

# Uporaba vmesnika za namensko programiranje z namenom razpoznavanja vzorcev

Avtor: Timotej Petrovčič

Mentorja: izr. prof. dr. Simon Dobrišek, as. dr. Klemen Grm



Predmet: Razpoznavanje vzorcev

Datum: 8 Januar 2023

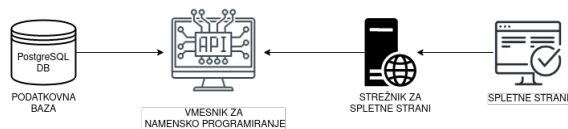
## Kazalo:

1	Namen . . . . .	1
2	Vmesnik za namensko programiranje . . . . .	1
3	Relacijska podatkovna baza . . . . .	3
4	Postavitev storitve . . . . .	3
5	Preizkušanje delovanja vmesnika . . . . .	4
6	Povzetek . . . . .	5
7	Viri in literatura . . . . .	5

**Ključne besede:** Vmesnik za namensko programiranje (API - Application programming interface), strežnik, odjemalec, podatkovna baza, transformacija Karhunea in Loeveja (PCA - analiza glavnih komponent), razvrščanje z enim najbližjim sosedom (1-NN)

## 1 Namen

Namen projekta je izdelava vmesnika za namensko programiranje na področju razpoznavanja vzorcev ter izdelava tako strežniške kot odjemalniške strani (backend/frontend). Prvotno smo se omejili na uvoz podatkov iz "UCI machine learning" repozitorija, saj so vzorčni podatki strukturirani tako, da vsak vzorec predstavlja novo vrstico v poljubni datoteki.



Slika 1: Struktura delovanja

## 2 Vmesnik za namensko programiranje

Vmesnik za namensko programiranje (Od sedaj naprej vmesnik) je bil napisan v programskem jeziku Python ter se v grobem deli na dva nivoja in sicer:

- Nivo spletnega dostopa ("recognition\_api")
- Nivo podatkovnega dostopa ("recognition.DAL")

Omenjena ločitev je pomembna, ker loči spletno stran vmesnika od podatkovne strani ter s tem omogoči neodvisno delovanje obeh nivojev, prepreči neželene napake v programu in izboljša varnost celotnega sistema.

Za izdelavo vmesnika smo uporabili asinhrono knjižnice, ki omogočajo sočasno/asinhrono izvajanje operacij (Concurrent execution).

### Nivo spletnega dostopa

Nivo spletnega dostopa je v grobem sestavljen iz 4 delov:

- Uvoz spremenljivk okolja
- Določitev različice izgradnje
- Določitev statičnih poti
- Inicializacija strežniške strani ter zagon vmesnika

Uvoz spremenljivk okolja omogoča dinamično upravljanje s programom preko .env datoteke v korenski mapi. Spremenljivke okolja v našem primeru določajo način delovanja, bazen povezav na podatkovno bazo, koren vmesniških poizvedb ter enolični identifikator vira za podatkovno bazo. Trenutna implementacija ne dovoljuje zagona vmesnika brez pravilno nastavljenih spremenljivk okolja, čemur se lahko izognemo ob dodatni nastavitvi privzetih vrednosti spremenljivkam okolja.

```
# Get environment variables
APP_CONFIG = os.getenv("APP_CONFIG")
URL_PREFIX = os.getenv("API_URL_PREFIX")
DB_URI = str(os.getenv("API_DB_CONNECTION_STRING"))
DB_MAX_CONNECTIONS = int(os.getenv("API_DB_MAX_CONNECTIONS"))
DB_POOL_RECYCLE = int(os.getenv("API_DB_POOL_RECYCLE"))
DB_MAX_OVERFLOW = int(os.getenv("API_DB_MAX_OVERFLOW"))
DB_POOL_TIMEOUT = int(os.getenv("API_DB_POOL_TIMEOUT"))
```

Slika 2: Spremenljivke okolja

Določitev različice izgradnje nastavi nivo beleženja podatkov med delovanjem vmesnika ter parametre bazena povezav na podatkovno bazo. Ločimo dve različici izgradnje: produkcijska različica ter razvojna različica. Različici izgradnje omogočata različne nastavitve delovanja vmesnika pred samo inicializacijo ter s tem odpravita neželene napake ali napačno beleženje delovanja vmesnika.

```
if __name__ == '__main__':
    # Set up operation mode for server
    if APP_CONFIG == 'dev':
        # Development build
        logging.basicConfig(level=logging.DEBUG)
        log = logging.getLogger()
        log.info("Running in development config")

    elif APP_CONFIG == 'prod':
        # Production build
        logging.basicConfig(level=logging.INFO)
        log = logging.getLogger()
        log.info("Running in production config")

    else:
        # If APP_CONFIG env variable is not set abort start
        logging.basicConfig(level=logging.INFO)
        log = logging.getLogger()
        log.info("Environment variable APP_CONFIG is not set")
        sys.exit(1)
```

