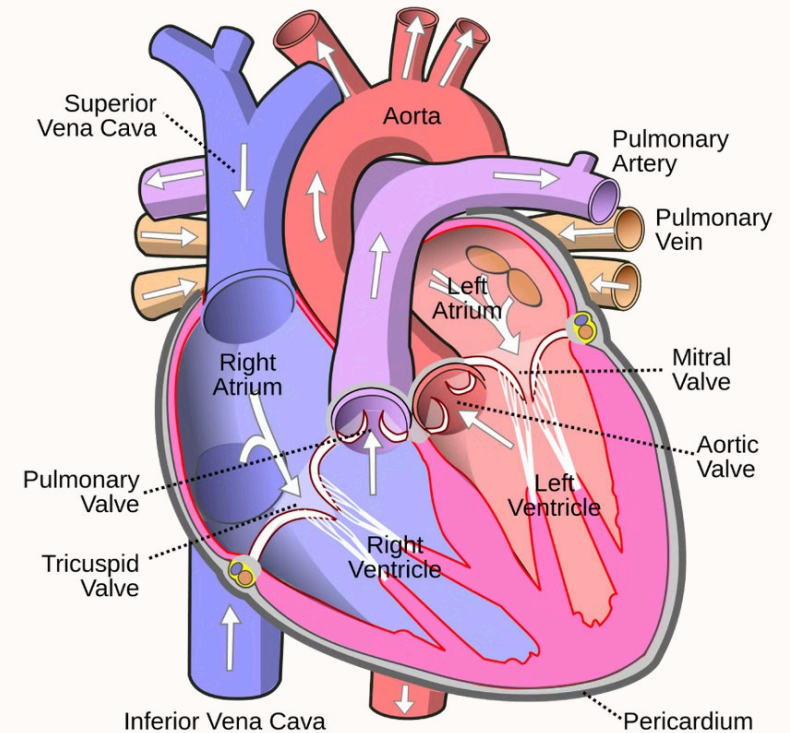


# CardioRanger

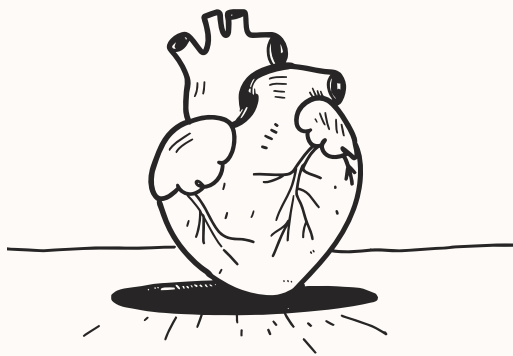
Aplikacja przewidująca potencjalną chorobą wieńcową serca (CHD) oparta na modelu uczenia maszynowego.



# Agenda

- 1 Opis danych oraz cel
- 2 Baza danych
- 3 Eksploracyjna analiza danych (EDA)
- 4 Model i kalibracja
- 5 Aplikacja

# Opis danych i cel



Badanie **Framingham Heart Study** to jedno z najbardziej przełomowych i długofalowych badań epidemiologicznych w historii medycyny. Rozpoczęte w 1948 roku w miasteczku Framingham w stanie Massachusetts, miało na celu zidentyfikowanie czynników ryzyka związanych z chorobami sercowo-naczyniowymi. Badanie objęło pierwotnie 5209 zdrowych mężczyzn i kobiet w wieku 30–62 lat, których stan zdrowia był monitorowany przez kolejne dekady.

Zbiór danych: <https://www.kaggle.com/datasets/aasheesh200/framingham-heart-study-dataset>

**Liczba obserwacji:** 4240

**Liczba kolumn:** 16 kolumn (płeć, wiek, wykształcenie, czy pali, czy przyjmuje leki na nadciśnienie, historia udaru, nadciśnienie tętnicze, czy występuję cukrzyca, poziom cholesterolu, skurczowe ciśnienie tętnicze, rozkurczowe ciśnienie tętnicze, wskaźnik masy ciała (BMI), tętno, poziom glukozy, czy uczestnik miał chorobę (target)).

