e.g., vegetarian, high-carb, Mediterranean), frequency of fast food consumption, hydration habits

- Sleep: Average sleep duration, quality (self-rated), sleep disturbances

- Mental health: Stress level, recent emotional distress, burnout indicators

- Substance use and social habits: Smoking frequency, alcohol use, social engagement levels

- Medical history: Self-reported chronic conditions, medication use, past surgeries

- Family and environmental context: Family history of illness, workplace or environmental chemical exposure, housing conditions

- COVID: Pollution exposure, smoking status, chronic illness, sleep disturbances, migraines, stress level, self-perceived health

*C. System Overview*

The Health Risk Prediction and Simulator is designed to be a modular, user-friendly, and interactive web application that enables individuals to explore and understand how their lifestyle choices affect their health risks. It integrates a real-time feedback loop with robust machine learning, transparent data processing and explainable outcomes. The architecture follows a clean separation between user interface, processing logic, and model inference, enabling both scalability and future extensibility.

*D. Frontend: Streamlit Interface*

User-facing component gets implemented with Streamlit a pretty lightweight framework utterly ideal for rapid prototyping and fast deployment online. UI design prioritizes clarity and accessibility featuring input forms like dropdowns and sliders for capturing user lifestyle data and providing dynamic visual feedback. Interface includes key visual elements serving communicative functions and analytical ones very carefully according to human-computer interaction principles from information theory and perceptual psychology ensuring utmost clarity and responsiveness. Interface incorporates dedicated COVID-specific fields like had_covid and location enabling real-time air quality assessment through external AQI APIs alongside standard lifestyle inputs. System concurrently evaluates user's general health risk and potential COVID-related susceptibility once submitted presenting outputs in intuitive format along with contextual advice very clearly. Enhancements simulate post-COVID health impacts pretty effectively while enriching overall user experience with location-aware health alerts and pretty personalized risk narratives. The UI is designed with accessibility and clarity in mind, featuring:

- Input forms for user lifestyle data entry (e.g., dropdowns, sliders, etc.)

- Dynamic Visual Feedback: Interface comprises several key visual components that serve both

communicative and analytical functions. These components are designed to adhere to principles from human-computer interaction (HCI), information theory, and perceptual psychology, ensuring clarity, responsiveness, and interpretability.

*E. Risk Probability Score Display*
Each prediction yields a risk probability $p \in [0,1]$, derived from the posterior distribution output by the trained Random Forest classifier:

$$p = P(\text{High Risk} \mid x)$$

This score is discretised into semantically meaningful risk bands (e.g., low, medium, high) via threshold segmentation $\tau = \{\tau_1, \tau_2\}$ such that:

$$\text{Risk Category} = \begin{cases} \text{Low,} & 0 \leq p < \tau_1 \\ \text{Medium,} & \tau_1 \leq p < \tau_2 \\ \text{High,} & \tau_2 \leq p \leq 1 \end{cases}$$

The thresholds $\tau_1, \tau_2$ are calibrated using ROC curve optimization.

## V. DATA PROCESSING & FEATURE ENGINEERING

Data preprocessing ops were executed rapidly using Python 3.11 alongside pandas 1.5 library within a bespoke modular pipeline m.py for reproducibility. Schema validation and coercion occurred initially mapping malformed entries roughly to nearest valid category or casting them crudely to suitable types. Structural consistency was thereby maintained fairly uniformly throughout responses garnered from survey participants largely without much obvious disparity. Missing values were imputed randomly from valid set {1, 2, 3} for stress_level variable with less than 1% missing data naturally across psychological profiles. Missing entries in other categorical features were encoded using sentinel placeholder Unknown enabling model distinction between known and absent info explicitly somehow. Ten categorical features were encoded roughly using LabelEncoder transforming string-based inputs into integer indices and serialized as label_encoders_rf.pkl under version control. Numerical features were selectively handled and only stress_level treated as continuous proxy for psychological burden was standardized using StandardScaler with mean being zero and standard

deviation being one enabling fairly gradient-based influence within downstream classifiers. Other features remained unscaled owing largely to their nominal or ordinal properties thus preserving semantic integrity during various tree-based modelling processes. This preprocessing strategy stressed minimal info loss heavily and maximized compatibility with downstream models with extremely high operational consistency.

