

# Dubinska analiza i otkrivanje znanja iz podataka

Korištenje Python paketa za statističku analizu i vizualizaciju podataka i analiza algoritma strojnog učenja i dubinske analize podataka

dr.sc. Damir Korenčić

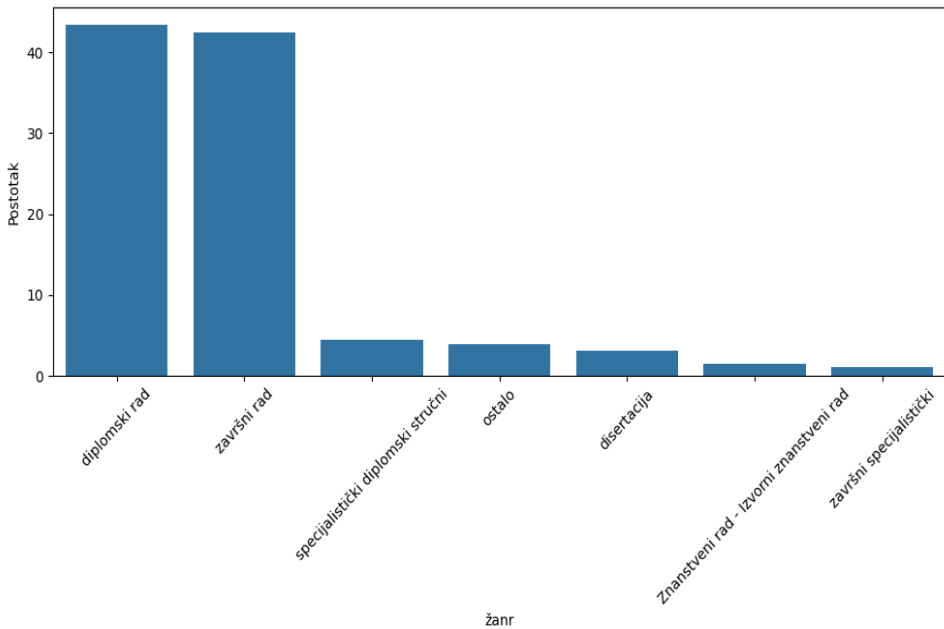
Institut Ruđer Bošković

Laboratorij za strojno učenje i reprezentacije znanja



10. lipnja 2026.

- Centralni repozitorij akademskih radova
  - Diplomski radovi
  - Doktorati
  - ...
- Dabar skup podataka
  - Sakupljanje podataka (crawling)
  - Uređivanje
  - **Čišćenje**
- 265.790 radova (kraj 2025.)
- <https://dabar.srce.hr>



- Matplotlib
  - Potrebna pozornost na detalje (mnogo koda)
  - Ne radi s Pandas DataFrame klasom
  - Zadane vrijednosti (defaults)
- Seaborn
  - API “iznad” Matplotlib-a
  - Jednostavne funkcije za standardne grafike
  - Integriran s DataFrame klasom
  - Razumne zadane vrijednosti (defaults)

# Pandas DataFrame i Series

```
1 def load_dframe_from_pkl(file_path: Path) -> DataFrame:
2     with open(file_path, 'rb') as f:
3         data = pickle.load(f)
4         assert isinstance(data, DataFrame)
5         return data
6     ...
7 df = load_dframe_from_pkl('dabar.pkl') # DataFrame - tablica
8 df['genre_hrv'] # Series - indeksirani niz, kolumna u tablici
9 df[df['genre_hrv']=='diplomski rad'] # DataFrame - pod-tablica
```

