



Novel Ensembled Approach To Healthcare Data Analysis



- Diabetes mellitus (DM) is a prevalent chronic disease worldwide, affecting millions of individuals.
- Our project aims to develop a predictive model that can assess the risk of diabetes based on various health parameters.

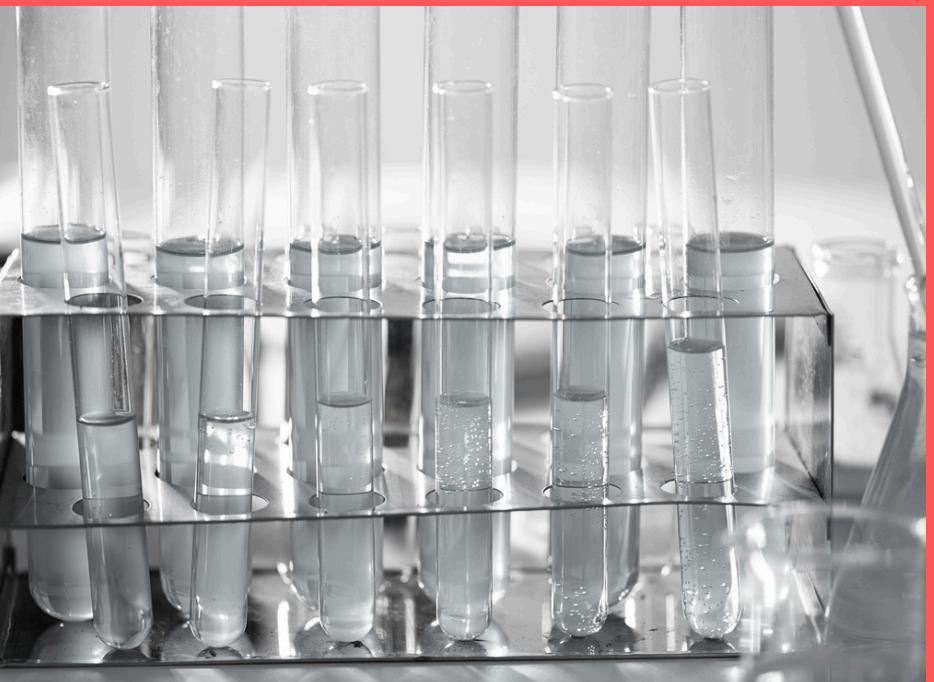
Introduction

- Leveraging machine learning algorithms, we seek to create a tool that aids in early detection and management of diabetes, ultimately improving healthcare outcomes.
- Improvement: Regular updates and detailed documentation to maintain effectiveness and support future research.



Problem Statement

"DEVELOPING A NOVEL ENSEMBLE APPROACH TO ENHANCE THE ACCURACY AND RELIABILITY OF HEALTHCARE DATA ANALYSIS TO PREDICT DIABETES DISEASE."



Motivation



- Enhancing Accessibility to Healthcare
- Enhancing Patient Outcomes
- Reducing Healthcare Costs
- Empowering Personalized Healthcare
- Research and Development



Objectives

- To increase the accuracy of the model using ensembled technique
- To provide more stable and reliable predictions by mitigating the weaknesses of individual models
- To performs well on unseen data, improving its applicability in real-world scenarios.
- To Validats model performance through techniques like k-fold cross-validation, ensuring robust evaluation.

