

# **Procjena financijskog rizika defaulta**

## Bayesovski optimizirani gradient boosting i objašnjivost pomoću SHAP-a

Matej Čiček

23. siječnja 2026.

**Sažetak.** U radu je razvijen reproducibilan sustav za procjenu rizika financijskog neispunjenja obveza (default/bankruptcy) koji integrira heterogene izvore podataka (CSV i XLSX), provodi naprednu predobradu (MICE imputacija, robustan tretman outliera), trenira i optimizira modele (logistička regresija, random forest i XGBoost) te uvodi Bayesovu optimizaciju hiperparametara pomoću Optuna okvira. Sustav uključuje rigoroznu validaciju (nested cross-validation), kalibraciju vjerojatnosti (izotonička kalibracija), interpretabilnost SHAP metodom te segmentaciju uzoraka (K-means). Integrirani podaci i metrike pohranjuju se u SQLite bazu, a rezultati i predikcije izlažu se kroz REST API (FastAPI), što omogućuje jednostavnu integraciju u druge aplikacije.

**Ključne riječi:** kreditni rizik, default, XGBoost, Optuna, SHAP, kalibracija, K-means, reproducibilnost, SQLite, FastAPI.

# Sadržaj

<b>Popis kratica</b>	<b>vi</b>
<b>1 Uvod</b>	<b>1</b>
1.1 Motivacija . . . . .	1
1.2 Ciljevi projekta . . . . .	1
1.3 Korišteni skupovi podataka i heterogenost . . . . .	2
1.4 Organizacija dokumenta . . . . .	2
<b>2 Glavni dio teme: formalizam, definicije i teorija</b>	<b>3</b>
2.1 Binarna klasifikacija i probabilističko predviđanje . . . . .	3
2.2 Metrike evaluacije . . . . .	3
2.3 Logistička regresija . . . . .	3
2.4 Random Forest . . . . .	4
2.5 XGBoost i gradijentno pojačavanje . . . . .	4
2.6 Bayesova optimizacija hiperparametara (Optuna) . . . . .	4
2.7 Nested cross-validation . . . . .	4
2.8 Imputacija: MICE . . . . .	4
2.9 Detekcija anomalija: Isolation Forest . . . . .	4
2.10 Kalibracija vjerojatnosti . . . . .	4
2.11 Interpretabilnost: SHAP . . . . .	5
2.12 Segmentacija: K-means . . . . .	5
<b>3 Opis implementacije (komponente sustava)</b>	<b>6</b>
3.1 Arhitektura rješenja . . . . .	6
3.2 Struktura repozitorija . . . . .	6
3.3 Reproducibilnost: Git i DVC . . . . .	7
3.4 Akvizicija podataka (Kaggle) . . . . .	7
3.5 Učitavanje i validacija integriteta . . . . .	8
3.6 EDA i statistička analiza . . . . .	8
3.7 Tretman outliera i anomalija . . . . .	10
3.8 Predobrada: MICE imputacija, skaliranje i kodiranje . . . . .	11
3.9 Integracija heterogenih skupova . . . . .	12
3.10 Selekcija značajki i baseline modeli . . . . .	13
3.11 Optuna + XGBoost optimizacija . . . . .	13
3.12 Rigorozna validacija: nested CV . . . . .	14
3.13 Odabir praga odluke . . . . .	15
3.14 Kalibracija vjerojatnosti . . . . .	16
3.15 SHAP interpretabilnost . . . . .	18
3.16 K-means segmentacija . . . . .	20

3.17	Pohrana integriranih podataka i metrika u SQLite . . . . .	21
3.17.1	Shema baze . . . . .	22
3.18	REST API sučelje (FastAPI) . . . . .	22
3.19	Pokretanje servisa . . . . .	23
<b>4</b>	<b>Prikaz rada aplikacije</b>	<b>24</b>
4.1	Pipeline: od podataka do modela . . . . .	24
4.2	Testiranje API-ja . . . . .	24
4.3	Primjer predikcije . . . . .	25
<b>5</b>	<b>Rezultati i vizualizacije</b>	<b>26</b>
5.1	Tablice metrika . . . . .	26
5.2	Nested CV rezultati . . . . .	26
5.3	Ključni grafovi . . . . .	26
<b>6</b>	<b>Kritički osvrt</b>	<b>29</b>
6.1	Praktična izvedivost . . . . .	29
6.2	Primjena i ograničenja . . . . .	29
6.3	Moguća poboljšanja . . . . .	30
<b>7</b>	<b>Zaključak</b>	<b>31</b>
<b>A</b>	<b>Dodatci</b>	<b>33</b>
A.1	Popis slika koje je potrebno umetnuti . . . . .	33

# Popis slika

3.1	Arhitektura sustava (predložak; korisnik umeće vlastitu sliku). . . . .	6
3.2	Histogrami numeričkih varijabli (Skup A). . . . .	9
3.3	Matrica korelacija (Skup A). . . . .	10
3.4	Kalibracijska krivulja prije kalibracije (predložak). . . . .	17
3.5	Kalibracijska krivulja nakon izotoničke kalibracije (predložak). . . . .	17
3.6	SHAP summary plot: globalna važnost značajki (predložak). . . . .	19
3.7	SHAP waterfall: lokalno objašnjenje za jedan uzorak (predložak). . . . .	20
3.8	K-means klasteri u PCA projekciji (predložak). . . . .	21
4.1	Primjer prikaza rada API-ja (predložak; korisnik umeće sliku). . . . .	25
5.1	ROC krivulja za XGBoost model (predložak). . . . .	27
5.2	Precision–Recall krivulja za XGBoost model (predložak). . . . .	28

# Popis tablica

5.1	Usporedba modela (primjer popunjavanja). . . . .	26
5.2	Nested CV rezultati (primjer). . . . .	26

# Popis kratica

AUC	Area Under ROC Curve
PR-AUC	Area Under Precision–Recall Curve
CV	Cross-Validation
MICE	Multiple Imputation by Chained Equations
SHAP	SHapley Additive exPlanations
API	Application Programming Interface
DVC	Data Version Control
PCA	Principal Component Analysis

# Poglavlje 1

## Uvod

### 1.1 Motivacija

Procjena financijskog rizika neispunjenja obveza (defaulta) predstavlja središnji problem u području kreditiranja, upravljanja portfeljem i regulatorne usklađenosti financijskih institucija. Odluke o odobravanju kredita, određivanju kamatnih stopa i alokaciji kapitala izravno ovise o sposobnosti modela da precizno i pouzdano procijeni vjerojatnost nastupa defaulta. Za razliku od općih klasifikacijskih problema, u kreditnom riziku nije dovoljna samo binarna odluka (odobriti/odbiti), već je nužna probabilistička predikcija: model mora generirati dobro kalibrirane vjerojatnosti koje se mogu izravno koristiti za definiranje pragova odlučivanja, izračun očekivanog gubitka (Expected Loss) i cjenovno upravljanje rizikom.

U industrijskoj praksi linearni modeli, poput logističke regresije, i dalje se koriste zbog svoje jednostavnosti i interpretabilnosti, no često ne uspijevaju uhvatiti složene nelinearne odnose prisutne u stvarnim financijskim podacima. Zbog toga ansambl metode, osobito gradijentno pojačavanje stabala odluke, predstavljaju dominantan pristup u modernim sustavima procjene kreditnog rizika. XGBoost se u tom kontekstu ističe kao de facto standard zahvaljujući visokoj prediktivnoj učinkovitosti, skalabilnosti i fleksibilnosti u radu s neravnotežnim skupovima podataka.

Međutim, povećanje modelne složenosti donosi i značajan izazov: smanjenu transparentnost odluka. Regulatorni okviri i poslovna praksa zahtijevaju da odluke o kreditiranju budu objašnjive i auditabilne, kako bi se osiguralo povjerenje korisnika i usklađenost s propisima. U tom smislu, metode objašnjivosti poput SHAP-a omogućuju formalizirano tumačenje doprinosa pojedinih značajki na razini cijelog modela i pojedinačnih predikcija. Time se premošćuje jaz između visoke prediktivne snage suvremenih modela i potrebe za razumijevanjem i nadzorom njihovih odluka.

Motivacija ovog projekta proizlazi upravo iz potrebe za integracijom snažnih modela strojnog učenja s metodama objašnjivosti i rigorozne validacije, kako bi se izgradio sustav koji je istovremeno učinkovit, pouzdan i primjenjiv u realnom financijskom okruženju.

### 1.2 Ciljevi projekta

Projekt je definiran kao inženjerski i istraživački pipeline s ciljevima:

- prikupljanje najmanje dva heterogena skupa podataka i integracija u zajednički okvir,



- implementacija predobrade (imputacija, tretman outliera, kodiranje i skaliranje),
- usporedba baseline modela s Bayesovski optimiziranim XGBoost modelom,
- rigorozna validacija (nested CV) i optimizacija praga odluke,
- kalibracija vjerojatnosti i procjena kvalitete probabilističkih predikcija,
- interpretabilnost (globalna i lokalna) pomoću SHAP-a,
- segmentacija uzoraka (K-means) i analiza klastera,
- pohrana integriranog skupa, metrika i artefakata u bazu te izrada API sučelja.

## 1.3 Korišteni skupovi podataka i heterogenost

Korištena su najmanje dva izvora podataka:

- **Skup A (CSV):** *Company Bankruptcy Prediction* [1] s ciljem **Bankrupt?**.
- **Skup B (XLSX):** *Credit Risk Modelling Dataset* [2] s binarnom ciljanom varijablom (npr. `default`).

