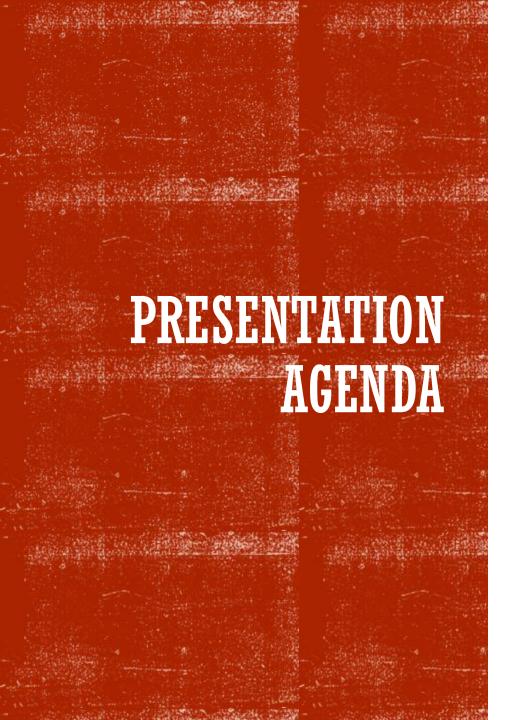


PREDICTION MODELS IN CYBERSECURITY AND HEALTHCARE

Assignment Project Presentation By Nikunj Gupta



- Overview of Prediction Models
 - Libraries used
 - Healthcare Prediction Models
 - Cybersecurity Prediction Models
 - Algorithm Analysis and Comparison
 - Input/Output Specifications
 - Implementation and Training Process
 - Model Performance and Results
 - Conclusions and Future Work

LIBRARIES USED













Streamlit:
An open-source
Python framework that
allows data scientists
and AI/ML engineers
to create interactive
web applications and
dashboards in just a
few lines of code
without requiring
front-end experience.

Pandas:
A Python library used for working with data sets, providing functions for analyzing, cleaning, exploring, and manipulating data to make it readable and relevant for data science applications.

Scikit-learn:
A machine learning library in Python that provides algorithms for applications like spam detection and image recognition, including gradient boosting, nearest neighbors, random forest, and logistic regression.

Seaborn:
A Python data
visualization library
that creates statistical
plots, such as line
plots with the ability to
show relationships
between variables
across different
subsets of data using
semantic groupings.

Matplotlib:
A comprehensive plotting library for Python that provides low-level control for creating static, animated, and interactive visualizations, including basic plots with customizable properties and styling options.

NumPy:
A Python library for working with arrays that provides functions for linear algebra, Fourier transforms, and matrices, offering array objects that are up to 50x faster than traditional Python lists.

Special mention - Perplexity AI: used as a research tool to find data quickly, perform multiple searches across the internet to get result required, and format presentation quickly





Disease Risk Prediction

 Algorithms: Random Forest, Logistic Regression, XGBoost

Patient Readmission Forecasting

 Algorithms: Random Forest, SVM, Neural Networks

Mortality Prediction

• Algorithms: Cox Regression, Random Forest, Deep Learning

Treatment Response Analysis

• Algorithms: Ensemble Methods, Gradient Boosting



Intrusion Detection Systems (IDS)

• Algorithms: Random Forest, SVM, Neural Networks, Ensemble Methods

Malware Detection

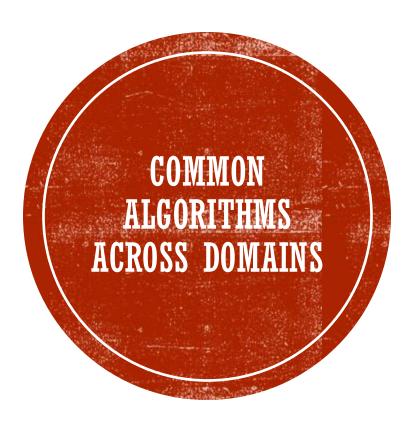
• Algorithms: CNN, Random Forest, Deep Learning