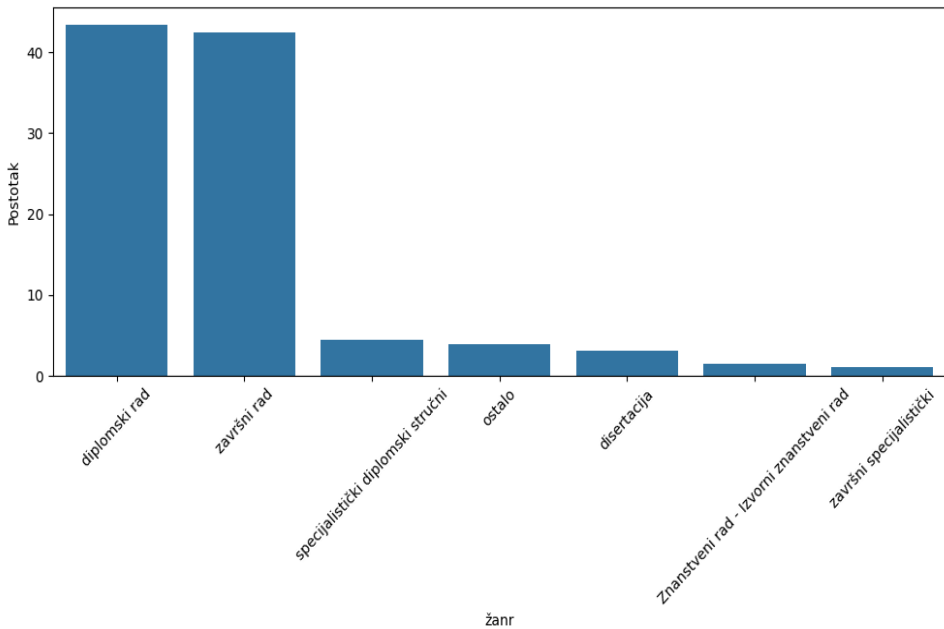
```
df['genre_hrv']
[(0, 'diplomski rad'), (1, 'završni rad'), ...]
```

```
df[df['genre_hrv']=='diplomski rad']
[(0, True), (1, False), ...]
```

# Pandas DataFrame i Series

```
1 df['genre_hrv'].value_counts() # Series!
```

```
diplomski rad, 115296  
završni rad, 112667  
specijalistički diplomski stručni, 11882
```



# Matplotlib i Seaborn

```
1 df = load_dframe_from_pkl(DABAR_RAW_DATASET)
2 genres = df['genre_hrv']
```

## Matplotlib:

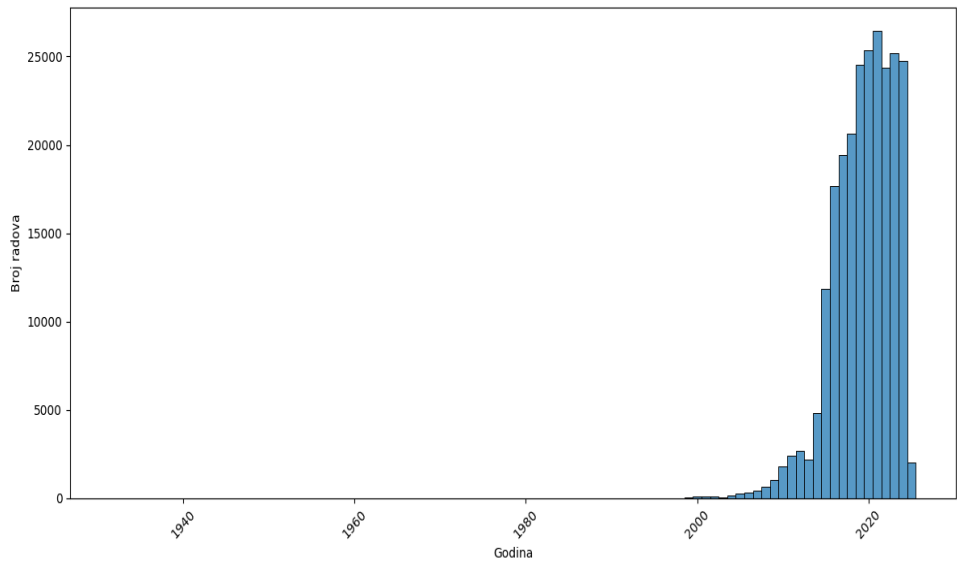
```
1 counts = genres.value_counts() # Series
2 heights = counts / counts.sum() * 100 # Series
3 x = range(len(counts))
4 plt.bar(x, heights.values)
```