Slika 3: Določitev različice izgradnje

Statične poti so sprva določene z usmerjevalno tabelo in dekoratorji poti, ki jim določijo krmilnike dogodkov. Vsak krmilnik sprejme različne argumente, a v splošnem smo podatke sprejemali v JSON formatu, saj je le-ta enostaven za uporabo ter omogoča dobro preverjanje napak v podatkih poizvedb.

```
# Configure routes table and all available methods
routes = web.RouteTableDef()

# Healthcheck
@routes.get('/healthz')
async def health_check(request):
    log.info("Api test running\n")
    return web.Response(text="## API test successfull ##\n")
```

Slika 4: Primer določitve statične poti

Inicializacija strežniške strani se začne s pridobitvijo zanke dogodkov. Nato se izvede inicializacija, ki zajame dodajanje dodatnih atributov objektu aplikacijske storitve. V našem primeru je to inicializacija modela podatkovne baze ("db\_model"), statičnih poti ter korena vmesniških poizvedb. Zagon vmesnika se zgodi po inicializaciji ter teče v neskončnost pod pogojem, da ga ne prekinemo.

```
# Get asyncio loop
loop = asyncio.get_event_loop()

# Create WebServer object and initialize it
server = API_Server()
loop.run_until_complete(server.initialize())

# Start the server
server.start_server(host='0.0.0.0', port=5000, loop=loop)
```

Slika 5: Pridobitev zanke dogodkov ter zagon vmesnika

```
# Initialization for API app object
async def initialize(self):
    log.info("API Initialization started")
    self.subapp = web.Application()

    log.info("Configuring API SQLAlchemy engine")
    self.subapp['engine'] = create_async_engine(DB_URI,
                                                echo = (False if APP_CONFIG == 'prod' else True),
                                                connect_args=({"server_settings": {"jit": "off"}},
                                                                pool_pre_ping=True,
                                                                pool_size=DB_MAX_CONNECTIONS,
                                                                pool_recycle=DB_POOL_RECYCLE,
                                                                max_overflow=DB_MAX_OVERFLOW,
                                                                pool_timeout=DB_POOL_TIMEOUT))

    log.info("Binding SQLAlchemy engine to application object")
    ahsa.setup(self.subapp, [ahsa.bind(DB_URI)])

    log.info("Creating DB tables")
    await ahsa.init_db(self.subapp, metadata)

    log.info("Adding routes to application object")
    self.subapp.router.add_routes(self.routes)

    # Add sub-app to set the IP/recognition-api request
    self.app = web.Application()
    self.app.add_subapp(URL_PREFIX, self.subapp)

    log.info("API initialization complete")

# Run API
def start_server(self, host, port, loop):
    log.info("Server starting on address: http://{}:{}".format(host, port))
    web.run_app(self.app, host=host, port=port, loop=loop)
```

Slika 6: Inicializacija vmesnika

## Nivo podatkovnega dostopa

Nivo podatkovnega dostopa določa logiko delovanja različnih poizvedb na vmesnik. V splošnem je sestavljen iz poizvedbe na podatkovno bazo, primerne obdelave podatkov in vrnitve ali je bila izbrana metoda uspešno oziroma neuspešno zaključena.

Za razpoznavanje vzorcev sta relevantni funkciji:

- *principalComponentAnalysis*
- *optimizedPrincipalComponentAnalysis*

Obe funkciji v splošnem izvajati operacijo analize glavnih komponent s pomočjo knjižnice *scikit-learn*.