**Cel :** Przewidzenie czy pacjent ma 10-letnie ryzyko przyszłej choroby wieńcowej (CHD).

# Baza danych

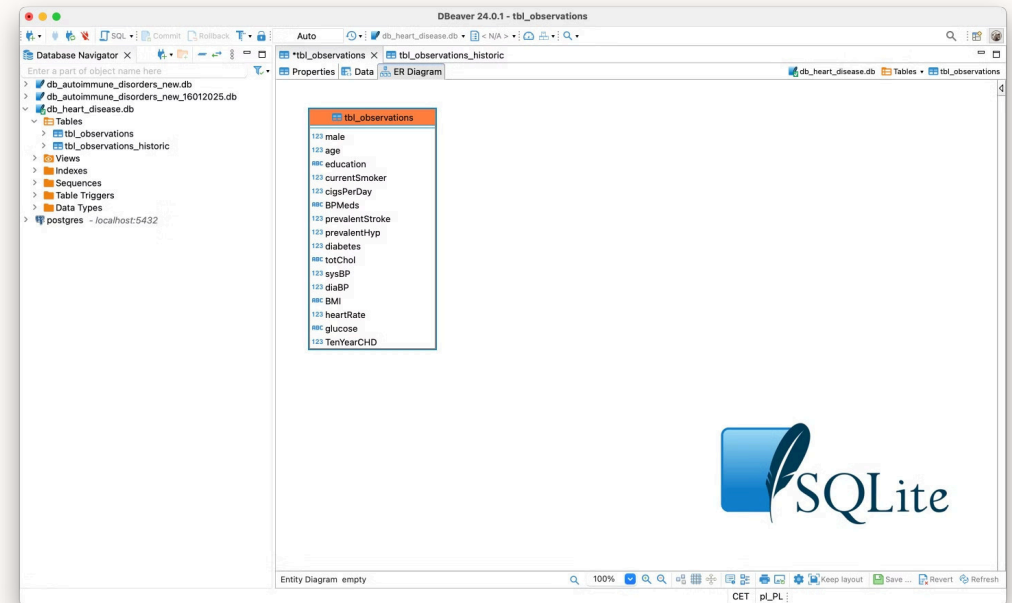
- **Technologia:** Baza danych została zbudowana w **SQLite**, lekkim i wydajnym systemie bazodanowym, który nie wymaga osobnego serwera i doskonale sprawdza się w aplikacjach lokalnych oraz analitycznych.

- **Struktura:** Baza zawiera **dwie tabele**:

**tbl\_observations** – główna tabela zawierająca **4240 rekordów**, wykorzystywana do budowy modelu predykcyjnego.

**tbl\_observations\_historic** – tabela przeznaczona do **archiwizacji danych historycznych**, nieużywana w modelu.

- **Zawartość:** Obie tabele posiadają **16 kolumn**, przechowujących kluczowe zmienne medyczne i demograficzne wykorzystywane w analizie ryzyka chorób serca.



# Eksploracyjna analiza danych (EDA)

## Przygotowanie danych:

- Zmieniono nazwy niektórych zmiennych i ich typy.
- Usunięto **15,21%** przypadków z brakami danych, pozostawiając **3 658** rekordów.
- Brak duplikatów i błędnych wartości.
- Zmienną kategoryczną (edukacja) zakodowano metodą **OneHotEncoder**.
- Cechy numeryczne poddano standaryzacji (**StandardScaler**).