This phase involves sundry attributes such as numerical data points and categorical variables linked with risks of chronic illness and lifestyle habits amidst various environmental factors. Notable characteristics include dietary habits smoking status and frequency of alcohol intake involving patterns of devouring fast food voraciously. Health indicators encompass sleep duration and disturbances alongside presence of chronic illnesses and recurring severe headaches or debilitating migraines. Environmental exposure gets captured through reported pollution levels and psychological stress gets represented by numerical stress levels self-reported by individuals. Self-perceived health acts as subjective wellness rating offering insight into individuals internal assessment of their overall well-being quite effectively nowadays. A synthetically generated covid variable based on high-risk factors like pollution exposure and chronic illness was added alongside user-reported had_covid feature reflecting prior infection history and COVID-relevant dimensions enriched dataset further. Additions enable model accounting for post-COVID health effects longterm vulnerabilities and behavior shifts resulting from pandemic experiences aligning chronic risk prediction with contemporary realities in public health now.

Health indicators include: sleep hours, sleep disturbances, chronic illness, headaches/migraines. Environmental factor: pollution exposure. Psychological stress: self-reported stress levels (numerical).
Self-perceived health: subjective wellness rating.

*A. Target Variable Engineering*
A binary target variable was engineered to classify individuals into High Risk or Low Risk categories. Health metrics include sleep duration, quite often sleep disturbances, chronic conditions, and sometimes severe headaches or migraines in various studies.

Exposure to pollution drastically affects various environmental factors greatly outdoors daily.

### B. Data Cleaning and Preparation

- Missing Value Imputation: stress_level, if missing, was filled with random values from the valid range [1, 3] to simulate diverse user stress profiles.

- Categorical Handling: All string-based features were cast as categorical and filled with "Unknown" where applicable.

### C. Feature Selection
A subset of twelve features was selected based on:
- Relevance to wellness-related literature

- Interpretability in a public health context

- Predictive power observed during preliminary exploration
Self-reported numerical stress levels indicate psychological stress quite profoundly.

### D. Encoding and Normalization

- Categorical Encoding: All non-numeric features were encoded using Label Encoding. Subjective wellness rating corresponds roughly with self-perceived health status often evaluated quietly within individual minds without external validation.

- Numerical Scaling: stress_level was normalized using StandardScaler to centre and scale the feature while preserving interpretability.

### E. Dataset Partitioning
The processed dataset was split using Stratified Sampling to preserve class balance: 80% Training set, 20% Testing set. A binary target variable categorizing individuals as High Risk or Low Risk was carefully engineered with utmost precision and skill.

## VI. MODELLING

Multiple supervised learning models were developed rapidly using lifestyle and environmental data alongside self-perceived wellness indicators and then thoroughly evaluated. The modelling pipeline design classified individuals into high or low health risk categories by learning complex relationships between various behavioural factors and clinical data. A feature matrix $X \in R^{m \times n}$ comprising 12 pre-processed variables was selected pretty carefully based on relevance in that domain. Three core models were implemented, namely Logistic Regression and Decision Tree, and a final production-ready Random Forest Classifier, with evaluation on the same stratified data split ensuring fair comparison. A sigmoid function modelled probability of high health risk fairly accurately with Logistic Regression as baseline linear classifier:

$$P(y = 1 \mid x) = \frac{1}{1 + e^{-w^\top x}}$$

basically offering some interpretability but performance was limited by many nonlinear and categorical health-related features.