Literture Survey:

Research Papers	Description	Author	Publication & Year of Publication
CatBoost Ensemble Approach for Diabetes Risk Prediction at Early Stages	In this paper, we propose an ensemble technique CatBoost which is a Gradient Boosting Decision Tree (GBDT) for diabetes prediction at early stages.	P. Suresh Kumar1 Anisha Kumari K Subhashree Mohapatra	IEEE 2021
Genome-wide association analysis of type 2 diabetes in the EPIC-InterAct study	The EPIC-InterAct project, centred in 8 countries in the European Prospective Investigations into Cancer and Nutrition study, is one of the largest prospective studies of T2D.	Lina Cai	2020
Genomic Prediction of 16 Complex Disease Risks Including Heart Attack, Diabetes, Breast and Prostate Cancer	Our results indicate that substantial improvements in predictive power are attainable using training sets with larger case populations. We anticipate rapid improvement in genomic prediction as more casecontrol data become available for analysis.	Louis Lello TimothyG. Raben SokeYuenYong	2019

Random Forest Algorithm for the Prediction of Diabetes

The literature introduces a weighted ensembling approach, where weights are determined based on the Area Under ROC Curve (AUC) of each ML model. AUC is chosen as the performance metric, optimized through hyperparameter tuning using the grid search technique.

Md. Kamrul Hasan

IEEE 2019

HYBRID PREDICTION MODEL FOR TYPE-2 DIABETES WITH CLASS IMBALANCE

This paper presents an efficient model for classifying type-2 diabetes using a hybrid approach. The dataset, obtained from Vanderbilt University, USA, focuses on 390 African-American individuals and includes parameters such as Glucose, Age, Cholesterol, Weight, and Waist/Hip ratio. To address imbalance, oversampling and undersampling methods are applied.

Rishi Kashyap
Surya Teja
CVN

IEEE 2020

Predictive Diabetes Mellitus by DNA sequences using Deep Learning

Deep learning algorithms and artificial intelligence (AI) are introduced. Based on DNA sequencing, a recent study used Long Short-Term Memory (LSTM) algorithms and Convolutional Neural Networks (CNNs) to identify diabetes types. The results were remarkable: on a labeled dataset, the suggested CNN-LSTM model reached 100% accuracy, highlighting AI's promise in healthcare

Lena abed
ALraheim
Hamza,
Hussein Attya
Lafta.
Sura Zaki Al-
Rashid,

IEEE 2023

Prediction of type 2 diabetes mellitus onset using logistic regressionbased scorecards

This study focused on analyzing data from 44,709 nondiabetic participants aged 40–69 in the UK Biobank to predict the risk of Type 2 Diabetes (T2D) onset within a mean timeframe of 7.3 years. Initially considering 798 potential predictors for T2D onset, the study employed gradient boosting decision trees, survival analysis, and logistic regression methods.

