

# **PROGNOSIS DIA**

A PROJECT REPORT

*Submitted by*

**RONITH T VINOD[Reg. No:RA2011003010883]**

**CHINMOYEE GOGOI[Reg No:RA2011003010884]**

*Under the guidance of*

**Dr. SIVAKUMAR B**

Associate Professor, Department of Computer Science and Engineering

*in partial fulfillment for the award of the degree*

*of*

**BACHELOR OF TECHNOLOGY**

**in**

**COMPUTER SCIENCE & ENGINEERING**



**DEPARTMENT OF COMPUTING TECHNOLOGIES  
COLLEGE OF ENGINEERING AND TECHNOLOGY  
SRM INSTITUTE OF SCIENCE AND TECHNOLOGY**

**KATTANKULATHUR- 603 203**

**NOVEMBER 2023**

# **SRM INSTITUTE OF SCIENCE AND TECHNOLOGY**

(Under Section 3 of UGC Act, 1956)

## **BONAFIDE CERTIFICATE**

Certified that 18CSP109L / I8CSP111L project report titled "**PROGNOSIS DIA**" is the bonafide work of **RONITH T VINOD[Reg No:RA201100301883]** and **CHINMOYEE GOGOI[Reg No:RA2011003010884]** who carried out the project work under my supervision. Certified further, that to the best of my knowledge the work reported here in does not form part of any other thesis or dissertation on the basis of which a degree or award was conferred on an earlier occasion for this or any other candidate.



Department of Computing Technologies  
**SRM Institute of Science and Technology**  
**Own Work Declaration Form**

**Degree/Course** : B.Tech in Computer Science and Engineering

**Student Names** : RONITH T VINOD, CHINMOYEE GOGOI

**Registration Number** : RA2011003010883,RA201103010884

**Title of Work** : PROGNOSIS DIA

I/We here by certify that this assessment complies with the University's Rules and Regulations relating to Academic misconduct and plagiarism, as listed in the University Website, Regulations, and the Education Committee guidelines.

I / We confirm that all the work contained in this assessment is our own except where indicated, and that we have met the following conditions:

- Clearly references / listed all sources as appropriate
- Referenced and put in inverted commas all quoted text (from books, web,etc.)
- Given the sources of all pictures, data etc that are not my own.
- Not made any use of the report(s) or essay(s) of any other student(s) either past or present
- Acknowledged in appropriate places any help that I have received from others (e.g fellow students, technicians, statisticians, external sources)
- Compiled with any other plagiarism criteria specified in the Course hand book / University website

I understand that any false claim for this work will be penalized in accordance with the University policies and regulations.

**DECLARATION:**

I am aware of and understand the University's policy on Academic misconduct and plagiarism and I certify that this assessment is my / our own work, except where indicated by referring, and that I have followed the good academic practices noted above.

**Student 1 Signature:**

**Student 2 Signature:**

**Date:**

If you are working in a group, please write your registration numbers and sign with the date for every student in your group.

## **ACKNOWLEDGEMENT**

We express our humble gratitude to **Dr. C. Muthamizhchelvan**, Vice-Chancellor, SRM Institute of Science and Technology, for the facilities extended for the project work and his continued support.

We extend our sincere thanks to Dean-CET, SRM Institute of Science and Technology, **Dr. T. V. Gopal**, for his invaluable support.

We wish to thank **Dr. Sivakumar B**, Assistant Professor, School of Computing, SRM Institute of Science and Technology, for her support throughout the project work.

We are incredibly grateful to our Head of the Department, **Dr. M. Pushpalatha**, Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for her suggestions and encouragement at all the stages of the project work.

We want to convey our thanks to our Project Coordinators, S. Godfrey Winster, Dr. M. Baskar, Dr. P. Murali, Dr. J. Selvin Paul Peter, Dr. C. Pretty Diana Cyril and Dr.G.Padmapriya, Panel Head, Dr. G. Usha, Associate Professor and Panel Members, Mrs. M. Ranjani, Assistant Professor, Dr. C. Pretty Diana Cyril, Assistant Professor and Dr. R. Anto Arockia Rosaline, Assistant Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for their inputs during the project reviews and support.

We register our immeasurable thanks to our Faculty Advisor, Dr. P Murali, Associate Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for leading and helping us to complete our course.

Our inexpressible respect and thanks to our guide, Dr. C. Sivakumar B, Associate Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for providing us with an opportunity to pursue our project under her mentorship. She provided us with the freedom and support to explore the research topics of our interest. Her passion for solving problems and making a difference in the world has always been inspiring.

We sincerely thank all the staff and students of Computing Technologies Department, School of Computing, S.R.M Institute of Science and Technology, for their help during our project. Finally, we would like to thank our parents, family members, and friends for their unconditional love, constant support and encouragement.

**RONITH T VINOD[Reg. No:RA2011003010883]**

**CHINMOYEE GOGOI[Reg No:RA2011003010884]**

## ABSTRACT

Dia Prognosis represents a revolutionary paradigm shift in medical diagnosis and patient care.

This innovative system harnesses the power of advanced data analytics, artificial intelligence and personalized medicine to revolutionize the way we approach disease prognosis.

By integrating large data sets, including genetic, clinical and environmental factors, Prognosis Dia provides a comprehensive and comprehensive view of patient health.

Using predictive algorithms, it can predict disease progression, identify potential complications, and tailor treatment plans with unprecedented accuracy.

This transformative technology not only improves patient outcomes but also provides healthcare professionals with the tools to make more informed decisions, ushering in a new era of traditional healthcare system.

Dynamic, personable. Dia Prognosis is more than just a diagnostic tool; it is a vision of a future where healthcare is more predictive, preventative and patient-centered.

Random forest is a machine learning algorithm that can be used for both regression and classification tasks. It is an ensemble learning method that combines multiple decision trees to create a more accurate and robust model.

In a random forest, each decision tree is built using a subset of training data and a random selection of features. This helps reduce overfitting and increases overall tree diversity.

When predicting, each tree in the forest generates a class prediction or continuous output independently, and the final prediction is obtained by aggregating the individual predictions of all the trees.

Random forests have several advantages over single decision trees, including higher accuracy, greater resistance to overfitting, and the ability to handle high-dimensional data.

They are also relatively easy to use and require minimal hyperparameter tuning.

Some applications of random forests include image classification, sentiment analysis, credit risk assessment, and disease diagnosis.

# TABLE OF CONTENTS

|  |           |
|--|-----------|
| <b>ABSTRACT</b>  | <b>6</b>  |
| <b>TABLE OF CONTENTS</b>                               | <b>7</b>  |
| <b>LIST OF FIGURES</b>                                 | <b>8</b>  |
| <b>1 INTRODUCTION</b>                                  | <b>9</b>  |
| <b>1.1 BACKGROUND</b>                                  | <b>9</b>  |
| <b>1.2 MOTIVATION</b>                                  | <b>10</b> |
| <b>1.3 OBJECTIVE</b>                                   | <b>12</b> |
| <b>2 LITERATURE SURVEY</b>                             | <b>16</b> |
| <b>3 ARCHITECTURE AND ANALYSIS OF PREDICTION MODEL</b> | <b>17</b> |
| <b>3.1 FRONT-END DESIGN</b>                            | <b>20</b> |
| <b>3.2 BACK-END DESIGN</b>                             | <b>23</b> |
| <b>4 METHODOLOGY</b>                                   | <b>26</b> |
| <b>4.1 BASIC INTRODUCTORY</b>                          | <b>26</b> |
| <b>4.2 ALGORITHM USED</b>                              | <b>30</b> |
| <b>4.3 MAIN CODE</b>                                   | <b>38</b> |
| <b>4.4 MAIN CODE EXPLAINATION</b>                      | <b>42</b> |
| <b>5 RESULTS AND DISCUSSION</b>                        | <b>45</b> |
| <b>5.1 PERFORMANCE ANALYSIS</b>                        | <b>45</b> |
| <b>5.2 PERFORMANCE WORKING SCREENSHOTS</b>             | <b>47</b> |
| <b>6 CONCLUSION AND FUTURE SCOPE</b>                   | <b>49</b> |
| <b>7 REFERENCES</b>                                    | <b>54</b> |

## **LIST OF FIGURES**

Fig.1 Proposed System

Fig.2 Working of Random Forest

Fig.3 Front-End UI

Fig.4 Backend(Excel API)

Fig.5 Importing SKLearn

Fig.6 Correlating the Model

Fig.7 Data Model

Fig.8 Model Training

Fig.9 Diff for Linear and Logistic

Fig.10 Formulating Regression

Fig.11 Working Model UI

Fig.12 Back-End Database

# **CHAPTER 1**

## **INTRODUCTION**

### **1.1 BACKGROUND**

Diabetes is a chronic disease that affects millions of people worldwide. It is caused by high levels of glucose (sugar) in the blood, which can lead to serious complications such as heart disease, kidney failure, and blindness. Early detection and management of diabetes are crucial to prevent these complications and improve patients' quality of life.

Machine learning algorithms can be used to predict the risk of diabetes in individuals based on their clinical and demographic data. Rf is one such algorithm that has been used for this purpose.

The Rf algorithm can analyze a large number of features and capture complex relationships between them, making it well-suited for diabetes prediction. In this approach, a dataset containing clinical and demographic information of patients, including age, gender, body mass index (BMI), blood pressure, and glucose levels, is used to train a Rf model. The model learns to identify patterns in the data that are associated with diabetes and uses this information to make predictions on new data.

During the training process, the Rf algorithm creates multiple decision trees, each using a subset of the features and a random selection of data samples. The final prediction is obtained by aggregating the predictions of all trees in the forest. This helps to reduce overfitting and improve the model's accuracy.

The diabetic predictor using the Rf algorithm can be a valuable tool for healthcare professionals to identify patients at risk of developing diabetes and provide early interventions to prevent or manage the disease. It can also aid in the development of personalized treatment plans and improve patient outcomes.

## 1.2 MOTIVATION

**Collecting and analyzing data:** The project involves collecting and analyzing large amounts of data related to the patient's medical history, lifestyle, and other factors that may influence their likelihood of developing diabetes.

**Data diversity:** Since diabetes is a complex disease and various factors affect the outcome, the data taken has multiple factors that change the result.