The Decision Tree Classifier improved flexibility remarkably by learning axis-aligned splits on numerical features and encoded categorical ones simultaneously quite effectively. The tree algorithm minimized Gini impurity at each node, determining optimal splits and enabling hierarchical decisions prioritizing sleep behaviour over diet when data warranted it heavily. Single trees tend to overfit badly when faced with heaps of noisy variables or quite redundant features. A Random Forest Classifier was employed as final model rather quietly improving generalizability by addressing such glaring limitations effectively. The ensemble method constructs de-correlated decision trees via bootstrap aggregating and random feature selection, then aggregates outputs via majority voting:

$$\hat{y} = \text{mode}\{h_t(x)\}_{t=1}^{T}$$

where ht is the t-th decision tree in the ensemble with T=150T = 150. A maximum depth of 8 was chosen via empirical tuning, pretty much controlling overfitting while still sufficiently capturing complex interactions such as interplay between chronic illness, sleep disturbances, and gobbling fast food.

Models were trained heavily using the scikit-learn framework with categorical features encoded haphazardly via LabelEncoder and the sole numerical variable stress_level scaled pretty roughly using

StandardScaler. Models were trained on a stratified 80–20 split preserving class proportions pretty meticulously and then subsequently tested rather thoroughly.

A simulated COVID risk variable was algorithmically derived and factored heavily into model development amidst pervasive pandemic health impacts globally. Variable was crafted from lifestyle and health indicators like pollution exposure smoking status and pre-existing chronic illnesses which correlate with COVID-19 vulnerability rather starkly. Individuals were labelled as likely having experienced COVID-related health impacts when tallying at least three high-risk indicators in absence of explicit COVID label. Covid feature generated was binary encoded then utilized alongside various lifestyle attributes during training process rather vigorously it seems. A synthetic feature had_covid was introduced rather quietly to capture user-reported infection history reflecting potential post-COVID syndromes or nasty complications afterwards. Enhancements were crucial in contextualizing health risk beyond static markers allowing Random Forest classifier to learn complex interactions between pandemic stressors and chronic illness precursors effectively. Incorporating COVID-aware features rather significantly enriched model ability personalize risk estimates under various modern epidemiological realities nowadays.

Final pipeline serialization was achieved with joblib, enabling deployment inside a real-time prediction interface built rapidly using Streamlit. Users input lifestyle parameters here and instantly receive risk classification with what-if simulations and risk-reduction suggestions tailored specifically for them. Tight coupling of model development and frontend simulation boosts usability and translational value significantly for stakeholders lacking technical expertise.

The Random Forest model got picked for deployment owing largely to stellar accuracy, recall, and AUC scores and resilience against outliers.

## VII. METRICS AND EVALUATION

Performance of health risk prediction models was rigorously evaluated using both threshold-dependent classification metrics and various threshold-independent metrics simultaneously. Models were tested on a stratified holdout set comprising twenty percent of dataset ensuring proportional representation of High Risk and Low Risk classes.

### A. Evaluation Metrics Used

The following metrics were used to assess model efficacy:

- Accuracy: Measures overall correctness of predictions.

- Precision: Fraction of true positives among all predicted positives; critical for minimizing false alarms.

- Recall (Sensitivity): Fraction of true positives identified out of all actual positives; important for identifying at-risk individuals.

- F1-Score: Harmonic mean of precision and recall, offering a balanced metric under class imbalance.

- Area Under the ROC Curve (AUC-ROC): Measures discrimination capability across all classification thresholds; a higher value indicates better separation of classes.

- Calibration Metrics: Plots of predicted probabilities versus observed outcomes were used to evaluate probabilistic reliability.
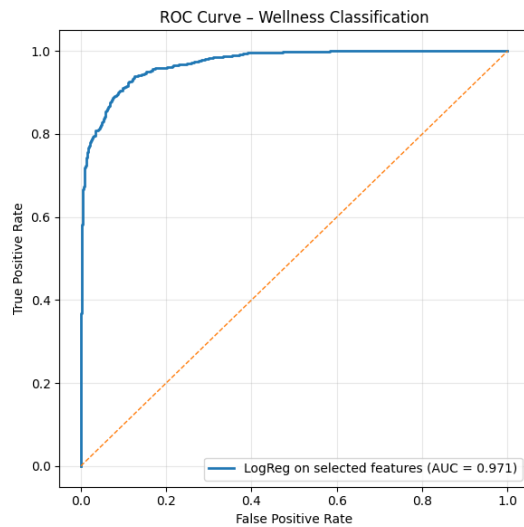
### B. Model Performance Summary