Skupovi su heterogeni po formatu (CSV i Excel), strukturi i semantici značajki.

## 1.4 Organizacija dokumenta

U nastavku se daje teorijska podloga, opis implementacije, prikaz rada aplikacije, kritički osvrt i zaključak, uz literaturu u IEEE stilu. Svaki relevantan isječak koda opisan je u tekstu prije ili poslije listinga.

# Poglavlje 2

## Glavni dio teme: formalizam, definicije i teorija

### 2.1 Binarna klasifikacija i probabilističko predviđanje

Neka su  $\mathbf{x} \in \mathbb{R}^d$  značajke, a  $y \in \{0, 1\}$  ciljana oznaka gdje  $y = 1$  označava default/bankrot. Model procjenjuje vjerojatnost

$$p(y = 1 \mid \mathbf{x}) \in [0, 1]. \quad (2.1)$$

U kreditnom riziku takva vjerojatnost se često mapira u odluku preko praga  $\tau$ :

$$\hat{y} = \mathbb{I}[p(y = 1 \mid \mathbf{x}) \geq \tau]. \quad (2.2)$$

### 2.2 Metrike evaluacije

Za neravnotežne probleme AUC i PR-AUC pružaju dopunske informacije. ROC-AUC mjeri rangiranje bez obzira na prag, dok PR-AUC bolje reflektira kvalitetu predikcije pozitivne klase. Za izabrani prag koristi se F1 mjera:

$$F1 = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}. \quad (2.3)$$

Za probabilističku kvalitetu koristi se Brier score [6]:

$$\text{Brier} = \frac{1}{n} \sum_{i=1}^n (p_i - y_i)^2. \quad (2.4)$$

### 2.3 Logistička regresija

Logistička regresija modelira:

$$p(y = 1 \mid \mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b), \quad (2.5)$$

gdje je  $\sigma$  logistička funkcija. Prednost je interpretabilnost i stabilnost, no ograničenje je slabija sposobnost modeliranja nelinearnih odnosa.

## 2.4 Random Forest

Random Forest (RF) trenira ansambl stabala na bootstrap uzorcima, uz slučajni odabir značajki po splitu. RF je robustan i često dobar baseline, no može imati slabiju kalibraciju i slabije performanse u odnosu na specijalizirane metode gradijentnog pojačavanja.

## 2.5 XGBoost i gradijentno pojačavanje

XGBoost [3] je učinkovita implementacija gradijentnog pojačavanja stabala odluke s regularizacijom. U neravnotežnim problemima tipično se koristi `scale_pos_weight` kako bi se penalizirale pogreške na manjinskoj klasi.

## 2.6 Bayesova optimizacija hiperparametara (Optuna)

Umjesto grid/random pretrage koristi se Bayesova optimizacija hiperparametara s Optuna okvirom [4]. Optuna iterativno predlaže hiperparametre (trialove) prema prethodnim rezultatima i maksimizira odabranu metriku (u ovom projektu PR-AUC).

## 2.7 Nested cross-validation

Nested CV razdvaja optimizaciju hiperparametara (inner loop) od procjene generalizacije (outer loop). Time se smanjuje optimistična pristranost pri procjeni performansi modela, posebno kada je prostor hiperparametara velik.

## 2.8 Imputacija: MICE

Za numeričke značajke koristi se IterativeImputer koji implementira MICE pristup [7]. Umjesto imputacije prosjekom/medijanom, MICE iterativno modelira svaku varijablu na temelju ostalih, često dajući realističnije imputacije.

## 2.9 Detekcija anomalija: Isolation Forest

Isolation Forest [8] izolira anomalije kroz slučajne particije prostora podataka. U projektu se koristi za izgradnju dodatne binarne značajke `outlier_flag` koja signalizira potencijalno atipične uzorke.

## 2.10 Kalibracija vjerojatnosti

Ansambl modeli mogu davati nekalibrirane vjerojatnosti. Izotonička kalibracija uči monotonu transformaciju koja usklađuje predikcije s empirijskim učestalostima [6]. Usporedba Brier score-a prije i poslije kalibracije daje jasan signal korisnosti.

## 2.11 Interpretabilnost: SHAP

SHAP [5] daje aditivna objašnjenja predikcija i povezuje doprinos značajki s formalizmom Shapleyjevih vrijednosti. Time se dobiva i globalna važnost značajki (summary plot) i lokalno objašnjenje za pojedini uzorak (waterfall).

## 2.12 Segmentacija: K-means

K-means [9] particionira uzorke u  $k$  klastera minimiziranjem unutar-klasterske varijance. Silhouette score se koristi za izbor  $k$ , a PCA projekcija omogućuje vizualnu provjeru separacije.

# Poglavlje 3

## Opis implementacije (komponente sustava)

### 3.1 Arhitektura rješenja

Sustav se sastoji od sljedećih komponenti:

1. Akvizicija podataka (Kaggle CLI) i validacija integriteta.
2. Predobrada i integracija heterogenih skupova u zajednički prostor značajki.
3. Modeliranje: baseline modeli i optimizirani XGBoost.
4. Validacija: standardna CV i nested CV.
5. Post-proces: threshold tuning, kalibracija, SHAP objašnjenja, K-means segmentacija.
6. Pohrana: SQLite baza za integrirane podatke i metrike.
7. Sučelje: REST API (FastAPI) za dohvat metrika i predikcije.

Slika 3.1: Arhitektura sustava (predložak; korisnik umeće vlastitu sliku).

### 3.2 Struktura repozitorija

Sljedeći isječak prikazuje preporučenu strukturu projekta. Odvajanje podataka, artefakata i izvornog koda olakšava održavanje i reproducibilnost.

```
1 default-risk-project/  
2   data/  
3     raw/  
4     interim/  
5     processed/  
6   artifacts/  
7     fig/  
8     models/  
9   src/  
10  api.py
```

```

11     features.py
12     train.py
13     default_risk.db
14     dvc.yaml
15     .gitignore
16     README.md

```

Listing 3.1: Primjer strukture projekta

### 3.3 Reproducibilnost: Git i DVC

Ovaj isječak demonstrira inicijalizaciju Git i DVC-a [10]. DVC prati velike datoteke i omogućuje vraćanje točne verzije skupa podataka korištene u eksperimentu.

```

1 git init
2 dvc init
3 dvc add data/raw/
4 git add .
5 git commit -m "Init: repo + DVC + raw data tracking"

```

Listing 3.2: Inicijalizacija Git i DVC

### 3.4 Akvizicija podataka (Kaggle)

U nastavku je kod koji provjerava dostupnost Kaggle tokena i preuzima dva skupa. Ovaj korak zadovoljava zahtjev projekta za dokumentiranje načina prikupljanja i citiranje izvora [1, 2].

```

1 from pathlib import Path
2 import os
3 import subprocess
4
5 def run(cmd):
6     r = subprocess.run(cmd, shell=True, capture_output=True, text
7                         =True)
8     return r.returncode, r.stdout.strip(), r.stderr.strip()
9
10 def kaggle_ready():
11     p = Path.home() / ".kaggle" / "kaggle.json"
12     if p.exists():
13         return True
14     if os.environ.get("KAGGLE_USERNAME") and os.environ.get("
15         KAGGLE_KEY"):
16         return True
17     return False
18
19 def kaggle_download(dataset_slug, out_dir):
20     out_dir = Path(out_dir)
21     out_dir.mkdir(parents=True, exist_ok=True)
22     cmd = f'kaggle datasets download -d "{dataset_slug}" -p "{
23         out_dir.as_posix()}" --unzip -q'

```

```

21     return run(cmd)
22
23 slug_a = "fedesoriano/company-bankruptcy-prediction"
24 slug_b = "abhirajmandal/credit-risk-modelling-dataset"
25
26 if kaggle_ready():
27     kaggle_download(slug_a, "data/raw/bankruptcy_csv")
28     kaggle_download(slug_b, "data/raw/credit_xlsx")

```

Listing 3.3: Preuzimanje heterogenih skupova (CSV i XLSX) s Kaggle-a

## 3.5 Učitavanje i validacija integriteta

Nakon preuzimanja, slijedi učitavanje i osnovna validacija (dimenzije, duplikati, udio nedostajućih). Ovaj korak služi kao kontrolna točka prije ulaska u predobradu.

```

1 import pandas as pd
2 import numpy as np
3
4 df_a_raw = pd.read_csv("data/raw/bankruptcy_csv/CompanyBankruptcy
5 .csv")
6 df_b_raw = pd.read_excel("data/raw/credit_xlsx/case_study1.xlsx")
7
8 def basic_validation(df):
9     return {
10         "redaka": int(df.shape[0]),
11         "stupaca": int(df.shape[1]),
12         "duplikata": int(df.duplicated().sum()),
13         "udio_nedostajucih": float(df.isna().mean().mean())
14     }
15
16 va = basic_validation(df_a_raw)
17 vb = basic_validation(df_b_raw)
18 va, vb

```

Listing 3.4: Učitavanje i osnovna provjera integriteta

## 3.6 EDA i statistička analiza

Ovaj dio pokazuje primjer EDA postupaka: histogrami numeričkih varijabli i prikaz korelacija. Takvi grafovi pomažu identificirati ekstremne distribucije, potencijalne transformacije i multikolinearnost.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 num_cols_a = df_a_raw.select_dtypes(include=[np.number]).columns.
5 tolist()
6 df_a_raw[num_cols_a].iloc[:, :12].hist(bins=30, figsize=(12,6))
7 plt.tight_layout()

