



Univerzitet „Džemal Bijedić“ u Mostaru

Fakultet Informacijskih Tehnologija

Predmet: Umjetna inteligencija

HealthPredict – AI agent za procjenu rizika od srčanih bolesti

Dokumentacija

Profesor:
prof.dr.sc. Nina Bijedić

Student:
Nejra Muminović, IB220043

Mostar, 2025

Svrha i cilj projekta

Cilj projekta HealthPredict je razvoj AI agenta koji omogućava procjenu rizika od srčanih bolesti na osnovu medicinskih i demografskih podataka pacijenata. Sistem je namijenjen kao pomoćni alat za brzu i informativnu procjenu zdravstvenog rizika, a ne kao zamjena za stručnu medicinsku dijagnozu.

Projekt je inspiriran činjenicom da su kardiovaskularne bolesti jedan od vodećih uzroka smrtnosti u svijetu, te da rana procjena rizika može imati ključnu ulogu u prevenciji i pravovremenom liječenju. Korištenjem historijskih medicinskih podataka i algoritama mašinskog učenja, HealthPredict omogućava korisnicima da dobiju uvid u vjerovatnoću postojanja srčanih oboljenja.

Svrha projekta je izrada web aplikacije koja koristi algoritme mašinskog učenja za analizu faktora kao što su starost, spol, tip angine, krvni pritisak, kolesterol, EKG rezultati i drugi relevantni parametri, te na osnovu njih predviđa rizik od srčanih bolesti.

Tehnologije

Za realizaciju AI agenta HealthPredict korištene su sljedeće tehnologije:

- Python – korišten za backend logiku, obradu podataka, treniranje modela i predikciju.
- Flask – korišten za izradu REST API-ja i komunikaciju između frontenda i backend-a.
- Random Forest Classifier – algoritam mašinskog učenja korišten za klasifikaciju rizika od srčanih bolesti.
- Logistic Regression – korišten kao dodatni model tokom treniranja i evaluacije.
- JavaScript – korišten za izradu frontend dijela aplikacije i komunikaciju s API-jem.
- HTML & CSS – korišteni za izradu korisničkog interfejsa.
- Pandas – korišten za obradu i manipulaciju medicinskih podataka.
- Scikit-learn – biblioteka korištena za treniranje, skaliranje i evaluaciju ML modela.
- CSV dataset (Heart Disease Dataset) – korišten kao izvor podataka za treniranje modela.

Zašto Random Forest Classifier?

Random Forest Classifier je odabran zbog svoje sposobnosti da:

- efikasno radi s većim brojem karakteristika,
- bude otporan na šum u podacima,
- pruži dobru tačnost bez potrebe za kompleksnim podešavanjem,
- omogućiti uvid u važnost pojedinih karakteristika koje utiču na predikciju.

Zahvaljujući ovim osobinama, Random Forest je pogodan za medicinske podatke koji često sadrže varijacije i nelinearne odnose.

Funkcionalnosti

AI agent HealthPredict nudi sljedeće funkcionalnosti:

- Predikcija rizika od srčanih bolesti
Korisnik unosi medicinske podatke (starost, spol, tip bola u prsima, krvni pritisak, holesterol, EKG rezultat, puls, itd.), a sistem vraća procjenu rizika (nizak ili visok rizik) zajedno sa vjerovatnoćom.
- Prikaz najvažnijih faktora
Sistem prikazuje koje karakteristike su najviše utjecale na odluku modela.
- Dodavanje novih pacijenata
Korisnici mogu dodavati nove podatke o pacijentima u bazu podataka.
- Ponovno treniranje modela (Retraining)
Nakon dodavanja novih podataka, model se može ponovno trenirati kako bi se poboljšala tačnost i prilagodio novim informacijama.
- Višejezična podrška (Bosanski / Engleski)
Korisnički interfejs omogućava promjenu jezika aplikacije.

Kod

Učitavanje podataka

Podaci se učitavaju iz CSV datoteke korištenjem Pandas biblioteke. Dataset sadrži medicinske i demografske podatke pacijenata relevantne za procjenu srčanih bolesti.

```
6 def load_data(path="data/heart.csv"):
7     df = pd.read_csv(path)
8     return df
```

Procesiranje podataka

Podaci se čiste i pripremaju za treniranje modela:

- uklanjaju se nedostajuće vrijednosti,
- kreira se ciljna varijabla (0 – nema bolesti, 1 – postoji rizik),
- vrši se enkodiranje kategorijskih varijabli,
- numeričke vrijednosti se skaliraju korištenjem StandardScaler.

```

10 def preprocess(df):
11     df = df.dropna()
12
13     df["target"] = df["num"].apply(lambda x: 1 if x > 0 else 0)
14     y = df["target"]
15
16     x = df.drop(["num", "target", "id", "dataset"], axis=1)
17
18     x["sex"] = x["sex"].map({"Male":1, "Female":0})
19     x["fbs"] = x["fbs"].map({True:1, False:0})
20     x["exang"] = x["exang"].map({True:1, False:0})
21
22     x = pd.get_dummies(x, columns=["cp", "restecg", "slope", "thal"], drop_first=True)
23
24     scaler = StandardScaler()
25     x_scaled = scaler.fit_transform(x)
26
27     return x_scaled, y, scaler, x.columns.tolist()

```

Podaci se prije treniranja modela prolaze kroz fazu procesiranja. Uklanjaju se redovi sa nedostajućim vrijednostima, zatim se kreira ciljna varijabla gdje vrijednost 1 označava prisustvo rizika od srčanih bolesti, dok 0 označava odsustvo rizika. Kategorijske varijable se enkodiraju korištenjem binarnog mapiranja i one-hot enkodiranja, dok se numeričke vrijednosti skaliraju pomoću StandardScaler metode radi poboljšanja performansi modela.

Treniranje modela

Model se trenira korištenjem Random Forest Classifier algoritma. Tokom treniranja čuvaju se:

- istrenirani model,
- scaler,
- lista feature-a.

Ovi podaci se koriste kasnije tokom predikcije kako bi ulazni podaci bili pravilno obrađeni.

```

78     print("Treniranje RandomForest sa GridSearch...")
79     rf = RandomForestClassifier(random_state=42)
80     param_grid = {
81         "n_estimators": [100, 200],
82         "max_depth": [None, 6, 10],
83         "min_samples_split": [2, 5]
84     }
85     grid = GridSearchCV(rf, param_grid, cv=5, scoring="roc_auc", n_jobs=-1)
86     grid.fit(X_train, y_train)
87     best_rf = grid.best_estimator_
88     print("Random Forest metrics:", evaluate_model(best_rf, X_test, y_test))
89     print("Best RF params:", grid.best_params_)
90
91     os.makedirs("models", exist_ok=True)
92     joblib.dump({"model": best_rf, "scaler": scaler, "features": feature_names}, MODEL_PATH)
93     print(f"Model spremljen u {MODEL_PATH}")

```

Ovaj dio koda služi za treniranje Random Forest Classifier modela i optimizaciju njegovih hiperparametara pomoću GridSearchCV metode. Najprije se inicijalizira Random Forest model uz korištenje parametra `random_state=42`, čime se osigurava ponovljivost rezultata. Zatim se definiše skup hiperparametara koje će GridSearch ispitivati, uključujući broj stabala u modelu (`n_estimators`), maksimalnu dubinu stabla (`max_depth`) i minimalan broj uzoraka potreban za dijeljenje čvora (`min_samples_split`). GridSearchCV se koristi za sistematsko testiranje svih kombinacija ovih parametara kroz 5-struku unakrsnu validaciju, pri čemu se kao kriterij za odabir najboljeg modela koristi ROC AUC metrika. Nakon treniranja, bira se najbolji model sa optimalnim parametrima, vrši se evaluacija njegovih performansi na testnom skupu podataka, a zatim se istrenirani model zajedno sa scalerom i listom feature-a trajno sprema na disk kako bi se kasnije mogao koristiti za predikciju rizika od srčanih bolesti.

Predikcija rizika srčanih bolesti

Na osnovu unesenih podataka, backend API obrađuje zahtjev i vraća:

- predikciju (nizak ili visok rizik),
- vjerovatnoću rizika,
- najutjecajnije karakteristike.

```
8  def predict(model, scaler, features, data: dict):
9      X = pd.DataFrame([data])
10
11     # poravnaj feature-e
12     X = X.reindex(columns=features, fill_value=0)
13
14     X_scaled = scaler.transform(X)
15     proba = model.predict_proba(X_scaled)[0][1]
16     pred = int(proba >= 0.5)
17
18     return {
19         "prediction": pred,
20         "probability": float(proba)
21     }
```

Funkcija `predict` prima istrenirani model, scaler, listu feature-a i podatke unesene od strane korisnika. Ulazni podaci se prvo pretvaraju u Pandas DataFrame, a zatim se poravnavaju prema feature-ima korištenim tokom treniranja modela, pri čemu se nedostajuće vrijednosti popunjavaju nulom. Nakon toga, podaci se skaliraju pomoću prethodno sačuvanog StandardScaler objekta. Random Forest Classifier zatim izračunava vjerovatnoću rizika od srčane bolesti, na osnovu koje se formira konačna predikcija (0 – nizak rizik, 1 – visok rizik). Funkcija vraća rezultat u JSON formatu koji sadrži predikciju i vjerovatnoću rizika.