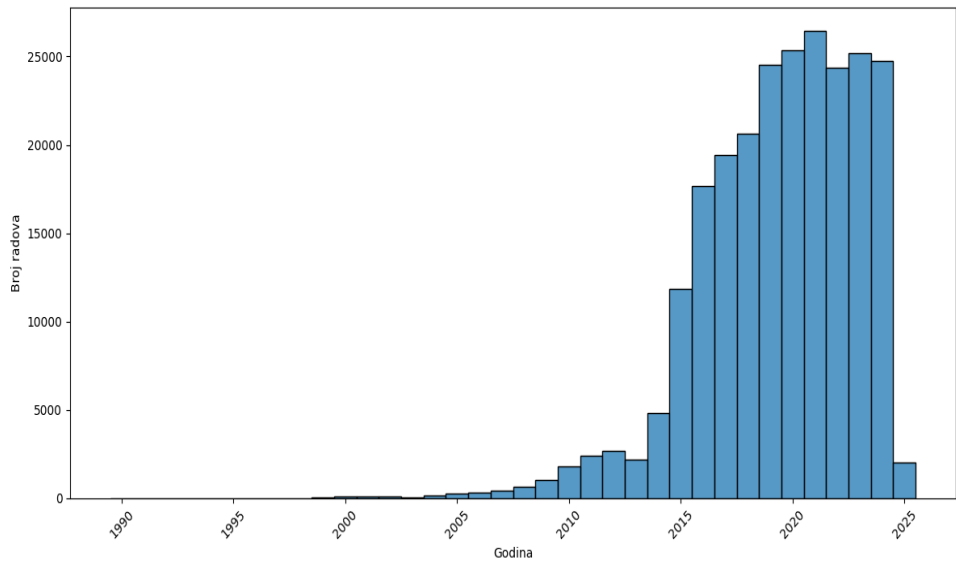
## Seaborn:

```
1 sns.countplot(x=genres, order=genres.value_counts().index,
↪ stat='percent')
```

## Seaborn+Matplotlib

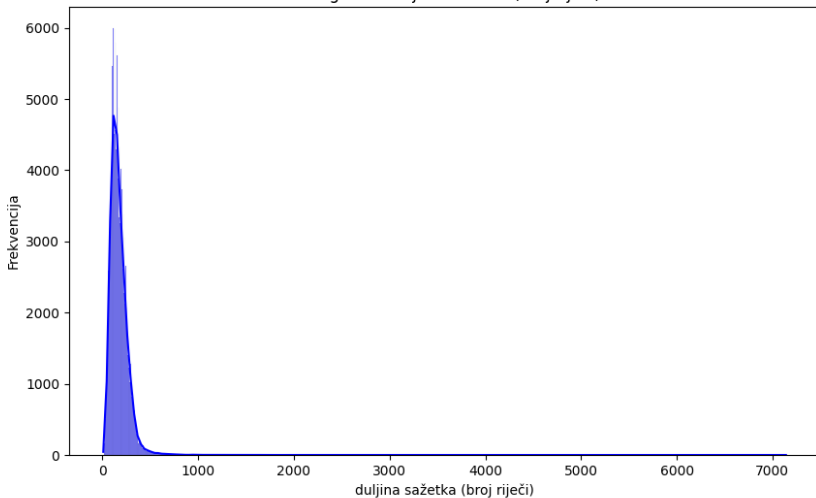
```
1  from matplotlib import pyplot as plt
2  plt.figure() # init. display
3  ...
4  sns.countplot(x=values, order=values.value_counts().index,
   ↪  stat='percent')
5  ...
6  plt.ylabel('postotak') # customize the plot
7  plt.xticks(rotation=45)
```



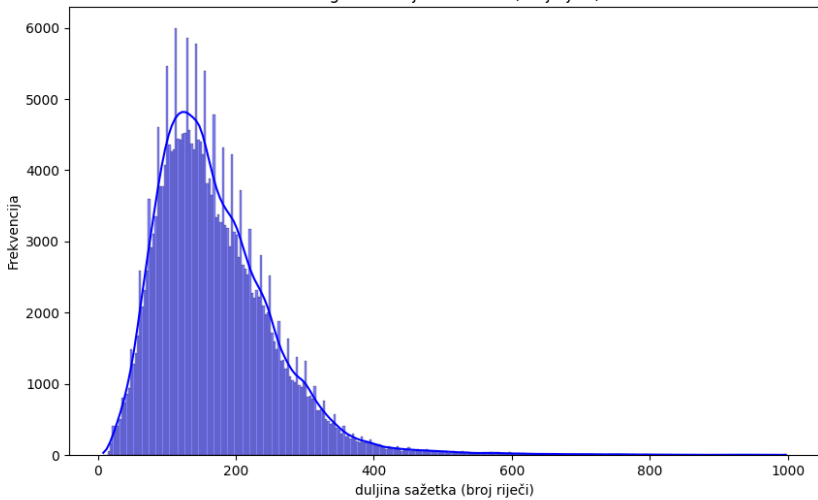


```
1  dates = pd.to_datetime(df['date_defended'], errors='coerce').dropna()
2  dates = dates[dates >= pd.to_datetime(cutoff_start)]
3  years = dates.dt.year
4  sns.histplot(x=years, discrete=True)
```