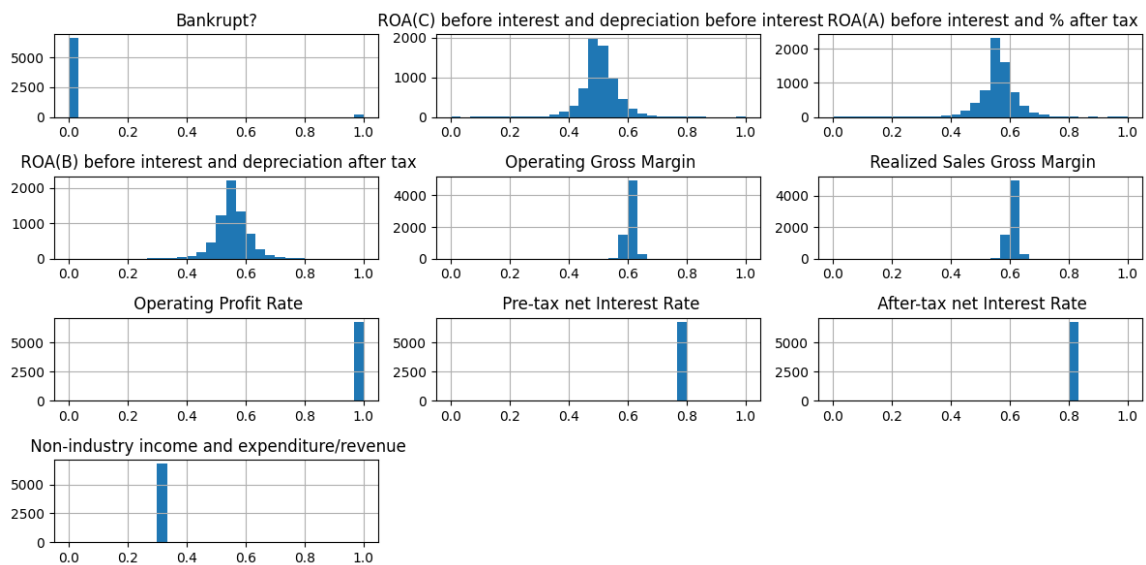
```

```

7 plt.show()
8
9 corr = df_a_raw[num_cols_a].corr(method="pearson")
10 plt.figure(figsize=(9,6))
11 plt.imshow(corr.values, aspect="auto")
12 plt.colorbar()
13 plt.title("Pearson korelacije (Skup A)")
14 plt.tight_layout()
15 plt.show()

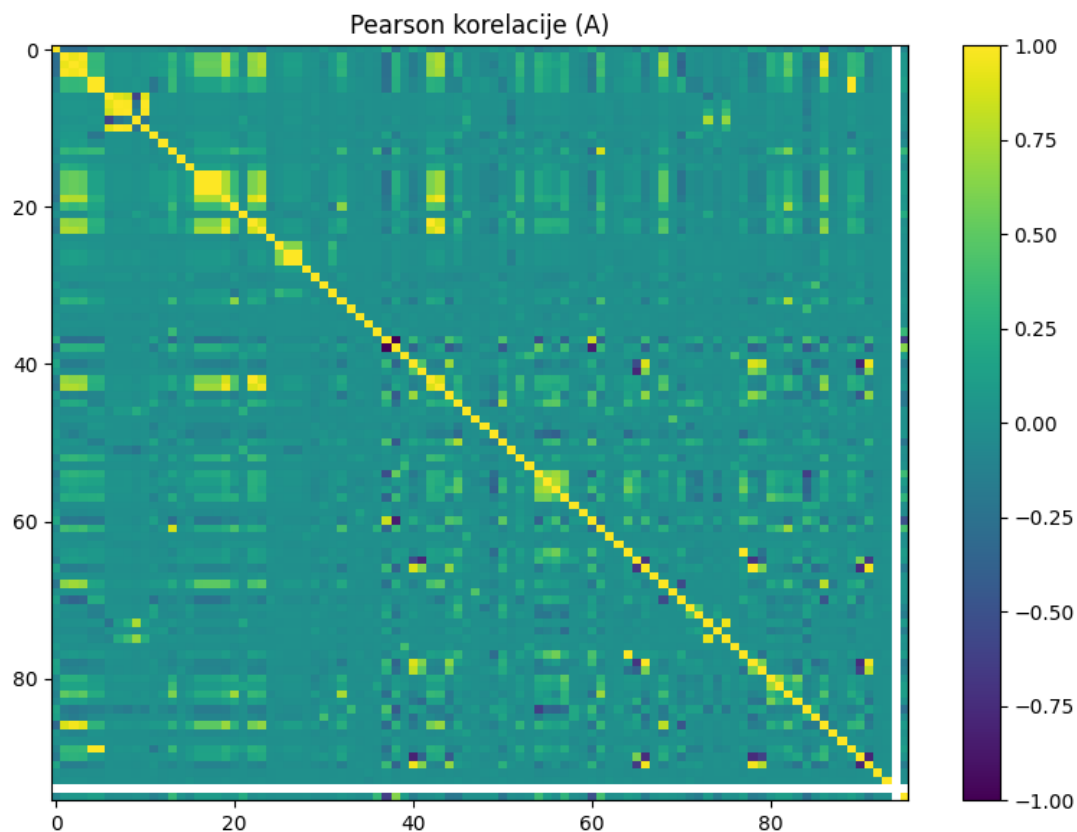
```

Listing 3.5: Distribucije i korelacije (primjer)



Slika 3.2: Histogrami numeričkih varijabli (Skup A).





Slika 3.3: Matrica korelacija (Skup A).

### 3.7 Tretman outliera i anomalija

Sljedeći isječak implementira IQR winsorization (klipanje ekstremnih vrijednosti) i izgradnju `outlier_flag` značajke Isolation Forest metodom [8]. Ovakav signal može pomoći modelu u prepoznavanju atipičnih obrazaca.

```

1 import numpy as np
2 from sklearn.ensemble import IsolationForest
3
4 def iqr_clip_fit(X, num_cols, k=1.5):
5     q1 = X[num_cols].quantile(0.25)
6     q3 = X[num_cols].quantile(0.75)
7     iqr = q3 - q1
8     lo = q1 - k * iqr
9     hi = q3 + k * iqr
10    return lo, hi
11
12 def iqr_clip_transform(X, num_cols, lo, hi):
13     X2 = X.copy()
14     X2[num_cols] = X2[num_cols].clip(lo, hi, axis=1)
15     return X2
16
17 def add_isoforest_flag(X, num_cols):
18     Xn = X[num_cols].replace([np.inf, -np.inf], np.nan)
19     Xn = Xn.fillna(Xn.median())

```

```

20     iso = IsolationForest(n_estimators=200, contamination="auto",
21                           random_state=42)
    return (iso.fit_predict(Xn.values) == -1).astype(int)

```

Listing 3.6: IQR winsorization + Isolation Forest outlier flag

## 3.8 Predobrada: MICE imputacija, skaliranje i kodiranje

Ovaj preprocesor provodi MICE imputaciju [7] na numeričkim varijablama, standardizaciju, a za kategorijske varijable (ako postoje) imputaciju najčešćom vrijednošću i one-hot kodiranje. Implementacija je robusna na slučaj kada u skupu nema kategorijskih varijabli.

```

1  import numpy as np
2  from sklearn.pipeline import Pipeline
3  from sklearn.compose import ColumnTransformer
4  from sklearn.preprocessing import OneHotEncoder, StandardScaler
5  from sklearn.impute import SimpleImputer
6  from sklearn.experimental import enable_iterative_imputer
7  from sklearn.impute import IterativeImputer
8
9  def make_preprocessor(num_cols, cat_cols, seed=42):
10     num_pipe = Pipeline([
11         ("imp", IterativeImputer(random_state=seed, max_iter=20,
12                                 sample_posterior=False)),
13         ("sc", StandardScaler())
14     ])
15     if len(cat_cols) == 0:
16         return ColumnTransformer([("num", num_pipe, num_cols)])
17     cat_pipe = Pipeline([
18         ("imp", SimpleImputer(strategy="most_frequent")),
19         ("oh", OneHotEncoder(handle_unknown="ignore"))
20     ])
21     return ColumnTransformer([
22         ("num", num_pipe, num_cols),
23         ("cat", cat_pipe, cat_cols)
24     ])
25
26 def to_feature_names(pre):
27     num_cols = list(pre.transformers_[0][2])
28     cat_cols = list(pre.transformers_[1][2]) if len(pre.
29     transformers_) > 1 else []
30     if len(cat_cols) == 0:
31         return num_cols
32     cat_pipe = pre.transformers_[1][1]
33     oh = cat_pipe.named_steps["oh"]
34     return num_cols + oh.get_feature_names_out(cat_cols).tolist()

```

Listing 3.7: Robustan ColumnTransformer preprocesor i dohvat naziva značajki

### 3.9 Integracija heterogenih skupova

Integracija se izvodi tako da se oba skupa zasebno preprocesiraju, a zatim se njihove matrice značajki spajaju u zajednički prostor značajki pomoću *outer* konkatencije uz popunjavanje nedostajućih značajki nulama. Dodatno se čuva informacija o izvoru uzorka kroz *dataset* atribut.

