

Disease Prediction Using Machine Learning

This notebook implements a predictive analytics model for early disease detection using a healthcare dataset. It includes preprocessing, EDA, model training, and evaluation.

```
# Step 1: Install Required Libraries
!pip install seaborn scikit-learn
```

```
Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packages (0.13.2)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: numpy!=1.24.0,>=1.20 in /usr/local/lib/python3.11/dist-packages (from seaborn) (2.0.2)
Requirement already satisfied: pandas>=1.2 in /usr/local/lib/python3.11/dist-packages (from seaborn) (2.2.2)
Requirement already satisfied: matplotlib!=3.6.1,>=3.4 in /usr/local/lib/python3.11/dist-packages (from seaborn) (3.10.0)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.15.2)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (24.1)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (11.2.1)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (2.9.0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.2->seaborn) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.2->seaborn) (2025.2)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.4->seaborn) (1.17.0)
```

```
# Step 2: Import Libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
```

```
# Step 3: Upload Dataset
from google.colab import files
uploaded = files.upload()
```

```
# Load the dataset
df = pd.read_csv(list(uploaded.keys())[0])
df.head()
```

Choose Files
healthcare_dataset.csv

- healthcare_dataset.csv(text/csv) - 8399221 bytes, last modified: 5/1/2025 - 100% done

Saving healthcare_dataset.csv to healthcare_dataset.csv

| | Name | Age | Gender | Blood Type | Medical Condition | Date of Admission | Doctor | Hospital | Insurance Provider | Billing Amount | Room Number | Admission Type | Discharge Date | Me |
|---|---------------|-----|--------|------------|-------------------|-------------------|------------------|----------------------------|--------------------|----------------|-------------|----------------|----------------|----|
| 0 | Bobby JacksOn | 30 | Male | B- | Cancer | 2024-01-31 | Matthew Smith | Sons and Miller | Blue Cross | 18856.281306 | 328 | Urgent | 2024-02-02 | P |
| 1 | LesLie TErRy | 62 | Male | A+ | Obesity | 2019-08-20 | Samantha Davies | Kim Inc | Medicare | 33643.327287 | 265 | Emergency | 2019-08-26 | |
| 2 | DaNnY sMitH | 76 | Female | A- | Obesity | 2022-09-22 | Tiffany Mitchell | Cook PLC | Aetna | 27955.096079 | 205 | Emergency | 2022-10-07 | |
| 3 | andrEw waTtS | 28 | Female | O+ | Diabetes | 2020-11-18 | Kevin Wells | Hernandez Rogers and Vang, | Medicare | 37909.782410 | 450 | Elective | 2020-12-18 | |
| 4 | adrIENNE bEIl | 43 | Female | AB+ | Cancer | 2022-09-19 | Kathleen Hanna | White-White | Aetna | 14238.317814 | 458 | Urgent | 2022-10-09 | |

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)



Step 4: Data Preprocessing

Check for missing values

print("Missing values per column:")

print(df.isnull().sum())

Fill or drop missing values based on strategy

df = df.dropna() # Or use df.fillna(method='ffill') if preferred

Remove duplicates

df = df.drop_duplicates()

Normalize column names

df.columns = df.columns.str.strip().str.lower().str.replace(' ', '_')

Encode categorical variables if any

for col in df.select_dtypes(include='object').columns:

le = LabelEncoder()

df[col] = le.fit_transform(df[col])

Missing values per column:

```

Name          0
Age           0
Gender        0
Blood Type    0
Medical Condition  0
Date of Admission  0
Doctor        0
Hospital      0
Insurance Provider  0
Billing Amount  0
Room Number   0
Admission Type  0

```

```
Discharge Date      0
Medication           0
Test Results         0
dtype: int64
```

```
# Step 5: Exploratory Data Analysis (EDA)
```

```
# Basic statistics
print(df.describe())
```

```
# Visualize correlations
plt.figure(figsize=(10, 6))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title("Feature Correlation Heatmap")
plt.show()
```

```

count    name    age    gender    blood_type \
mean    24991.842703    51.535185    0.500237    3.492013
std     14430.654340    19.605661    0.500004    2.289690
min       0.000000    13.000000    0.000000    0.000000
25%     12402.250000    25.000000    0.000000    1.000000

```

Step 6: Feature Selection and Splitting

target_column_name = 'disease' # Replace 'disease' with the actual name

X = df.drop(target_column_name, axis=1)

y = df[target_column_name]

Split the dataset

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

Feature scaling

scaler = StandardScaler()

X_train = scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

```

count    name    age    gender    blood_type \
mean    24991.842703    51.535185    0.500237    3.492013
std     14430.654340    19.605661    0.500004    2.289690
min       0.000000    13.000000    0.000000    0.000000
25%     12402.250000    25.000000    0.000000    1.000000

```

Step 7: Model Building

Using Random Forest Classifier

model = RandomForestClassifier(n_estimators=100, random_state=42)

model.fit(X_train, y_train)

Make predictions

y_pred = model.predict(X_test)

```

mean    24991.842703    51.535185    0.500237    3.492013
std     14430.654340    19.605661    0.500004    2.289690
min       0.000000    13.000000    0.000000    0.000000
25%     12402.250000    25.000000    0.000000    1.000000

```

Step 8: Model Evaluation

Accuracy and report

print("Accuracy:", accuracy_score(y_test, y_pred))

print("\nClassification Report:")

print(classification_report(y_test, y_pred))

Confusion Matrix

conf_matrix = confusion_matrix(y_test, y_pred)

sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues")

plt.title("Confusion Matrix")

plt.xlabel("Predicted")

plt.ylabel("Actual")

plt.show()

Step 9: Optional Deployment with Streamlit (code not run in Colab)

Save the model (if needed)

import joblib

joblib.dump(model, "disease_model.pkl")

print("Model saved. You can deploy it using Streamlit or Flask.")