Network Anomaly Detection

 Algorithms: Isolation Forest, Autoencoders, LSTM, Clustering

Phishing Identification

 Algorithms: Logistic Regression, SVM, Random Forest



Random Forest

- Most widely adopted approach
- Handles missing data, provides feature importance

Neural Networks & Deep Learning

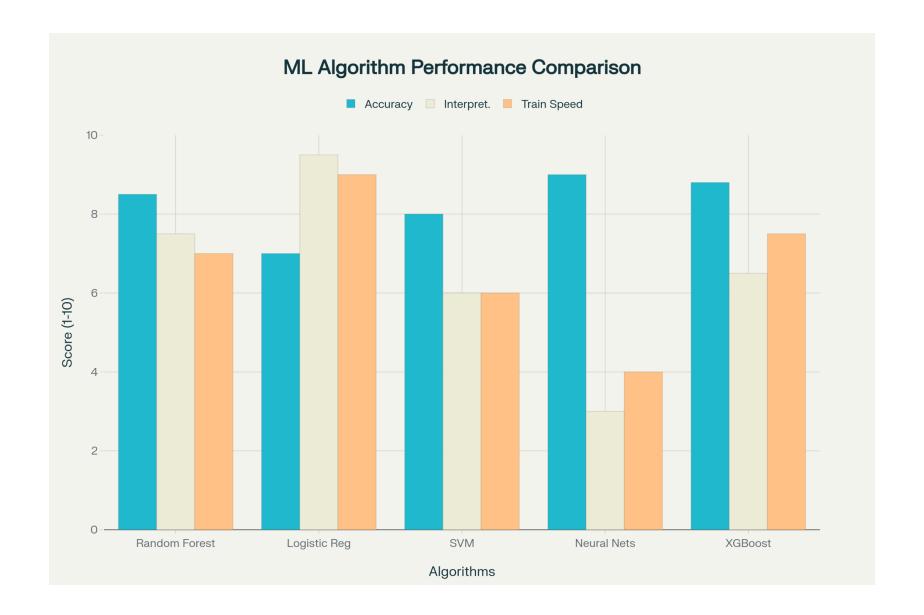
• Exceptional performance in complex pattern recognition

Support Vector Machines

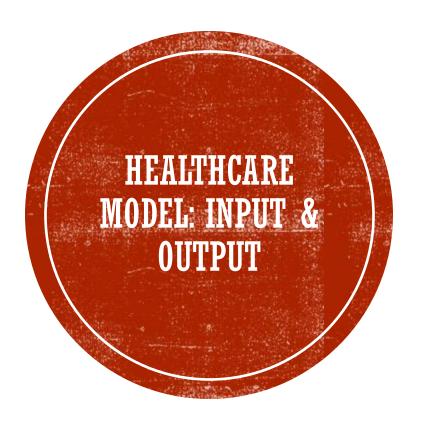
• Excel in high-dimensional data scenarios

Logistic Regression

• Interpretable with probabilistic outputs



Model Name	Input Types	Output Types	
Linear Regression	Numerical features Continuous numerical values		
Logistic Regression	Numerical and categorical features	Binary/multiclass probabilities, class labels	
Decision Tree	Numerical, categorical, ordinal features Class labels, continuous values (regression trees)		
Random Forest	umerical, categorical, ordinal features Class predictions, probability scores, continuo values		
XGBoost	Numerical and categorical features	Class predictions, probability scores, regression values	
Support Vector Machine (SVM)	Numerical features, categorical (after encoding)	Binary/multiclass classifications, regression values	
K-Nearest Neighbors (KNN)	Numerical features	Class labels, continuous values (regression)	
Neural Network/MLP	Numerical, categorical (encoded), text (vectorized), image data	Classifications, regression values, complex predictions	
Naive Bayes	Numerical and categorical features	Class predictions with probability scores	
K-Means Clustering	Numerical features	Cluster assignments, centroids	
DBSCAN	Numerical features	Cluster assignments, core points identification	
Principal Component Analysis (PCA)	Numerical features	Reduced-dimension numerical features	
LSTM	Sequential data, time series, text sequences	Sequential predictions, next value prediction	
Convolutional Neural Network (CNN)	Image data (pixel values), 2D/3D arrays	Image classifications, object detection, feature maps	

