

Information extraction from contracts

Thi Hien Ha

May 2020

Smlouva o poskytování služeb

uzavřená dle § 1746 a násl. zákona č. 89/2012 Sb., občanský zákoník ve znění pozdějších předpisů a na základě veřejné zakázky s názvem „Dodávka a instalace aplikace pro využití revizí a pasportizace“

mezi

TESCO SW a.s.

se sídlem:

tr. Kosmonautů 1288/1, Hodolany, Olomouc, PSČ 779 00

IČO:

258925333

DIC:

CZ699000785

zastoupená:

RNDr. **Josefem Tesáříkem**, předsedou představenstva

zapsána v obchodním rejstříku vedeném Krajským soudem v Ostravě, oddíl B, vložka 2530

(dále jako „**poskytovatel**“)

a

Krajská zdravotní, a.s.

se sídlem:

Sociální péče 3316/12A, Ústí nad Labem, PSČ 401 13

IČO:

25488627

DIC:

CZ25488627

zastoupená:

Ing. **Petrem Fialou**, generálním ředitelem společnosti na základě pověření představenstvem společnosti ze dne 17. 12. 2015

zapsána v obchodním rejstříku vedeném Krajským soudem v Ústí nad Labem, oddíl B, vložka 1550

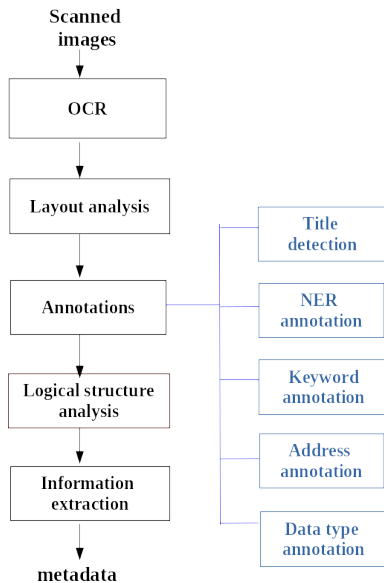
(dále jako „**uživatel**“)

tuto
smlouvu
(dále jen „smlouva“)

Fields to be extracted

- ▶ General info: contract number, contract date, legislation
- ▶ Parties: ORG, address, company id, vat number, representative, bank info, role

Pipeline



what is new in the annotation

- ▶ Process all pages of the contract
- ▶ Title: based on 3 criteria: font size, alignment and keyword
- ▶ NER: using Slavic-Bert-Ner:
<https://github.com/deepmipt/Slavic-BERT-NER>

Data source: <https://smlouvy.gov.cz>

The screenshot displays the 'Publikující smluvní strana' (Publishing contract party) section for 'DOMTEL s.r.o.'. The page is divided into several columns: 'Publikující smluvní strana', 'Smluvní strany' (Contract parties), 'Verze záznamu' (Record version), and 'Soubory' (Files). The 'Publikující smluvní strana' column contains the following information:

- Název subjektu:** Mateřská škola, Praha 3, Sudoměřská 54/1137
- Číslo smlouvy:** 11773596
- ID verze:** 12645996
- Číslo verze:** 1
- Zveřejnění:** 16.05.2020 00:12:34
- Zařetěžení:** Mateřská škola, Praha 3, Sudoměřská 54/1137
- Datová schránka:** qpyfa8

The 'Smluvní strany' column shows details for 'DOMTEL s.r.o.':

- Název:** DOMTEL s.r.o.
- Číslo smlouvy:** 04802925
- Datová schránka:** qpyfa8
- Adresa:** Kottarova 562/24, Praha 8-Úžehův, 180 00 Praha 8
- Účastník / Odbor:**

The 'Verze záznamu' column indicates 'Verze smlouvy: 1' and 'Datum publikace: 16.05.2020'. A button labeled 'Zobrazit detail verze smlouvy' is visible.

The 'Soubory' column lists two files:

- regisr_smlouva_12645996.pdf (825.42 kB, 16.05.2020 00:13:03)
- regisr_smlouva_12645996.pdf (825.42 kB, 16.05.2020 00:13:03)

The 'Adresa záznamu' column provides the URL: <https://smlouvy.gov.cz/smlouva/12645996>.

Data preparation

Prepare training set and test set

- ▶ Down load data from website
- ▶ Apply 2-step filters on downloaded data
- ▶ Randomly select 60 files for the experiment, than random 10 of 60 is used as training set and 50 others as test set.

Create GT from metadata

- ▶ Read html files to map each contract with metadata file
- ▶ convert metadata file into desired format

Problem with the metadata

- ▶ Info is not in the pdf file but in the metadata
- ▶ Info in different formats, e.g: date, address
- ▶ Lack of info: VAT, representative, role, bank info, legislation

Parties' information

- ▶ All potential parties are extracted
- ▶ The more info, the higher confidence.
- ▶ Evaluation: Find the party which has highest number of common fields and compare with that party only.

Result on the training set

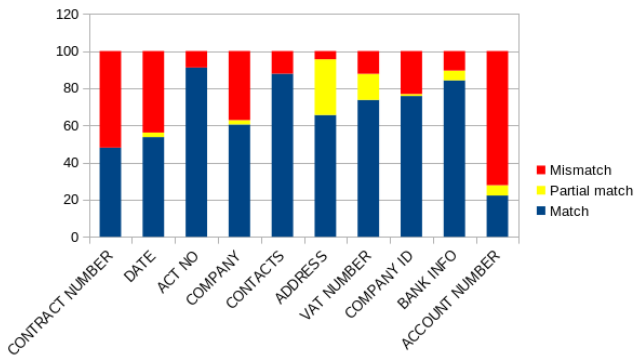
Match	104	83.87%
Patital match	8	6.45%
Mismatch	12	9.68%
Total	124	100%

Mismatch analysis:

- ▶ ORGs are marked as LOCs by both NER and libpostal (4).
E.g: "HLAVNÍ MĚSTO PRAHA"
- ▶ OCR errors: e.g dates are hand-written (2)
- ▶ info is spread more than 1 line, e.g law (1)
- ▶ Party's info is in different blocks (role: 4, bank: 1)

Result on the test set

Match	433	66.31%
Patital match	41	6.28%
Mismatch	179	27.41%
Total	653	100%



Error analysis on the test set

- ▶ OCR errors: hand-written dates, low quality images, covered areas
- ▶ Keywords are not appeared in the training set, e.g "Smlouva č." (8/13 missing contract number), 13/13 missing account number (Bankovní spojen, č.ú)
- ▶ ORG and LOC mis-classified
- ▶ Party info in different blocks