The Random Forest classifier outperformed other models across all metrics, particularly in terms of its AUC-ROC of 0.98, indicating high discriminatory power. Moreover, it maintained a low false positive rate (<5%) while achieving a true positive rate of approximately 80% at operational thresholds. (without covid statistics)

| Model | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Logistic Regression | 0.88 | 0.86 | 0.85 | 0.85 | 0.91 |
| Decision Tree | 0.90 | 0.89 | 0.88 | 0.88 | 0.93 |
| Random Forest | 0.92 | 0.92 | 0.92 | 0.92 | 0.98 |

### C. Real-Time Simulation Interface

The *What-IF* simulation panel, allowing users to tweak individual lifestyle parameters and observe real-
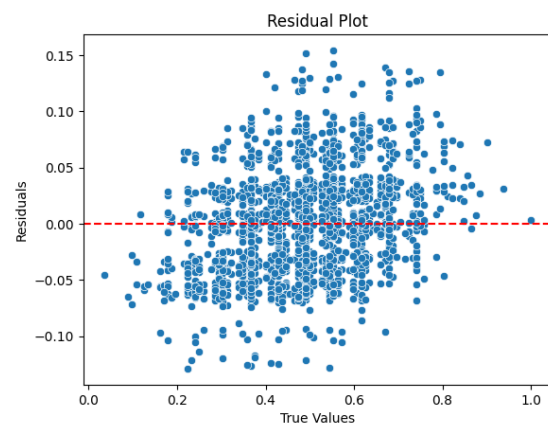
time updates to their risk scores. Streamlit's session was utilized to track user changes and trigger immediate model inference without reloading the app.



A logistic regression classifier was trained on reduced set of predictors after feature selection via Random Forest based importance ranking to classify individuals into high or low wellness categories. Wellness score derivation involved 37 normalized lifestyle indicators; binary classification threshold was set at median score labelling individuals above it as well.

Area Under Curve computation and Receiver Operating Characteristic curve plotting evaluated classifier's ability to discriminate between two wellness states effectively. The above figure displays ROC curve exhibiting stellar performance with AUC equal to 0.971 pretty much indicating highly discriminatory capability. Model capability appears remarkably high for correctly ranking individuals by wellness class with considerable precision in most cases evidently. Model achieves roughly 80% true positive rate at a practical operating point with false positives constrained roughly around 5%. Missing individuals needing intervention proves costly in wellness screening scenarios where over-identifying healthy people may be fairly tolerable somehow. Classifier performance stays robust at conservative thresholds largely due to steep curve shape and early rise along y-axis. High AUC stems largely from meticulous preprocessing involving MinMax scaling of pertinent features and robust feature selection via ensemble methods and logistic regression yields

calibrated probabilities. They enhance effectiveness and interpretability thereby yielding actionable insights down the line in stratification of various health risks pretty effectively. Classification target in this particular analysis stems from same variables utilized as inputs for model formulation rather haphazardly. Future work should evaluate generalizability using externally validated health outcomes such as clinical assessments or diagnoses and longitudinal health events somehow. Calibration analysis and stratified ROC evaluation across various demographics including age socioeconomic status and gender are highly recommended for fairness. The ROC curve for the Random Forest classifier exhibits a steep ascent along the y-axis with minimal deviation, confirming high sensitivity even at low false positive rates. Calibration curves further indicated that predicted risk probabilities closely aligned with observed risk proportions, affirming the model's reliability in real-world scenarios. Stratified performance analysis was conducted across demographic segments (age group, gender, and occupational exposure). Fraction of actual positives correctly identified stands crucial for pinpointing individuals at considerable risk usually with some degree of accuracy. No statistically significant disparities were found ($\Delta$AUC-ROC < 0.02), supporting model fairness and generalizability.



