

PHASE 5

PROJECT : AI Based Diabetes Prediction System

DOMAIN : Artificial Intelligence

DATE : 02-11-2023

PROJECT STATEMENT

The problem at hand is to develop an AI-based system for predicting the likelihood of an individual developing diabetes. Diabetes is a chronic health condition that affects millions of people worldwide, and early detection is crucial for effective management and prevention of complications. The AI model should take various health-related features as input and output a prediction regarding the risk of diabetes for an individual. The objective is to create a system that can assist healthcare professionals in identifying at-risk individuals and provide personalized recommendations for prevention and management.

DESIGN THINKING PROCESS

1. **Empathize:** Understand the needs and concerns of potential users and stakeholders, such as healthcare professionals, patients, and researchers. Conduct interviews, surveys, and gather insights to define the problem accurately.
2. **Define:** Clearly define the problem statement, scope, and objectives of the AI-based diabetes prediction system. Create user personas and identify the key features and functionalities required.
3. **Ideate:** Brainstorm potential solutions and features. Consider innovative approaches to data collection, model development, and user interaction. Encourage creative thinking to address the problem effectively.
4. **Prototype:** Develop a prototype or proof of concept to visualize the system's functionality. This may include a user interface for data input and result presentation.
5. **Test:** Test the prototype with potential users and gather feedback. Refine the system based on user input and iterate on the design and functionality.
6. **Develop:** Once the prototype is validated, proceed with the development of the AI model, data pipeline, and user interface.
7. **Test & Iterate:** Continuously test and refine the system throughout development to ensure it meets the needs of the end users.
8. **Deploy:** Deploy the AI-based diabetes prediction system in a real-world healthcare setting, ensuring data security and compliance with regulations.

SOLUTION ARCHITECTURE

Solution architecture is a structured approach to designing complex systems or projects, outlining the components, relationships, and processes to achieve specific goals or solve problems efficiently. It provides a high-level blueprint for project development and implementation.

1. Gather and preprocess data from various healthcare sources, ensuring data quality.
2. Train AI models for diabetes prediction and evaluate their performance.
3. Deploy models in a secure, scalable environment with a user-friendly interface.
4. Continuously monitor and update the system while complying with regulations.
5. Collaborate with healthcare professionals and gather user feedback for improvements.