Histogram - duljina sažetka (broj riječi)



Histogram - duljina sažetka (broj riječi)

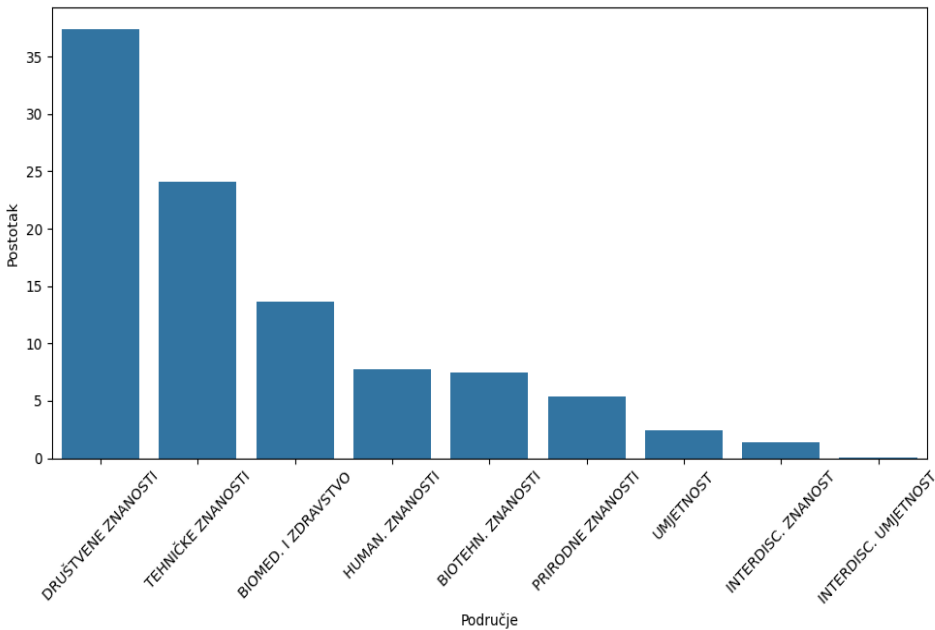


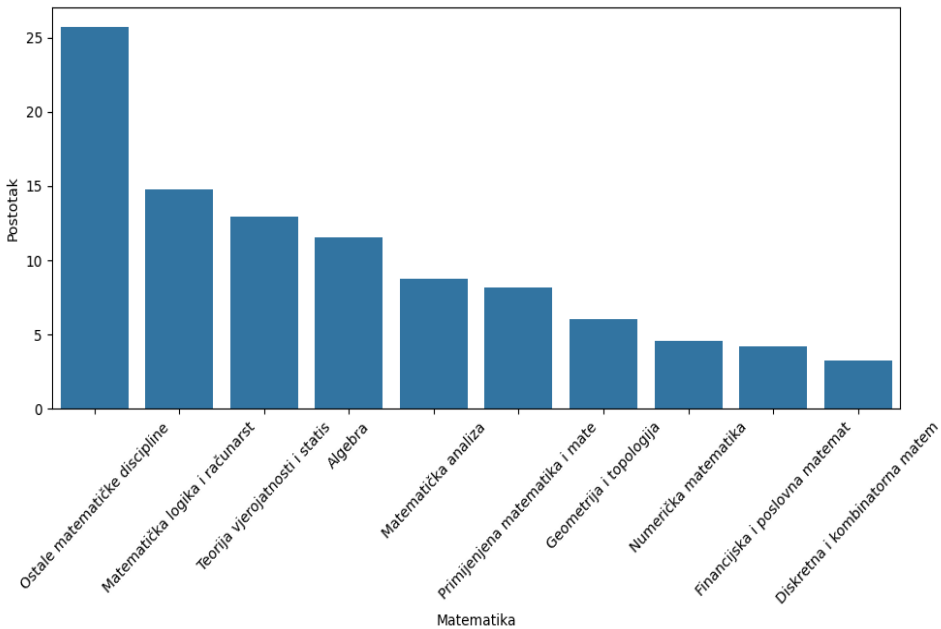
# Pandas i Seaborn

```
1 summary = numbers.describe() # Series
2 print(summary) # Series!
3 ...
4 sns.histplot(numbers, bins="auto", color='blue', kde=True)
```