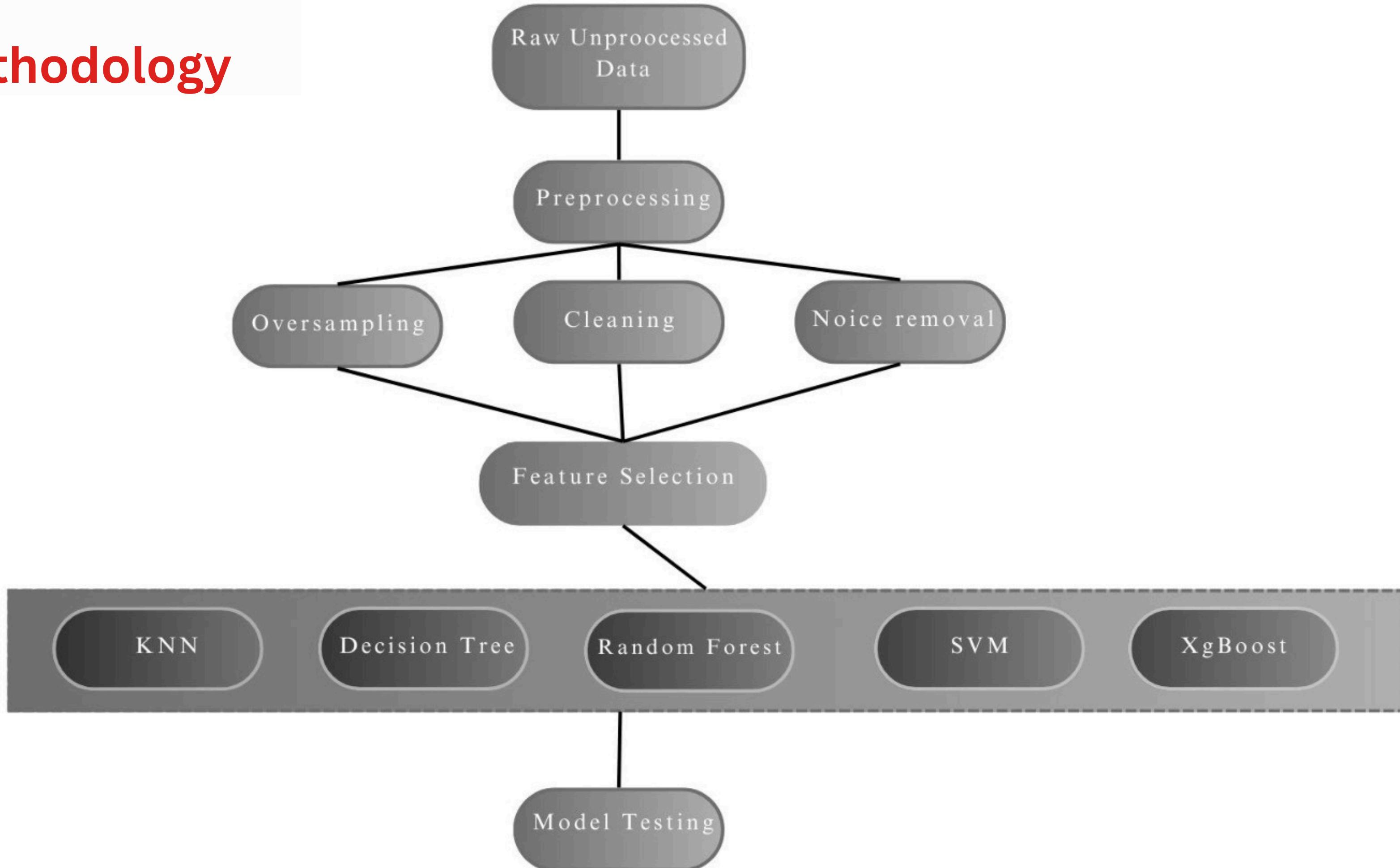
Yochai Edlitz
Eran Segal

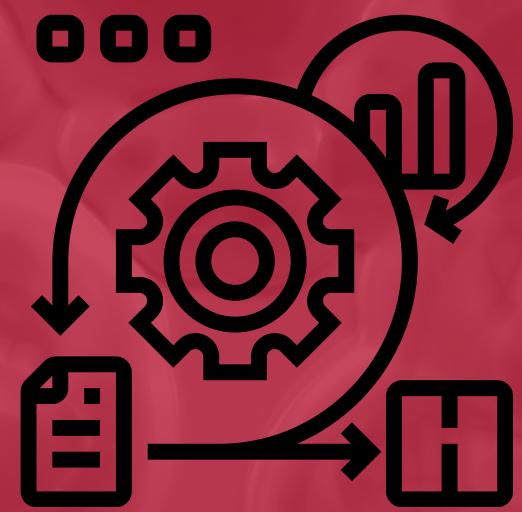
2020

<p>Machine Learning Tools for Long-Term Type 2 Diabetes Risk Prediction</p>	<p>A novel ensemble WeightedVotingLRRFs ML model is introduced to enhance diabetes prediction, achieving an Area Under the ROC Curve (AUC) of 0.884. The model incorporates optimal weights determined through a bi-objective genetic algorithm based on Sensitivity and AUC.</p>	<p>NIKOS FAZAKIS OTILIA KOCSIS</p>	<p>IEEE 2021</p>
<p>Prediction Model for Type 2 Diabetes using Stacked Ensemble Classifiers</p>	<p>The study employs 10-fold cross-validation and assesses performance using precision, recall, F-measure, and accuracy metrics. Experimental results demonstrate the superiority of the proposed model, achieving high accuracy rates of up to 93.18%, 98.87%, and 96.09% for datasets I, II, and III, respectively.</p>	<p>Norma Latif Fitriyani Muhammad Syafrudin</p>	<p>2020</p>