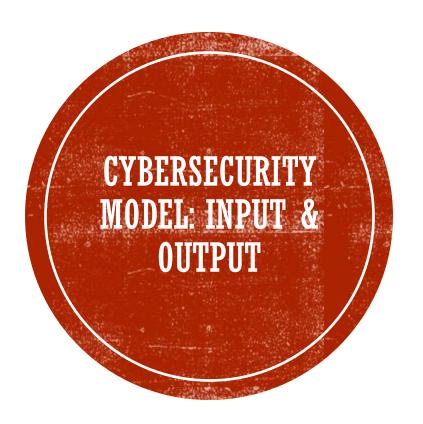


Input Features:

- Patient Demographics: Age (18-90), Gender, Blood Type
- Medical Data: Diabetes, Hypertension, Asthma, Test Results
- Administrative: Admission Type, Insurance, Billing (\$1,000-\$50,000)
- Hospital Stay Duration: 1-30 days

Output:

- Binary Classification: High Risk (1) or Low Risk (0)
- Probability Scores: 0.0 to 1.0 likelihood of adverse outcomes

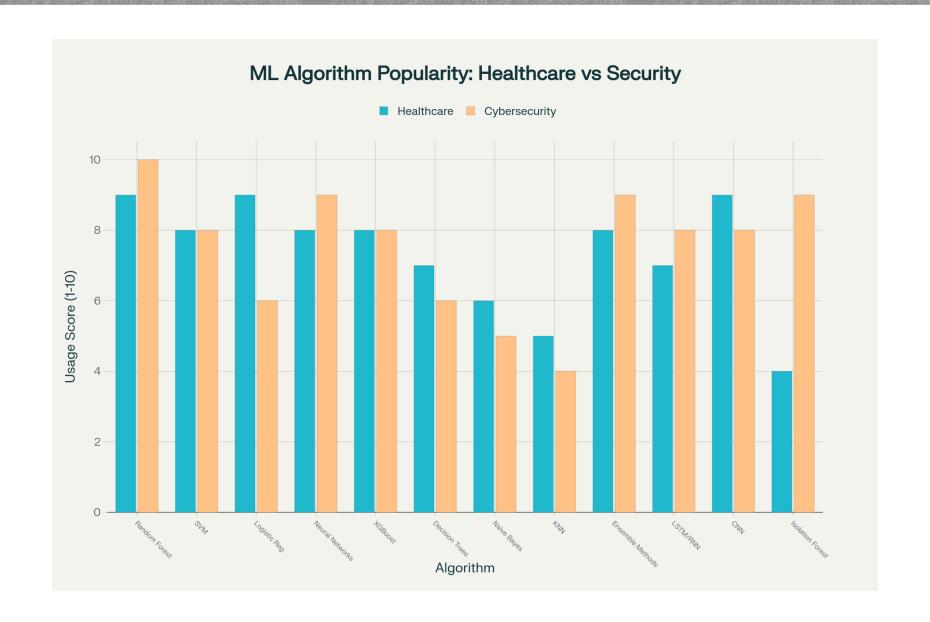


Input Features:

- Network Traffic: Flow duration, bytes/second, packets/second
- Flow Characteristics: Forward/backward packet counts, total lengths
- Packet Analysis: IAT means, PSH/URG flag counts, average sizes

Output:

- Binary Classification: Attack (1) or Benign (0)
- Probability Scores: 0.0 to 1.0 likelihood of malicious activity





Algorithm Chosen:

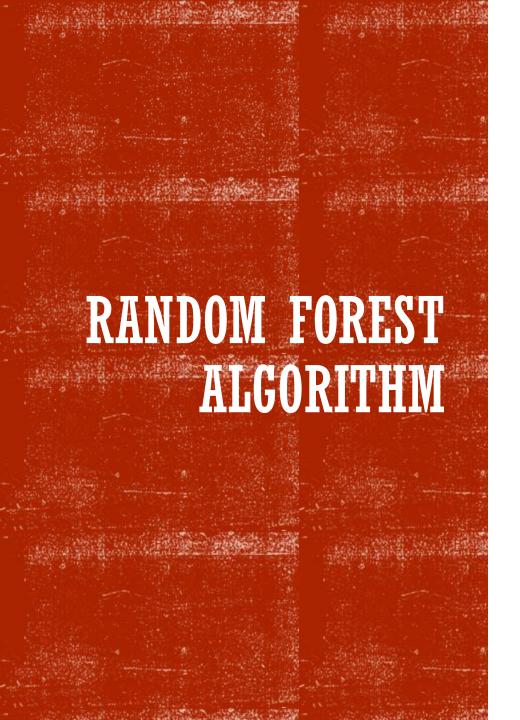
- Random Forest Classifier with 100 estimators
- Fixed random state (42) for reproducible results

Data Preprocessing:

- Healthcare: Label Encoding, StandardScaler, SimpleImputer
- Cybersecurity: StandardScaler normalization only

Feature Ranking:

Feature importance rankings for interpretability



- What is Random Forest?
 An ensemble machine learning algorithm that combines multiple decision trees to make more accurate predictions by leveraging the "wisdom of crowds" principle.
- Why It Works: The randomization creates diverse, uncorrelated trees that collectively make better decisions than any single tree, similar to how a group of experts often outperforms one expert.

HOW IT WORKS:



Sampling - Creates multiple random subsets of training data as decision trees for each feature



Random Feature Selection - Each tree uses only a random subset of features at each split



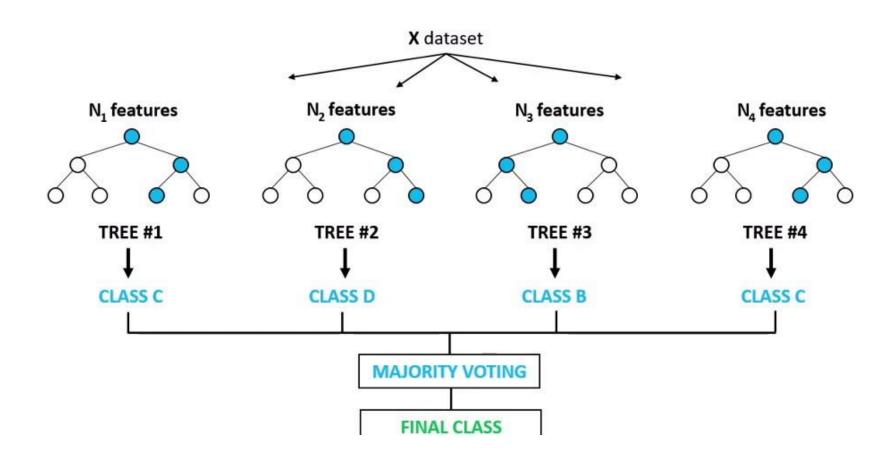
Build Multiple
Trees(estimators) - Constructs
many decision trees using
different data/feature
combinations for single feature

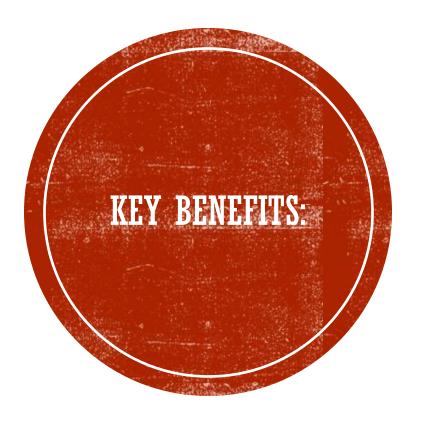


Aggregate Predictions -Combines results via majority voting (classification) or averaging (regression)



Random Forest Classifier





- Reduces Overfitting Multiple trees prevent memorizing training data
- Higher Accuracy Ensemble performs better than individual trees
- Handles Missing Data Robust in case of incomplete datasets
- Works for Both -Cybersecurity and Healthcare datasets