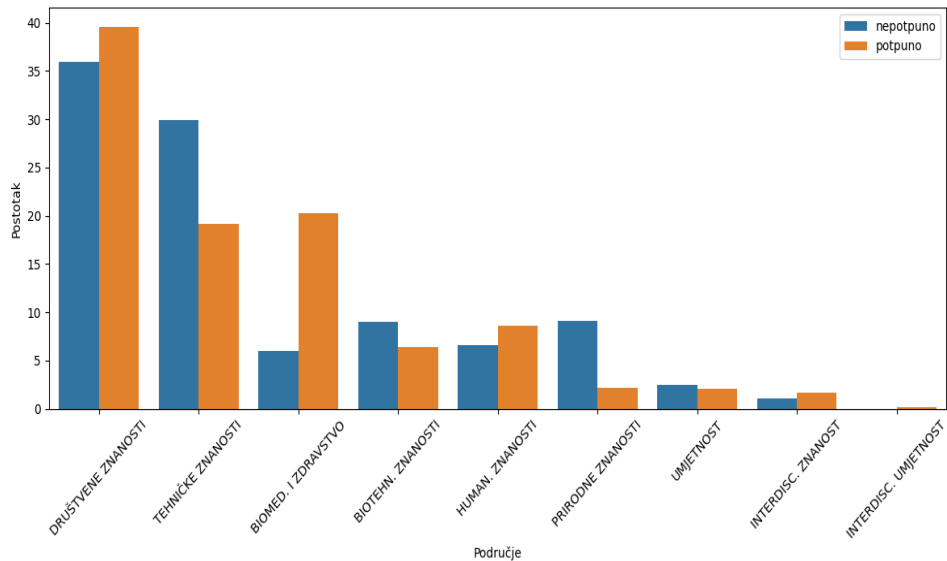
```
count 262287.000000
mean 171.845753
std 104.487170
min 8.000000
25% 110.000000
50% 154.000000
75% 214.000000
max 7141.000000
```

- Područje → Polje → Grana
- Prirodne znanosti → Matematika → Algebra
- Nacionalno vijeće za znanost: “Pravilnik o znanstvenim i umjetničkim područjima, poljima, i granama” (2009)





- Potpuni unos:  
Prirodne znanosti → Matematika → Algebra
- Nepotpuni unos:  
Prirodne znanosti → Matematika



# Pregled tema u skupu podataka

- ključne riječi: 'dramsko pismo', 'filmsko pismo'
- normalizacija: 'dramskoPismo', 'filmskoPismo'

```
1 from wordcloud import WordCloud
2 ...
3 for subjects in df[keyword_column]:
4     subjects_list.extend([normalize(subject) for subject in subjects
5         ↪ if subject])
6     ...
7 wordcloud = WordCloud(width=3200, height=1600,
8     ↪ background_color='white', normalize_plurals=False)
9 wordcloud.generate(' '.join(subjects_list)) # jedan dokument
10 plt.figure(figsize=(20, 10))
11 plt.imshow(wordcloud, interpolation='bilinear')
```