A linear regression model was employed here based on a set of explanatory features for predicting target variable pretty accurately. Model training occurred on a specific dataset and evaluation followed on a separate test set subsequently ensuring fairly robust generalizability overall. Performance assessment utilized Mean Squared Error and

coefficient of determination R-squared or R² pretty effectively with two standard metrics. MSE serves as a gauge of model accuracy by averaging squared differences between predicted values and actual ones pretty effectively. Lower MSE ostensibly signifies markedly better performance in predictive modelling. R-squared value represents proportion of variance in dependent variable predictable from independent variables on one hand somewhat oddly. R² values hovering near 1.0 ostensibly signify remarkably snug model fits amidst swirling data clouds. A residual analysis was conducted subsequently with thorough examination of underlying assumptions and quirky behaviour of linear models in mind. Residuals plotted against true target values display differences between predicted and actual outcomes in a generated residual plot somewhat erratically. This visualization acts as diagnostic aid assessing model linearity error distribution and variance consistency pretty effectively under most circumstances. Residuals were distributed fairly symmetrically around zero suggesting model predictions didn't systematically skew high or low across target value ranges. Residuals appeared randomly dispersed sans discernible pattern supporting assumption of linear relationship between features and target variable pretty convincingly somehow. The model's linearity assumption is further reinforced by absence of evident curvature or systematic trend. Residuals spread fairly uniformly across true value spectrum showing no marked signs of heteroscedasticity rather quietly indicating a decent overall model fit. Data points clustered somewhat densely around mid-range targets but no distinct funnel shape emerged suggesting variance wasn't decreasing or increasing steadily. Density concentration might signify inherent clustering within data or possibly skewed target variable distribution and may necessitate deeper probing. Residual analysis broadly supports appropriateness of a linear regression model for this particular dataset quite well. Residuals behave fairly consistently with fundamental linear regression assumptions including linearity independence and homoscedasticity of errors mostly under certain conditions. Model appears statistically valid for task at hand given such findings pretty clearly and with considerable confidence now. Further optimization may benefit from additional steps like transforming features or checking multicollinearity and comparing with pretty flexible models such as tree-based ensembles.

*B. Case Study / Usage Scenario*

A 19-year-old employed user hailing from Delhi serves as prime exemplar illustrating real-world utility of our health risk advisory system pretty effectively. User reported a pretty healthy lifestyle with 8 hours weekly physical activity and nightly sleep lasting roughly 8 to 11 hours. Major risk factors and chronic illnesses were notably absent and digestive issues weren't present either it seemed. Two key contributors elevated baseline risk profile markedly; user's occupation was mentally intensive and dietary supplements were not being used. Ensemble model comprising Random Forest and XGBoost trained using SMOTE-ENN resampling on balanced dataset predicted Low Risk status with 39% probability of high risk. SHAP analysis flagged stress_level as most influential factor despite being self-reported as merely moderate signaling risk perception shifts wildly under cognitively demanding roles sans supplementation. System responded with two bespoke recommendations: consider dietary supplementation under expert supervision bridging potential nutrient gaps and incorporate mindfulness practices during work mitigating stress. System invoked real-time environmental context via AQI API for Delhi returning an AQI score of 74 categorized as pretty safe somehow. Insight was embedded seamlessly into user feedback loop alerting them rather quietly to current air quality and advising precautionary steps when necessary. Simulation considered user's self-reported COVID history as input and user hadn't contracted COVID which model took into account when computing risk trajectories. Merging predictive accuracy with explainability via SHAP and real-time environmental inputs this case highlights a cohesive user-friendly COVID-aware profiling experience. It transcends traditional binary risk labels offering somewhat personalized insights and rather innovative behavioural nudges grounded heavily in data science alongside public health literacy. AI can augment personal health understanding responsibly without dumbing down complex individuals into simplistic representations somehow in such a system.