```
1 import pandas as pd
2 import numpy as np
3
4 def split_types(df):
5     num = df.select_dtypes(include=[np.number]).columns.tolist()
6     cat = [c for c in df.columns if c not in num]
7     return num, cat
8
9 Xa = df_a_raw.drop(columns=["Bankrupt?"])
10 ya = df_a_raw["Bankrupt?"].astype(int)
11
12 Xb = df_b_raw.drop(columns=["default"])
13 yb = df_b_raw["default"].astype(int)
14
15 num_a, cat_a = split_types(Xa)
16 num_b, cat_b = split_types(Xb)
17
18 prep_a = make_preprocessor(num_a, cat_a, seed=42)
19 prep_b = make_preprocessor(num_b, cat_b, seed=42)
20
21 Za = prep_a.fit_transform(Xa)
22 Zb = prep_b.fit_transform(Xb)
23
24 fa = to_feature_names(prepare_a)
25 fb = to_feature_names(prepare_b)
26
27 Xa_enc = pd.DataFrame(Za.toarray() if hasattr(Za, "toarray") else
28                       Za, columns=fa)
29
30 Xb_enc = pd.DataFrame(Zb.toarray() if hasattr(Zb, "toarray") else
31                       Zb, columns=fb)
32
33 Xa_enc["dataset"] = "A"
34 Xb_enc["dataset"] = "B"
35
36 X_all = pd.concat([Xa_enc, Xb_enc], axis=0, join="outer",
37                   ignore_index=True).fillna(0.0)
38 y_all = pd.concat([ya.reset_index(drop=True), yb.reset_index(drop=True)],
39                   ignore_index=True)
40 X_all.shape, y_all.mean()
```

Listing 3.8: Integracija A i B skupa u zajednički prostor značajki

## 3.10 Selekcija značajki i baseline modeli

Zbog potencijalno visoke dimenzionalnosti uvodi se selekcija značajki. RFE postupno uklanja najmanje važne značajke prema logističkoj regresiji te ostavlja podskup za treniranje modela. Nakon toga se treniraju baseline modeli (LR i RF) radi usporedbe.

```
1 from sklearn.feature_selection import RFE
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import roc_auc_score,
   average_precision_score
5
6 X = X_all.drop(columns=["dataset"]).astype(float).values
7 y = y_all.values
8
9 lr_base = LogisticRegression(max_iter=1000, solver="liblinear",
   class_weight="balanced", random_state=42)
10 rfe = RFE(estimator=lr_base, n_features_to_select=min(60, X.shape
   [1]//2), step=0.1)
11 rfe.fit(X, y)
12
13 mask = rfe.support_
14 X_sel = X[:, mask]
15
16 lr = LogisticRegression(max_iter=1000, solver="liblinear",
   class_weight="balanced", random_state=42)
17 rf = RandomForestClassifier(n_estimators=600, class_weight="
   balanced_subsample", random_state=42)
18
19 lr.fit(X_sel, y)
20 p_lr = lr.predict_proba(X_sel)[:,1]
21
22 rf.fit(X_sel, y)
23 p_rf = rf.predict_proba(X_sel)[:,1]
24
25 roc_auc_score(y, p_lr), average_precision_score(y, p_lr),
   roc_auc_score(y, p_rf)
```

Listing 3.9: RFE selekcija + baseline modeli (LR i RF)

## 3.11 Optuna + XGBoost optimizacija

Sljedeći isječak definira objective funkciju za Optuna: za svaki trial predlaže hiperparametre XGBoost modela, evaluira ih stratificiranom CV-om i vraća prosječni PR-AUC.

```
1 import optuna
2 from xgboost import XGBClassifier
3 from sklearn.model_selection import StratifiedKFold
4 from sklearn.metrics import average_precision_score
5 import numpy as np
6
```

```

7 def objective(trial, X, y, seed=42, n_splits=5):
8     params = {
9         "n_estimators": trial.suggest_int("n_estimators", 200,
10            1200),
11         "max_depth": trial.suggest_int("max_depth", 2, 8),
12         "learning_rate": trial.suggest_float("learning_rate",
13            0.01, 0.25, log=True),
14         "subsample": trial.suggest_float("subsample", 0.6, 1.0),
15         "colsample_bytree": trial.suggest_float("colsample_bytree",
16            0.6, 1.0),
17         "min_child_weight": trial.suggest_float("min_child_weight",
18            1e-2, 20.0, log=True),
19         "reg_lambda": trial.suggest_float("reg_lambda", 1e-3,
20            50.0, log=True),
21         "reg_alpha": trial.suggest_float("reg_alpha", 1e-4, 10.0,
22            log=True),
23         "gamma": trial.suggest_float("gamma", 0.0, 5.0),
24         "objective": "binary:logistic",
25         "eval_metric": "aucpr",
26         "random_state": seed,
27         "n_jobs": -1
28     }
29     posw = float((y==0).sum() / max(1, (y==1).sum()))
30     params["scale_pos_weight"] = posw
31
32     cv = StratifiedKFold(n_splits=n_splits, shuffle=True,
33        random_state=seed)
34     scores = []
35     for tr_idx, va_idx in cv.split(X, y):
36         m = XGBClassifier(**params)
37         m.fit(X[tr_idx], y[tr_idx])
38         p = m.predict_proba(X[va_idx])[:,1]
39         scores.append(average_precision_score(y[va_idx], p))
40     return float(np.mean(scores))
41
42 study = optuna.create_study(direction="maximize")
43 study.optimize(lambda t: objective(t, X_sel, y, seed=42, n_splits
44    =5), n_trials=40, show_progress_bar=False)
45
46 study.best_value, study.best_params

```

Listing 3.10: Optuna objective funkcija za XGBoost (cilj: PR-AUC)

## 3.12 Rigorozna validacija: nested CV

Nested CV implementacija osigurava da se hiperparametri odabiru isključivo unutar train dijela vanjskog folda, a evaluacija se radi na neviđenim vanjskim test foldovima.

```

1 from sklearn.model_selection import StratifiedKFold
2 from sklearn.metrics import roc_auc_score,
   average_precision_score

```

```

3 from xgboost import XGBClassifier
4 import numpy as np
5 import optuna
6
7 outer = StratifiedKFold(n_splits=5, shuffle=True, random_state
8     =42)
9 outer_rows = []
10 for fold, (tr_idx, te_idx) in enumerate(outer.split(X_sel, y),
11     start=1):
12     Xo_tr, Xo_te = X_sel[tr_idx], X_sel[te_idx]
13     yo_tr, yo_te = y[tr_idx], y[te_idx]
14
15     st = optuna.create_study(direction="maximize")
16     st.optimize(lambda t: objective(t, Xo_tr, yo_tr, seed=42,
17         n_splits=3), n_trials=15, show_progress_bar=False)
18
19     bp = dict(st.best_params)
20     bp.update({
21         "objective": "binary:logistic",
22         "eval_metric": "aucpr",
23         "random_state": 42,
24         "n_jobs": -1,
25         "scale_pos_weight": float((yo_tr==0).sum() / max(1,(yo_tr
26             ==1).sum()))
27     })
28
29     m = XGBClassifier(**bp)
30     m.fit(Xo_tr, yo_tr)
31     p = m.predict_proba(Xo_te)[: ,1]
32
33     outer_rows.append({
34         "fold": fold,
35         "AUC": float(roc_auc_score(yo_te, p)),
36         "PR_AUC": float(average_precision_score(yo_te, p))
37     })
38
39 outer_rows, float(np.mean([r["AUC"] for r in outer_rows])), float
40     (np.mean([r["PR_AUC"] for r in outer_rows]))

```

Listing 3.11: Nested CV (outer + inner Optuna)

### 3.13 Odabir praga odluke

Ovaj dio pokazuje kako se prag  $\tau$  može optimizirati za F1 ili prema jednostavnoj profit funkciji (cost-sensitive pristup). U stvarnoj primjeni težine se određuju prema poslovnim troškovima FP i FN.

```

1 import numpy as np
2 from sklearn.metrics import precision_recall_curve, f1_score
3

```