komunikacija analiza organizacija  
turizam  
hrvatska  
ključne riječi  
sigurnost  
kvaliteta  
život  
kvaliteta  
organizacija  
medicinska sestra  
automatizacija  
model  
kvaliteta života  
film  
organizacija  
umjetna inteligencija  
medicinska sestra  
tehničar  
trudnoća  
implementacija  
znanje  
ljudska prava  
suradnja  
istra  
poduzeće  
potrošači  
sigurnost  
digitalna transformacija  
kvaliteta  
metoda konačnih elemenata  
kulturalna baština  
konkurentna prednost  
projekt  
globalizacija  
promet  
poduzetništvo  
aplikacija  
ljudski potencijali  
metoda konačnih elemenata  
kulturalna baština  
konkurentna prednost  
oglašavanje  
energija  
inovacije  
simulacija  
tržište  
računovodstvo  
0 0  
identitet  
električna energija  
zaštita  
arhitektura  
stavovi  
strategija  
motivacija  
prevencija  
fizioterapija  
baza podataka  
internet  
upravljanje  
energetika  
učinkovitost  
prehrana  
hrvatski jezik  
kultura  
liječenje  
adolescenti  
poljoprivreda  
usporedba  
tjelesna aktivnost  
obnovljivi izvori energije  
financijalno  
inovacije  
senzori  
diskriminacija  
mobilna aplikacija  
kvaliteta života  
medicinska sestra  
sustav  
republica hrvatska  
razvoj  
djeca  
republica hrvatska  
zaštita okoliša  
ponašanje  
potrošača  
zdravstvenaljeva  
web aplikacija  
razvoj  
djeca  
matlab  
zaštita okoliša  
ponašanje  
potrošača  
zdravstvenaljeva  
web aplikacija  
razvoj  
djeca  
mladi  
bol  
inkluzija  
sport  
glazba  
diplomski ispit  
koncert  
klasifikacija  
umjetnost  
arduino  
gospodarstvo  
programiranje  
boja  
strojno učenje  
percepcija  
temperatura  
metode  
nogomet  
roditelji  
terapija  
brand  
održavanje  
mediji  
konstrukcija  
održavanje  
mediji  
konstrukcija  
studenti  
html  
css  
geminizam  
proračun  
ambalaza  
ambalaza  
c  
stres društveno odgovorno poslovanje  
financiranje  
fotografija  
glazba  
diplomski ispit  
koncert  
klasifikacija  
umjetnost  
arduino  
gospodarstvo  
programiranje  
boja  
strojno učenje  
percepcija  
temperatura  
metode  
nogomet  
roditelji  
terapija  
brand  
održavanje  
mediji  
konstrukcija  
održavanje  
mediji  
konstrukcija  
studenti  
html  
css  
geminizam  
proračun  
ambalaza  
ambalaza  
c  
marketing  
zdravlje  
android  
društvene mreže  
prodaja zagreb  
tehnologija  
odnosi javnošću  
igra  
europska unija  
kulturalni turizam  
republica hrvatska  
zaštita okoliša  
ponašanje  
potrošača  
zdravstvenaljeva  
web aplikacija  
razvoj  
djeca  
mladi  
bol  
inkluzija  
sport  
glazba  
diplomski ispit  
koncert  
klasifikacija  
umjetnost  
arduino  
gospodarstvo  
programiranje  
boja  
strojno učenje  
percepcija  
temperatura  
metode  
nogomet  
roditelji  
terapija  
brand  
održavanje  
mediji  
konstrukcija  
održavanje  
mediji  
konstrukcija  
studenti  
html  
css  
geminizam  
proračun  
ambalaza  
ambalaza  
c  
marketing  
zdravlje  
android  
društvene mreže  
prodaja zagreb  
tehnologija  
odnosi javnošću  
igra  
europska unija  
kulturalni turizam  
republica hrvatska  
zaštita okoliša  
ponašanje  
potrošača  
zdravstvenaljeva  
web aplikacija  
razvoj  
djeca  
mladi  
bol  
inkluzija  
sport  
glazba  
diplomski ispit  
koncert  
klasifikacija  
umjetnost  
arduino  
gospodarstvo  
programiranje  
boja  
strojno učenje  
percepcija  
temperatura  
metode  
nogomet  
roditelji  
terapija  
brand  
održavanje  
mediji  
konstrukcija  
održavanje  
mediji  
konstrukcija  
studenti  
html  
css  
geminizam  
proračun  
ambalaza  
ambalaza  
c  
marketing  
zdravlje  
android  
društvene mreže  
prodaja zagreb  
tehnologija  
odnosi javnošću  
igra  
europska unija  
kulturalni turizam