Example

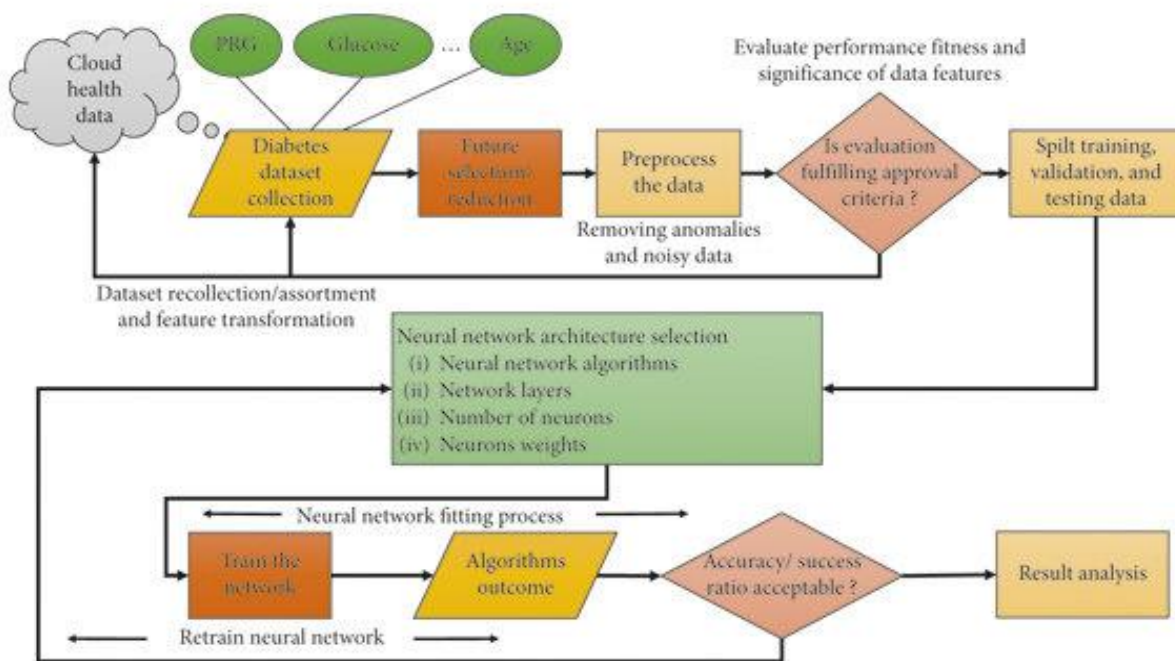


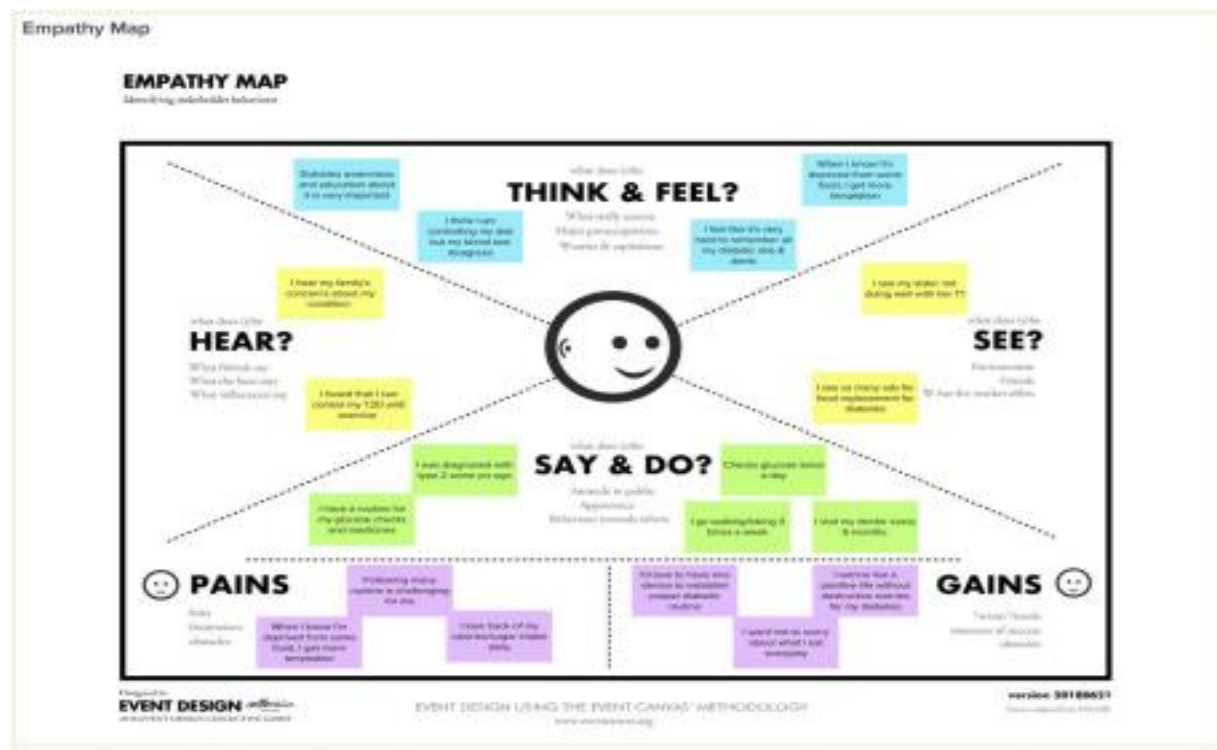
Figure: AI diabetes based prediction system

Reference: <https://www.xenonstack.com/blog/artificial-intelligence-diabetes-detection>

EMPATHY MAP

Empathy Map Canvas:

In creating an empathy map canvas for an AI-based diabetes prediction system, we focus on understanding the typical user, their expressed needs, thoughts, emotions, and actions. For instance, users might vocalize their struggles with blood sugar control, harbor concerns about long-term complications, experience frustration from constant monitoring, and engage in behaviors like blood sugar tracking and medication adherence. This canvas provides a holistic view of the user's perspective, helping guide the development of a more user-centric and effective diabetes prediction system.



1.User Persona: Start by defining a representative user persona, such as a middle aged individual with type 2 diabetes.

2.Label "Says": Create a section labeled "Says" where you note down direct statement or quotes from the user regarding their diabetes management. For example, "I struggle with managing my blood sugar levels."

3.Label "Thinks": In this section, jot down the thoughts, concerns, or worries that the user may have about their diabetes. For instance, "I'm concerned about the long term effects of diabetes."