Dodavanje novih podataka

Novi podaci o pacijentima se validiraju, pohranjuju i koriste za buduće treniranje modela.

```
26 @app.route("/add", methods=["POST"])
27 def add_data():
28     try:
29         df = load_data()
30         data = request.json
31
32         numeric_fields = ["age", "sex", "trestbps", "chol", "fbs", "thalch", "exang", "oldpeak", "ca"]
33         for field in numeric_fields:
34             if field not in data or data[field] in ["", None]:
35                 return jsonify({"status": "error",
36                                "message_bs": "Molimo popunite sva obavezna polja.",
37                                "message_en": "Please fill all required fields."}), 400
38                 data[field] = float(data[field])
39
40         string_fields = ["cp", "restecg", "slope", "thal"]
41         for field in string_fields:
42             if field not in data or data[field] in ["", None]:
43                 return jsonify({"status": "error",
44                                "message_bs": "Molimo popunite sva obavezna polja.",
45                                "message_en": "Please fill all required fields."}), 400
46
47         data.setdefault("id", len(df) + 1)
48         data.setdefault("dataset", "new")
49         data.setdefault("num", 0)
50
51         df = pd.concat([df, pd.DataFrame([data])], ignore_index=True)
52         df.to_csv("data/heart.csv", index=False)
53
54         return jsonify({"status": "success",
55                        "message_bs": "Pacijent uspješno dodan.",
56                        "message_en": "Patient added successfully."})
57
58     except Exception as e:
59         print("Add Data Error:", e)
60         return jsonify({"status": "error",
61                        "message_bs": f"Greška: {str(e)}",
62                        "message_en": f"Error: {str(e)}"}, 500
```

Ova funkcija `/add` omogućava dodavanje novih podataka o pacijentima u backend sistem. Prvo učitava postojeći skup podataka, a zatim prima JSON podatke sa zahtjeva. Funkcija provjerava da li su sva obavezna numerička i kategorijska polja popunjena i vrši konverziju numeričkih vrijednosti u odgovarajući tip. Nakon validacije, dodaju se interne kolone poput `id`, `dataset` i `num` kako bi novi zapis bio kompatibilan sa postojećim podacima. Na kraju, novi pacijent se dodaje u `DataFrame` i trajno pohranjuje u CSV datoteku, čime podaci postaju dostupni za buduće treniranje modela. U slučaju greške, funkcija vraća odgovarajući JSON odgovor sa opisom problema.

Ponovno treniranje modela

Korisnik može pokrenuti ponovno treniranje modela koristeći prošireni skup podataka, čime se osigurava kontinuirano poboljšanje tačnosti predikcija.

```
66 @app.route("/retrain", methods=["POST"])
67 def retrain():
68     try:
69
70         df = load_data()
71         df = df.dropna()
72         if df.empty:
73             return jsonify({
74                 "status": "error",
75                 "message_bs": "Nije moguće retrenirati model jer nema validnih korisnika.",
76                 "message_en": "Cannot retrain the model because there are no valid users."
77             })
78
79
80         os.system("python train_model.py")
81         global model, scaler, features
82         model, scaler, features = load_model()
83
84         return jsonify({
85             "status": "success",
86             "message_bs": "Model je uspješno retreniran i ponovno učitano.",
87             "message_en": "Model retrained and reloaded successfully."
88         })
89
90     except Exception as e:
91         print("Retrain Error:", e)
92         return jsonify({
93             "status": "error",
94             "message_bs": f"Greška prilikom retreniranja: {str(e)}",
95             "message_en": f"Error retraining model: {str(e)}"
96         }), 500
```

Ova funkcija `/retrain` u Flask API-ju omogućava ponovno treniranje modela na ažuriranom skupu podataka. Kada korisnik pošalje POST zahtjev na ovaj endpoint:

- Funkcija učitava postojeći dataset pozivom `load_data()` i uklanja sve redove sa nedostajućim vrijednostima (`NaN`).
- Provjerava se da li dataset sadrži validne podatke; ako ne, vraća se greška u JSON formatu.
- Ako su podaci validni, pokreće se skripta `train_model.py` koja trenira model koristeći Random Forest algoritam.
- Nakon što se treniranje završi, funkcija ponovno učitava model, scaler i listu feature-a u memoriju pomoću `load_model()`.
- Na kraju, funkcija vraća JSON odgovor s informacijom da je model uspješno retreniran i ponovno učitano.