## Wnioski z EDA:

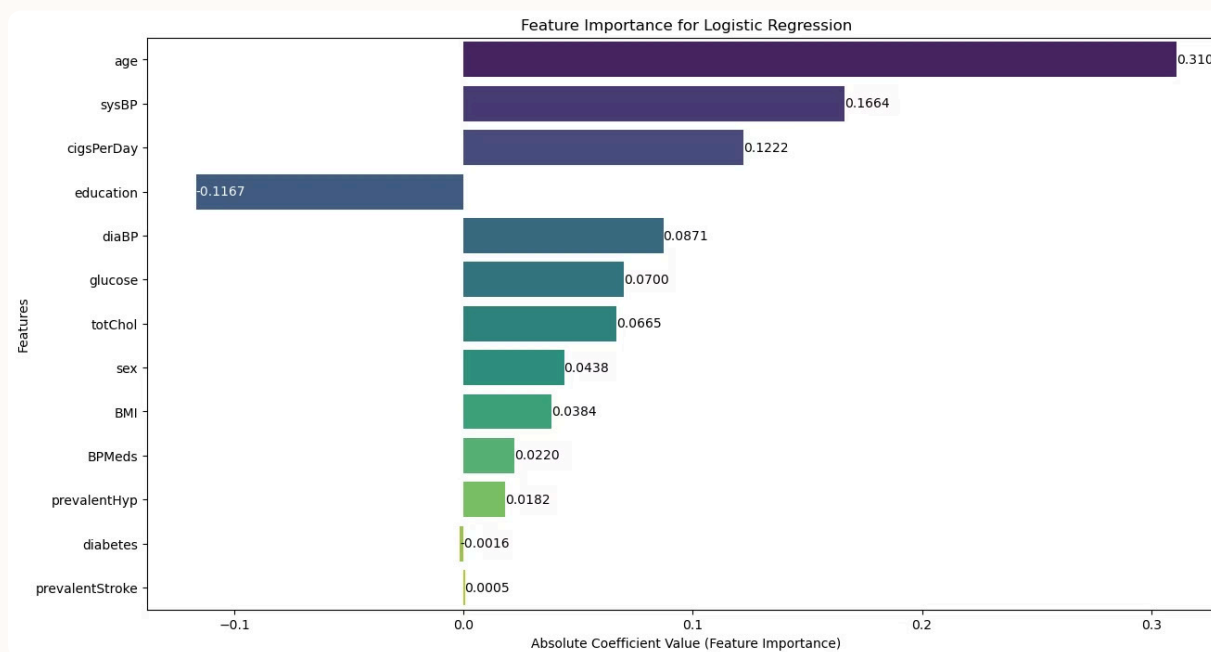
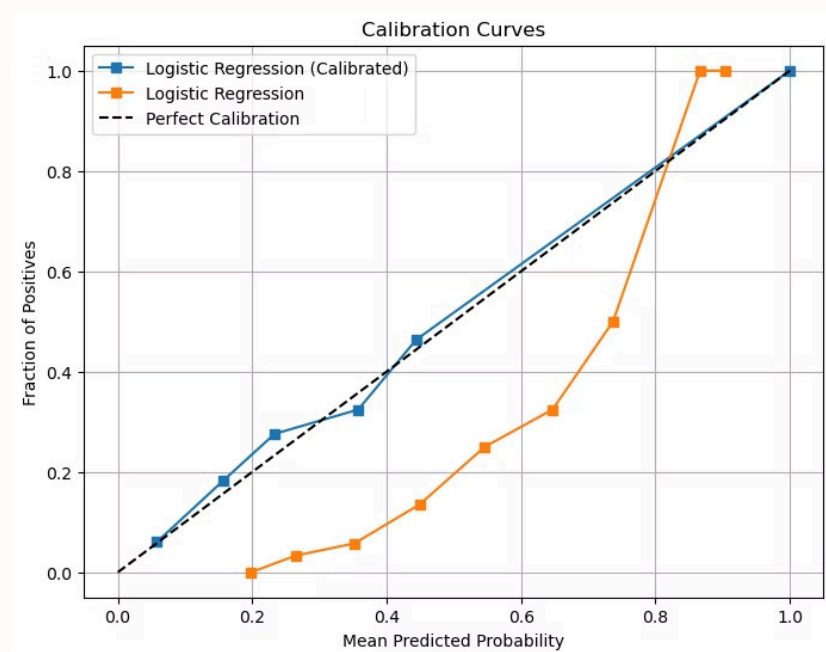
- Usunięto zmienne **currentSmoker** (palenie tytoniu) i **heartRate** (tętno) ze względu na brak istotnego wpływu na zmienną zależną.
- Największy wpływ na ryzyko choroby serca mają: **age** (wiek), **sysBP** (ciśnienie skurczowe), **prevalentHyp** (nadciśnienie tętnicze) oraz **CigsPerDay** (ilość papierosów spalanych dziennie).

# Model

W procesie tworzenia aplikacji przetestowanych zostało łącznie 4 algorytmy uczenia maszynowego na zbiorze danych repróbrowanych metodą SMOTE: **Regresja Logistyczna**, **KNN**, **Las Losowe** oraz **CatBoost**. Kluczową metryką ewaluacyjną, na podstawie której porównywana była skuteczność modeli, był wskaźnik **AUC**.

Najlepsze wyniki osiągnął model **Regresji Logistycznej**, uzyskując **AUC** na poziomie **75%**. Optymalne parametry obejmowały parametr kary L2, współczynnik regularyzacji C równy 0.001 oraz algorytm optymalizacji funkcji kosztu **liblinear**.