4. **Label "Feels":** Describe the emotions and feelings that the user experiences in relation to their diabetes. For example, "I feel overwhelmed by the constant monitoring and lifestyle changes."

5. **Label "Does":** Outline the actions and behaviors the user engages in concerning diabetes management. This could include activities like checking blood sugar, taking medications, and making dietary choices.

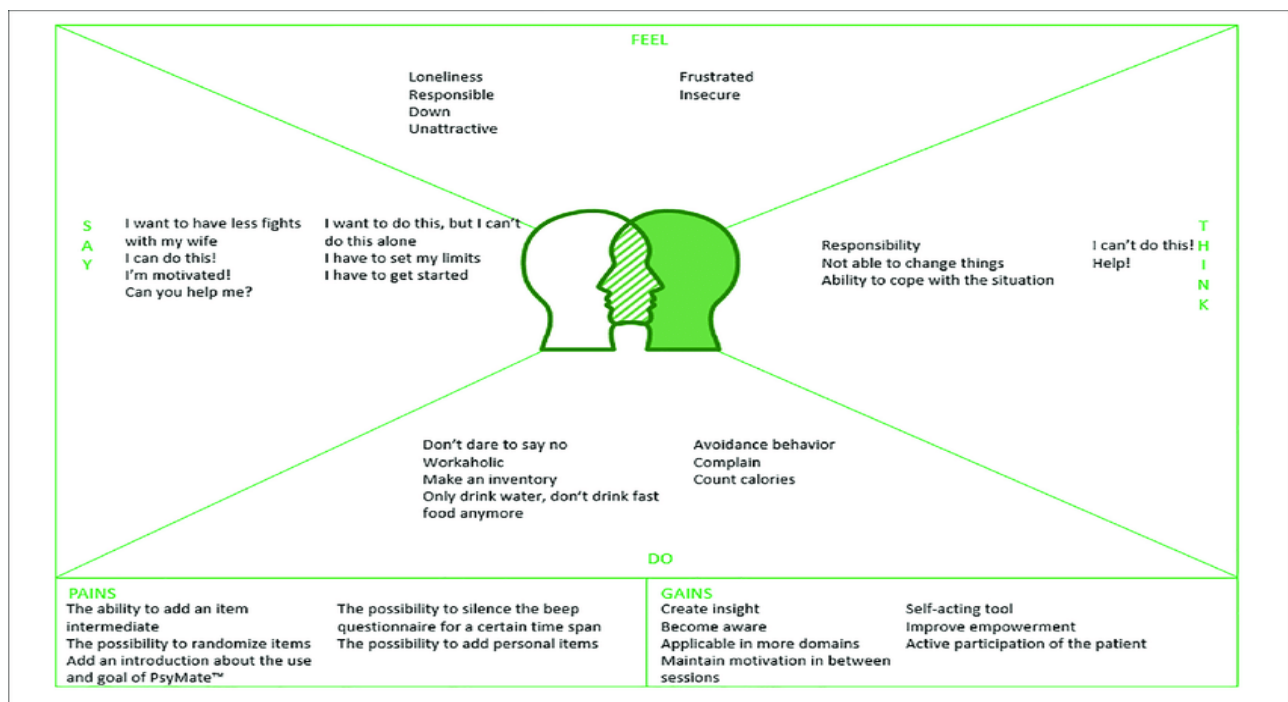
6. **Environment:** Consider the user's physical environment, including where they live and work, as it can influence their diabetes management.

7. **Influences:** Identify the factors and people that influence the user's decisions and behaviors, such as healthcare professionals, family, or friends.

8. **Pain Points:** List the challenges and difficulties the user faces in managing their diabetes, such as the complexity of tracking data or the fear of complications.

9. **Gains:** Note the user's goals, aspirations, and desired outcomes related to their diabetes, such as better control and improved quality of life.

10. **Takeaways:** Summarize the key insights gained from the empathy map that can inform the design and development of the AI-based diabetes prediction system, ensuring it addresses the user's needs, emotions, and actions effectively.