<p>Diabetes Prediction Using Ensembling of Different Machine Learning Classifiers</p>	<p>In this literature, we are proposing a robust framework for diabetes prediction where the outlier rejection, filling the missing values, data standardization, feature selection, K-fold cross-validation, and different Machine Learning classifiers were employed. T</p>	<p>MD. KAMRUL HASAN MD. ASHRAFUL ALAM</p>	<p>2020</p>

Methodology

Algorithms





Proposed Methodology

- Data Collection and Preprocessing
- Model Selection and Ensemble Techniques
- Model Evaluation
- Cross-Validation and Generalization
- Implementation and Deployment

Algorithms Used



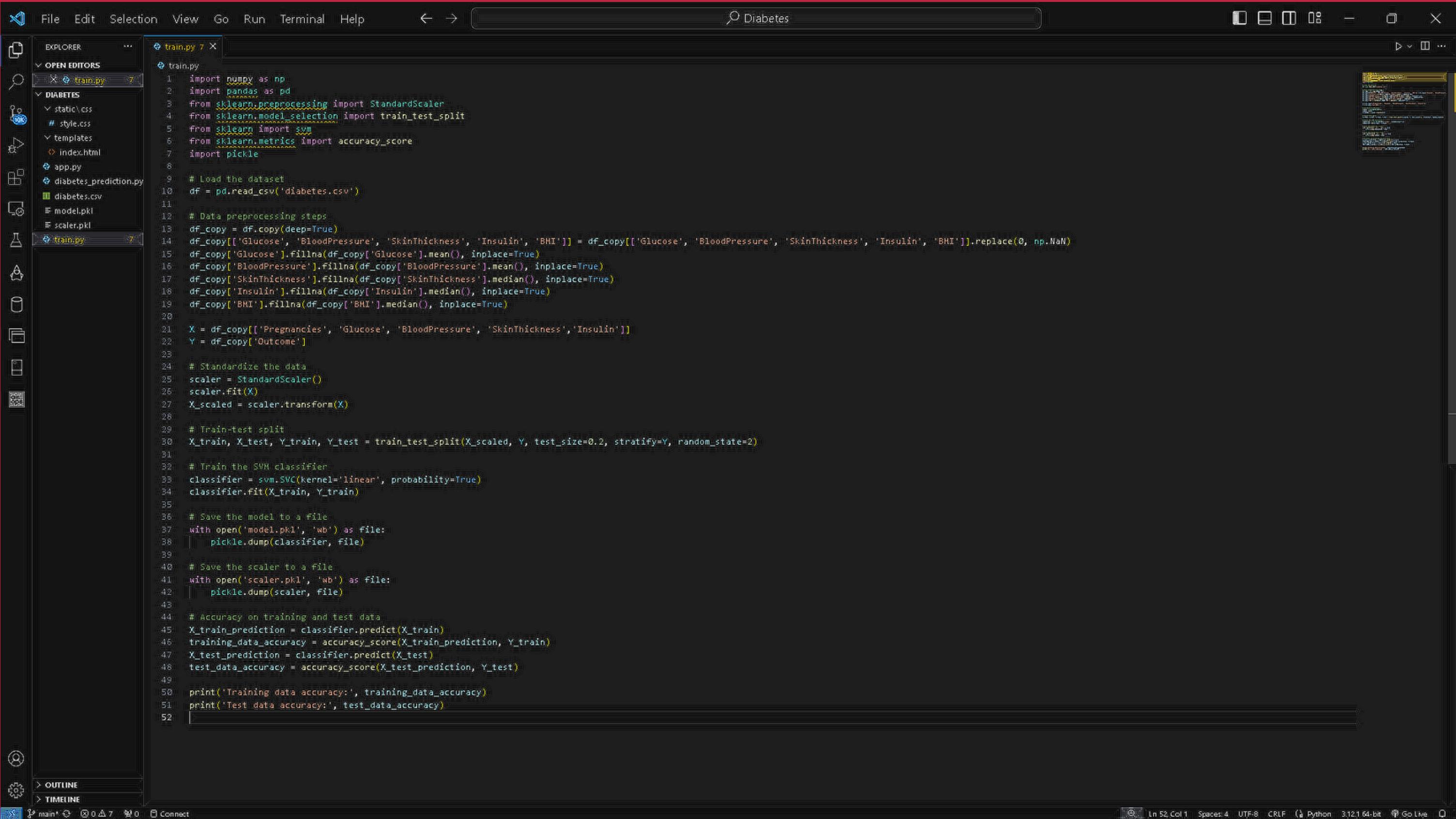
1. K-Nearest Neighbors (KNN)
2. Support Vector Machines (SVM)
3. Decision Trees
4. Random Forest
5. XGBoost (Extreme Gradient Boosting)