```

4 p = m.predict_proba(X_sel)[: ,1]
5 pr, rc, th = precision_recall_curve(y, p)
6
7 f1s = 2*pr*rc/(pr+rc+1e-12)
8 best_i = int(np.nanargmax(f1s[: -1])) if len(th) else 0
9 thr_f1 = float(th[best_i]) if len(th) else 0.5
10 yhat = (p >= thr_f1).astype(int)
11
12 thr_f1, f1_score(y, yhat)
13
14 C_TP, C_FP, C_FN, C_TN = 1.0, -0.2, -1.0, 0.0
15 def profit(y_true, y_pred):
16     y_true = np.asarray(y_true); y_pred = np.asarray(y_pred)
17     tp = ((y_true==1)&(y_pred==1)).sum()
18     fp = ((y_true==0)&(y_pred==1)).sum()
19     fn = ((y_true==1)&(y_pred==0)).sum()
20     tn = ((y_true==0)&(y_pred==0)).sum()
21     return C_TP*tp + C_FP*fp + C_FN*fn + C_TN*tn
22
23 profits = [(t, profit(y, (p>=t).astype(int))) for t in np.
24             linspace(0.01,0.99,99)]
25 best_t_profit = max(profits, key=lambda x: x[1])[0]
26 best_t_profit

```

Listing 3.12: Threshold tuning: F1 i profit funkcija

## 3.14 Kalibracija vjerojatnosti

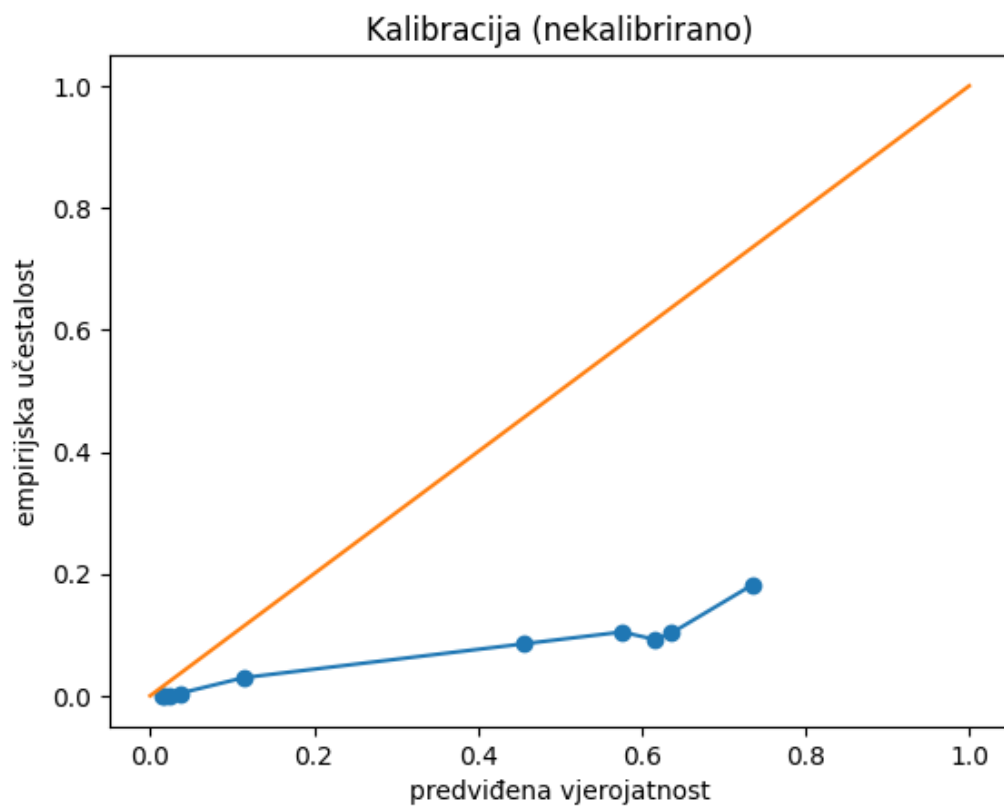
Sljedeći isječak provodi izotoničku kalibraciju i uspoređuje Brier score prije i poslije. U dokumentaciji se preporučuje uključiti i kalibracijske krivulje.

```

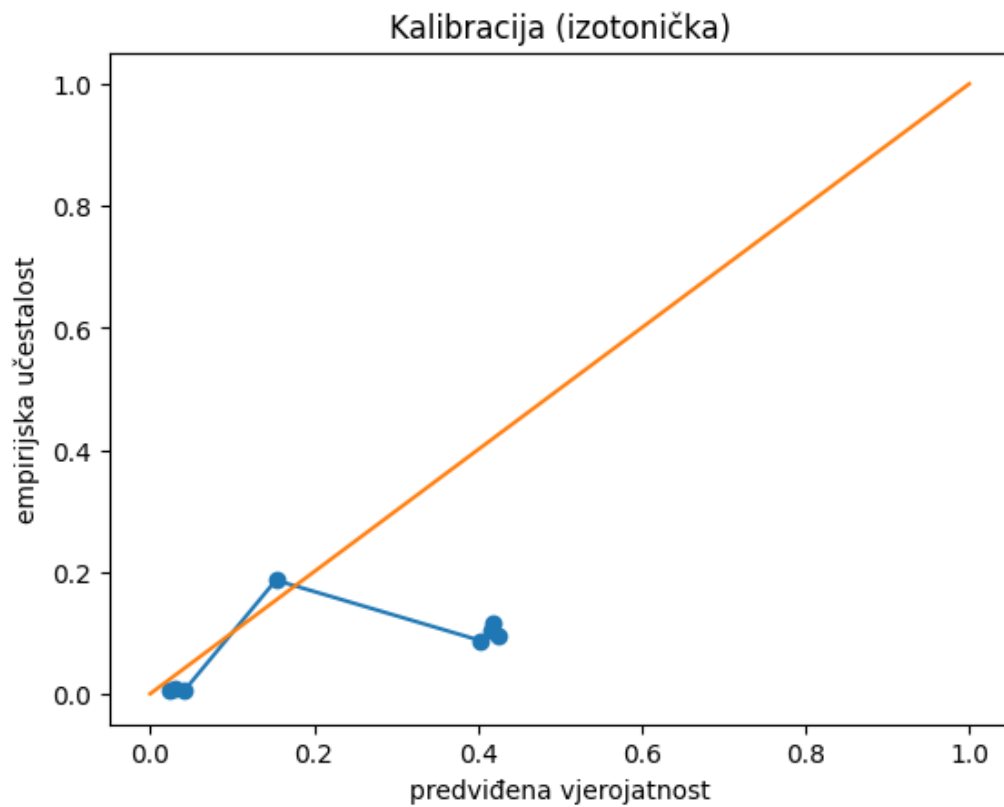
1 from sklearn.calibration import CalibratedClassifierCV
2 from sklearn.metrics import brier_score_loss
3
4 cal = CalibratedClassifierCV(m, method="isotonic", cv=3)
5 cal.fit(X_sel, y)
6
7 p_uncal = m.predict_proba(X_sel)[: ,1]
8 p_cal = cal.predict_proba(X_sel)[: ,1]
9
10 brier_score_loss(y, p_uncal), brier_score_loss(y, p_cal)

```

Listing 3.13: Izotonička kalibracija i Brier score



Slika 3.4: Kalibracijska krivulja prije kalibracije (predložak).



Slika 3.5: Kalibracijska krivulja nakon izotoničke kalibracije (predložak).

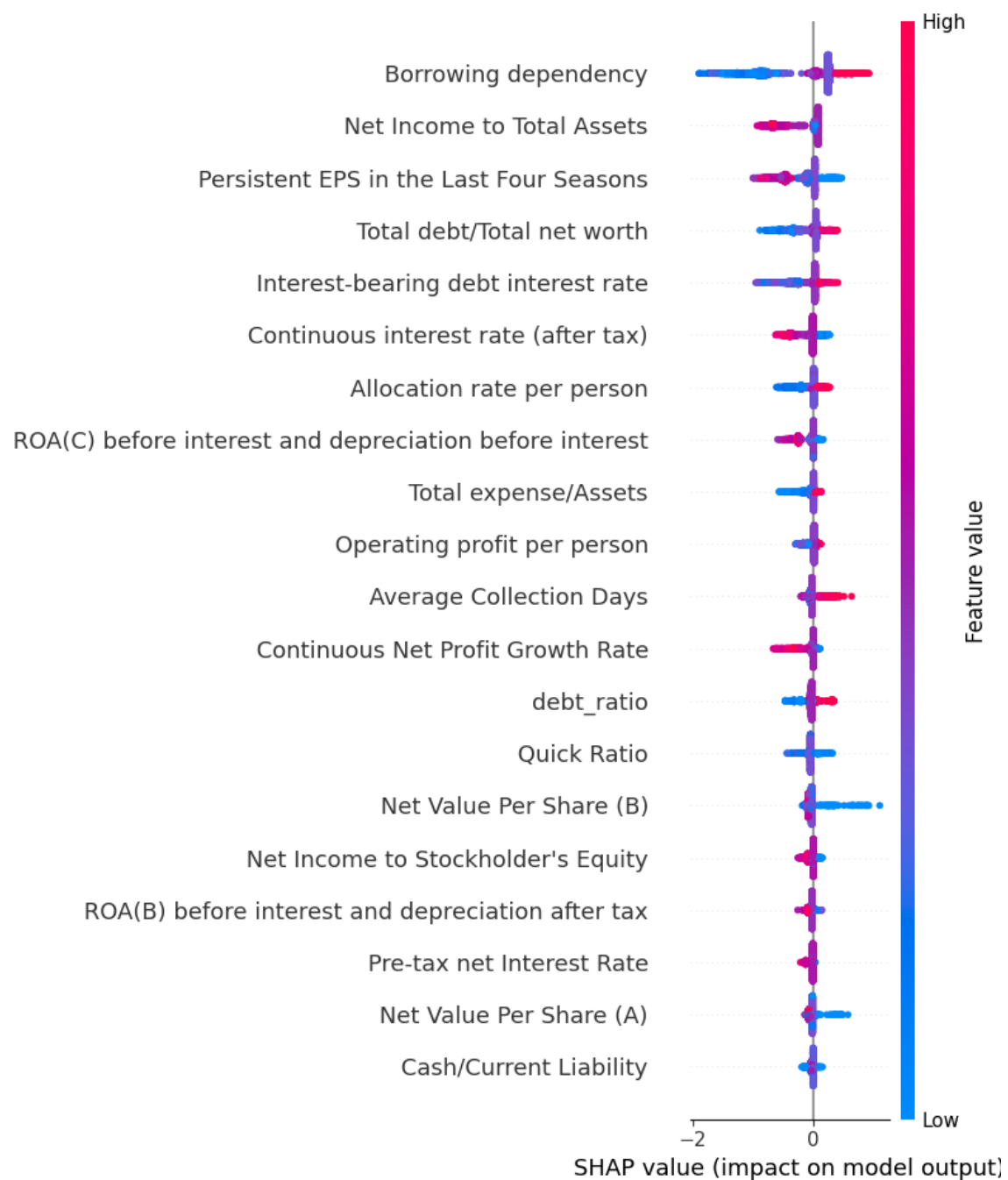


## 3.15 SHAP interpretabilnost

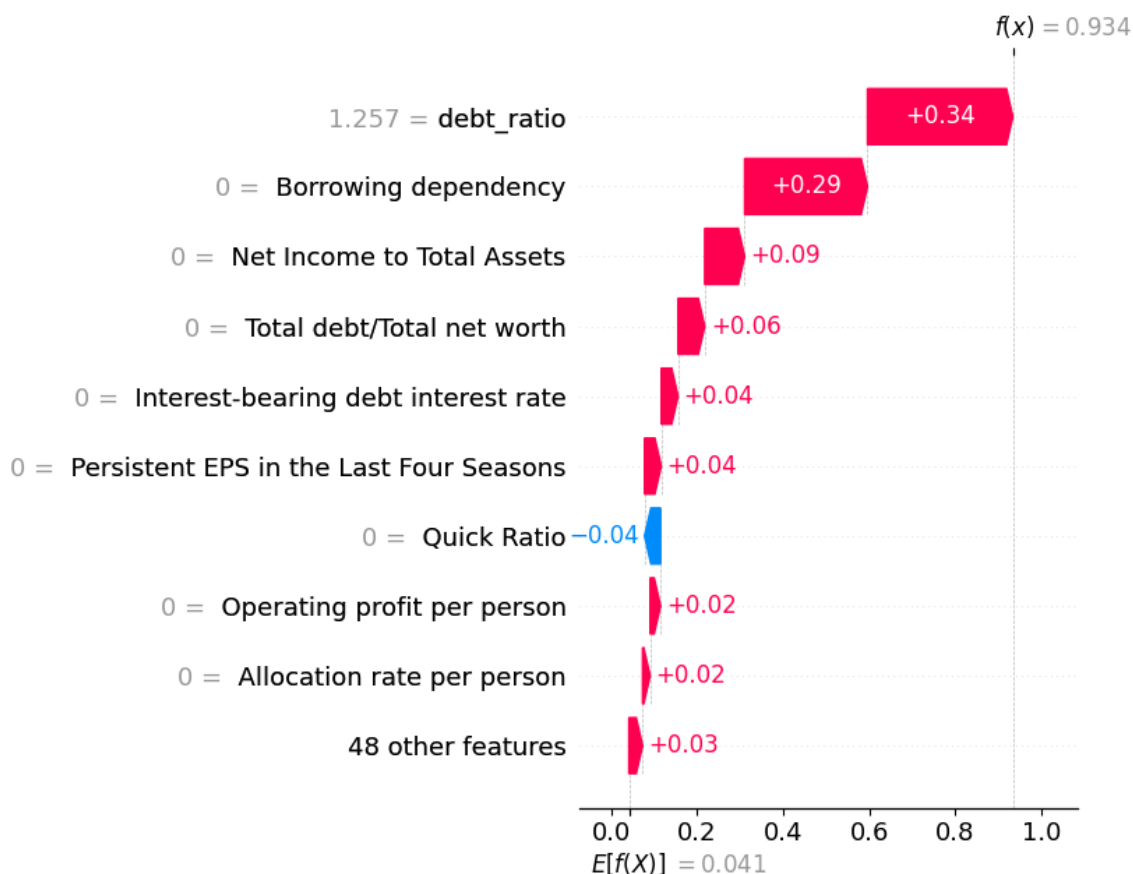
Ovaj dio računa SHAP vrijednosti i sprema grafove. Preporuka je koristiti uzorak podataka radi vremenske učinkovitosti, te navesti barem jedan globalni i jedan lokalni graf.