PROPOSED SOLUTION

| | |
|-------------------------------------|--|
| 1. Problem Statement | Diabetes is a prevalent and serious health issue affecting millions of people worldwide. Early detection and management of diabetes are crucial to prevent complications. However, many individuals are unaware of their risk factors or hesitant to get tested. The problem statement should elaborate on the prevalence, impact, and challenges associated with diabetes diagnosis and management. |
| 2. Idea / Solution Description: | Our proposed solution is an AI-based diabetes prediction system. Leveraging machine learning algorithms, this system will analyze relevant medical and lifestyle data to predict the likelihood of an individual developing diabetes. The system will offer personalized risk assessments and recommendations for preventive measures. |
| 3. Novelty / Uniqueness: | The uniqueness of our solution lies in its ability to harness the power of AI and machine learning to predict diabetes risk accurately. It will continuously adapt and improve its predictions as more data becomes available. Additionally, it will provide personalized insights to users, promoting better health choices. |
| 4. Social Impact / Health Benefits: | <ul style="list-style-type: none">- The AI-based diabetes prediction system aims to address several key social and health-related aspects:- Early Detection: By identifying individuals at risk, it promotes early diagnosis and intervention, reducing the likelihood of complications.- Improved Quality of Life: People can make informed lifestyle choices to manage their health and reduce the risk of diabetes. |
| 5. Business Model (Revenue Model): | <ul style="list-style-type: none">- Our revenue model could include:- Subscription-Based Access: Charging users a fee for ongoing access to personalized diabetes risk assessments and recommendations.- Data Monetization: Collaborating with healthcare organizations, research institutions, or pharmaceutical companies to provide anonymized, aggregated data for research purposes. |
| 6. Scalability of the Solution: | <ul style="list-style-type: none">- Our AI-based diabetes prediction system is highly scalable:- Data Integration: It can integrate with various data sources, such as electronic health records, wearable devices, and lifestyle apps.- User Base: The system can accommodate a growing user base without compromising prediction accuracy.- Continuous Improvement: The system will continuously update its algorithms with new data and research, ensuring scalability and adaptability. |

BRAINSTORMING

1. Data Collection & Integration:

- Collect medical records, lifestyle data, and genetic information for comprehensive risk assessment.
- Integrate with wearable devices and apps to provide real-time data.

2. Machine Learning Algorithms:

- Explore different machine learning models (e.g., logistic regression, neural networks) for predicting diabetes.
- Investigate feature engineering techniques to improve model accuracy.

3. User Interface & Experience:

- Develop a user-friendly app or web platform for data input, risk assessment, and recommendations.
- Consider interactive dashboards, visualizations, and alerts.

4. Personalized Recommendations:

- Provide tailored dietary, exercise, and lifestyle recommendations based on risk assessment.
- Offer medication and healthcare provider suggestions when needed.

5. Data Security & Privacy:

- Implement robust security measures to protect user data and maintain compliance with healthcare regulations.
- Ensure transparent data handling and user consent mechanisms.

6. Research Collaboration:

- Collaborate with healthcare institutions and research organizations for data sharing and validation.
- Explore opportunities for contributing to diabetes research.

7. Continuous Learning and Updates:

- Enable the system to learn from user data and research advancements to improve prediction accuracy. Develop a system for providing regular updates and improvements.

DEVELOPMENT PHASE

DESCRIPTION OF THE DATASET:

This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective is to predict based on diagnostic measurements whether a patient has diabetes.

Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.

- Pregnancies: Number of times pregnant
- Glucose: Plasma glucose concentration a 2 hours in an oral glucose tolerance test
- BloodPressure: Diastolic blood pressure (mm Hg)
- SkinThickness: Triceps skin fold thickness (mm)
- Insulin: 2-Hour serum insulin (μ U/ml)
- BMI: Body mass index ($\text{weight in kg}/(\text{height in m})^2$)
- DiabetesPedigreeFunction: Diabetes pedigree function
- Age: Age (years)
- Outcome: Class variable (0 or 1)

DATASET LINK:

<https://www.kaggle.com/datasets/mathchi/diabetes-data-set>

PROJECT OVERVIEW

The AI-Based Diabetes Prediction System is a data-driven healthcare solution designed to predict the likelihood of an individual developing diabetes based on a set of key medical and demographic features. The project's primary objectives include:

- ❖ Data Preprocessing and Loading
- ❖ Data Exploration
- ❖ Model Selection
- ❖ Model Evaluation

1.DATA COLLECTION AND LOADING

The very first step is to choose the dataset for our model .This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective is to predict based on diagnostic measurements whether a patient has diabetes.

Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.

Now we have to set the development environment to build our project. Import the necessary libraries to built the model. We are ensembleing the random forest ,SVM and logistic regression model .

Code:

```
import numpy as np
import pandas as pd
import pickle
from sklearn.ensemble import RandomForestClassifier, VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
```

Code Explanation:

➤numpy (import numpy as np):

NumPy is a fundamental library for scientific computing in Python.