**Multiple data sources:** Diabetes is a personal issue, therefore data is not available to the public, to resolve this issue the data collected has been taken from multiple sources.

**Mixed Data:** Data from multiple ethnic and racial backgrounds has been taken to better predict the result.

**Algorithms used:** Since machine learning involves many different algorithms, we have done a comparative study involving random forest and logistic regression.

Motivation for Developing a Diabetic Predictor:

**Improving patient health:** Diabetes is a chronic disease that affects millions of people worldwide. The development of diabetes prediction tools is driven by the desire to improve the health and well-being of people with diabetes.

By accurately predicting their risk of diabetes and its complications, we can enable individuals to take proactive steps to manage their health and reduce the impact of the disease.

**Early intervention:** Early detection and intervention are essential to control diabetes. Diabetes prediction tools can identify people at risk before symptoms worsen, allowing for early medical intervention, lifestyle changes and preventative measures. This can lead to better health outcomes and reduced health care costs.

**Personalized Medicine:** Diabetes is a complex disease with many different risk factors, including genetics, lifestyle and environmental factors. Diabetes prediction tools can apply the concept of personalized medicine, tailoring recommendations and treatments to each individual's unique risk profile. This approach can lead to more efficient and effective health care.

**Reduces the healthcare burden:** Diabetes places a significant burden on healthcare systems, both in terms of costs and resources. The development of diabetes prediction tools can help optimize resource allocation, reduce emergencies, and prevent diabetes-related complications, thereby easing the burden on the healthcare infrastructure.

**Data-Driven Insights:** In the era of Big Data and advanced analytics, diabetes prediction tool development can harness the power of data to better understand the progression of diabetes. This knowledge can guide ongoing research and lead to better diabetes management strategies.

**Patient Empowerment:** Diabetes prediction tools can enable individuals to make informed decisions about their health. By providing them with information about risks and potential interventions, it encourages active participation in their own health care, leading to better compliance with treatment plans and healthier lifestyle.

**Prevent complications:** Diabetes can cause serious complications, including heart disease, kidney failure and vision loss. The motivation for diabetes prediction tools stems from the ability to prevent these complications by identifying and treating risks early, thereby improving the quality of life of those affected.

**Research and innovation:** Diabetes prediction tool development stimulates research and innovation in healthcare, machine learning and artificial intelligence. It paves the way for creating innovative tools and technology to better manage diabetes and ultimately find a cure.

**Global health impact:** Diabetes is a global health problem. Diabetes prediction tools have the potential to positively impact communities and healthcare systems worldwide by addressing the prevention and management of this common disease. In summary, the motivation for developing a diabetic predictor is deeply rooted in the desire to enhance the lives of individuals at risk of or living with diabetes, reduce the burden on healthcare systems, and leverage cutting-edge technology to promote a healthier and more informed society burden on healthcare systems, and leverage cutting-edge technology to promote a healthier and more informed society.

## 1.3 OBJECTIVES

**Preventing Diabetes:** By identifying people at risk of diabetes, interventions can be provided to help prevent or delay the onset of the disease. This can save lives and money because diabetes is an expensive disease to treat.

**Improved health outcomes:** People diagnosed with diabetes early can get the treatment they need to control blood sugar levels and prevent complications. This can improve their quality of life and reduce the risk of death.

**Reduce health care costs:** By preventing or delaying the onset of diabetes, health care costs can be reduced. In fact, diabetes is a chronic disease that requires expensive treatment.

**Improve public health:** By reducing the number of people with diabetes, the overall health of the population can be improved. This could lead to a reduction in deaths from diabetes-related complications.

**Early detection and risk assessment:** Developing an accurate and reliable diabetes prediction tool that can identify people at risk of diabetes before symptoms appear, thereby allowing early intervention.

**Identifying risk factors:** Identify and analyze key risk factors associated with diabetes, including genetic, lifestyle and environmental factors, to provide a prediction model .

**Personalized Predictions:** Create a system that provides personalized risk assessments to individuals, taking into account their unique risk profile, thereby providing personalized recommendations impersonal.

**Data integration:** Aggregates and integrates diverse data sources, such as electronic health records, genetic data, lifestyle information, and environmental data, to improve accuracy accuracy of prediction.

**Machine learning algorithms:** Develop and improve machine learning algorithms and prediction models to continuously improve the accuracy and reliability of diabetes prediction tools.

**User-Friendly Interface:** Design a user-friendly interface for healthcare professionals and individuals to access and interpret the predictions and recommendations, making the tool accessible and understandable.

**Healthcare Integration:** Facilitate the integration of the diabetic predictor into healthcare systems to support healthcare providers in patient management and early intervention.

**Preventive Education:** Incorporate educational components to empower individuals with knowledge about diabetes risk and preventive measures they can take.

**Long-Term Monitoring:** Implement a system for continuous monitoring and follow-up to assess the effectiveness of interventions and track changes in an individual's risk profile over time.

**Research and Innovation:** Encourage ongoing research and innovation by providing a platform for data analysis, insights, and the development of new preventive and treatment strategies.

**Reducing health care costs:** Assessing the impact of diabetes predictors on health care costs by reducing the number of diabetes-related complications, number of emergency department visits rescue and hospitalization.

**Global Healthcare Coverage:** Aims to make diabetes prediction tools accessible and adaptable to different healthcare systems and regions for global impact in preventing and managing diabetes.

**Ethical considerations:** Ensure that the project takes into account ethical and privacy concerns related to data collection, storage, and use, and complies with applicable regulations onion.

**Patient Empowerment:** Promote active participation of patients in their own health care by providing them with the tools and knowledge to make informed decisions about their health.

**Collaboration and Collaboration:** Foster collaboration with healthcare organizations, research organizations, and stakeholders to leverage the expertise and resources needed for project success.

**Continuous Improvement:** Committed to continuously evaluating and improving the Diabetes Predictor, incorporating feedback from users and healthcare providers to improve its effectiveness.

**Sustainability:** Develop a plan for long-term sustainability and scalability of the project to ensure its ongoing impact on diabetes prevention and management.

## CHAPTER 2

### LITERATURE SURVEY

| S.No | Paper Name  | Journal Name   | Year       |
|------|---|--|------------|
| 1    | "A comparative study of machine learning algorithms for diabetes prediction" by Hadi Zare                   | Journal of Medical Systems                                     | May 2020   |
| 2    | "Prediction of diabetes mellitus using machine learning algorithms" by Saurabh Pal                          | Journal of Medical Systems                                     | June 2019  |
| 3    | "Predicting diabetes with machine learning techniques: A review" by Asmaa Abbas                             | Journal of Diabetes Research                                   | March 2020 |
| 4    | "Predicting Type 2 Diabetes Mellitus using Machine Learning Techniques" by R. Kavitha and S. Kulothungan    | International journal of Engineering and Technology            | 2018       |
| 5    | "A Machine Learning-based Predictive Model for Type 2 Diabetes Mellitus" by H. V. Shinde and S. R. Kulkarni | International Journal of Advanced Research in Computer Science | 2021       |
| 6    | "Early Diabetes Prediction using Machine Learning Algorithms" by K. Gunasekaran and S. Muthukumar           | International Journal of Advanced Science and Technology       | 2019       |

## CHAPTER 3

### ARCHITECTURE AND ANALYSIS OF PREDICTION MODEL

Architecture diagram for diabetes prediction project: Data sources: This is the starting point of the architecture, where data is collected from various sources. These sources may include electronic health records, genetic databases, lifestyle surveys, environmental sensors, and more.

**Data entry:** Data from various sources is entered into the system. This may involve pre-processing and data transformation to ensure consistent data format and quality.

**Data storage:** Data is stored in a secure and scalable database or data store. This can be a combination of traditional databases and data lakes for structured and unstructured data, respectively.

**Machine Learning Model:** The predictive model is a central component of the architecture. Machine learning algorithms and predictive models are trained on historical data to develop the diabetic prediction capabilities.

**Personalization Engine:** A personalization engine takes into account an individual's unique data, such as genetic information, lifestyle choices, and medical history, to fine-tune the predictions.

**User Interface:** There are two primary user interfaces: Healthcare Professional Interface: This is designed for healthcare providers and allows them to access patient data, predictions, and recommendations for treatment and interventions.

**Patient Interface:** Designed for individuals at risk or living with diabetes, this interface provides them with access to their risk assessments and personalized recommendations.

**Continuous Monitoring:** Data from patients, including health metrics, lifestyle changes, and treatment adherence, are continuously monitored and integrated into the system. This information is used to update predictions and recommendations over time.

**Educational Content:** The system may provide educational content to patients to help them understand their risk factors and make informed decisions about their health.

**Integration with Healthcare Systems:** The system is integrated into healthcare institutions' existing systems, ensuring that healthcare professionals can seamlessly access and interact with the diabetic predictor.

**Evaluation and Feedback Loop:** Continuous evaluation and feedback mechanisms are in place to assess the effectiveness of the system. This includes feedback from healthcare professionals and patients, as well as the analysis of healthcare costs and outcomes.