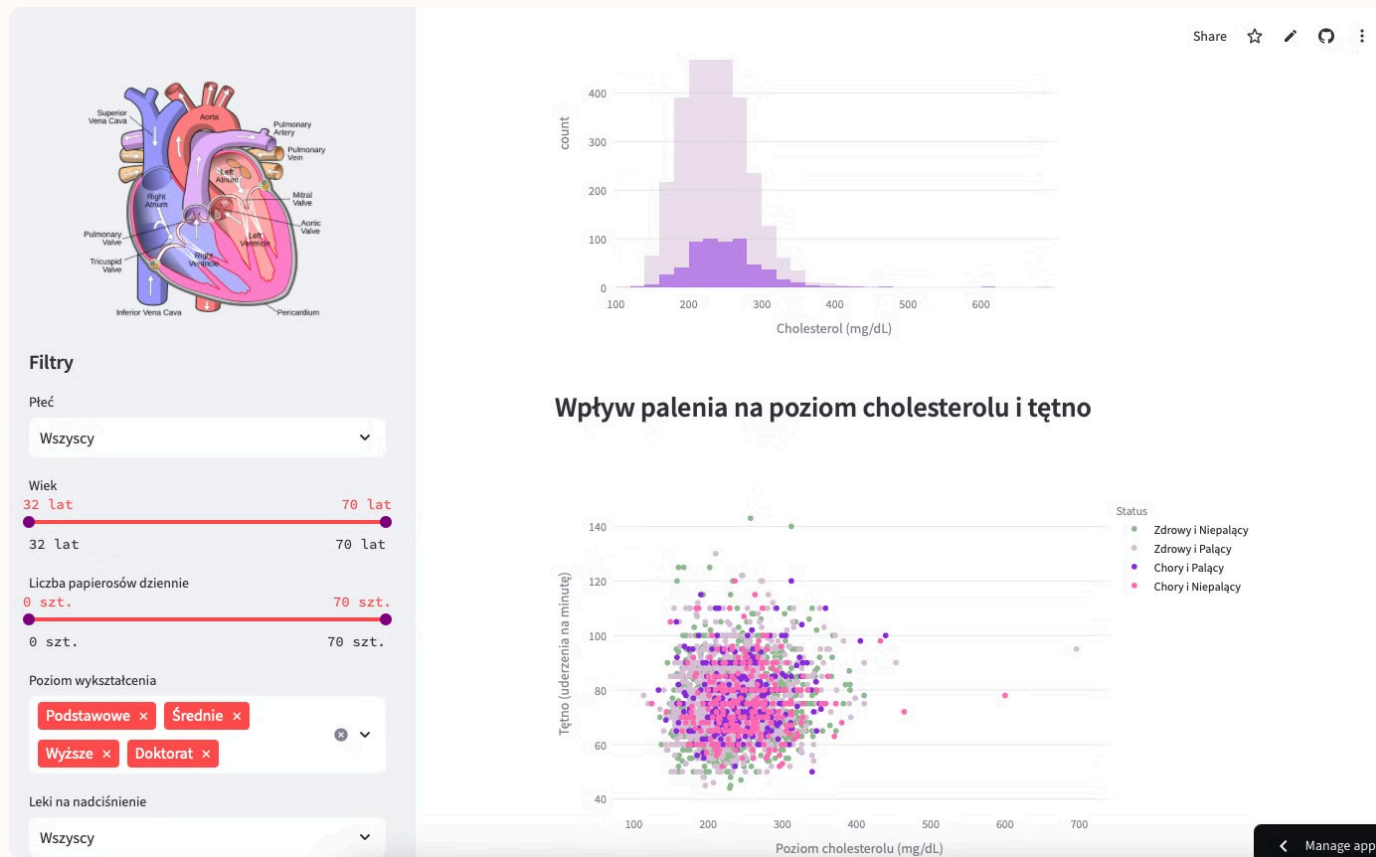
Aby zwiększyć stabilność i wiarygodność predykcji, finalny model został dodatkowo **skalibrowany** przy użyciu **regresji izotonicznej**. Istotność zmiennych w modelu została szczegółowo przeanalizowana, co pozwoliło na lepsze zrozumienie czynników wpływających na przewidywanie ryzyka choroby wieńcowej serca.