➤pandas (import pandas as pd):

Pandas is a widely used library for data manipulation and analysis. It provides data structures like DataFrame .

➤ import pickle:

Used for saving and loading Python objects, including machine learning models, to and from files.

➤ **from sklearn.ensemble import RandomForestClassifier, VotingClassifier:**

Imports machine learning models, RandomForestClassifier for handling complex data and VotingClassifier for combining multiple classifier predictions.

➤ **from sklearn.linear_model import LogisticRegression:**

Imports LogisticRegression, a simple yet effective linear classification algorithm commonly used in binary classification.

➤ **from sklearn.svm import SVC:**

Imports SVC (Support Vector Classifier), a powerful classification algorithm for finding optimal decision boundaries in binary and multi-class classification tasks.

Code:

```
# Load the dataset
data = pd.read_csv("/content/sample_data/diabetes.csv")
```

Code Explanation:

➤ The primary task is to read the data from the "diabetes.csv" file. This can be done using the read_csv function provided by pandas. The loaded data is stored in a pandas DataFrame, which is a two-dimensional, size-mutable, and potentially heterogeneous tabular data structure with labeled axes (rows and columns). In this case, the DataFrame is assigned to the variable "data."

➤ After running this code, the "data" variable will contain the contents of the "diabetes.csv" file as a DataFrame, and we can use various pandas methods and functions to analyze, manipulate, and visualize the data as needed for our project or analysis.

2.DATA EXPLORATION:

Exploring the data is a crucial step in understanding the dataset, identifying patterns, and gaining insights that can guide our modeling process. Here are some common data exploration techniques :

Code:

```
data.head()
# To get the number of rows and columns in the dataset
data.shape
# To get the statistical measures of the data
data.describe()# To get the statistical measures of the data
data.columns
```

Code Explanation:

- **data.head():** It gives you a quick preview of what your data looks like, including the first five rows by default.
- **data.shape:** This will return the number of rows and columns in your dataset. The output will be in the format (number of rows, number of columns).
- **data.describe():** This code provides statistical summary measures for your dataset.
- **data.columns:** This code returns the column names (features) of your dataset.

Output:

```
data = pd.read_csv('Content/sample_data/diabetes.csv')
data.head()
```

| | Pregnancies | Glucose | BloodPressure | SkintThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|-------------|---------|---------------|----------------|---------|------|--------------------------|-----|---------|
| 0 | 6 | 128 | 72 | 36 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 29.6 | 0.351 | 31 | 0 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 33 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 30 | 168 | 43.1 | 2.260 | 33 | 1 |

```
# To get the number of rows and columns in the dataset
data.shape
```

(768, 9)

```
data.describe()
```

To get the statistical measures of the data

| | Pregnancies | Glucose | BloodPressure | SkintThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|-------|-------------|------------|---------------|----------------|------------|------------|--------------------------|------------|------------|
| count | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 |
| mean | 3.845152 | 120.854531 | 68.155408 | 20.536458 | 79.799479 | 31.992579 | 0.471875 | 33.240885 | 0.345508 |
| std | 3.368778 | 31.870610 | 18.360607 | 15.952118 | 95.204002 | 7.680150 | 0.331329 | 11.760232 | 0.476951 |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.078606 | 21.000000 | 0.000000 |
| 25% | 1.000000 | 99.000000 | 62.000000 | 0.000000 | 0.000000 | 27.300000 | 0.143750 | 24.000000 | 0.000000 |
| 50% | 3.000000 | 117.000000 | 72.000000 | 21.000000 | 38.000000 | 32.000000 | 0.172500 | 29.000000 | 0.000000 |
| 75% | 8.000000 | 140.250000 | 68.000000 | 31.000000 | 127.250000 | 36.000000 | 0.626250 | 41.000000 | 1.000000 |
| max | 17.000000 | 199.000000 | 120.000000 | 99.000000 | 846.000000 | 67.900000 | 2.425000 | 81.000000 | 1.000000 |

Handling Missing Values:

Check for missing values in your dataset and decide how to handle them. Checking for missing values is an essential step in data preprocessing to ensure that the dataset is clean and ready for analysis or modeling.