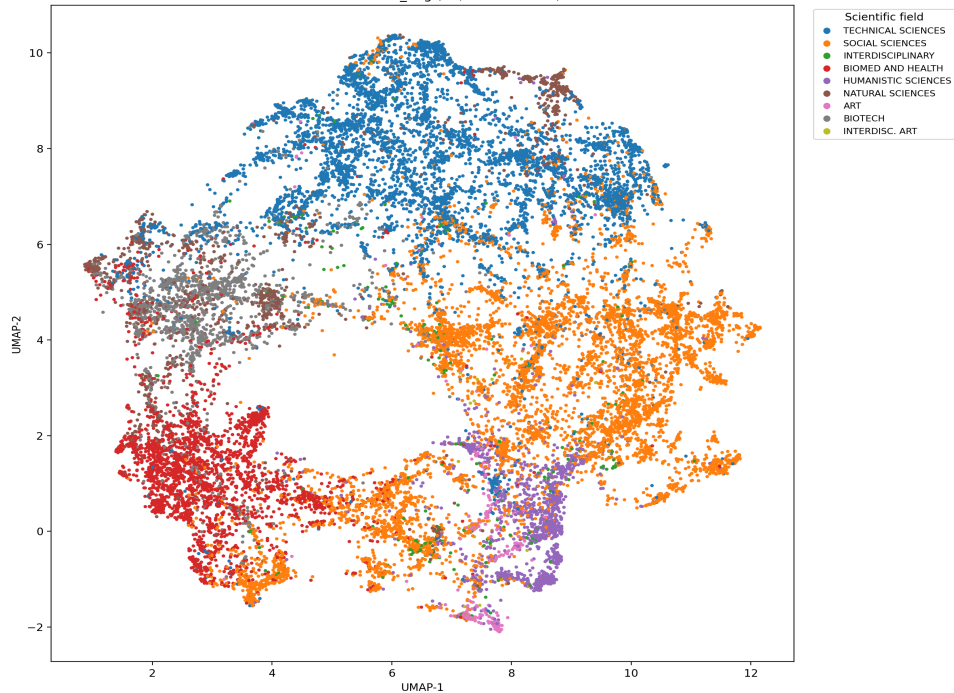
financijskiIzveštaji marketing studenti društvenoOdgovornoPoslovanje percepcija emocije hotelijerstvo algoritam menadžment  
medicina medicina troškovi eu trudnoća transport temperatura pretilost programiranje proizvodnja  
društveneMreže tehnologija promovija zdravlje  
komunikacija razvoj dječja  
održiviRazvoj razvoja dječja  
republikaHrvatska  
kvalitetaživotasigurnost  
turizam  
obrazovanje edukacija  
organizacija  
prevenција  
upravljanje  
mediji  
prevencija  
obrazovanje  
edukacija  
tjelesnaAktivnost globalizacijaoptimizacija automatizacija medicinskaSestra

# UMAP 2D vizualizacija

- “Semantic Embedding”
  - Sažeci tekstova → visokodim. vektor
  - `sentence-transformers/all-mpnet-base-v2`
- UMAP algoritam
  - Konstrukcija grafa u visokodim. prostoru
    - ▶ broj susjeda
    - ▶ mjera udaljenosti
  - Aproksimacija grafa u 2D prostoru

```
1 reducer = umap.UMAP(n_components=2, n_neighbors=30,  
↳ min_dist=min_dist, metric='cosine', random_state=random_state)  
2 return reducer.fit_transform(embeddings)
```

UMAP of abstract\_eng (20,000 abstracts)



- **Alat za istraživanje**
- Alat za zaključivanje (?!)
  - Validacija rezultata
  - U kombinaciji s drugim argumentima
- Varijacija parametara (i random seed-a)
  - Uočavanje novih obrazaca
  - Robustnost zaključaka

- Visual Studio Code
  - `sudo snap install --classic code`
- Jupyter ekstenzija
  - "jupyter.notebookFileRoot": "\$workspaceFolder"
- Data Wrangler ekstenzija
- Jupyter integriran sa GitHub Copilot-om

- Jednostavan kod (u pravilu)
- Lakša provjera rezultata
- ...
- Katalog alata i tehnika (čovjek)  
-> Implementacija (AI)

Hvala na pažnji!