



Jelena Matković
Novi Sad, 2024.

Evaluacija *ChatGPT* modela za ekstrakciju metapodataka iz sudskih odluka

Motivacija

- Pravosuđe se odlikuje konstantnim **povećanjem broja pravnih dokumenata** (zakoni, sudske odluke)
- Potreba za sistemom sa **efikasnom navigacijom** kroz pravne dokumente
- **Transparentnost** pravosudnih odluka
- Organizacija koja omogućava **pretragu** sudskih odluka **po zadatim atributima**
- Transformacija u **mašinski čitljive** dokumente radi obrade i analize

Problem i cilj

- **Ručna ekstrakcija metapodataka** je vremenski zahtevna i podložna greškama
- Potreba za efikasnim načinom ekstrakcije metapodataka iz **nestrukturiranih tekstualnih dokumenata**
- **Tradicionalne metode:**
 - **Regularni izrazi:** Veliki broj šablona, teški za održavanje, slabo skalabilni
 - **Neuronske mreže:** Zahtevaju mnogo obeleženih podataka za treniranje
- **Cilj:** Primena **ChatGPT-a** za jednostavnu i preciznu ekstrakciju metapodataka iz sudskih odluka

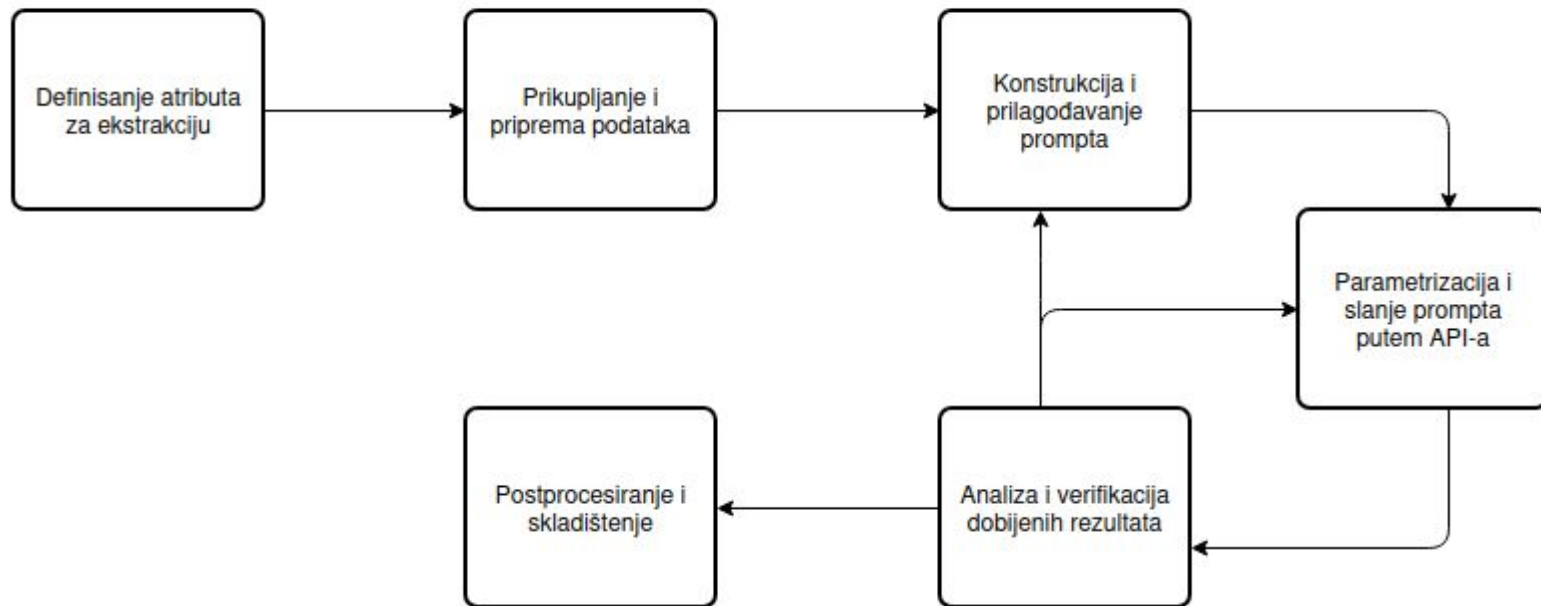
Srodna istraživanja

- **Ekstrakcija metapodataka u e-trgovini:**
 - Dobijanje vrednosti za karakteristike proizvoda
 - *ChatGPT* pokazao uspeh sa *F1* skorom 87.86 u izdvajanju ključnih podataka o proizvodima
 - Testirani različiti tipovi upita
- **Ekstrakcija podataka iz novozelandskih odluka:**
 - Korišćenje *OCR* tehnologije i regularnih izraza za obradu *PDF*-ova sudskih odluka
 - Uspešno ekstrahovano preko 87% podataka

Prikupljanje i priprema podataka

- **Izvor podataka:** 4000 *PDF* dokumenata sa sudskim odlukama iz Novog Zelanda objavljenih od strane okružnih sudova
- **Razlike u formatiranju:** Varijabilna struktura dokumenata otežava konzistentnu ekstrakciju metapodataka
- **Prepoznavanje važnih informacija:** metapodaci navođeni na početku dokumenta
- **Proces pripreme:**
 - Primena *OCR* tehnike za prevođenje odluka iz *PDF*-a u tekstualni format
 - Preuzimanje teksta prve dve stranice dokumenata za analizu modelom *ChatGPT*-a