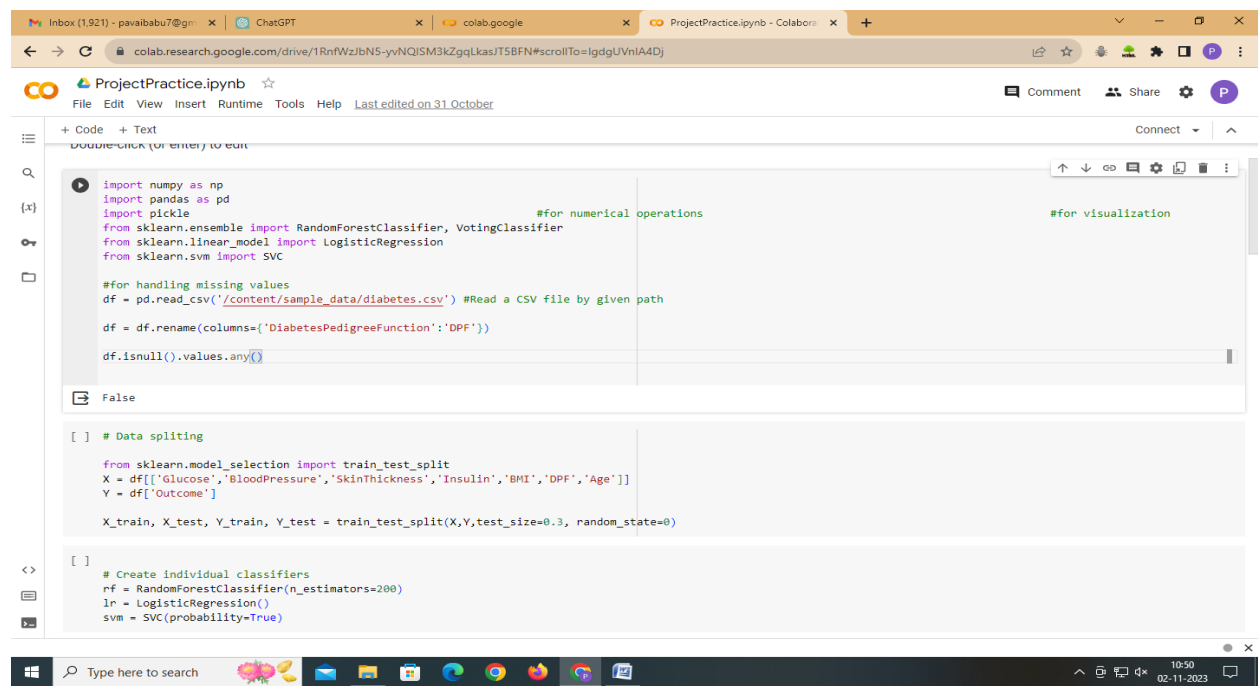
Code:

```
data.isnull().values.any()
```

Code Explanation:

If the code returns **True**, it means that there are missing values in the DataFrame **data**. If it returns **False**, it means that there are no missing values in the DataFrame.

Output:



The screenshot shows a Google Colab notebook titled 'ProjectPractice.ipynb'. The code is as follows:

```
import numpy as np
import pandas as pd
import pickle                                     #for numerical operations
from sklearn.ensemble import RandomForestClassifier, VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC                       #for visualization

#for handling missing values
df = pd.read_csv('/content/sample_data/diabetes.csv') #Read a CSV file by given path

df = df.rename(columns={'DiabetesPedigreeFunction': 'DPF'})

df.isnull().values.any()

False

[ ] # Data splitting

from sklearn.model_selection import train_test_split
X = df[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DPF', 'Age']]
Y = df['Outcome']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3, random_state=0)

[ ]

# Create individual classifiers
rf = RandomForestClassifier(n_estimators=200)
lr = LogisticRegression()
svm = SVC(probability=True)
```

Code:

```
from sklearn.model_selection import train_test_split  
X = df[['Glucose','BloodPressure','SkinThickness','Insulin','BMI','DPF','Age']]  
Y = df['Outcome']
```

Code Explanation:

1. X: This contains the independent variables (features) and is used to store all the data attributes that will be used as input for your machine learning model.
2. Y: This contains the dependent variable (outcome or label) and is used to store the target you want to predict or classify.

3.MODEL SELECTION:

Using an ensemble approach for an AI-based Diabetes Prediction System is a wise choice, as it can enhance prediction accuracy and robustness by combining the strengths of multiple machine learning models. Here's how you can implement an ensemble approach for diabetes prediction:

The goal of using an ensemble approach is to leverage the strengths of each individual classifier to improve the overall prediction accuracy of the system. Here's a brief overview of each of these classifiers:

1.Random Forest Classifier:

- RandomForest is an ensemble learning method based on decision trees. It creates multiple decision trees during training and combines their predictions. It is known for its robustness, ability to handle complex data, and resistance to overfitting.

