

## **Basic Projects**

### **1. Project: Diabetes Prediction**

Description: Develop a predictive model to identify individuals at high risk of developing diabetes based on various health indicators.

Dataset: Pima Indians Diabetes Database on Kaggle.

#### ***Pima Indians Diabetes Database***

Content: This dataset comprises diagnostic measurements from 768 female patients of Pima Indian heritage. It includes data points like the number of pregnancies, glucose concentration, blood pressure, skin thickness, insulin levels, BMI, diabetes pedigree function, and age.

Usage: The dataset is commonly used to predict the onset of diabetes based on these health indicators.

### **2. Project: Heart Disease Detection**

Description: Use machine learning to predict the presence of heart disease in patients based on attributes like cholesterol levels, resting blood pressure, and other heart-related metrics.

Dataset: Heart Disease UCI on UCI Machine Learning Repository.

#### ***Heart Disease UCI***

Content: The dataset contains 76 attributes, but all published experiments refer to using a subset of 14 of them. It includes attributes like age, sex, chest pain type, resting blood pressure, cholesterol levels, fasting blood sugar, and more.

Usage: It's utilized for predicting the presence of heart disease in patients, making it a standard dataset for classification problems in healthcare.

### **3. Project: Brain Cancer Classification**

Description: Create a classification model to classify types of brain cancer.

Dataset:

The Cancer Imaging Archive (TCIA): TCIA is a service that de-identifies and hosts a large archive of medical images of cancer accessible for public download. The data are organized as “collections”; typically, patients' imaging related by a common disease (e.g., lung cancer), image modality, or research focus. Check for datasets specifically related to brain cancer, such as the Brain Tumor collections, which can include MRI, CT scans, and more.

#### **4. Project: Liver Disease Prediction**

Description: Predict the onset of liver disease in patients by analyzing blood and demographic data.

Dataset: Indian Liver Patient Dataset on Kaggle.

##### ***Indian Liver Patient Dataset***

Content: It comprises 583 data points collected from patients in North East of Andhra Pradesh, India. The dataset has 11 variables, including indicators of liver function tests, age, gender, total bilirubin, direct bilirubin, alkaline phosphatase, and more.

Usage: This dataset is utilized to predict whether a patient is likely to have liver disease based on these biochemical and demographic attributes.

#### **5. Project: Medical Cost Personalization**

Description: Analyze patient data to predict healthcare costs based on factors like age, BMI, smoking status, and region.

Dataset: Medical Cost Personal Datasets on Kaggle.

*Content: This dataset includes information on age, sex, BMI, children, smoker status, region, and charges for 1,338 individuals.*

*Usage: It's used to analyze and predict personal healthcare costs based on demographic and health-related factors, which can be beneficial for insurance cost prediction and personal health finance management.*

# Innovative Projects

## 1. Project: AI-Driven Telemedicine Chatbot

Description: Develop an NLP model that provides medical consultation, diagnosis suggestions, and general health advice through a chat interface.

Dataset: Medical Dialogue Dataset is available for research purposes, focusing on COVID-19 dialogues but can be adapted for broader medical use.

### ***Medical Dialogue Dataset***

Content: This dataset contains dialogue data focusing on COVID-19, providing real doctor-patient conversations, which are de-identified. It includes various medical conditions and queries addressed in the dialogues.

Usage: It's suited for developing healthcare chatbots, specifically for NLP applications that require understanding and generating medical dialogue.

## 2. Project: Real-Time Anomaly Detection in EHR

Description: Implement a system that uses machine learning to detect anomalies in real-time electronic health record (EHR) data, identifying potential data entry errors or patient health issues.

Dataset: MIMIC-III offers a rich, freely available dataset derived from EHR data.

### ***MIMIC-III (Medical Information Mart for Intensive Care III)***

Content: MIMIC-III is a large database comprising de-identified health-related data associated with over forty thousand patients who stayed in critical care units. It includes information like vital signs, medications, laboratory measurements, observations, and notes.

Usage: This dataset is widely used for research in various aspects of healthcare analytics, including anomaly detection, predictive modeling, and patient care optimization.

## 3. Project: Mental Health Condition Monitoring from Social Media

Description: Use natural language processing to analyze social media posts for signs of mental health conditions, aiming to provide early warnings or insights into population mental health trends.

Dataset: Multi-Task Learning for Mental Health Conditions dataset on Kaggle can be a starting point.

### ***Mental Health Conditions Dataset from Social Media***

Content: The dataset consists of posts from Reddit, labeled with various mental health conditions. It's a collection meant for multi-task learning, providing insights into how mental health conditions are discussed in social media contexts.

Usage: This dataset is valuable for analyzing social media language to detect signs of mental health issues, aiding in mental health research and early intervention strategies.

## **4. Project: Predictive Modeling for Epidemic Outbreaks**

Description: Develop a model to predict the spread of infectious diseases, such as COVID-19, using time-series data and demographic information.

Dataset: COVID-19 Data Repository by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University provides a comprehensive, regularly updated dataset.

### ***COVID-19 Data Repository by CSSE at Johns Hopkins University***

Content: This repository provides daily updated time series data on confirmed cases, deaths, and recovery from COVID-19 globally, broken down by country and region.

Usage: It's extensively used for predictive modeling, epidemiological studies, and understanding the spread and impact of COVID-19 across different geographies.