Experimental Results

Algorithms	Accuracy	Precision	Recall	F1 Score
KNeighboursClassifier()	69%	60%	39%	47%
SVM	76%	75%	50%	60%
Decision Tree	68%	57%	46%	51%
Random Forest	72%	63%	50%	53%
XgBoost	75%	65%	65%	65%

Implementation

Implementation



The screenshot shows a code editor interface with a dark theme. The top bar includes standard menu items: File, Edit, Selection, View, Go, Run, Terminal, Help, and a search bar labeled "Diabetes". The left sidebar contains an "EXPLORER" panel showing project files: "train.py" (7), "static\css\style.css", "templates\index.html", "app.py", "diabetes_prediction.py", "diabetes.csv", "model.pkl", and "scaler.pkl". The main editor area displays the "train.py" file content:

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.model_selection import train_test_split
5 from sklearn import svm
6 from sklearn.metrics import accuracy_score
7 import pickle
8
9 # Load the dataset
10 df = pd.read_csv('diabetes.csv')
11
12 # Data preprocessing steps
13 df_copy = df.copy(deep=True)
14 df_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] = df_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']].replace(0, np.NaN)
15 df_copy['Glucose'].fillna(df_copy['Glucose'].mean(), inplace=True)
16 df_copy['BloodPressure'].fillna(df_copy['BloodPressure'].mean(), inplace=True)
17 df_copy['SkinThickness'].fillna(df_copy['SkinThickness'].median(), inplace=True)
18 df_copy['Insulin'].fillna(df_copy['Insulin'].median(), inplace=True)
19 df_copy['BMI'].fillna(df_copy['BMI'].median(), inplace=True)
20
21 X = df_copy[['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin']]
22 Y = df_copy['Outcome']
23
24 # Standardize the data
25 scaler = StandardScaler()
26 scaler.fit(X)
27 X_scaled = scaler.transform(X)
28
29 # Train-test split
30 X_train, X_test, Y_train, Y_test = train_test_split(X_scaled, Y, test_size=0.2, stratify=Y, random_state=2)
31
32 # Train the SVM classifier
33 classifier = svm.SVC(kernel='linear', probability=True)
34 classifier.fit(X_train, Y_train)
35
36 # Save the model to a file
37 with open('model.pkl', 'wb') as file:
38     pickle.dump(classifier, file)
39
40 # Save the scaler to a file
41 with open('scaler.pkl', 'wb') as file:
42     pickle.dump(scaler, file)
43
44 # Accuracy on training and test data
45 X_train_prediction = classifier.predict(X_train)
46 training_data_accuracy = accuracy_score(X_train_prediction, Y_train)
47 X_test_prediction = classifier.predict(X_test)
48 test_data_accuracy = accuracy_score(X_test_prediction, Y_test)
49
50 print('Training data accuracy:', training_data_accuracy)
51 print('Test data accuracy:', test_data_accuracy)
52
```

The bottom status bar shows: Ln 52, Col 1, Spaces: 4, UTF-8, CRLF, Python 3.12.1 64-bit, Go Live.

Implementation

The screenshot shows a code editor interface with a dark theme. The top bar includes standard menu items: File, Edit, Selection, View, Go, Run, Terminal, Help, and a search bar labeled "Diabetes". The left sidebar contains an "EXPLORER" view showing a file tree for a project named "DIABETES". The "index.html" file is selected in the tree and is also the active editor tab. The main editor area displays the HTML code for the user interface:

```
<!DOCTYPE html>
<html >
<!--From https://codepen.io/frytyler/pen/EGdtg-->
<head>
    <meta charset="UTF-8">
    <title>Diabetes Predictor</title>
    <link href='https://fonts.googleapis.com/css?family=Pacifico' rel='stylesheet' type='text/css'>
    <link href='https://fonts.googleapis.com/css?family=Arimo' rel='stylesheet' type='text/css'>
    <link href='https://fonts.googleapis.com/css?family=Hind:300' rel='stylesheet' type='text/css'>
    <link href='https://fonts.googleapis.com/css?family=Open+Sans+Condensed:300' rel='stylesheet' type='text/css'>
    <link rel="stylesheet" href="{{ url_for('static', filename='css/style.css') }}>
</head>
<body>
    <div class="login">
        <h1>Diabetes Predictor</h1>
        <form action="/" method="post">
            <input type="text" name="pregnancies" placeholder="Pregnancies" required="required" />
            <input type="text" name="glucose" placeholder="Glucose" required="required" />
            <input type="text" name="blood_pressure" placeholder="Blood Pressure" required="required" />
            <input type="text" name="skin_thickness" placeholder="Skin Thickness" required="required" />
            <input type="text" name="insulin" placeholder="Insulin" required="required" />
            <br>
            <br>
            {{ prediction_text }}
        </div>
    </body>
</html>
```

The code defines a simple HTML form for inputting pregnancy, glucose, blood pressure, skin thickness, and insulin levels. It also includes a "Predict" button and a placeholder for predicted output. The "prediction_text" variable is likely a template variable being rendered by the application.

User interface implementaion

Our Model

The image displays three screenshots of a "Diabetes Predictor" web application, arranged horizontally against a red background.

Home Page: The first screenshot shows the application's main interface with the title "Diabetes Predictor". It features five input fields for "Pregnancies", "Glucose", "Blood Pressure", "Skin Thickness", and "Insulin", each with a dropdown menu showing values like 0, 80, 120, 60, and 70 respectively. Below these is a blue "Predict" button, and at the bottom, the text "None".

User Input: The second screenshot shows the same interface after user interaction. The "Glucose" field now has the value "70" selected. The "Predict" button is highlighted in blue, and the text "None" is displayed below the input fields.

Output: The third screenshot shows the results of the prediction. The "Predict" button is now greyed out. The text "The person is not diabetic" is displayed prominently below the input fields.

Future Scope

1. Integration of Genetic Data
2. Real-time Monitoring Systems
3. Expansion to Other Chronic Diseases
4. Explore Additional Algorithms

Conclusion



- Ensemble methods significantly improve the accuracy of diabetes prediction by combining the strengths of multiple models.
- Ensemble techniques effectively reduce both bias and variance, leading to better generalization on unseen data.
- The flexibility to combine different algorithms allows for the creation of highly customized models tailored to specific prediction tasks.

References

1. Hamza, L.A. et al. (2023). Predictive Diabetes Mellitus From DNA Sequences Using Deep Learning. *Al-Bahir Journal for Engineering and Pure Sciences*, 3(2). <https://doi.org/10.55810/2313-0083.1042>