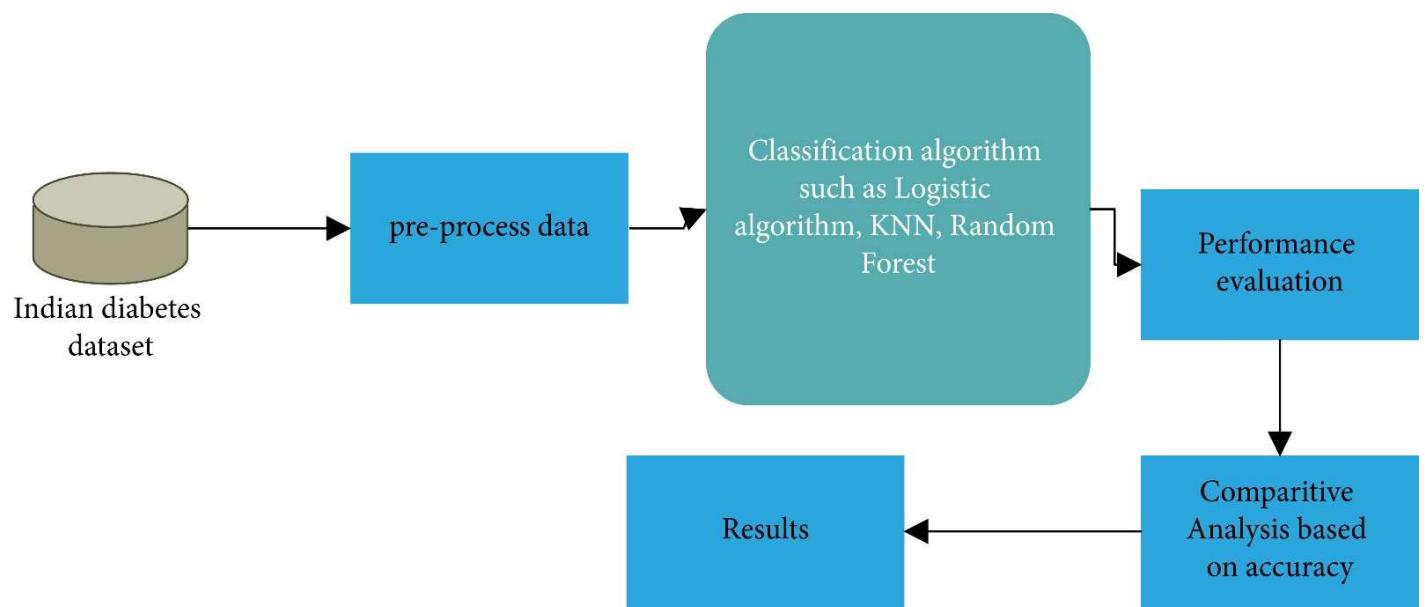
**Research and Innovation:** Data and insights generated by the system are made available for research and innovation purposes, contributing to the advancement of diabetes prevention and management strategies.

**Scalability and Sustainability:** The architecture is designed to be scalable and sustainable, ensuring that it can handle growing data volumes and adapt to evolving healthcare needs over time.

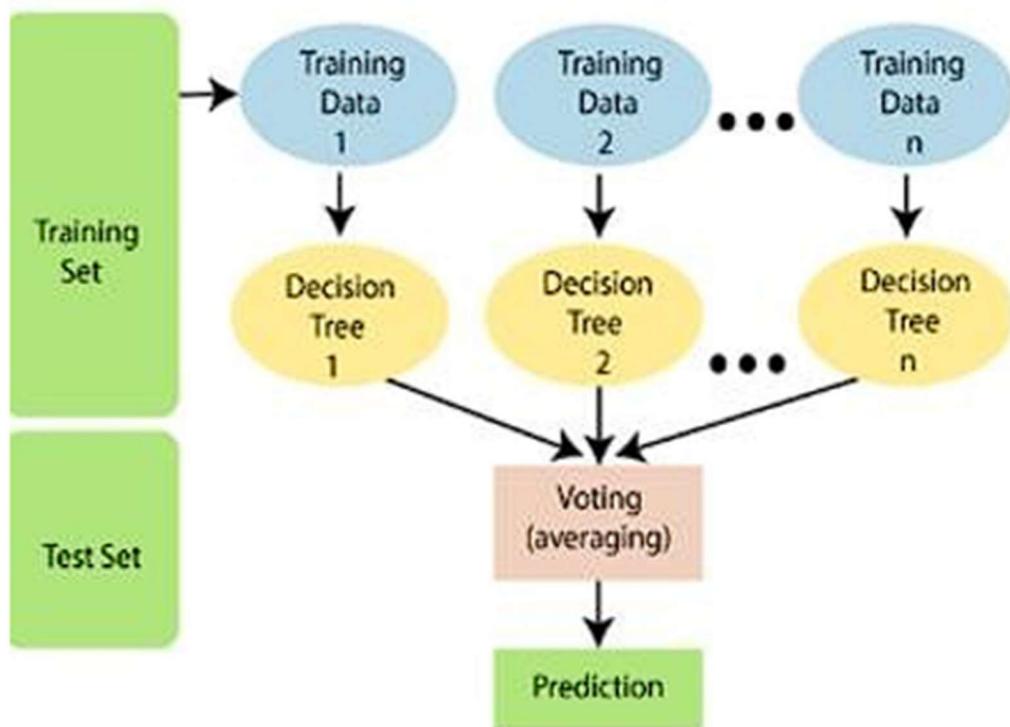
**Security and Compliance:** Security measures, including data encryption, access controls, and compliance with relevant healthcare regulations (e.g., HIPAA), are integrated into the entire architecture to protect patient data.

**Global reach:** The architecture is designed to be flexible to adapt to different health systems and regions, supporting global impact on diabetes prevention and management.

This architecture diagram illustrates the different components and their interactions in the Diabetes Predictor project, focusing on data-driven predictions, personalization, human interface use and continuous feedback loop to ensure the effectiveness of the system in improving diabetes prevention and management.



**Fig.1 Proposed System**



**Fig.2 Working of Random Forest**

### 3.1 FRONT END DESIGN

#### Landing page:

- Header: Project name and logo.
- Introduction: Briefly describe the purpose of data collection and the importance of predicting diabetes.
- Start Button: A clearly visible “Start” button allows the user to begin data collection.

## **Data Collection Form:**

- Property Name: These property names can be displayed as labels for input fields, making it clear to the user what information is being collected.
  1. "Pregnancy" - Box to enter the number of pregnancies.
  2. "Glucose" - Plasma glucose concentration input field.
  3. "Blood pressure" - Diastolic blood pressure input field.
  4. "Skin thickness" - Input field for the thickness of the triceps skin fold (mm).
  5. "Insulin" - Input field for 2-h serum insulin.
  6. "BMI" - Input field for body mass index.
  7. "Diabetes Pedigree Function" - Input field for the diabetes pedigree function.
  8. "Age" - Input field for the age of the patient.
- Input Validation: Ensure that input fields have validation to accept only numeric or relevant values, and display error messages if invalid data is entered.

## **Submit Button**

After users have entered their data, provide a "Submit" button to save the information.

## **Confirmation Page**

After data submission, display a confirmation page that thanks the user for their input and explains the purpose of data collection.

## Data Privacy Notice

Includes an explanation of how collected data will be used, securely stored, and any applicable privacy policies.

## Footer

- Contact details.
- Link to privacy policy and terms of use.
- Links to social networks.

This front-end user interface is designed to collect the data attributes needed to predict diabetes in a user-friendly manner, with clear labels and input fields.

It also emphasizes data privacy and provides options for new and returning users through login/registration features.

Additionally, the interface must include appropriate data validation to ensure users provide accurate information.

The screenshot shows a mobile-style application interface titled "DIABETES PREDICTOR". At the top, there is a small instruction: "HOLD MOUSE ON TITLES WHICH ARE CONFUSING". Below the title, there are eight input fields, each with a label and a corresponding text input box:

- Blood Sugar Level
- Diastolic Blood Pressure
- Body Mass Index (xx.x)
- Insulin
- Age (in Years)
- Number Of Pregnancies
- Skin Thickness
- Diabetes Pedigree Function (0.xxx)

**Fig.3 Front End UI**

## **3.2 BACK-END DESIGN**

To create a backend for a Diabetes Predictor project where input values are processed and calculated values are stored using the Excel API, you can use a web application framework like Flask (Python) as an example. Here is an overview of the backend components and processes:

1. **Set up the framework and backend environment:** Set up the Python environment and install the necessary libraries and dependencies, including Flask, Pandas, and Excel API libraries (e.g., openpyxl).
2. **Input and process data:** Create an API endpoint to receive input from the user interface, such as values for "Pregnancy", "Glucose", "Blood Pressure", and other attributes. Validate and clean input data to ensure it is in the correct format and within acceptable ranges.
3. **Calculate data:** Use input data to calculate predicted outcomes (e.g., diabetes risk) based on a pre-trained model or machine learning algorithm.
4. **Store data in Excel:**
  - Use the Excel API library to open Excel spreadsheets.
  - Create a new spreadsheet or use an existing spreadsheet to store collected data.
  - Add input values (e.g., "Pregnancy", "Glucose", etc.) and calculation results to the appropriate cells of the spreadsheet.

5. **Save Excel file:** Save the edited Excel file to a specified location or cloud storage (for example, Google Drive or Dropbox) for data storage.
6. **API Endpoints:** Define API endpoints for various actions, such as data submission, data retrieval, and user management.
7. **Error Handling:** Implement error handling to manage exceptions, data validation errors, and other issues that may arise during data input and processing.
8. **Security:** Implement security measures to protect data and user privacy. This may include user authentication and authorization to access and modify data.
9. **Testing:** Develop and execute test cases to ensure the API behaves as expected and handles various scenarios effectively.
10. **Documentation:** Provide clear and comprehensive documentation for the API, including endpoint descriptions and example requests and responses.
11. **Deployment:** Deploy the backend API to a web server or cloud platform to make it accessible over the internet.
12. **Ongoing Maintenance:** Regularly monitor and maintain the backend to ensure it functions smoothly, address any issues, and update it as needed. Additionally, it's important to consider data security and privacy when storing sensitive medical data in Excel or any other format. Compliance with relevant regulations (e.g., HIPAA) is crucial when handling healthcare data.