```
1 import shap
2 import matplotlib.pyplot as plt
3
4 expl = shap.TreeExplainer(m)
5 X_sample = X_sel[:2000]
6 sv = expl.shap_values(X_sample)
7
8 shap.summary_plot(sv, X_sample, show=False)
9 plt.savefig("artifacts/fig/shap_summary.png", dpi=200,
10             bbox_inches="tight")
11 plt.close()
12
13 idx = 0
14 shap.plots._waterfall.waterfall_legacy(expl.expected_value, sv[
15     idx], features=X_sample[idx], show=False)
16 plt.savefig("artifacts/fig/shap_waterfall.png", dpi=200,
17             bbox_inches="tight")
18 plt.close()
```

Listing 3.14: SHAP TreeExplainer: globalno i lokalno objašnjenje



Slika 3.6: SHAP summary plot: globalna važnost značajki (predložak).



Slika 3.7: SHAP waterfall: lokalno objašnjenje za jedan uzorak (predložak).

### 3.16 K-means segmentacija

U nastavku se bira broj klastera prema silhouette score-u i vizualizira segmentacija u PCA prostoru. Nakon toga se tipično izvodi analiza klastera: stopa defaulta po klasteru i prosjeci ključnih značajki.

```

1 from sklearn.cluster import KMeans
2 from sklearn.metrics import silhouette_score
3 from sklearn.decomposition import PCA
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 k_vals = range(2, 9)
8 sil = []
9 for k in k_vals:
10     km = KMeans(n_clusters=k, random_state=42, n_init=10)
11     lab = km.fit_predict(X_sel)
12     sil.append(silhouette_score(X_sel, lab))
13
14 k_best = int(list(k_vals)[int(np.argmax(sil))])
15
16 km = KMeans(n_clusters=k_best, random_state=42, n_init=10)
17 lab = km.fit_predict(X_sel)

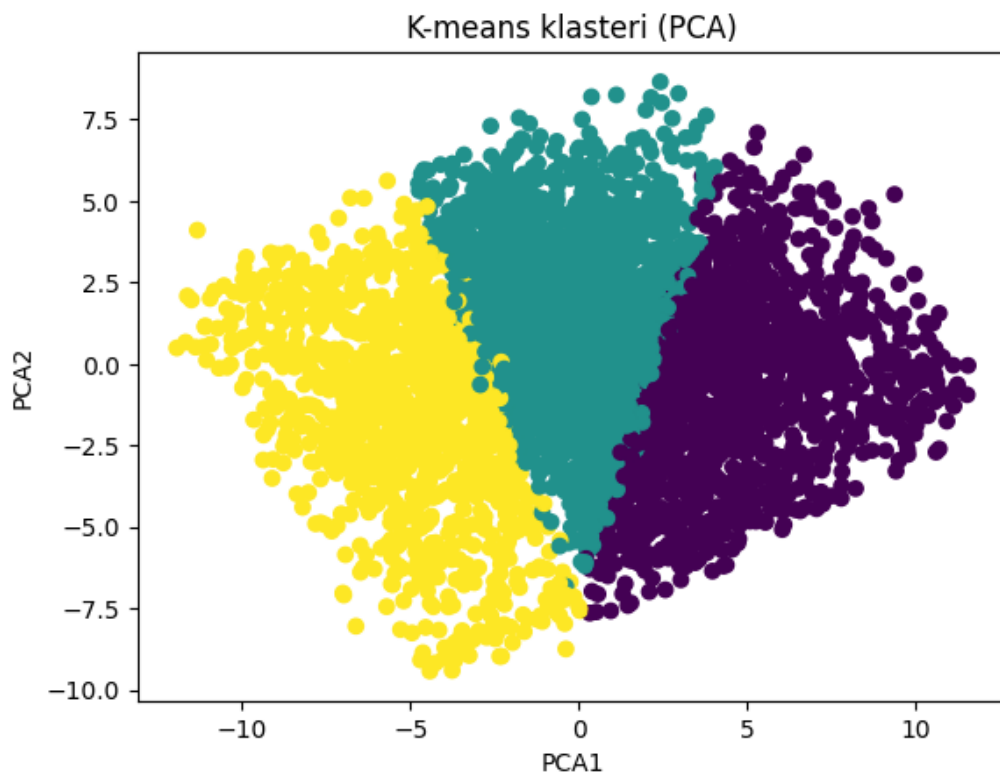
```

```

18
19 pca = PCA(n_components=2, random_state=42)
20 Z = pca.fit_transform(X_sel)
21
22 plt.figure()
23 plt.scatter(Z[:,0], Z[:,1], c=lab)
24 plt.savefig("artifacts/fig/kmeans_pca.png", dpi=200, bbox_inches=
    "tight")
25 plt.close()
26
27 k_best, float(np.max(sil))

```

Listing 3.15: K-means: izbor k, učenje i PCA vizualizacija



Slika 3.8: K-means klasteri u PCA projekciji (predložak).

### 3.17 Pohrana integriranih podataka i metrika u SQLite

Pohrana u SQLite omogućuje ponovnu uporabu rezultata bez retreniranja te jednostavan dohvat kroz API sloj. Sljedeći isječak prema odabrane značajke i ciljne oznake, te primjer tablice metrika.

```

1 import sqlite3
2 import pandas as pd
3

```

```

4 DB = "default_risk.db"
5 con = sqlite3.connect(DB)
6
7 pd.DataFrame(X_sel).assign(y=y).to_sql("dataset_selected", con,
    index=False, if_exists="replace")
8
9 metrics = pd.DataFrame([
10     {"model": "LR", "AUC": None, "PR_AUC": None, "F1": None, "thr"
        : None},
11     {"model": "RF", "AUC": None, "PR_AUC": None, "F1": None, "thr"
        : None},
12     {"model": "XGB", "AUC": None, "PR_AUC": None, "F1": None, "thr"
        : None}
13 ])
14 metrics.to_sql("metrics", con, index=False, if_exists="replace")
15
16 con.close()

```

Listing 3.16: SQLite: spremanje procesiranih podataka i metrika

### 3.17.1 Shema baze

Sljedeći SQL primjer opisuje minimalnu shemu tablica. U praksi se tablice generiraju iz pandas DataFrame objekata, ali eksplicitna shema je korisna za dokumentaciju.