Tok procesa ekstrakcije



Prompt inženjering

- **Prompt inženjering:**
 - Kreiranje i prilagođavanje upita za tačnu identifikaciju i ekstrakciju
 - Korišćeni **zatvoreni promptovi** koji jasno definišu informacije koje treba izdvojiti (naziv suda, datum itd.)
- **Zero-shot prompting:** Oslanjanje na unapred stečeno znanje modela bez dodatnih primera
- **Few-shot prompting:** Upotreba nekoliko primera za usmeravanje modela ka preciznijoj ekstrakciji

Prompt inženjering

- Primer prompta:

Extract the following attributes for every court decision:

1. Court name
2. Case numbers
3. Neutral citation
4. Applicants, plaintiffs or prosecutors
5. Defendants or respondents
6. Date of hearing
7. Date of judgement
8. Appereances
9. Name of the judge

Organize data by following example:

- 1: <value>
- 2: <value>

Ekstrakcija metapodataka

- **Primarni metapodaci:** naziv suda, identifikaciona oznaka odluke, datum saslušanja, datum presude, ime sudije, zastupnici, neutralni citat, imena učesnika u sporu
- **Proces ekstrakcije:**
 - Kreiranje instrukcije za ekstrakciju atributa
 - Pridruživanje teksta sudske odluke instrukciji
 - Slanje prompta modelu putem API-ja
 - Dobijanje odgovora u formatu <atribut> : <vrednost>

Ekstrakcija metapodataka

- Ekstrahovani podaci čuvani u CSV datoteci:

	Case_name	Case_ID	Neutral_citation	Court	Hearing	Judgment	Plaintiff	Defendant	Appearances	Judge
0	Police v Samsudeen [2021] NZDC 17615 (4 Septem...	CRI-2020-004-006912	[2021] NZDC 17615	District Court at Auckland	4 September 2021	4 September 2021	NEW ZEALAND POLICE	AHAMED AATHIL MOHAMED SAMSUDEEN	B Dickey and H Steele for the Prosecution, H Le...	P Winter
1	R v Moore [2021] NZDC 8794 (7 May 2021)	CRI-2020-090-001698	[2021] NZDC 8794	District Court at Auckland	7 May 2021	7 May 2021	THE QUEEN	JESSIE MOORE	L Oh for the Crown, M Mortimer for the Defendant	N R Dawson
2	R v Gordon Stables [2023] NZDC 1384 (27 Januar...	CRI-2021-043-000626	[2023] NZDC 1384	District Court at New Plymouth	27 January 2023	27 January 2023	The King	Matthew Paul John Gordon-Stables	H Bullock for the Crown, P Keegan for the Defen...	A S Greig
3	R v Lamb [2018] NZDC 14857 (20 July 2018)	CRI-2017-069-000892	[2018] NZDC 14857	District Court at Rotorua	20 July 2018	20 July 2018	THE QUEEN	Liam Lamb	A Gordan for the Crown, L Te Kani for the Defen...	L M Bidois

Evaluacija rezultata

- **Merene metrike:**

- **Preciznost (*Precision*):** Procenat tačno ekstrahovanih metapodataka u odnosu na ukupan broj ekstrahovanih

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Osetljivost (*Recall*):** Procenat tačno identifikovanih metapodataka u odnosu na sve koji su stvarno prisutni

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **F1 skor:** Kombinacija preciznosti i osetljivosti za evaluaciju uspešnosti

Evaluacija rezultata

Naziv atributa	Preciznost	Osetljivost	F1 skor
Naziv suda	1	1	1
ID odluke	1	1	1
Neutralni citat	0.9	1	0.95
Tužilac	1	1	1
Optuženi	1	1	1
Datum saslušanja	0.92	1	0.96
Datum odluke	0.98	1	0.99
Zastupnici	0.95	1	0.97
Ime sudije	0.8	1	0.89
Ukupno	0.91	1	0.95

Analiza grešaka

- Vrednosti neutralnih citata i datuma ekstrahovani i tamo gde nisu navedeni
- Format imena sudija i zastupnika nije konzistentan

Naziv atributa	Stvarna vrednost	Dobijena vrednost
Neutralni citat	Nije navedena	R V LAMELANGI DC TAU CRI-2010-070-002390
Neutralni citat	Nije navedena	R v Goodwin
Ime sudije	S M Harrop	Judge S M Harrop
Datum saslušanja	Nije navedena	23. January 2022.

Prednosti ChatGPT-a

- **Robustnost:**
 - Prilagodljiv različitim formatima i stilovima teksta
- **Razumevanje konteksta:**
 - Otporan na promenu redosleda informacija
 - Razume tražene informacije bez dodatnog objašnjenja
 - Treniran nad velikim brojem podataka

Mane ChatGPT-a

- **Halucinacije:**
 - Tendencija da “izmišlja” vrednosti ukoliko ne može da ih pronađe u tekstu
- **Cena:**
 - Upotreba modela nad velikoj količini dokumenata može biti skupa
- **Vreme:**
 - Analiza i ekstrakcija iziskuje više vremena od tradicionalnih metoda zbog čekanja odgovora modela

Zaključak

- **ChatGPT** predstavlja moćan alat za ekstrakciju metapodataka iz pravnih dokumenata sa visokim stepenom tačnosti
- **Potencijal za širu primenu** u pravnim sistemima i drugim oblastima koje se oslanjaju na nestrukturirane tekstove
- Za preciznije ekstrakovanje potrebno **dodatno obučavanje**
- **Buduća istraživanja:**
 - Proširenje prompta i za složenije attribute
 - Prilagođavanje različitim pravnim sistemima i tipovima dokumenata
 - Obuhvatanje većeg broja dokumenata

Hvala na pažnji!
Pitanja?