|    | A           | B       | C             | D             | E       | F    | G                | H   | I       | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z | AA | AB | AC |
|----|-------------|---------|---------------|---------------|---------|------|------------------|-----|---------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|----|----|----|
| 1  | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI  | DiabetesPedersen | Age | Outcome |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 2  | 6           | 148     | 72            | 35            | 0       | 33.6 | 0.627            | 50  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 3  | 1           | 85      | 66            | 29            | 0       | 26.6 | 0.351            | 31  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 4  | 8           | 183     | 64            | 0             | 0       | 23.3 | 0.672            | 32  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 5  | 1           | 89      | 66            | 23            | 94      | 28.1 | 0.167            | 21  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 6  | 0           | 137     | 40            | 35            | 168     | 43.1 | 2.288            | 33  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 7  | 5           | 116     | 74            | 0             | 0       | 25.6 | 0.201            | 30  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 8  | 3           | 78      | 50            | 32            | 88      | 31   | 0.248            | 26  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 9  | 10          | 115     | 0             | 0             | 0       | 35.3 | 0.134            | 29  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 10 | 2           | 197     | 70            | 45            | 543     | 30.5 | 0.158            | 53  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 11 | 8           | 125     | 96            | 0             | 0       | 0    | 0.232            | 54  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 12 | 4           | 110     | 92            | 0             | 0       | 37.6 | 0.191            | 30  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 13 | 10          | 168     | 74            | 0             | 0       | 38   | 0.537            | 34  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 14 | 10          | 139     | 80            | 0             | 0       | 27.1 | 1.441            | 57  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 15 | 1           | 189     | 60            | 23            | 846     | 30.1 | 0.398            | 59  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 16 | 5           | 166     | 72            | 19            | 175     | 25.8 | 0.587            | 51  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 17 | 7           | 100     | 0             | 0             | 0       | 30   | 0.484            | 32  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 18 | 0           | 118     | 84            | 47            | 230     | 45.8 | 0.551            | 31  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 19 | 7           | 107     | 74            | 0             | 0       | 29.6 | 0.254            | 31  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 20 | 1           | 103     | 30            | 38            | 83      | 43.3 | 0.183            | 33  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 21 | 1           | 115     | 70            | 30            | 9       | 34.6 | 0.529            | 32  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 22 | 3           | 126     | 88            | 41            | 235     | 39.3 | 0.709            | 27  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 23 | 8           | 99      | 84            | 0             | 0       | 35.4 | 0.388            | 50  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 24 | 7           | 196     | 90            | 0             | 0       | 39.8 | 0.451            | 41  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 25 | 9           | 119     | 80            | 35            | 0       | 29   | 0.263            | 29  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 26 | 11          | 143     | 94            | 33            | 146     | 36.6 | 0.254            | 51  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 27 | 10          | 125     | 70            | 26            | 115     | 31.1 | 0.205            | 41  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 28 | 7           | 147     | 76            | 0             | 0       | 39.4 | 0.257            | 43  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 29 | 1           | 97      | 66            | 15            | 140     | 23.2 | 0.487            | 22  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 30 | 13          | 145     | 82            | 19            | 110     | 22.2 | 0.245            | 57  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 31 | 5           | 117     | 92            | 0             | 0       | 34.1 | 0.337            | 38  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 32 | 5           | 109     | 75            | 26            | 0       | 36   | 0.546            | 60  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 33 | 3           | 158     | 76            | 36            | 245     | 31.6 | 0.851            | 28  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 34 | 3           | 88      | 58            | 11            | 54      | 24.8 | 0.267            | 22  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 35 | 6           | 92      | 92            | 0             | 0       | 19.9 | 0.188            | 28  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 36 | 10          | 122     | 78            | 31            | 0       | 27.6 | 0.512            | 45  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 37 | 4           | 103     | 60            | 33            | 192     | 24   | 0.966            | 33  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 38 | 11          | 138     | 76            | 0             | 0       | 33.2 | 0.42             | 35  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 39 | 9           | 102     | 76            | 37            | 0       | 32.9 | 0.665            | 46  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 40 | 2           | 90      | 68            | 42            | 0       | 38.2 | 0.503            | 27  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 41 | 4           | 111     | 72            | 47            | 207     | 37.1 | 1.39             | 56  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 42 | 3           | 180     | 64            | 25            | 70      | 34   | 0.271            | 26  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 43 | 7           | 133     | 84            | 0             | 0       | 40.2 | 0.696            | 37  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 44 | 7           | 106     | 92            | 18            | 0       | 22.7 | 0.235            | 48  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 45 | 9           | 171     | 110           | 24            | 240     | 45.4 | 0.721            | 54  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 46 | 7           | 159     | 64            | 0             | 0       | 27.4 | 0.294            | 40  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 47 | 0           | 180     | 66            | 39            | 0       | 42   | 1.893            | 25  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 48 | 1           | 146     | 56            | 0             | 0       | 29.7 | 0.564            | 29  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 49 | 2           | 71      | 70            | 27            | 0       | 28   | 0.586            | 22  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 50 | 7           | 103     | 66            | 32            | 0       | 39.1 | 0.344            | 31  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 51 | 7           | 105     | 0             | 0             | 0       | 0    | 0.305            | 24  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 52 | 1           | 103     | 80            | 11            | 82      | 19.4 | 0.491            | 22  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 53 | 1           | 101     | 50            | 15            | 36      | 24.2 | 0.526            | 26  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 54 | 5           | 88      | 66            | 21            | 23      | 24.4 | 0.342            | 30  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 55 | 8           | 176     | 90            | 34            | 300     | 33.7 | 0.467            | 58  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 56 | 7           | 150     | 66            | 42            | 342     | 34.7 | 0.718            | 42  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 57 | 1           | 73      | 50            | 10            | 0       | 23   | 0.248            | 21  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 58 | 7           | 187     | 68            | 39            | 304     | 37.7 | 0.254            | 41  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 59 | 0           | 100     | 88            | 60            | 110     | 46.8 | 0.962            | 31  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 60 | 0           | 146     | 82            | 0             | 0       | 40.5 | 1.781            | 44  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 61 | 0           | 105     | 64            | 41            | 142     | 41.5 | 0.173            | 22  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 62 | 2           | 84      | 0             | 0             | 0       | 0    | 0.304            | 21  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 63 | 8           | 133     | 72            | 0             | 0       | 32.9 | 0.27             | 39  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 64 | 5           | 44      | 62            | 0             | 0       | 25   | 0.587            | 36  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 65 | 2           | 141     | 58            | 34            | 128     | 25.4 | 0.699            | 24  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 66 | 7           | 114     | 66            | 0             | 0       | 32.8 | 0.258            | 42  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 67 | 5           | 99      | 74            | 27            | 0       | 29   | 0.203            | 32  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 68 | 0           | 109     | 88            | 30            | 0       | 32.5 | 0.855            | 38  | 1       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |
| 69 | 2           | 109     | 92            | 0             | 0       | 42.7 | 0.845            | 54  | 0       |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |

**Fig4. Backend(EXCEL API)**

# CHAPTER 4

## METHODOLOGY

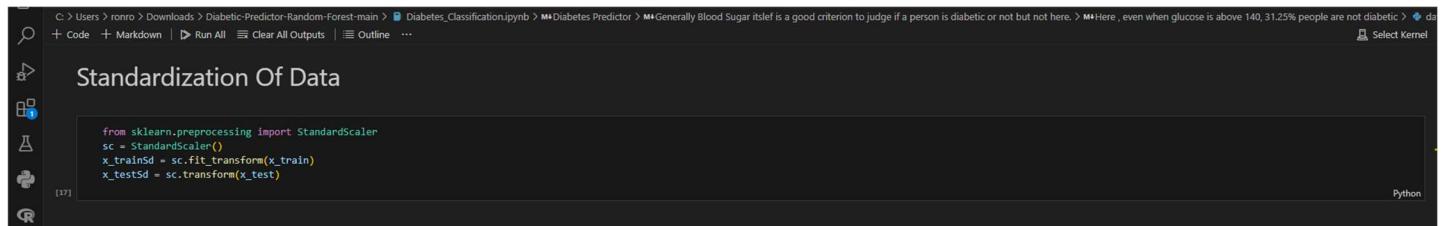
### 4.1 BASIC INTRODUCTORY

The system architecture for the diabetes prediction tool using the Rf algorithm can be divided into several components:

1. **Data collection and preprocessing:** The first component involves the collection collects clinical and demographic data from a variety of sources, such as electronic health records, medical devices, and patient self-assessments. The data are then preprocessed to remove missing values, outliers, and redundant features.

| Serial no | Attribute Names            | Description                      |
|-----------|----------------------------|----------------------------------|
| 1         | Pregnancies                | Number of times pregnant         |
| 2         | Glucose                    | Plasma glucose concentration     |
| 3         | Blood Pressure             | Diastolic blood pressure         |
| 4         | Skin Thickness             | Triceps skin fold thickness (mm) |
| 5         | Insulin                    | 2-h serum insulin                |
| 6         | BMI                        | Body mass index                  |
| 7         | Diabetes pedigree function | Diabetes pedigree function       |
| 8         | Outcome                    | Class variable (0 or 1)          |
| 9         | Age                        | Age of patient                   |

**2. Model Training:** The preprocessed data is used to train the Rf model using a machine learning framework like scikit-learn in Python. The model is optimized using techniques such as hyperparameter tuning and cross-validation to improve model accuracy.



A screenshot of a Jupyter Notebook interface. The title bar shows the path: C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > Diabetes\_Classification.ipynb. The main area has a dark theme and displays the following code:

```
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
x_trainSc = sc.fit_transform(x_train)
x_testSc = sc.transform(x_test)
```

The code imports the StandardScaler from sklearn.preprocessing, creates an instance of it, and then applies it to both the training and testing datasets. The cell number [17] is visible at the bottom left of the code block.

**Fig5. Importing SKLearn**

**3. Model deployment:** The trained model is deployed to a web or mobile application where users can enter their clinical and demographic data and receive risk prediction's chance of diabetes. The model can be deployed on cloud platforms such as AWS or Azure to ensure scalability and availability.

app.py diabetes.csv Diabetes\_Classification.ipynb index.html profile README.md requirements.txt

C:\Users\niraj>Downloads\Diabetic-Predictor-Random-Forest-main> Diabetes\_Classification.ipynb > Diabetes Predictor > Generally Blood Sugar itself is a good criterion to judge if a person is diabetic or not but not here. > Here, even when glucose is above 140, 31.25% people are not diabetic > Select Kernel

+ Code + Markdown | Run All | Clear All Outputs | Outline ...