Analiza glavnih komponent je statistična metoda obdelave podatkov, ki omogoča zmanjšanje dimenzije vzorcev ob ohranjanju največje količine splošne informacije za razvrščanje le-teh. Postopek je v grobem sestavljen iz naslednjih korakov:

- Normalizacija vzorcev
- Izračun kovariančne matrike normaliziranih vzorcev
- Izračun lastnih vrednosti in lastnih vektorjev kovariančne matrike
- Izbira končnega števila komponent vzorcev, ki je enako ali manjše številu izbranih lastnih vrednosti/vektorjev razporejenih od največjega do najmanjšega (Najprej izberemo večje nato manjše vrednosti)
- Pretvorba normaliziranih vzorcev v nove vzorce z manjšo dimenzijo

V našem primeru normalizacijo podatkov izvede razred *StandardScaler*, ki osnovnim vzorcem odšteje povprečno vrednost ter rezultat deli s standardnim odklonom vzorcev.

Sledi uporaba *PCA* razreda, ki izvede izračun kovariančne matrike vzorcev, določi lastne vrednosti oziroma vektorje preko SVD algoritma ter izvede pretvorbo v izbrani N-dimenzijski prostor z uporabo N-največjih lastnih vrednosti. Sledi izračun variance, ki določi kako velik vpliv ima na osnovne normalizirane podatke posamezna komponenta. Izračuna se kot razmerje med vsoto izbranih lastnih vrednosti za zmanjšanje dimenzije ter vsoto vseh lastnih vrednosti.

Po zgornjem koraku se funkciji *principalComponentAnalysis* in *optimizedPrincipalComponentAnalysis* razčlenita točno zaradi variance. Ne-optimizirana funkcija izbira komponente po vrsti, medtem ko optimizirana izbere tisto kombinacijo komponent, ki predstavlja največjo varianco. Večja varianca torej omogoča bolj natančno razpoznavanje po izvedbi analize.

Sledi preizkus med osnovnimi normaliziranimi podatki ter dimenzijsko reduciranimi normaliziranimi podatki s pomočjo razpoznavnika enega najbližjega sosedu (1-NN). Za razpoznavnik uporabimo razred *KNeighborsClassifier*, kateremu podamo argument *n\_neighbors* enak ena, saj s tem določimo število najbližjih sosedov, ki jih bo razpoznavnik upošteval.

```
# Principal Component Analysis
numpyFeaturesList = np.array(listOfFeatures)
sc = StandardScaler()
X_normalised = sc.fit_transform(numpyFeaturesList)

# Loop over all possible number of components for PCA and find the least number of components with requested accuracy
pca = PCA()
pca.fit(X_normalised)
cumsum = np.cumsum(pca.explained_variance_ratio_)
optimalNumOfComponents = np.argmax(cumsum >= json_data.get("RequestedVariance")) + 1

# Execute PCA for optimal number of components for requested accuracy
pca = PCA(n_components = optimalNumOfComponents)
X_optimal_pca = pca.fit_transform(X_normalised)
explained_variance = pca.explained_variance
logging.info("Explained variance for {} components: {}".format(optimalNumOfComponents, explained_variance))

# k-NN model definition and cross-correlation comparison score
knn = KNeighborsClassifier(n_neighbors = 1)
scores_original = cross_val_score(knn, X_normalised, listOfFeatureNames)
scores_pca = cross_val_score(knn, X_optimal_pca, listOfFeatureNames)
```

Slika 6: Izsek funkcije *principalComponentAnalysis*

### 3 Relacijska podatkovna baza

Strukturo relacijske podatkovne baze ("db\_model") določijo sami podatki, ki jih uporabljamo. Vsaka podatkovna zbirka vsebuje več vzorcev, torej lahko sklepamo, da bo izbrana shema v obliki eden proti mnogim (*one-to-many relationship database*). V našem primeru ustvarimo znotraj podatkovne baze dve tabeli, prvo za imena različnih podatkovnih zbirk ter drugo za posamezne vzorce v zbirkah. Vsakemu elementu v drugi tabeli moramo torej ustvariti referenco na element v prvi tabeli, saj bomo le-tako lahko vedeli, kateri vzorec pripada posamezni zbirki v podatkovni bazi. Prav tako je potrebno upoštevati kaskadno brisanje vzorcev, saj bo v primeru izbrisa podatkovne zbirke iz prve tabele potrebno izbrisati vse vzorce iz druge tabele, ki se nanašajo na izbrano podatkovno zbirko.

```
metadata = Metadata()
Base = declarative_base(metadata = metadata)

# One feature set to have multiple children which represent separate features
class Feature_set(Base):
    __tablename__ = 'Feature_set_table'

    ID = Column(Integer, primary_key = True, unique = True, autoincrement = True)
    created = Column(DateTime, default = datetime.now)
    last_updated = Column(DateTime, default = datetime.now, onupdate = datetime.now)
    name = Column(String(30), unique = True)

    feature_children = relationship('Features', cascade='all,delete', backref='Feature_set_table', )

class Features(Base):
    __tablename__ = 'Features_table'

    ID = Column(Integer, primary_key = True, autoincrement = True)
    feature = Column(String(10000))

    parent_id = Column(Integer, ForeignKey('Feature_set_table.ID', ondelete='CASCADE'))
    parent = relationship('Feature_set', backref = backref('Features_table', cascade = 'all, delete
```