2.Logistic Regression:

- Logistic Regression is a simple yet effective linear classification algorithm. It is widely used in binary classification problems, such as predicting the probability of a binary outcome (e.g., diabetes vs. non-diabetes). It's interpretable and computationally efficient.

3.Support Vector Machine (SVM):

- SVM is a powerful classification algorithm that finds a hyperplane that best separates data points belonging to different classes. It can handle both linear and non-linear classification tasks through kernel functions. SVMs are known for their ability to handle high-dimensional data and find optimal decision boundaries.

Code:

```
# Create individual classifiers
rf = RandomForestClassifier(n_estimators=200)
lr = LogisticRegression()
svm = SVC(probability=True)
```

Code Explanation:

- **rf = RandomForestClassifier(n_estimators=200)**: This line creates an instance of the **RandomForestClassifier** model. It specifies that the random forest should consist of 200 decision trees (you can adjust this number as needed). This model is suitable for handling complex data and is robust against overfitting.
- **lr = LogisticRegression()**: Here, you're creating an instance of the **LogisticRegression** model. This model is a straightforward yet effective choice for linear classification tasks, particularly in binary classification. It's known for its interpretability and computational efficiency.
- **svm = SVC(probability=True)**: This line creates an instance of the **SVC** (Support Vector Classifier) model with the **probability** parameter set to **True**. The **probability** parameter allows the model to predict probabilities, which is often useful in classification tasks. SVMs are known for their ability to find optimal decision boundaries and handle high-dimensional data.

4.MODEL EVALUATION:

The Voting Classifier combines the predictions of these three classifiers by taking a weighted average (soft voting) of their predicted probabilities. This ensemble approach can often lead to better performance than any single classifier on its own, as it leverages their individual strengths and mitigates their weaknesses.

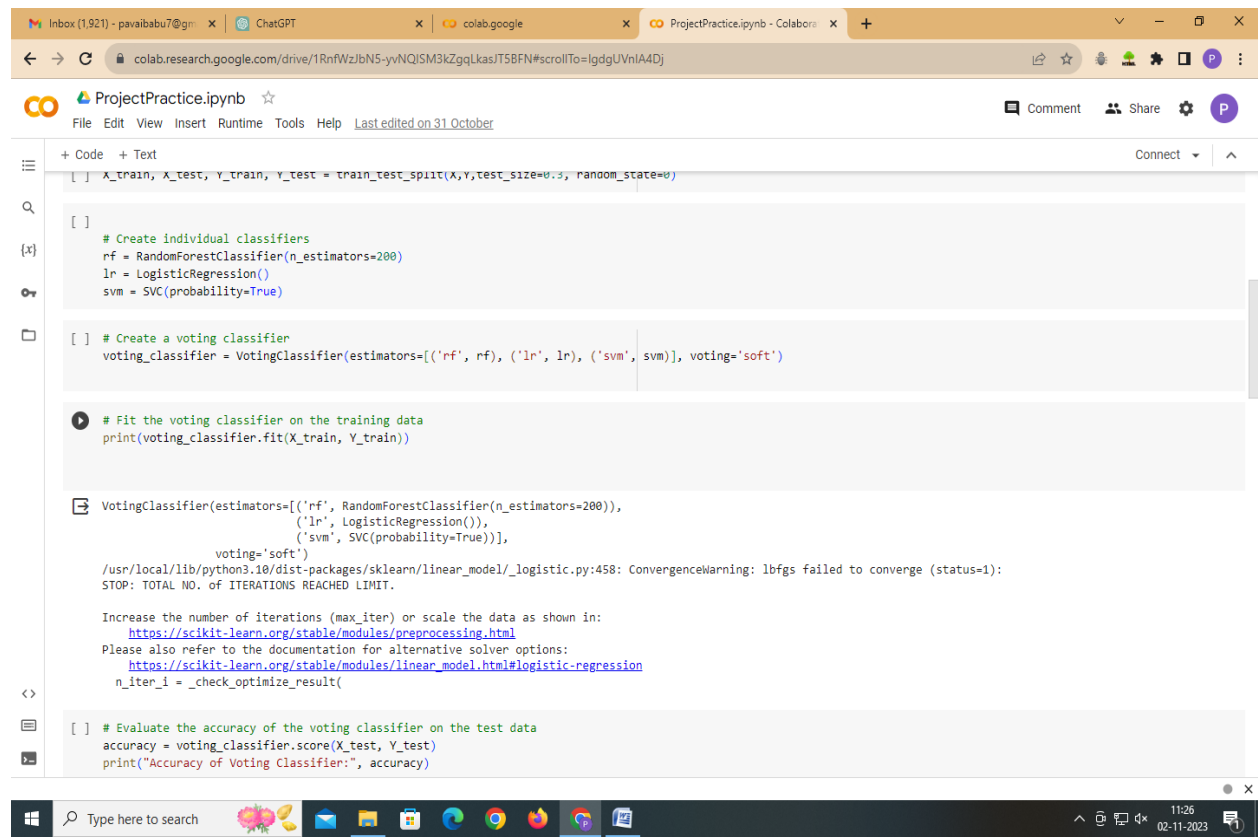
The choice of using this ensemble of classifiers is a common practice in machine learning, especially when the dataset is complex and diverse. It can provide a more robust and accurate prediction system by aggregating the diverse perspectives of multiple models.