## Correlation

```

import seaborn as sns
import matplotlib.pyplot as plt

corrmat = data.corr()
top_corr_features = corrmat.index
plt.figure(figsize=(20,20))

#heat map

g= sns.heatmap(data[top_corr_features].corr(),annot=True,cmap="RdYlGn")

```

[7] Python

```

data2 = data[data['glucose'] > 140]

diabetes = len(data2.loc[data['result'] == 1])
no_diabetes = len(data2.loc[data['result'] == 0])

```

[4] Python

```

(diabetes,no_diabetes)

```

[5] Python

```

... (132, 60)

```

```

data.isnull().values.any()

```

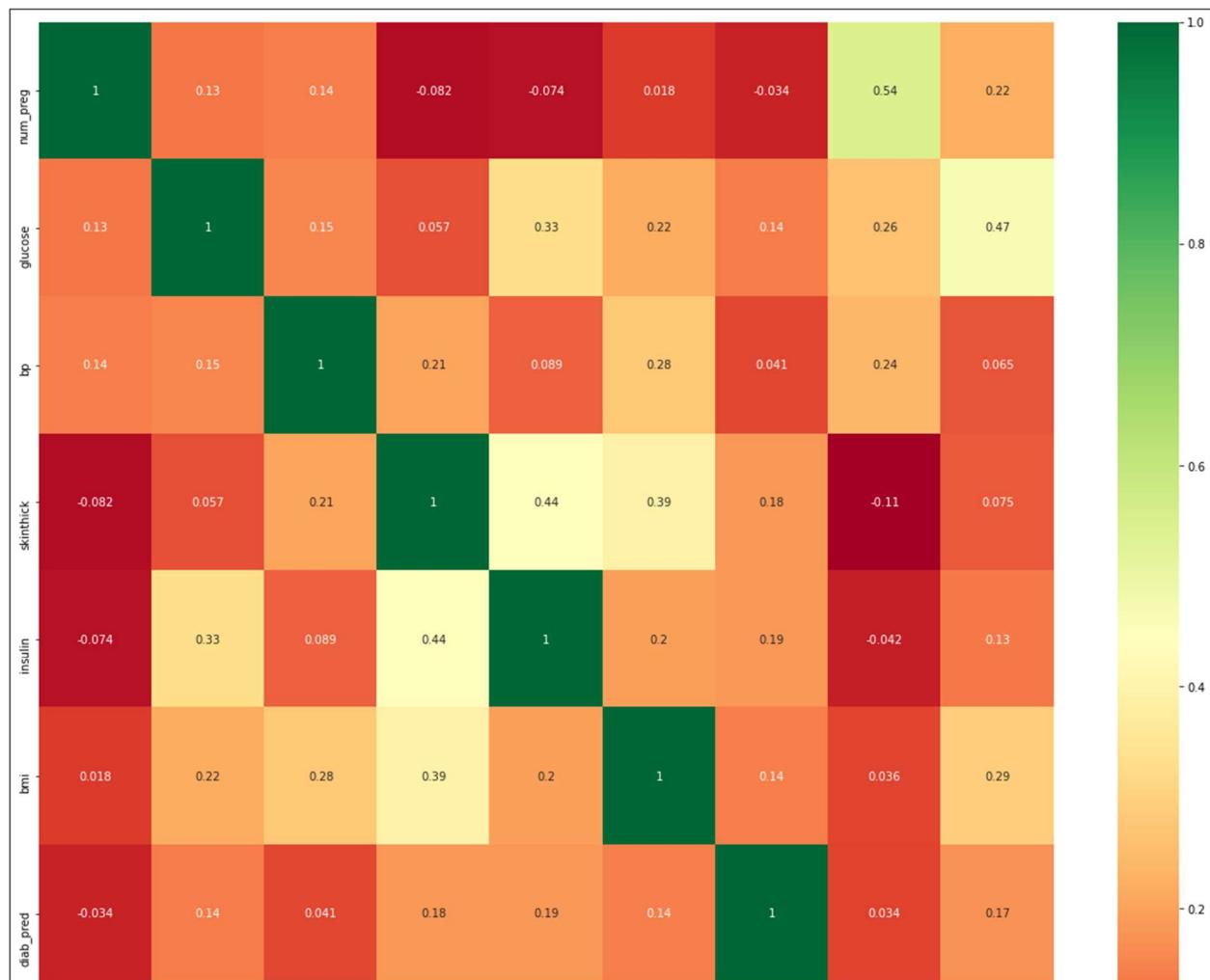
[6] Python

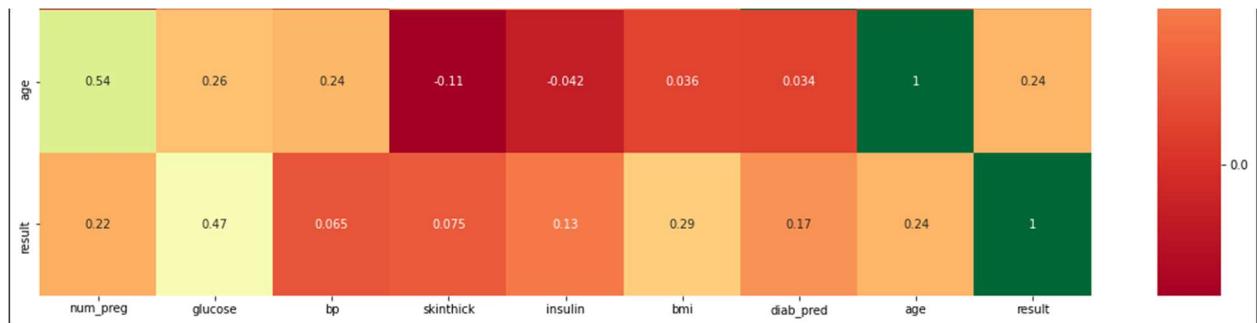
```

... False

```

**Fig6. Correlating the Model**





**Fig7. Data Model**

4. **User interface:** The user interface component provides a user-friendly interface for users to input their data and receive the prediction. The interface can be designed using web development tools such as React, Angular, or Vue.js, or mobile app development tools such as Flutter or React Native.

5. **Data privacy and security:** The system should ensure the privacy and security of patient data by implementing secure data storage, encryption, and access controls. It should also comply with relevant regulations such as HIPAA or GDPR. In general, the system architecture of the diabetes prediction machine using the Rf algorithm includes data collection and preprocessing, training the machine learning model, deploying the model to the application web or mobile application, providing a user interface and ensuring data privacy and security.

## 4.2 ALGORITHMS USED

**Random forest algorithm in the Prognosis Dia:** The random forest algorithm is a flexible and powerful machine learning method that can be especially useful in predicting the risk and outcomes of diabetes.

In project, the key element is to process the input data and calculate the predicted diabetes risk.

### Here's how it works:

- **Ensemble learning:** The Random Forest algorithm is an ensemble learning technique, meaning it combines predictions from multiple decision trees to produce a single correct result more and stronger. Each decision tree in the ensemble is trained on a random subset of data and features.
- **Input Data:** The algorithm takes user-supplied input data attributes from the user interface, including "Pregnancy", "Glucose", "Blood Pressure" and others, as well as historical data used to train the model.

- **Training:** The initial Random Forest model is trained on a diverse dataset consisting of known outcomes (e.g. diabetes or non-diabetes) associated with multiple sets of attributes calculate different inputs. This historical data is used to create a set of decision trees, each of which learns different patterns and relationships in the data.

## Fig8. Model Training

- **Predictions:** When users submit their data through the user interface, the algorithm uses the entire decision tree to make predictions. It evaluates how well the input data fits the models and relationships identified during training and calculates the diabetes risk probability.

- **Aggregate:** The predictions from each decision tree are combined through a voting mechanism. In the context of project, this means that the algorithm predicts the most likely outcome (diabetes or non-diabetes) based on the collective decision of the trees. This ensemble approach reduces the risk of overfitting and enhances the model's generalization capabilities.
- **Scoring:** The Random Forest algorithm can also provide a probability score, indicating the confidence level of the prediction. This score can help users understand the reliability of the prediction.
- **Personalization:** One of the key advantages of the Random Forest algorithm is its ability to accommodate personalized data. In project, this means it can take into account each user's unique attributes, such as age, blood glucose level, etc., to make personalized predictions .
- **Continuous Learning:** The model can be periodically updated with new data, ensuring that it remains accurate and relevant over time.

This is important for monitoring changes in a patient's health status and adjusting predictions as more information becomes available.

The Random Forest algorithm is a valuable tool in Diabetes Prediction project because it can handle complex, multidimensional data and make accurate predictions about diabetes risk road. Its global nature, ability to handle personal data and adaptability make it well suited to the dynamic nature of healthcare and personalized medicine.

It enables the project to deliver meaningful and useful insights to healthcare professionals and patients, facilitating informed decision-making and mainstream healthcare management dynamic.

## **Karnaugh Map**

A Karnaugh map, commonly known as a K-map, is a graphical representation of a truth table used in the design and simplification of digital logic.

K-maps are a valuable tool for minimizing Boolean expressions, optimizing logic circuits, and simplifying complex Boolean functions.

### **Here is an overview of K maps and their applications:**

- **Grid representation:** Karnaugh maps consist of a grid, where each grid cell represents a possible combination of input variable in binary format. The number of rows and columns in the grid depends on the number of input variables. In general, for  $n$  input variables, a  $2^n$  grid is created.
- **Truth table mapping:** The grid values are populated by mapping the corresponding binary values to the output of the truth table. Each cell contains the binary input values as well as the output for that input combination.
- **Clustering and simplification:** K-map is used to identify groups of adjacent cells with the same output value. These groups are called “min terms” or “max terms”. By grouping these cells, you can simplify Boolean expressions. The size and arrangement of the groups depends on the logic function to be minimized. The goal is to reduce the number of terms and/or variables in a Boolean expression while maintaining its logical behavior

- **Minimize Boolean functions:** K-maps are mainly used to minimize Boolean functions, making logic circuits more efficient and cost-effective. By grouping cells into K-maps, you can identify common factors and create simplified expressions for logic functions. They are especially useful for minimizing the number of gates and input variables, reducing power consumption, and optimizing circuit speed.
- **Advantages:** K-maps offer a visual, intuitive, and systematic approach to simplifying Boolean functions. They help designers quickly grasp the relationships between different input combinations. They ensure that simplification is done optimally by finding the largest groups possible to minimize the expression.
- **Limit:** K-maps are most effective for small to moderately complex Boolean functions. For very large functions, manual mapping can become tedious and software tools may be preferred.

The effectiveness of K-map depends on the choice of grouping variables, which may require some expertise.

## **Logistic Regression**

A Logistic Regression algorithm is used for prediction.