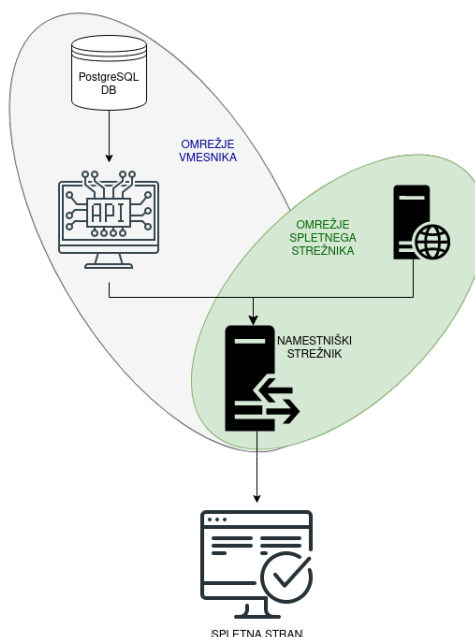
Slika 7: Model podatkovne baze

### 4 Postavitev storitve

Za postavitev storitve smo uporabili orodje Docker, ki omogoči optimizirano delovanje aplikacij glede na obremenjenost delujoče storitve. Prav tako smo za ažurnost knjižnic v Pythonu uporabili orodje Poetry, ki omogoči določitev specifičnih izdaj knjižnic ter tako izloči možnost kasnejše programske nekompatibilnosti. Znotraj orodja Docker smo storitev delili na sledeče podenote:

- Vmesnik za namensko programiranje
- Podatkovna baza
- Spletni strežnik (Narejena zasnova)
- Namestniški strežnik

Vsaka izmed podenot ima svoj zagonski cikel ter je odvisna od ostalih. Poleg osnovnih treh enot (vmesnik, podatkovna baza in spletni strežnik) smo dodali tudi namestniški strežnik, ki predstavlja dodaten nivo varnosti v primeru izpostavitve vmesnika internetu ter hitrejše nalaganje spletnega vmesnika s pomočjo statične hrambe določenih skript. Prav tako je bilo potrebno postaviti pravilno omrežje med aplikacijami samimi ter omejiti dostop do določenih.



Slika 8: Omrežje postavitev storitve



*Slika 11: Primer delovanja strežnika ob zgornjih treh proizvodbah*

## 6 Povzetek

Trenutno je izdelan prototip delovanja storitve (Vmesnika ter spletnega strežnika), ki ga je kasneje mišljeno izboljšati in razširiti na poljubno velikost (Zavzame 8 izbirnih projektov). Glede same analize glavnih komponent pa se je izkazalo, da je z reduciranjem komponent mogoče doseči le omejeno količino razpoznanih podatkov. Najboljši rezultat smo dobili pri 2 komponentah, kjer je najnižja razpoznana vrednost po reduciranju bila 80% z uporabo osnovnih podatkov pa smo dosegli najnižjo vrednost 90%.

## 7 Viri in literatura

- S. Dobrišek: Razpoznavanje vzorcev, gradivo za predavanja, FE 2022
- K. Grm: Razpoznavanje vzorcev, gradivo za laboratorijske vaje, FE 2022
- Dokumentacija orodja Docker
- Dokumentacija namestniškega strežnika NGINX
- Dokumentacija o jezikih HTML5, CSS ter JS
- Dokumentacija uporabljenih knjižnic v Pythonu:
  - Vmesnik - "recognition\_api": *logging, os, sys, asyncio, asyncpg, aiohttp, aiohttp\_sqlalchemy, sqlalchemy*
  - Vmesnik - "recognition\_DAL": *logging, numpy, sqlalchemy, sklearn*
  - Vmesnik - "db\_model": *datetime, sqlalchemy*
  - Spletni strežnik - "webserver": *logging, os, sys, asyncio, aiohttp*