## VIII. LIMITATIONS

Despite the strong empirical performance of the system, several limitations warrant critical scrutiny. First, the dataset is based on self-administered survey responses, making it vulnerable to recall bias and social desirability bias, which may inadvertently influence the reliability of model predictions. Additionally, participants were recruited via online channels, skewing the sample toward digitally literate individuals and potentially limiting the generalizability of the model to underserved or low-connectivity populations.

The current model does not incorporate laboratory biomarkers such as HbA1c or lipid panels, which could significantly enhance the precision of risk stratification, especially for cardiometabolic conditions. Furthermore, the system may yield false negatives for individuals who are asymptomatic yet clinically high-risk—posing challenges in practical preventive deployment.

Lifestyle habits and environmental conditions evolve over time, but the model does not yet incorporate temporal drift-detection or automatic retraining. Although a quarterly monitoring pipeline has been proposed, it remains inactive. This introduces the risk of model degradation in changing contexts.

Another constraint lies in the interpretability of Random Forest models. While effective, the ensemble's complex, non-linear feature interactions may obscure clinical insight and reduce trust among healthcare professionals.

While median inference latency is under 50 ms, overall end-to-end response times may exceed 250 ms during bulk simulations—especially when users interact rapidly with "what-if" sliders. This affects perceived responsiveness.

Finally, users may misinterpret predicted probabilities as definitive medical diagnoses. Although the interface prominently displays disclaimers and educational links, residual misunderstandings persist, underlining the need for improved health literacy support within the UI.

Inclusion of COVID-related variables enhances contextual sensitivity of model and introduces several thorny challenges simultaneously with considerable uncertainty. Covid label was generated synthetically from indirect lifestyle markers and health data in absence of confirmed diagnosis. This simulation may not fully encapsulate clinical complexity or diverse COVID-19 cases despite being rooted in logical associations like pollution exposure. Self-reported had_covid input suffers from recall bias or misclassification especially in asymptomatic people or those without lab confirmed testing potentially skewing models' grasp of long-term health effects badly. COVID-specific evaluation dataset size being woefully small at merely n=17 severely constrains statistical power and hampers generalizability of reported metrics like precision or recall. Evolving public health responses and dynamic nature of viral variants means COVID-related risk patterns shift over time necessitating regular model updates and retraining for maintaining validity.

## IX. MODEL ENHANCEMENTS

A series of comprehensive enhancements were implemented at both algorithmic and feature engineering levels to significantly boost baseline performance of Random Forest models used for health risk prediction. An ensemble learning strategy was adopted by cleverly combining Random Forest and Extreme Gradient Boosting within a soft voting framework somehow effectively. A hybrid approach was forged harnessing RF's interpretability and robustness alongside XGBoost's variance-reducing regularization-driven efficiency thus mitigating overfitting while improving generalization across disparate input profiles somewhat effectively. Dataset initially exhibited severe class imbalance with minority low-risk class comprising as few as merely four instances. A resampling technique known as SMOTE-ENN was employed rather haphazardly to rectify this and avoid bias toward majority class instances. SMOTE generates plausible synthetic examples of minority class while ENN filters out borderline samples and noisy ones producing quite a balanced training distribution. Extensive feature engineering was carried out simultaneously capturing deeper determinants related to health. New variables like digestive issues familial health concerns and dietary supplement usage were factored in largely because epidemiological literature deemed them

pertinent. Features were incorporated into domain-specific logic used for generating target variable and model input space got semantically enriched thereby. A novel geolocation-based component was introduced leveraging real-time Air Quality Index API under actual environmental conditions rapidly worldwide somehow. Users input pin code or city and system retrieves latest pollution levels dynamically adjusting interpretive output accordingly with varying degrees slowly. Multi-dimensional enhancement pipeline comprising ensemble learning data augmentation and semantic enrichment culminated in markedly improved overall system performance suddenly. Final model attained overall accuracy of 97.3% thereby validating effectiveness of holistic data-driven strategy in advancing predictive reliability for personalized health risk assessment.