Here's a description of how the Logistic Regression algorithm can be applied to this project:

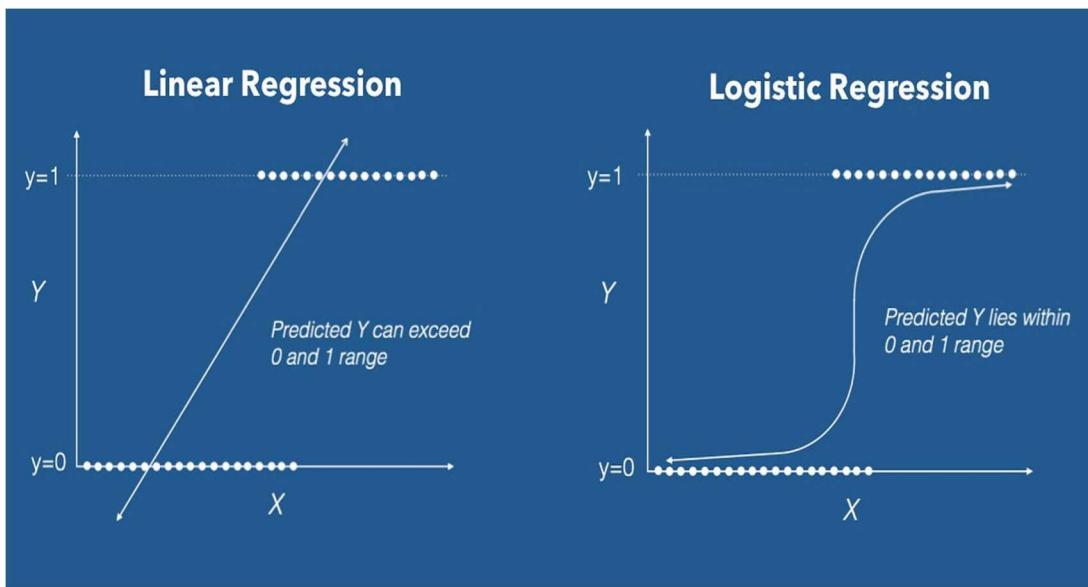
**Logistic Regression Algorithm in the Diabetic Predictor Project:** Logistic Regression is a widely used classification algorithm in machine learning, and it's particularly useful for binary classification problems like predicting whether a person is diabetic or not in project.

**Here's how it works in this context:**

- **Binary Classification:** Logistic Regression is well-suited for binary classification problems where the goal is to categorize data into one of two classes. In case, it's used to predict whether a person is diabetic (class 1) or not (class 0).
- **Input Data:** The algorithm takes input data, including various features such as the number of pregnancies, glucose levels, blood pressure, and other attributes, as independent variables. These attributes act as predictors in the model.
- **Logistic Function:** Logistic regression uses a logistic function (also known as a sigmoid function) to transform a linear combination of input features into probability values. This function assigns a value between 0 and 1 to a linear combination. This represents the probability of belonging to the positive class (diabetes).
- **Training the model:** In the training phase, the algorithm adjusts the parameters (weights and biases) to find the best fit for the training data. Use optimization techniques such as gradient descent to minimize logistic loss (also known as log loss or cross-entropy loss) by adjusting model parameters.
- **Decision Boundaries:** Logistic regression uses decision boundaries to classify data points. The project sets decision boundaries based on a predefined threshold (typically 0.5). If the predicted probability is greater than the threshold, the model

classifies the person as having diabetes. Otherwise, it is classified as non-diabetic.

- **Model evaluation:** Various metrics are used to evaluate the performance of a logistic regression model, such as precision, recall, F1 score, and ROC curve. These metrics can be used to assess how well the model can distinguish between diabetic and non-diabetic patients.
- **Probability values:** Logistic regression also provides probability values for predictions. The project can interpret these values as the probability that a person has diabetes. For example, a value of 0.8 means you have an 80% chance of developing diabetes.
- **Feature Importance:** Logistic regression provides insight into the importance of each feature when making predictions. This information can help you understand which input variables have the greatest impact on determining diabetes risk.
- **Real-time predictions:** Once the logistic regression model is trained, it can generate real-time predictions based on new input data provided by the user through the web application.



**Fig9. Diff for Linear & Logistic**

$$\text{Accuracy} = \frac{(TP+TN)}{(TP+TN+FN+FP)}$$

$$\text{Recall} = \frac{(TP)}{(TP+FN)}$$

$$\text{Precision} = \frac{(TP)}{(TP+FP)}$$

$$F - \text{measure} = \frac{2 \times (\text{Precision} \times \text{Recall})}{(\text{Precision} + \text{Recall})}$$

| Classification | Precision | Recall | F-measure | Accuracy |
|----------------|-----------|--------|-----------|----------|
| RF (K-fold)    | 0.744     | 0.750  | 0.746     | 74.69%   |
| RF (Splitting) | 0.779     | 0.771  | 0.774     | 77.14%   |
| LR (K-fold)    | 0.761     | 0.768  | 0.761     | 76.82%   |
| LR (Splitting) | 0.788     | 0.789  | 0.788     | 78.85%   |

**Fig10. Formulating Regression**

## 4.3 MAIN CODE

- PyFlask

```
C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > app.py > predict
 1  from flask import Flask, render_template, request
 2  import jsonify
 3  import requests
 4  import pickle
 5  import numpy as np
 6  import sklearn
 7  from sklearn.preprocessing import StandardScaler
 8  app = Flask(__name__)
 9  model = pickle.load(open('diabetes.pkl', 'rb'))
10  @app.route('/', methods=['GET'])
11  def Home():
12      return render_template('index.html')
13
14
15  standard_to = StandardScaler()
16  @app.route("/predict", methods=['POST'])
17  def predict():
18      #we read all the data required and we fit this data into the pickle we
19      #have made
20      if request.method == 'POST':
21          num_preg = int(request.form['num_preg'])
22          glucose = int(request.form['glucose'])
23          bp = int(request.form['bp'])
24          skinthick = int(request.form['skinthick'])
25          insulin = int(request.form['insulin'])
26          bmi = float(request.form['bmi'])
27          diab_pred = float(request.form['diab_pred'])
28          age = int(request.form['age'])
29
30          features = np.array([[num_preg,glucose,bp,skinthick,insulin,bmi,diab_pred,age]])
31          util_output = model.predict(features)
32          output = util_output[0];
33          if output==0:
34              return render_template('index.html',prediction_text="Not Diabetic")
35          elif output==1:
36              return render_template('index.html',prediction_text="Diabetic (81 percent chance)")
37          else:
38              return render_template('index.html',prediction_text = "Data given cannot determine result")
39      else:
40          return render_template('index.html')
41
42  if __name__=="__main__":
43      app.run(debug=True)
44
```

## Headers for the PyFlask

```
C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > requirements.txt
 1 appdirs==1.4.4
 2 argh==0.26.2
 3 async-generator==1.10
 4 atomicwrites==1.4.0
 5 black==19.10b0
 6 brotlipy==0.7.0
 7 certifi==2021.5.30
 8 cycler==0.10.0
 9 cytoolz==0.11.0
10 dask-searchcv==0.2.0
11 Flask==2.0.1
12 future==0.18.2
13 gunicorn==20.1.0
14 HeapDict==1.0.1
15 ipython-genutils==0.2.0
16 itsdangerous==2.0.1
17 jsonify==0.5
18 locket==0.2.1
19 mccabe==0.6.1
20 mypy-extensions==0.4.3
21 olefile==0.46
22 pandocfilters==1.4.2
23 parso==0.7.0
24 pathspec==0.7.0
25 pyparsing==2.4.7
26 python-dateutil==2.8.1
27 pywin32-ctypes==0.2.0
28 pywinpty==0.5.7
29 PyYAML==5.4.1
30 QDarkStyle==2.8.1
31 QtPy==1.9.0
32 scikit-image==0.18.1
33 Send2Trash==1.5.0
34 terminado==0.9.4
35 webencodings==0.5.1
36 Werkzeug==2.0.1
37 wincertstore==0.2
38 wrapt==1.12.1
39 zict==2.0.0
40 requests
41 sklearn
42 numpy
```

- HTML Code

```
C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > index.html > html > body > style > .result
1  <!DOCTYPE html>
2  <html lang="en">
3
4  <head>
5      <meta charset="UTF-8">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7      <link rel="preconnect" href="https://fonts.gstatic.com" crossorigin>
8      <link href="https://fonts.googleapis.com/css2?family=Montserrat:wght@500&family=Roboto+Condensed&display=swap" rel="stylesheet">
9      <link href="data:image/x-icon;base64,AAABAAEAEBAQAAEABAAoQAAFgAACgAAAAQAAAIAAAAEBAAAAAAgAAAAAAGAAAAAAEAAAAAAAwP/AAA
10     <title>Diabetes Predictor</title>
11 </head>
12
13 <body>
14
15     <div>
16         <h3 class="result">{{ prediction_text }}</h3>
17         <form action="{{ url_for('predict') }}" method="post" class="form">
18             <h2>DIABETES PREDICTOR</h2>
19             <h5>HOLD MOUSE ON TITLES WHICH ARE CONFUSING</h5>
20             <h3 title="Plasma glucose concentration a 2 hours in an oral glucose tolerance test">Blood Sugar Level</h3>
21                 <input name="glucose" required = "true">
22             <h3 title="If Blood Pressure is 120/80, enter 80 here">Diastolic Blood Pressure</h3>
23                 <input name="bp" required = "true">
24             <h3 title ="Body mass index (weight in kg/(height in m)^2)">Body Mass Index (xx.x)</h3>
25                 <input name="bmi" required = "true">
26             <h3 title="2-Hour serum insulin (mu U/ml)">Insulin</h3>
27                 <input name="insulin" required = "true">
28             <h3>Age (in Years)</h3>
29                 <input name="age" required="true">
30             <h3>Number Of Pregnancies</h3>
31                 <input name="num_preg" required = "true">
32             <h3 title="Triceps skin fold thickness (mm)">Skin Thickness</h3>
33                 <input name="skinthick" required = "true">
34             <h3 title = "A function which stores the likelihood of diabetes in family history">Diabetes Pedigree Function (0.xxx)</h3>
35                 <input name="diab_pred" required="true">
36
37                 <button class="submitButton" type="submit" >Check For Diabetes Now</button>
38             </form>
39             <br><br>
40         </div>
41
```

```
C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > ◇ index.html > ◇ html > ◇ body > ◇ style > ◇ .result
42
43
44
45     <style>
46         body {
47             height: max-height;
48             background-color: ■#ffffff;
49             background-image: url("data:image/svg+xml,%3Csvg xmlns='http://www.w3.org/2000/svg' width='100%25'%3E%3Cdefs%3E%3Clinear
50             background-attachment: fixed;
51             background-size: cover;
52             display: flex;
53             justify-content: center;
54             align-items: center;
55             font-family: 'Roboto Condensed', sans-serif;
56         }
57         input {
58             color: □#333;
59             font-size: 1.2rem;
60             margin: 0 auto;
61             padding: 1rem 1rem;
62             border-radius: 0.3rem;
63             background-color: ■#FFCCCB;
64             border: none;
65             width: 90%;
66             display: block;
67             border-bottom: 0.3rem solid transparent;
68             transition: all 0.3s;
69             outline: none;
70             margin-bottom: 25px;
71             font-family: 'Roboto Condensed', sans-serif;
72         }