2. Zou, Q. et al. (2018). Predicting Diabetes Mellitus With Machine Learning Techniques. *Front. Genet.* 9:515. <https://doi.org/10.3389/fgene.2018.00515>
3. Das, B. (2022). A deep learning model for identification of diabetes type 2 based on nucleotide signals. *Neural Comput & Applic*, 34, 12587–12599. <https://doi.org/10.1007/s00521-022-07121-8>
4. Kim, J. et al. (2018). Genetic prediction of type 2 diabetes using deep neural network. *Clin Genet*, 93(4), 822–829. doi: 10.1111/cge.13175.
5. Alshamlan, H. et al. (2020). A Gene Prediction Function for Type 2 Diabetes Mellitus using Logistic Regression. doi: 10.1109/ICICS49469.2020.239549.
6. Kumar, A. et al. (2017). SVMRFE based approach for prediction of most discriminatory gene target for type II diabetes. *Health Cares Data*, 12, 28–37. <https://doi.org/10.1016/j.gdata.2017.02.008>.
7. Lello, L. et al. (2019). Genomic Prediction of 16 Complex Disease Risks Including Heart Attack, Diabetes, Breast and Prostate Cancer. *Sci Rep*, 9(1):15286. <https://doi.org/10.1038/s41598-019-51258-x>.
8. Loh, M. et al. (2022). Identification of genetic effects underlying type 2 diabetes in South Asian and European populations. *Commun Biol*, 5(1):329. <https://doi.org/10.1038/s42003-022-03248-5>.

9. Diabetes Genetics Initiative of Broad Institute of Harvard and MIT, & others. (2007). Genome-wide association analysis identifies loci for type 2 diabetes and triglyceride levels. *Science*, 316(5829), 1331–6. <https://doi.org/10.1126/science.1142358>.
10. Cai, L. et al. (2020). Genome-wide association analysis of type 2 diabetes in the EPIC-InterAct study. *Sci Data*, 7(1):393. <https://doi.org/10.1038/s41597-020-00716-7>.
11. Cole, J.B. & Florez, J.C. (2020). Genetics of diabetes mellitus and diabetes complications. *Nat Rev Nephrol*, 16(7), 377–390. <https://doi.org/10.1038/s41581-020-0278-5>.
12. Cho, N.H. & others. (2018). IDF Diabetes Atlas: Global estimates of diabetes prevalence for 2017 and projections for 2045. *Diabetes Res Clin Pract*, 138, 271–281. <https://doi.org/10.1016/j.diabres.2018.02.023>.
13. Butt, U.M. et al. (2021). Machine Learning Based Diabetes Classification and Prediction for Healthcare Applications. *J Healthc Eng*, 2021, 9930985. <https://doi.org/10.1155/2021/9930985>.
14. Oliullah, K. et al. (2023). A stacked ensemble machine learning approach for the prediction of diabetes. *J Diabetes Metab Disord*. <https://doi.org/10.1007/s40200-023-01321-2>.
15. Tripathi, D. et al. (2022). Diabetes Prediction Using Machine Learning Analytics: Ensemble Learning Techniques.
16. Nagpal, A. et al. (2023). A novel ensemble machine learning framework for early stage diabetes mellitus prediction. *Multidisciplinary Science Journal*.
17. National Diabetes Statistics Report | Diabetes | Centers for Disease Control and Prevention. 2022. <https://www.cdc.gov/diabetes/data/statistics-report/index.html>.
18. Hosseini Sarkhosh, S.M. et al. (2022). Predicting diabetic nephropathy in type 2 diabetic patients using machine learning algorithms. *J Diabetes Metab Disord*, 21(2), 1433–41.
19. Hemanth, S. & Alagarsamy, S. (2023). Hybrid adaptive deep learning classifier for early detection of diabetic retinopathy using optimal feature extraction and classification. *J Diabetes Metab Disord*, 1–15.