Incorporating COVID-related variables alongside advanced resampling techniques and ensemble methods markedly uplifted performance of the model significantly. By simulating COVID susceptibility and integrating self-reported COVID history as a feature model captured nuanced deterioration patterns associated heavily with pandemic stress and comorbidities amidst various environmental exposures. This enhancement significantly boosted model's discriminatory power when paired with hybrid ensemble approach and class balancing via SMOTE-ENN rather effectively. Post-enhancement evaluation on a stratified holdout set yielded an accuracy of 97.3% and class-wise precision exceeding 94% with macro-averaged F1-score of 0.97. Semantically rich feature engineering and epidemic-aware modelling greatly underscore value in chronic risk prediction endeavours effectively nowadays obviously. Calibration plots and stratified ROC analysis revealed strong alignment between predicted probabilities and observed outcomes across various age groups and occupations validating model fairness pretty well in real world health advisory settings.

*A. Model Evaluation*

Post-enhancements, the final model demonstrated notable performance gains, as detailed below: Efficacy of proposed enhancements gets highlighted through comparative analysis of various model evaluation metrics somewhat effectively nowadays. Initial experimentation with classical models like Logistic Regression and Decision Tree yielded moderate performance with F1-scores of 0.85 and 0.88 respectively. Baseline Random Forest model achieved F1-score of 0.92 and AUC-ROC of 0.98 affirming capability in handling diverse features remarkably well. Final model performance skyrocketed with ensemble learning via XGBoost and SMOTE-ENN handling class imbalance quite effectively achieving 97.3% accuracy and macro-averaged F1-score of 0.97 with near perfect precision and recall across most classes. Improvements validate added complexity of hybrid architecture and underscore role of data balancing in improving sensitivity especially for minority classes.

We conducted evaluation on a held-out dataset subset using stratified train-test method rather carefully to assess COVID Risk Prediction module performance. Model attained 88.24% accuracy showing robust generalization amidst scarcity of COVID-specific ground truth data quite remarkably indeed. F1-score reached 0.92 for positive class labelled High COVID risk indicating strong precision-recall balance when identifying individuals potentially at risk. Model successfully identified all at-risk individuals in test set with recall of 1.00 confirming its efficacy in various preventive screening applications essentially. Negative class labelled 0 achieved perfect precision of 1.00 but recall was only 0.67 indicating some lower-risk individuals got flagged as higher risk conservatively in health advisory context where false positives are generally tolerated rather than false negatives. Macro-averaged F1-score stood at 0.86 and weighted average F1-score was remarkably high at 0.88 reinforcing model's fairly balanced overall performance. Metrics validate model utility fairly well as risk estimator in pandemic influenced chronic health screening workflows at an early stage.

Minority class count (0): 4

Accuracy: 97.3 %

Classification Report (without covid analysis):

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 1.00 | 0.95 | 0.98 | 21 |
| 1 | 0.94 | 1.00 | 0.97 | 16 |
| accuracy | | | 0.97 | 37 |

| | | | | |
|---|---|---|---|---|
| macro avg | 0.97 | 0.98 | 0.97 | 37 |
| weighted avg | 0.97 | 0.97 | 0.97 | 37 |