```

```
C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > ◇ index.html > ◇ html > ◇ body > ◇ style > ◇ .result
73     h2{
74         color : □#3B0918;
75         font-size : 1.75rem;
76         margin : auto;
77         text-align: center;
78         font-weight : strong;
79         text-shadow: 0 2px 2px □black;
80     }
81     h3{
82         color : □#3B0918;
83         font-size : 1.3 rem;
84         display : block;
85         margin-bottom: 10px;
86         margin-left: 5px;
87     }
88     h5{
89         color : □#3B0918;
90         text-align: center;
91     }
92     .submitButton{
93         border-radius: 4px;
94         border: none;
95         background-color: ■#FFCCCB;
96         color: □#3B0918;
97         text-align: center;
98         text-transform: uppercase;
99         font-size: 22px;
100        padding: 1.5rem 1.5rem;
101        transition: all 0.4s;
102        cursor: pointer;
103        width: 100%;
104        display: block;
105        margin-top : 20px;
106        font-family: 'Roboto Condensed', sans-serif;
107    }

```

```
C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > <> index.html > ⌂ html
107         }
108     .submitButton:hover {
109         background-color : □#228b22;
110         color : ■white;
111     }
112     .form {
113         background-color: ■#C85250;
114         box-shadow: 0 20px 20px □black;
115         width: 400px;
116         padding : 1rem 1rem;
117     }
118     .result {
119         color : ■white;
120         text-align : center;
121         font-size : 2rem;
122     }
123
124     </style>
125 </body>
126
127 </html>
```

## Creating Pickle File

```
import pickle
filename = 'diabetes.pkl'
pickle.dump(final_model, open(filename, 'wb'))
```

[24]

### 4.4 MAIN CODE EXPLANATION

The web application allows users to enter certain medical and health data and the model will predict whether the user has diabetes based on this data.

**Here is an explanation of the code:**

- **Import libraries:** The code starts by importing the necessary libraries and modules, including Flask to create the web application, select to load the machine learning model pre-trained, neat for number of operations and scikit-learn (sklearn) for pre-processing the data.
- **Initializing Flask application:** app = Flask(\_\_name\_\_) initializes Flask web application.
- **Loading machine learning model:** model = pickles.load(open('diabetes.pkl', 'rb')) loads pre-trained machine learning model to predict diabetes from pickles file typically created during model training.
- **Create route home:** @app.route('/',methods=['GET']) set route to root URL ("/"). When the user visits the root URL, the Home() function will be called. The Home() function displays an HTML template named "index.html", which is the main web page where users can enter data.
- **Data preprocessing:** Standard\_to = StandardScaler() creates a StandardScaler object to normalize the input data. Normalization ensures that the input data is converted to a mean of 0 and a standard deviation of 1. Generate a predicted route: @app.route("/predict", METHODS=[' POST']) defines a route where users can send data to prediction using the HTTP POST method. This route is typically triggered when a user submits form data on the website. The predict() function was called when accessing this route.

- **Predict function:** In `predict()` function, it checks whether the HTTP method used is POST or not. If so, the function retrieves data from the form fields in the HTML template. The function then preprocesses the input data by normalizing it using a `StandardScaler` object. Standardized data is fed to the machine learning model and predictions are derived using `model.predict(feature)`. The prediction is then used to determine whether the user has diabetes.
- **Prediction returned:** Function that uses `prediction` to display a message on the web page indicating whether the user has diabetes or not. Result message displayed on template "index.html".
- **Running the application:** The `if __name__ == "__main__"` code block ensures that the Flask application only runs when the script is run directly (not when it is imported as a module).
- **Run the app in debug mode:** `app.run(debug=True)` runs the Flask app in debug mode, allowing real-time code changes and debugging during development. In short, this code configures a web application that uses Flask to deploy a diabetes prediction model. Users can input their medical data, which is then processed and fed into a machine learning model to predict whether they have diabetes or not. The results are displayed on the website.

# CHAPTER 5

## RESULTS AND DISCUSSION

### 5.1 PERFORMANCE ANALYSIS

Performance analysis of a project, especially a machine learning project, typically involves evaluating key metrics to assess the effectiveness and efficiency of the system.

In the context of a diabetes prediction project, here are some aspects of performance analysis to consider.

- **Accuracy and Precision:** Evaluate the accuracy of a diabetes prediction model by comparing its predictions to known results. Calculate the percentage of true positive and true negative predictions. Accuracy measures the accuracy of positive predictions. If a patient claims to have diabetes, it is important to assess how often the model accurately predicts diabetes.
- **Recall and Sensitivity:** Recall (sensitivity) measures the ability of a model to correctly identify all cases of diabetes. It is important to evaluate whether the model has missing cases (false negatives).
- **Specificity:** Specificity measures the ability of a model to correctly identify non-diabetic cases. It is important to assess the extent to which the model does not misclassify non-diabetic cases as diabetic.
- **F1 score:** The F1 score combines precision and recall into one metric to provide a balanced assessment of a model's performance.

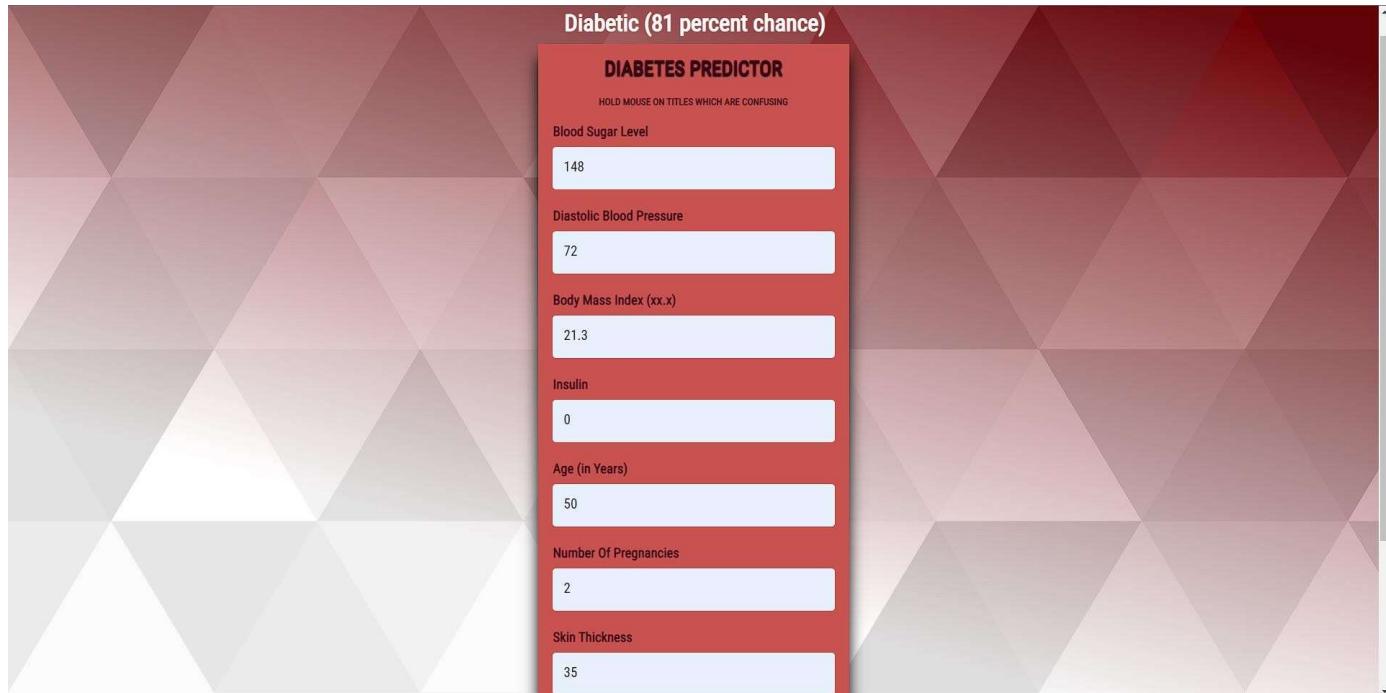
- **ROC curves and AUC:** Receiver operating characteristic (ROC) curves are useful for evaluating the tradeoff between true positive rate and false positive rate at different classification thresholds. The area under the ROC curve (AUC) summarizes the overall performance of the model. A higher AUC indicates better performance.
- **Confusion Matrix:** The confusion matrix is a useful tool for understanding model performance. Displays true positives, true negatives, false positives, and false negatives, providing insight into the strengths and weaknesses of the model.
- **Cross-validation:** Evaluate model performance across different subsets of data using techniques such as k-fold cross-validation. This helps check the stability and generalizability of the model.
- **Data Split:** Split the dataset into training, validation, and test sets. Evaluate the model's performance on data not seen during training to ensure that the model generalizes well.
- **Performance over time:** If the project is deployed in a live environment, monitor its performance over time. Check whether the model's accuracy is stable or decreasing, and update the model if necessary.
- **Efficiency:** Evaluate the computational efficiency of models and web applications. Consider metrics such as response time, resource usage, and memory usage.
- **User Feedback:** Collect feedback from users such as: Healthcare professionals and patients understand satisfaction and the real-world impact of the project.

- **Privacy and Security:** Assess the project's ability to maintain privacy and security, especially when sensitive health information is involved.
- **Regulatory Compliance:** Ensure that the project complies with relevant health regulations, including HIPAA and all other applicable legal standards.