```

1 CREATE TABLE dataset_selected (
2     f0 REAL, f1 REAL, ...,
3     y INTEGER
4 );
5
6 CREATE TABLE metrics (
7     model TEXT,
8     AUC REAL,
9     PR_AUC REAL,
10    F1 REAL,
11    thr REAL
12 );

```

Listing 3.17: Primjer sheme (konceptualno)

## 3.18 REST API sučelje (FastAPI)

API sloj omogućuje standardiziran pristup metrikama i predikcijama. U nastavku je minimalna implementacija: `/health` za provjeru stanja, `/metrics` za dohvat metrika i `/predict` za predikciju vjerojatnosti defaulta.

```

1 from fastapi import FastAPI
2 from pydantic import BaseModel
3 import numpy as np
4 import sqlite3
5 import pandas as pd

```

```

6
7 app = FastAPI()
8
9 class PredictIn(BaseModel):
10     features: dict
11
12 def read_sql(table, db_path="default_risk.db"):
13     con = sqlite3.connect(db_path)
14     try:
15         return pd.read_sql_query(f"SELECT * FROM {table}", con)
16     finally:
17         con.close()
18
19 @app.get("/health")
20 def health():
21     return {"status": "ok"}
22
23 @app.get("/metrics")
24 def metrics():
25     df = read_sql("metrics")
26     return df.to_dict(orient="records")
27
28 @app.post("/predict")
29 def predict(inp: PredictIn):
30     x = np.zeros(X_sel.shape[1], dtype=float)
31     for k, v in inp.features.items():
32         idx = int(k)
33         x[idx] = float(v)
34     p = float(cal.predict_proba(x.reshape(1, -1))[:, 1][0])
35     return {"p_default": p}

```

Listing 3.18: FastAPI: endpointi health, metrics i predict

## 3.19 Pokretanje servisa

Ovaj isječak prikazuje tipični način pokretanja FastAPI aplikacije pomoću Uvicorn poslužitelja.

```

1 uvicorn src.api:app --host 127.0.0.1 --port 8000 --reload

```

Listing 3.19: Pokretanje FastAPI aplikacije

# Poglavlje 4

## Prikaz rada aplikacije

### 4.1 Pipeline: od podataka do modela

Prikaz rada aplikacije organiziran je kroz korake:

1. preuzimanje podataka i validacija,
2. EDA i statistička analiza,
3. predobrada i integracija,
4. treniranje baseline modela,
5. Optuna optimizacija XGBoost-a,
6. nested CV evaluacija,
7. threshold tuning i kalibracija,
8. SHAP i K-means analize,
9. pohrana u SQLite,
10. izlaganje kroz API i testiranje endpointa.

### 4.2 Testiranje API-ja

Sljedeći primjeri pokazuju provjeru stanja servisa i dohvat metrika.

```
1 curl http://127.0.0.1:8000/health
2 curl http://127.0.0.1:8000/metrics
```

Listing 4.1: Testiranje /health i /metrics endpointa

## 4.3 Primjer predikcije

U praksi se šalje JSON koji sadrži vrijednosti značajki. Ovaj primjer pretpostavlja da se u `features` šalju indeksi značajki (jednostavniji oblik), dok je produkcijski pristup tipično mapiranje po nazivima.