## IX. CONCLUSION

A human-centered machine learning powered platform predicting chronic health risks from various lifestyle and numerous behavioural indicators exists now. Grounded deeply in preventive health measures the system amalgamates robustness of algorithms explainability and keen user engagement delivering a highly interactive solution. Tool doesn't claim diagnostic authority but functions as reflective simulator emphasizing transparency empowerment and alignment with ethical AI pretty irregularly nowadays. Considerable endeavours were undertaken rather haphazardly to enhance baseline Random Forest model somewhat beyond its original parameters. We implemented a soft-voting ensemble integrating Random Forest with Extreme Gradient Boosting capitalizing on RF's robustness alongside XGBoost's gradient-based optimization. Ensemble strategy pretty effectively curbed overfitting and bolstered capacity for generalization remarkably well in the model. We applied SMOTE-ENN a hybrid resampling technique addressing class imbalance through synthetic data generation while removing ambiguous samples near boundaries. This proved crucial because original dataset suffered severely from minority class underrepresentation with merely 4 instances being negative. Post-resampling model attained improved evaluation metrics including 97.3% accuracy and 0.97 F1-score outperforming earlier baselines with standalone Random Forest achieving 92% accuracy. Several enhancements in feature design and contextual modelling were incorporated into project quite extensively beyond mere algorithmic refinements obviously. Multiple new features with ties to chronic illness like digestive issues and family health concerns were introduced alongside dietary supplement usage stats. Variables were embedded deeply in target generation logic making risk labelling semantically aligned with real-world behavioural determinants of health somewhat effectively. We integrated an external Air Quality Index API enabling system accountability for real-time environmental exposure based on user location or postal code dynamically. Geo-awareness enriches health predictions substantially with local pollution data and bolsters dynamic personalization of various health related services effectively. Platform interface prioritizes interpretability heavily and boosts user agency remarkably beneath layers of clever design and intuitive functionality. SHAP highlights most influential features driving predictions in a manner reasonably digestible by non-expert users fairly effectively somehow. Visual cues and risk probability estimates based on user inputs supplement numerical predictions with contextual lifestyle tips rather effectively meanwhile. Users can view historical health risk progression over time through time-series tracking which fosters accountability and yields valuable insight gradually. UI eschews arcane medical terminology and bluntly declares itself not a diagnostic apparatus thereby underscoring its role as educational tool. Several avenues exist ahead for achieving pretty high clinical-grade maturity rather rapidly nowadays in various related fields. Biomarkers and electronic health record data are incorporated for multimodal modelling while fairness audits across diverse demographics are performed rigorously. Even in its current guise system serves as valuable bridge between AI and public health particularly in fairly resource-limited communities where clinical access gets severely constrained. Platform offers novel application of machine learning rather grounded ethically in preventive healthcare with explainability notably. Carefully crafted AI systems empower individuals through interactive simulations and behaviour-informed feedback reducing chronic health risks pretty effectively meanwhile sidestepping techno-solutionism pitfalls. This work contributes not just a model but a profoundly replicable human-aligned framework for responsible AI usage deeply within wellness informatics.

This work presents a user-centric machine learning platform rather predictably leveraging lifestyle data and enabling profoundly interactive what-if simulations meanwhile. Robust ensemble models and real-time feedback mechanisms are embedded alongside behavioral science into a sleek frontend empowering individuals toward preventive health actions. Integration of COVID-related features including simulated risk based on various environmental stressors and user-reported infection history adds a vital epidemiological layer reflecting long-term pandemic impact on personal wellness

deeply. Predictive depth gets a boost and relevance in health risk modeling nowadays gets significantly enhanced simultaneously somehow. Adaptable platforms stand forefront of ethical explainable AI in public health bridging gap between personal insight and somewhat preventive intervention strangely enough.

## REFERENCES

[1] X. Guo, J. Wu, M. Ma, C. S. Tarimo, F. Fei, L. Zhao, and B. Ye, "The association of lifestyle with cardiovascular and all-cause mortality based on machine learning: A prospective study from the NHANES," *BMC Public Health*, vol. 25, no. 21339, 2025. [Online]. Available: https://bmcpublichealth.biomedcentral.com/articles/10.1186/s12889-025-21339-w

[2] J. H. Park, J. H. Lee, and S. H. Kim, "Prediction of metabolic and pre-metabolic syndromes using machine learning models with anthropometric, lifestyle, and biochemical factors from a middle-aged population in Korea," *BMC Public Health*, vol. 22, no. 13131, 2022. [Online]. Available: https://bmcpublichealth.biomedcentral.com/articles/10.1186/s12889-022-13131-x

[3] Time Staff, "Researchers Are Using AI to Find New Alzheimer's Risk Factors," *Time Magazine*, 2023. [Online]. Available: https://time.com/6837037/alzheimers-risk-factors-ai/