Code:

```
# Create a voting classifier
voting_classifier = VotingClassifier(estimators=[('rf', rf), ('lr', lr), ('svm', svm)],
voting='soft')

# Fit the voting classifier on the training data
print(voting_classifier.fit(X_train, Y_train))
```

Output:



```
[ ] X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3, random_state=0)

[ ]
# Create individual classifiers
rf = RandomForestClassifier(n_estimators=200)
lr = LogisticRegression()
svm = SVC(probability=True)

[ ] # Create a voting classifier
voting_classifier = VotingClassifier(estimators=[('rf', rf), ('lr', lr), ('svm', svm)], voting='soft')

# Fit the voting classifier on the training data
print(voting_classifier.fit(X_train, Y_train))

VotingClassifier(estimators=[('rf', RandomForestClassifier(n_estimators=200)),
                           ('lr', LogisticRegression()),
                           ('svm', SVC(probability=True))],
                  voting='soft')
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear\_model.html#logistic-regression
n_iter_1 = _check_optimize_result(

[ ] # Evaluate the accuracy of the voting classifier on the test data
accuracy = voting_classifier.score(X_test, Y_test)
print("Accuracy of Voting Classifier:", accuracy)
```

Moving on to model accuracy is an important step to evaluate how well your machine learning models are performing. We can assess the model accuracy using various metrics. The accuracy score is a common metric to measure the overall correctness of your model's predictions.

Code:

```
# Evaluate the accuracy of the voting classifier on the test data
```

```
accuracy = voting_classifier.score(X_test, Y_test)  
print("Accuracy of Voting Classifier:", accuracy)
```

Code Explanation:

- **voting_classifier.score(X_test, Y_test):** The **score** method of the **voting_classifier** object is used to evaluate the model's accuracy on the provided test data. **X_test** contains the features, and **Y_test** contains the true target labels.
- **accuracy = ...:** The result of the accuracy calculation is assigned to the variable **accuracy**.
- **print("Accuracy of Voting Classifier:", accuracy):** This line prints the accuracy of the Voting Classifier on the test data to the console.

The accuracy score represents the proportion of correct predictions made by our model on the test data. It's a common metric to assess the model's performance in classification tasks. The accuracy value will be between 0 (no correct predictions) and 1 (all predictions are correct).

In our case, the accuracy value will tell you how well the Voting Classifier is performing in predicting diabetes based on the test data.

Output:



```
# Evaluate the accuracy of the voting classifier on the test data  
accuracy = voting_classifier.score(X_test, Y_test)  
print("Accuracy of Voting Classifier:", accuracy)
```

Accuracy of Voting Classifier: 0.7575757575757576


Code:

```
# Save the voting classifier to a file
filename = 'model/voting_diabetes.pkl'
pickle.dump(voting_classifier, open('/content/sample_data/code.py', 'wb'))
print("SUCCESS")
```

Code Explanation:

- It specifies the file path where the model will be saved using the **filename** variable.
- It saves the **voting_classifier** model to the specified file path using **pickle.dump()**.
- It prints a success message indicating that the model has been saved.

Output:



```
✓ [22] # Save the voting classifier to a file
In filename = 'model/voting_diabetes.pkl'
pickle.dump(voting_classifier, open('/content/sample_data/code.py', 'wb'))
print("SUCCESS")

SUCCESS
```

CHECK OUT THE CODE:

<https://colab.research.google.com/drive/1RnfWzJbN5yvNQISM3kZggLkasJT5BFN#scrollTo=lgdguVnIA4Dj>