Performance analysis is a continuous process that requires continuous monitoring and evaluation to ensure that the project is achieving its intended goals and delivering value to users.

Depending on the results of analysis, you may need to refine model, update web application, or make other improvements to improve performance.

## 5.2 PERFORMANCE WORKING SCREENSHOTS



The screenshot displays a mobile application interface titled "DIABETES PREDICTOR". At the top, it shows a prediction result: "Diabetic (81 percent chance)". Below this, there is a note: "HOLD MOUSE ON TITLES WHICH ARE CONFUSING". The interface consists of several input fields with placeholder values:

| Parameter                | Value |
|--------------------------|-------|
| Blood Sugar Level        | 148   |
| Diastolic Blood Pressure | 72    |
| Body Mass Index (xx.x)   | 21.3  |
| Insulin                  | 0     |
| Age (in Years)           | 50    |
| Number Of Pregnancies    | 2     |
| Skin Thickness           | 35    |

**Fig11. Working Model(UI)**

```
C: > Users > ronro > Downloads > Diabetic-Predictor-Random-Forest-main > diabetes.csv
1   Pregnancies,Glucose,BloodPressure,SkinThickness,Insulin,BMI,DiabetesPedigreeFunction,Age,Outcome
2   6,148,72,35,0,33.6,0.627,50,1
3   1,85,66,29,0,26.6,0.351,31,0
4   8,183,64,0,0,23.3,0.672,32,1
5   1,89,66,23,94,28.1,0.167,21,0
6   0,137,40,35,168,43.1,2.288,33,1
7   5,116,74,0,0,25.6,0.201,30,0
8   3,78,50,32,88,31,0.248,26,1
9   10,115,0,0,0,35.3,0.134,29,0
10  2,197,70,45,543,30.5,0.158,53,1
11  8,125,96,0,0,0,0.232,54,1
12  4,110,92,0,0,37.6,0.191,30,0
13  10,168,74,0,0,38,0.537,34,1
14  10,139,80,0,0,27.1,1.441,57,0
15  1,189,60,23,846,30.1,0.398,59,1
16  5,166,72,19,175,25.8,0.587,51,1
17  7,100,0,0,0,30,0.484,32,1
18  0,118,84,47,230,45.8,0.551,31,1
19  7,187,74,0,0,29.6,0.254,31,1
20  1,183,30,38,83,43.3,0.183,33,0
21  1,115,70,30,96,34.6,0.529,32,1
22  3,126,88,41,235,39.3,0.704,27,0
23  8,99,84,0,0,35.4,0.388,50,0
24  7,196,90,0,0,39.8,0.451,41,1
25  9,119,80,35,0,29,0.263,29,1
26  11,143,94,33,146,36.6,0.254,51,1
27  10,125,70,26,115,31.1,0.205,41,1
28  7,147,76,0,0,39.4,0.257,43,1
29  1,97,66,15,140,23.2,0.487,22,0
30  13,145,82,19,110,22.2,0.245,57,0
```

**Fig12. Backend Database**

## CHAPTER 6

### CONCLUSION AND FUTURE SCOPE

In summary, PrognosisDia with Rf algorithm can be an effective tool for early detection and treatment of diabetes. By analyzing clinical and demographic data, this algorithm can predict an individual's risk of diabetes and inform personalized treatment plans.

However, there is always room for improvement and future expansion of the system.

Possible future improvements include:

- **Integration with wearable devices:** This system can be improved by integrating with wearable devices that can collect real-time data such as heart rate, blood sugar levels, and physical activity. This provides more accurate predictions and allows patients to monitor their health status in real time
- **Inclusion of genetic data:** Genetic data can provide valuable information about an individual's predisposition to diabetes. Incorporating genetic data into the analysis can extend the system and improve prediction accuracy.
- **Integration with electronic health records:** The system can be enhanced by integrating it with electronic health records to access a patient's complete medical history. This can provide a more comprehensive view of the patient's health and help identify other health conditions that may affect their risk of developing diabetes.
- **Personalized recommendations:** The system can be enhanced by providing personalized recommendations based on an individual's risk factors and health history. This can help patients take preventative measures to reduce their risk of developing diabetes.

- **Continual learning:** The system can be enhanced by implementing continual learning techniques to improve the model's accuracy over time.

This involves updating the model with new data and retraining it periodically to capture any changes in the patient population or the risk factors associated with diabetes.

In summary, a diabetic predictor using the Rf algorithm has the potential to improve diabetes prevention and management.

However, future enhancements such as integration with wearable devices and electronic health records, incorporation of genetic data, personalized recommendations, and continual learning can further improve the accuracy and effectiveness of the system.

## Future Scope

The future scope of your Diabetic Predictor project can be quite broad, as it involves healthcare, data analysis, and machine learning. Here are some potential areas for future development and expansion:

- **Enhanced Predictive Models:**

- Continuously improve and fine-tune the predictive model. Incorporate more features and data sources for more accurate predictions.

- **Personalized Health Recommendations:**

- Develop a feature that provides personalized health recommendations based on the user's health data and risk assessment. These recommendations can include dietary advice, exercise routines, and lifestyle modifications.

- **Integration with Health Devices:**

- Integrate the application with wearable devices and health monitoring tools to automate data input and provide real-time feedback and alerts to users.

- **Telehealth Services:**

- Expand the project to include telehealth services, allowing users to connect with healthcare professionals for virtual consultations and monitoring.

- **Mobile Application:**

- Develop a mobile application version of the project to increase accessibility and convenience for users.

- **Data Privacy and Security:**

- Continuously update and enhance data privacy and security measures to comply with evolving healthcare regulations and standards.

- **Multilingual Support:**

- Offer the application in multiple languages to make it accessible to a broader audience.

- **Data Visualization and Reporting:**

- Implement interactive data visualization tools and generate comprehensive health reports for users.

- **Machine Learning Updates:**

- Stay up-to-date with the latest machine learning and artificial intelligence techniques to further improve prediction accuracy and model performance.

- **Population Health Analysis:**
  - Analyze aggregated and anonymized user data to gain insights into population health trends and contribute to public health research.
- **Integration with Electronic Health Records (EHR):**
  - Explore the integration of your project with EHR systems to provide healthcare professionals with a comprehensive view of a patient's health history.
- **Collaboration with Healthcare Institutions:**
  - Collaborate with healthcare providers, clinics, and hospitals to incorporate the project as part of their healthcare services.
- **Research and Publications:**
  - Contribute to medical research by using the project's data for scientific studies and publications.
- **User Education and Engagement:**
  - Develop educational content and engagement features to empower users with knowledge and encourage them to take control of their health.
- **Scalability and High Availability:**
  - Ensure the system can handle increased user loads as it grows, with high availability and robust infrastructure.
- **Regulatory Compliance:**
  - Stay updated with healthcare regulations and standards to ensure continued compliance with data protection and privacy laws.

- **AI Chatbots:**

- Implement AI-powered chatbots to provide immediate responses to user queries and facilitate health-related conversations.

- **Geospatial Analysis:**

- Use location-based data to offer localized health recommendations, clinics, and resources.

- **AI for Early Detection:**

- Explore the use of advanced AI techniques for early detection of diabetes and related health conditions.

- **Global Expansion:**

- Consider expanding the project to serve users in different regions or countries, accounting for variations in healthcare practices and data sources.

The future scope of your Diabetic Predictor project is not only promising but also aligns with the broader trends in healthcare technology and data-driven personalized health management. Regularly assessing user needs and staying current with advancements in healthcare technology and data analytics will be essential for the project's long-term success.

## **CHAPTER 7**

### **REFERENCES**

- [1] Gauri D. Kalyankar, Shivananda R. Poojara and Nagaraj V. Dharwadkar,” Predictive Analysis of Diabetic Patient Data Using Machine Learning and Hadoop”, International Conference On I-SMAC,978-1-5090-3243-3,2017.
- [2] Ayush Anand and Divya Shakti,” Prediction of Diabetes Based on Personal Lifestyle Indicators”, 1st International Conference on Next Generation Computing Technologies, 978-1-4673-6809-4, September 2015.
- [3] B. Nithya and Dr. V. Ilango,” Predictive Analytics in Health Care Using Machine Learning Tools and Techniques”, International Conference on Intelligent Computing and Control Systems, 978-1-5386-2745-7,2017.
- [4] Dr Saravana kumar N M, Eswari T, Sampath P and Lavanya S,” Predictive Methodology for Diabetic Data Analysis in Big Data”, 2nd International Symposium on Big Data and Cloud Computing,2015.
- [5] Aiswarya Iyer, S. Jeyalatha and Ronak Sumbaly,” Diagnosis of Diabetes Using Classification Mining Techniques”, International Journal of Data Mining & Knowledge Management Process (IJDKP) Vol.5, No.1, January 2015.
- [6] P. Suresh Kumar and S. Pranavi “Performance Analysis of Machine Learning Algorithms on Diabetes Dataset using Big Data Analytics”, International Conference on Infocom Technologies and Unmanned Systems, 978-1-5386-0514-1, Dec. 18-20, 2017.

[7] Mani Butwall and Shraddha Kumar," A Data Mining Approach for the Diagnosis of Diabetes Mellitus using Random Forest Classifier", International Journal of Computer Applications, Volume 120 - Number 8,2015.

d

[8] K. Rajesh and V. Sangeetha, "Application of Data Mining Methods and Techniques for Diabetes Diagnosis", International Journal of Engineering and Innovative Technology (IJEIT) Volume 2, Issue 3, September 2012.

[9]Humar Kahramanli and Novruz Allahverdi,"Design of a Hybrid System for the Diabetes and Heart Disease", Expert Systems with Applications: An International Journal, Volume 35 Issue 1-2, July, 2008.

[10] B.M. Patil, R.C. Joshi and Durga Toshniwal,"Association Rule for Classification of Type-2 Diabetic Patients", ICMLC '10 Proceedings of the 2010 Second International Conference on Machine Learning and Computing, February 09 - 11, 2010.

[11] Dost Muhammad Khan<sup>1</sup>, Nawaz Mohamudally<sup>2</sup>, "An Integration of K-means and Decision Tree (ID3) towards a more Efficient Data Mining Algorithm ", Journal Of Computing, Volume 3, Issue 12, December 2011.