U slučaju greške tokom procesa, funkcija hvata iznimku i vraća JSON odgovor sa detaljima o problemu.

Korisnički interfejs

Interfejs za predikciju rizika od srčanih bolesti:

The image displays three sequential screenshots of the HealthPredict web application, showing the process of entering patient data for heart disease risk prediction. The interface features a dark blue header with the 'HealthPredict' logo and a language selector set to 'Bosanski'. The background of the form area is a blue DNA double helix pattern.

Screenshot 1: 'Predikcija bolesti srca' form

- Starost: 40-70
- Spol: Muški
- Tip angine: tipična angina
- Krvni pritisak u mirovanju: 100-130
- Holesterol: (input field)

Screenshot 2: Continuation of the form

- (input field): 150-300
- Šećer u krvi na prazan stomak: Ne
- EKG u mirovanju: normalno
- Maksimalni puls: 60-120
- Angina izazvana vježbom: Ne
- Oldpeak: 0-6
- Nagib ST segmenta: ravan

Screenshot 3: Final form and prediction button

- Oldpeak: 0-6
- Nagib ST segmenta: ravan
- CA: 0-3
- Thal: normalno
- Predvidi rizik** (button)

Ovdje će biti prikazan rezultat procjene rizika za srčane bolesti.

Interfejs za dodavanje novih podataka:

HealthPredict

Jezik: Bosanski

Dodaj novog pacijenta

Starost

40-70

Spol

Odaberi

Tip angine

Odaberi

Krvni pritisak u mirovanju

100-130

Holesterol

150-300

HealthPredict

Jezik: Bosanski

Šećer u krvi na prazan stomak

Odaberi

EKG u mirovanju

Odaberi

Maksimalni puls

60-120

Angina izazvana vježbom

Odaberi

Oldpeak

0-6

Nagib ST segmenta

Odaberi

CA

0-3

HealthPredict

Jezik: Bosanski

Odaberi

Oldpeak

0-6

Nagib ST segmenta

Odaberi

CA

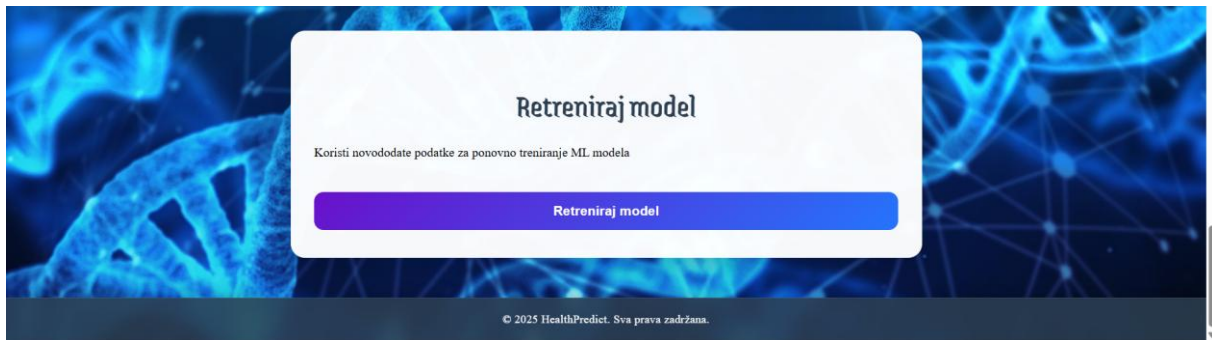
0-3

Thal

Odaberi

Dodaj pacijenta

Interfejs za ponovno treniranje modela:



Zaključak

Aplikacija HealthPredict predstavlja moćan alat za procjenu rizika od srčanih bolesti koristeći podatke o pacijentima i mašinsko učenje. Korištenjem Random Forest algoritma, aplikacija omogućava preciznu predikciju rizika, pruža uvid u najutjecajnije karakteristike svakog pacijenta i olakšava dodavanje novih podataka radi kontinuiranog poboljšanja modela.

Integracija jednostavnog korisničkog sučelja, mogućnosti višjejezične podrške i funkcionalnosti ponovnog treniranja modela osigurava da korisnici – bilo da su zdravstveni stručnjaci ili sami pacijenti – mogu donositi informirane odluke temeljem pouzdanih predikcija.

Ova aplikacija pokazuje značajan potencijal u primjeni AI tehnologija u zdravstvu, čineći procjenu rizika dostupnijom, bržom i preciznijom.