Data Generation:

- Synthetic data with configurable parameters
- Realistic patient profiles and network traffic patterns

Training Workflow:

- 80-20 data split with stratification
- Comprehensive metrics calculation
- Real-time prediction interface



⟨/⟩

The project includes a comprehensive Streamlit web application featuring:



User Interface: Intuitive usage pages for both healthcare and cybersecurity models



Data Synthesis: Configurable custom synthetic data generation for training and testing



Model Training: Automated training with one click Prediction making



Results Visualization: Performance metrics, confusion matrices, and feature importance plots



Prediction Interface: Real-time prediction capabilities with user input forms

MODEL PERFORMANCE & EVALUATION

- Evaluation Metrics Framework:
 - Accuracy: Overall correctness of predictions
 - AUC: Discrimination ability across thresholds
 - Precision: Quality of positive predictions
 - Recall: Completeness of positive case identification
 - F1-Score: Balanced measure of precision and recall
 - Confusion Matrix: Detailed classification breakdown

Metric	Description	Healthcare Application	Cybersecurity Application
Accuracy	Percentage of correct predictions out of total predictions	Overall accuracy in predicting patient risk levels	Overall accuracy in detecting attacks vs benign traffic
AUC Score	Area Under ROC Curve - measures discrimination ability	How well the model distinguishes high-risk from low-risk patients	How well the model distinguishes attacks from normal traffic
Precision	True Positives / (True Positives + False Positives)	Of patients predicted as high-risk, how many actually are high-risk	Of traffic predicted as attacks, how much is actually malicious
Recall	True Positives / (True Positives + False Negatives)	Of all actual high-risk patients, how many are correctly identified	Of all actual attacks, how many are correctly detected
F1-Score	Harmonic mean of Precision and Recall	Balance between identifying high-risk patients and avoiding false alarms	Balance between detecting attacks and minimizing false positives
Confusion Matrix	2x2 matrix showing TP, TN, FP, FN counts	Shows correct vs incorrect risk classifications	Shows correct vs incorrect attack/benign classifications
Classification Report	Detailed per-class precision, recall, F1-score	Detailed performance for high-risk and low-risk classes	Detailed performance for attack and benign classes
Feature Importance	Ranking of features by their contribution to predictions	Which patient factors most influence risk predictions	Which network features most indicate malicious activity
Cross-Validation	K-fold cross-validation capability built-in	Validates model stability across different patient populations	Validates model stability across different network conditions
Model Interpretability	Feature importance provides model explainability	Clinicians can understand why a patient is flagged as high-risk	Security analysts can understand why traffic is flagged as malicious



Implementation Architecture:

- Streamlit web application with intuitive interfaces
- Professional-grade visualization components

Performance Outcomes:

- High accuracy suitable for critical decision-making
- Clear feature importance rankings for interpretability
- Scalable architecture for largescale data processing
- Real-time prediction capabilities

Component	Description
Healthcare Prediction Model	Random Forest model for healthcare risk prediction
Cybersecurity Prediction Model	Random Forest model for cybersecurity threat detection
Features(columns) - Healthcare	Age, Gender, Medical Conditions, Billing, Hospital Days, Test Results, etc.
Features(columns) - Cybersecurity	Flow Duration, Packet Counts, Bytes, Network Traffic Features, etc.
Evaluation Metrics	Accuracy, AUC Score, Confusion Matrix, Classification Report
User Interface	Streamlit web application with multiple pages for easy usage
Training Process	Train-test split (80-20), StandardScaler, Label Encoding



Random Forest provides optimal balance:

- Accuracy, interpretability, computational efficiency
- Cross-domain applicability demonstrated by using same model for different purposes and getting optimal results

Critical Success Factors:

- Comprehensive evaluation metrics beyond accuracy
- Feature importance for model explainability
- User-friendly interfaces



GitHub link to project file: https://github.com/Nikunj-Gupta-1/Model_trainer.git

To view the Web-app : https://nikunj-s--ml-trainermodel-final.streamlit.app