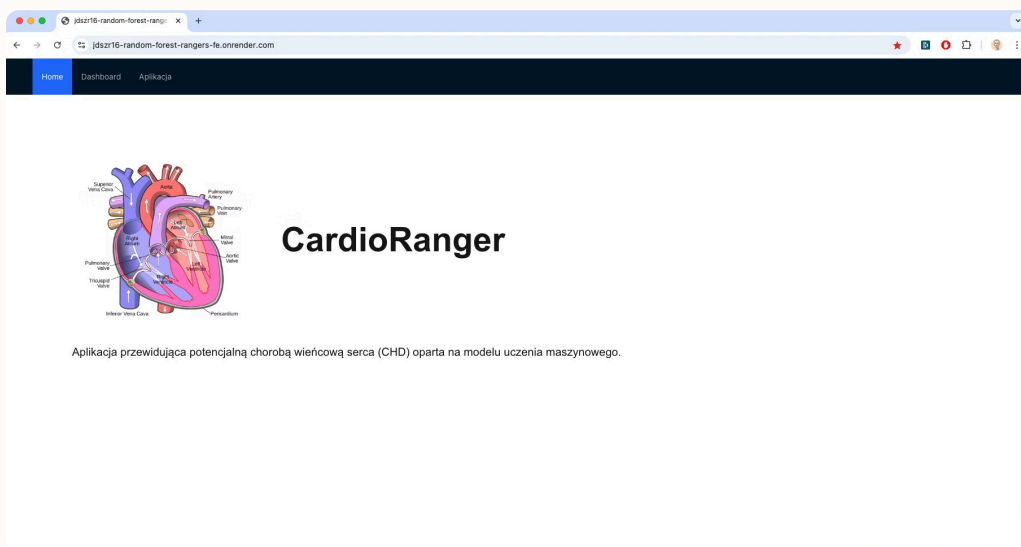
# Aplikacja

**Dashboard** umożliwia interaktywne analizowanie czynników wpływających na zdrowie serca.

**CardioRanger** to dashboard stworzony w **Streamlit**, który wizualizuje kluczowe wskaźniki zdrowotne i pozwala na eksplorację danych zapisanych w bazie **SQLite**.



# Aplikacja



The screenshot shows the form page of the CardioRanger application. The form is organized into three columns of input fields. The first column contains fields for Age (Wiek: 22), Height (Wzrost: 180 cm), Number of cigarettes per day (Liczba papierosów dziennie: 30), History of stroke (Przebyty udar: Nie), and Cholesterol (Cholesterol: 230 mg/dL). The second column contains fields for Blood pressure (Wykształcenie: Średnie), Weight (Waga: 90 kg), Medication for blood pressure (Leki na ciśnienie: Nie), Blood pressure (Ciśnienie: Nie), and Blood pressure (Ciśnienie skurczowe: 100 mmHg). The third column contains fields for Gender (Płeć: Męska), BMI (BMI: 27.7), Glucose (Glukoza: 80 mg/dL), Diabetes (Cukrzyca: Nie), and Blood pressure (Ciśnienie rozkurczowe: 70 mmHg). A "Wyślij" button is located at the bottom of the form. Below the form, there is a checkbox labeled "Jestem lekarzem". A text box below the checkbox contains the following text: "Na podstawie podanych informacji, Twoje ryzyko zachorowania na choroby serca w ciągu następnych lat wynosi 1.31 %." "Masz 22 lata i jesteś mężczyzną. Ciekawe jest, że palisz 30 papierosów dziennie. To znacznie zwiększa ryzyko problemów z sercem i innymi chorobami, dlatego warto przemyśleć rzucenie palenia." "Nie bierzesz leków na nadciśnienie, nie miałeś udaru, nie masz nadciśnienia ani cukrzycy, co jest pozytywne. Twoje ciśnienie krwi (100/70 mmHg) jest w normie, choć dolna granica może być bliska dolnej normy (60-90 mmHg dla skurczowego i 40-80 mmHg dla rozkurczowego)." "Twój poziom cholesterolu wynosi 230 mg/dL. Górna granica normy to 200 mg/dL, więc warto skonsultować się z lekarzem, aby omówić zdrowe nawyki. Twój wskaźnik masy ciała (BMI) to 27.7, co oznacza, że możesz być lekko nadmiernie ciężki. Normą jest 18.5-24.9."



# CardioRanger - Twój przewodnik po zdrowiu serca!