```
1 curl -X POST "http://127.0.0.1:8000/predict" \  
2   -H "Content-Type: application/json" \  
3   -d "{\"features\": {\"0\": 0.12, \"1\": -1.03, \"2\": 0.44}}"
```

Listing 4.2: Primjer poziva `/predict` endpointa

Slika 4.1: Primjer prikaza rada API-ja (predložak; korisnik umeće sliku).



# Poglavlje 5

## Rezultati i vizualizacije

### 5.1 Tablice metrika

U tablicu se unose rezultati iz izvršenog notebooka. Ako se metrike razlikuju od navedenih, potrebno je ažurirati vrijednosti u tablici.

Tablica 5.1: Usporedba modela (primjer popunjavanja).

Model	AUC	PR-AUC	F1	Prag
Logistička regresija	0.7909	0.2164	0.2566	0.7294
Random Forest	0.6256	0.1189	0.2564	0.1417
XGBoost (Optuna)	0.7921	0.2464	0.2673	0.7202

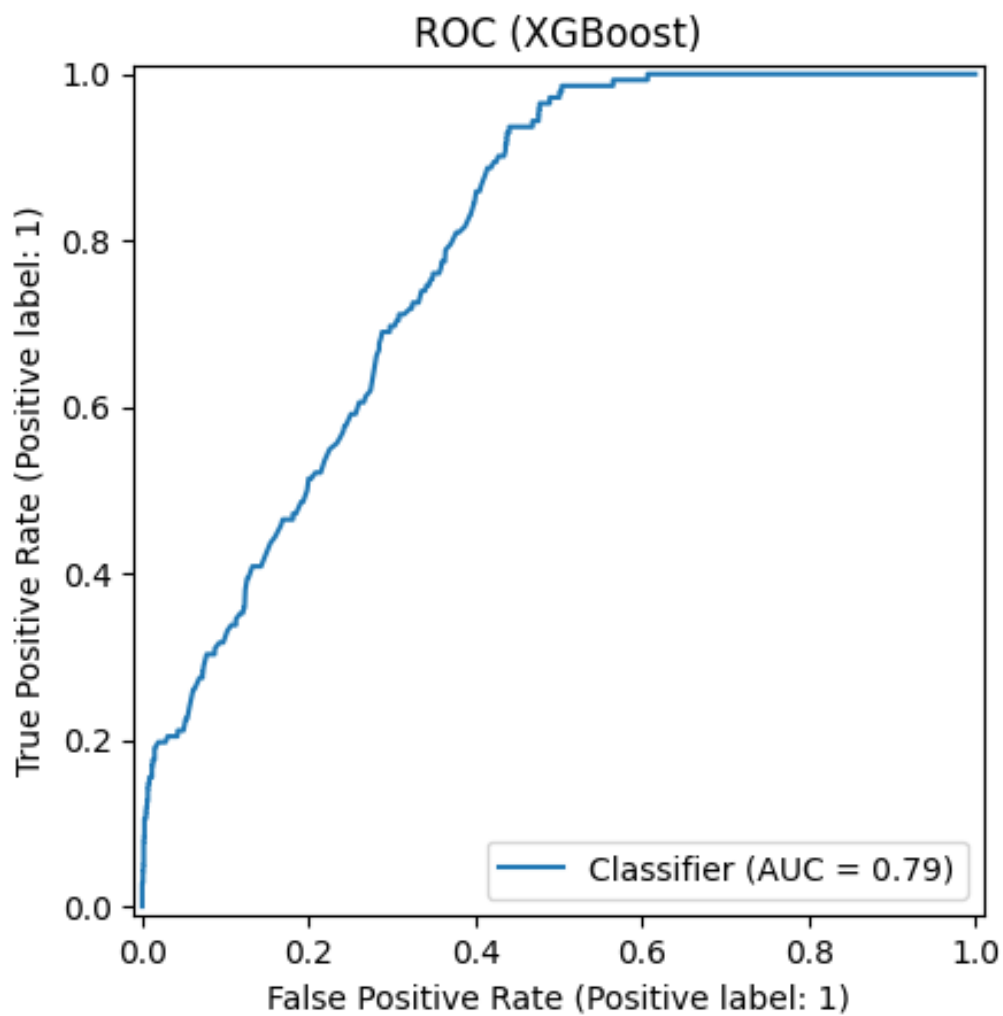
### 5.2 Nested CV rezultati

Tablica 5.2: Nested CV rezultati (primjer).

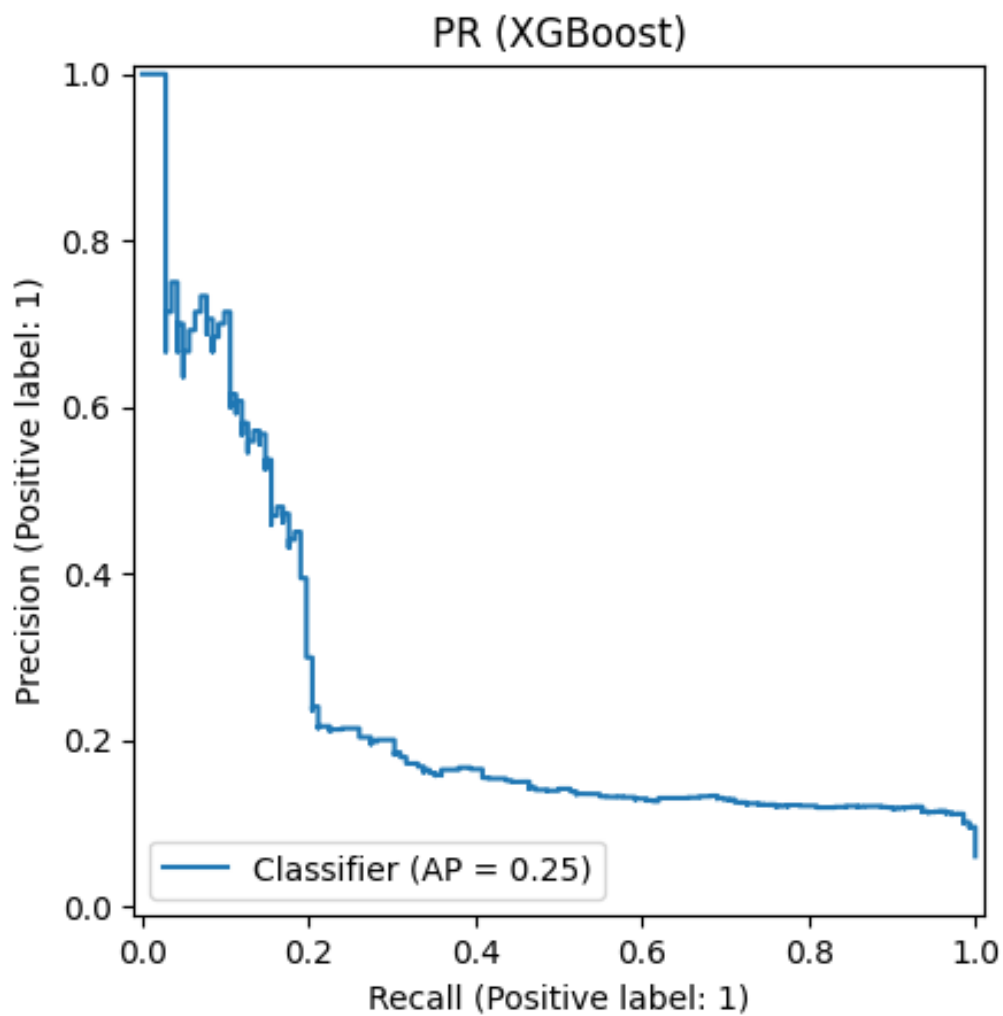
Fold	AUC	PR-AUC
1	0.8171	0.2520
2	0.8097	0.2014
3	0.7742	0.1702
4	0.8159	0.2890
5	0.8059	0.2671
Prosjek	0.8046	0.2359

### 5.3 Ključni grafovi

U nastavku su predlošci za umetanje grafova iz notebooka. Korisnik treba spremiti slike u mapu `fig/`.



Slika 5.1: ROC krivulja za XGBoost model (predložak).



Slika 5.2: Precision–Recall krivulja za XGBoost model (predložak).

# Poglavlje 6

## Kritički osvrt

### 6.1 Praktična izvedivost

Razvijeni sustav pokazuje visoku razinu praktične izvedivosti zahvaljujući modularnoj arhitekturi i oslanjanju na široko prihvaćene, stabilne Python biblioteke kao što su scikit-learn, xgboost, optuna, shap i fastapi. Takav tehnološki odabir omogućuje jednostavno održavanje, nadogradnju pojedinih komponenti i razumijevanje sustava od strane drugih inženjera ili analitičara.

Posebna prednost sustava jest korištenje SQLite baze podataka za pohranu integriranih podataka, metrika i rezultata modeliranja. Time se značajno smanjuje infrastrukturna složenost u usporedbi s distribucijskim bazama podataka, što je osobito pogodno za prototipe, istraživačke projekte i manje produkcijske sustave. Unatoč svojoj jednostavnosti, SQLite omogućuje pouzdanu i konzistentnu pohranu rezultata te njihovu ponovnu uporabu bez potrebe za ponovnim treniranjem modela.

Izlaganje funkcionalnosti sustava putem REST API-ja dodatno povećava njegovu praktičnu vrijednost. API sloj omogućuje jednostavnu integraciju modela procjene default rizika u postojeće informacijske sustave, web aplikacije ili automatizirane procese odlučivanja. Takav pristup odvaja modeliranje od korisničkog sučelja i poslovne logike, što je u skladu s modernim principima dizajna softverskih sustava.

Konačno, naglasak na reproducibilnosti kroz jasno definirani pipeline, verzioniranje podataka i strogu validaciju (nested cross-validation) čini sustav pogodnim ne samo za demonstraciju koncepta, već i kao temelj za daljnji razvoj prema produkcijskom okruženju. Time se potvrđuje da predloženo rješenje nije isključivo akademska vježba, već realno primjenjiv okvir za procjenu kreditnog rizika.

### 6.2 Primjena i ograničenja

Prednosti:

- XGBoost uz Bayesovu optimizaciju često daje dobar kompromis između performansi i stabilnosti [3, 4].
- Kalibracija poboljšava uporabljivost vjerojatnosti za poslovna pravila [6].
- SHAP omogućuje auditabilnost i objašnjenja prema regulatornim i poslovnim potrebama [5].

- Segmentacija K-means može otkriti profile rizika i podržati strategije upravljanja portfeljem [9].

Ograničenja:

- Integracija heterogenih skupova kroz union prostora značajki može rezultirati sparsnom i visokodimenzionalnom matricom, što povećava potrebu za selekcijom i regularizacijom.
- Nested CV i Optuna povećavaju računalni trošak te zahtijevaju balans između rigoroznosti i vremena.
- K-means pretpostavlja sferne klastere i osjetljiv je na skaliranje; interpretacija klastera zahtijeva dodatnu analizu.
- SHAP može biti računalno zahtjevan na velikim uzorcima pa se preporučuje uzorkovanje.

## 6.3 Moguća poboljšanja

- Uvođenje monitoringa (data drift, concept drift) i periodičnog retreniranja.
- Uvođenje cost-sensitive optimizacije praga prema realnim troškovima i prihodima.
- Povećanje transparentnosti kroz modelne kartice i automatsku generaciju izvještaja.
- Razmatranje alternative segmentacije (Gaussian Mixture Models, HDBSCAN) i usporedba.

# Poglavlje 7

## Zaključak

U ovom radu razvijen je cjelovit i reproducibilan pipeline za procjenu financijskog rizika neispunjenja obveza koji obuhvaća sve ključne faze modernog sustava strojnog učenja: od heterogene akvizicije i integracije podataka, preko napredne predobrade i usporedbe baseline modela, do Bayesovski optimiziranog XGBoost modela uz rigoroznu nested cross-validation evaluaciju. Dodatne komponente, poput optimizacije praga odluke i kalibracije vjerojatnosti, osiguravaju da su izlazi modela praktično uporabivi u kontekstu stvarnih poslovnih odluka.

Poseban doprinos rada očituje se u integraciji metoda interpretabilnosti i analize podataka. Primjena SHAP metode omogućuje transparentno i formalno objašnjenje odluka složenog ansambl modela, čime se postiže auditabilnost nužna za regulatorne i poslovne zahtjeve. Segmentacija uzoraka K-means metodom dodatno proširuje analitičku vrijednost sustava, omogućujući identifikaciju različitih profila rizika unutar populacije.

Rezultati rada pokazuju da kombinacija snažnih modela strojnog učenja, rigorozne validacije i metoda objašnjivosti može rezultirati sustavom koji je istovremeno prediktivno učinkovit, transparentan i reproducibilan. Takav sustav predstavlja čvrst temelj za daljnji razvoj prema produkcijskim rješenjima u području kreditnog rizika te potvrđuje da suvremeni pristupi strojnog učenja mogu zadovoljiti i tehničke i regulatorne zahtjeve financijske industrije.

# Bibliografija

- [1] Kaggle, “Company Bankruptcy Prediction,” dataset, pristupljeno: 23. siječnja 2026..
- [2] Kaggle, “Credit Risk Modelling Dataset,” dataset, pristupljeno: 23. siječnja 2026..
- [3] T. Chen and C. Guestrin, “XGBoost: A Scalable Tree Boosting System,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2016, pp. 785–794.
- [4] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A Next-generation Hyperparameter Optimization Framework,” in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.
- [5] S. M. Lundberg and S.-I. Lee, “A Unified Approach to Interpreting Model Predictions,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [6] A. Niculescu-Mizil and R. Caruana, “Predicting Good Probabilities with Supervised Learning,” in *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, 2005, pp. 625–632.
- [7] S. van Buuren and K. Groothuis-Oudshoorn, “mice: Multivariate Imputation by Chained Equations in R,” *Journal of Statistical Software*, vol. 45, no. 3, pp. 1–67, 2011.
- [8] F. T. Liu, K. M. Ting, and Z.-H. Zhou, “Isolation Forest,” in *2008 Eighth IEEE International Conference on Data Mining*, 2008, pp. 413–422.
- [9] J. MacQueen, “Some Methods for Classification and Analysis of Multivariate Observations,” in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1967, pp. 281–297.
- [10] DVC, “Data Version Control (DVC) Documentation,” službena dokumentacija, pristupljeno: 23. siječnja 2026..

# Dodatak A

## Dodatci

### A.1 Popis slika koje je potrebno umetnuti

U nastavku je popis predloženih slika koje ovaj dokument referencira. Potrebno ih je spremiti u mapu `fig/`:

- `architecture.png` – arhitektura sustava.
- `eda_hist_a.png` – histogrami skupa A.
- `eda_corr_a.png` – korelacije skupa A.
- `roc_xgb.png` – ROC krivulja.
- `pr_xgb.png` – PR krivulja.
- `calibration_uncalibrated.png` – kalibracija prije.
- `calibration_isotonic.png` – kalibracija poslije.
- `shap_summary.png` – SHAP summary.
- `shap_waterfall.png` – SHAP waterfall.
- `silhouette.png` – silhouette graf.
- `kmeans_pca.png` – PCA vizualizacija klastera.
- `api_demo.png` – prikaz rada API-ja (screenshot).