

Information extraction challenges and their evaluation

Filip Graliński

Applica.ai / Adam Mickewicz University

February, 22nd

Information Extraction Challenges

- ▶ Applica Kleister
 - ▶ Kleister NDA
 - ▶ Kleister Charity
- ▶ PolEval 2020 Annual Reports

Information Extraction Challenges

- ▶ Applica Kleister
 - ▶ Kleister NDA
 - ▶ Kleister Charity
- ▶ PolEval 2020 Annual Reports

1. INFORMACJE O SPÓŁKACH WCHODZĄCYCH W SKŁAD GRUPY KAPITAŁOWEJ

JEDNOSTKA DOMINIUJĄCA – PREZENTACJA SPÓŁKI



Zakłady Magnezytowe „ROP CZYCE” S.A. (ZMR S.A.)

Siedziba: Ropczyce, woj. podkarpackie

Adres: ul. Przemysłowa 1, 39-100 Ropczyce

Regon: 690026060

NIP: 818-00-02-127

www.ropczyce.com.pl

PRZEDMIOT DZIAŁALNOŚCI

Przedmiot działalności ZMR S.A. obejmuje produkcję i sprzedaż zasadowych wyrobów og które są niezbędnym elementem konstrukcji wyłożeń pieców i urządzeń cieplnych pr wysokich temperaturach, głównie w hutnictwie żelaza i stali, hutnictwie metali nieżelaz przemysle cementowo-wapienniczym, odlewniczym.

Spółka świadczy także usługi w zakresie nawęglania i ulepszania ciepłego wyrobów oraz pn badawczo-rozwojowe w dziedzinie związanej z przedmiotem jej działalności.

period_from ?

period_to ?

postal_code ?

city ?

...

1. INFORMACJE O SPÓŁKACH WCHODZĄCYCH W SKŁAD GRUPY KAPITAŁOWEJ

JEDNOSTKA DOMINIUJĄCA – PREZENTACJA SPÓŁKI



Zakłady Magnezytowe „ROPCZYCE” S.A. (ZMR S.A.)

Siedziba: Ropczyce, woj. podkarpackie

Adres: ul. Przemysłowa 1, 39-100 Ropczyce

Regon: 690026060

NIP: 818-00-02-127

www.ropczyce.com.pl

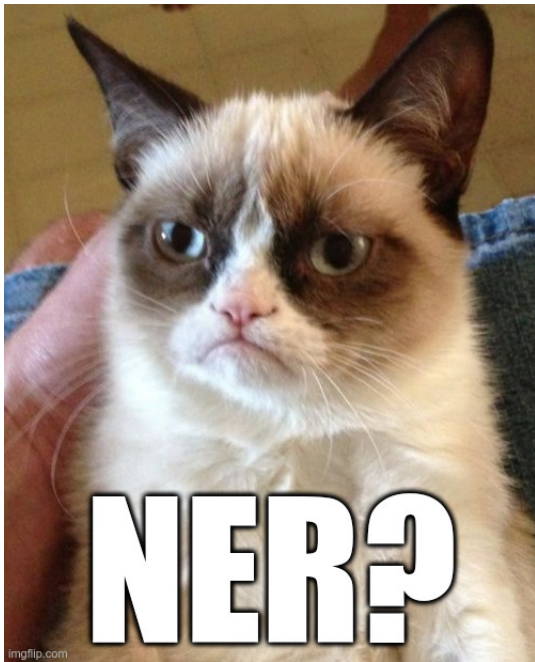
PRZEDMIOT DZIAŁALNOŚCI

Przedmiot działalności ZMR S.A. obejmuje produkcję i sprzedaż zasadowych wyrobów og które są niezbędnym elementem konstrukcji wyłożen pieców i urządzeń cieplnych pr wysokich temperaturach, głównie w hutnictwie żelaza i stali, hutnictwie metali nieżelaz przemysle cementowo-wapienniczym, odlewniczym.

Spółka świadczy także usługi w zakresie nawęglania i ulepszania ciepłego wyrobów oraz pn badawczo-rozwojowe w dziedzinie związanej z przedmiotem jej działalności.

period_from 2012-01-01
period_to 2012-06-30
postal_code 39-100
city Ropczyce
...

company, drawing_date,
period_from, period_to,
postal_code, city, street,
street_no, people



... no!

This is an **Information Extraction** task, not NER*.

- ▶ we are interested in the information not where it is
- ▶ not just any person, but CEO, etc.

* But of course you could use NER as a part of the pipeline

Entities	Description
<i>NDA dataset</i>	
party	parties appearing in the agreement (each of them is treated as a separate entity)
jurisdiction	state or country whose law governs the agreement
effective_date	date on which the contract becomes legally binding
term	duration of the agreement
<i>Charity dataset</i>	
address__post_town	post town (part of a charity address)
address__postcode	postcode (part of a charity address)
address__street_line	street with the house number (part of a charity address)
charity_name	name of the charitable organization
charity_number	identification number in the charity register
report_date	date of reporting
income_annually	annual income in British pounds (GBP)
spending_annually	annual spending in British pounds (GBP)

Evaluation metric

F1-score will be used as the evaluation metric

```
--metric MultiLabel-F1
```

Possible approaches

- ▶ handcrafted rule, e.g. regexps (as a baseline)

Possible approaches

- ▶ handcrafted rule, e.g. regexps (as a baseline)
- ▶ standard NER (for general entities) + role classification

Possible approaches

- ▶ handcrafted rule, e.g. regexps (as a baseline)
- ▶ standard NER (for general entities) + role classification
- ▶ specialized NER (but you need to **autotag** entities)

Possible approaches

- ▶ handcrafted rule, e.g. regexps (as a baseline)
- ▶ standard NER (for general entities) + role classification
- ▶ specialized NER (but you need to **autotag** entities)
- ▶ end-to-end (generative models)

Possible approaches

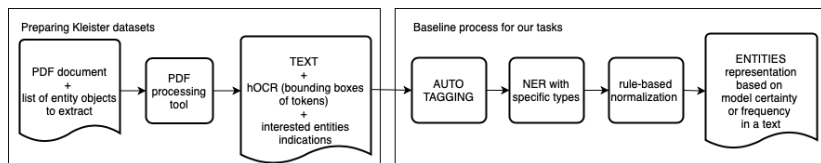
- ▶ handcrafted rule, e.g. regexps (as a baseline)
- ▶ standard NER (for general entities) + role classification
- ▶ specialized NER (but you need to **autotag** entities)
- ▶ end-to-end (generative models)
- ▶ ensembles

Possible approaches

- ▶ handcrafted rule, e.g. regexps (as a baseline)
- ▶ standard NER (for general entities) + role classification
- ▶ specialized NER (but you need to **autotag** entities)
- ▶ end-to-end (generative models)
- ▶ ensembles

Łukasz Garncarek, Rafał Powalski, Tomasz Stanisławek, Bartosz Topolski, Piotr Halama, Filip Graliński, *LAMBERT: Layout-Aware (Language) Modeling using BERT for information extraction*, <https://arxiv.org/abs/2002.08087>

Autotagging approach



Kleister-NDA dataset (pdf2djvu)

Entity name	Flair	BERT	RoBERTa	LayoutLM	LAMBERT	Autotag.	Human
effective_date	79.37	80.20	81.50	82.08	85.27	79.00	100 %
party	70.13	71.60	80.83	75.28	78.70	33.15	98 %
jurisdiction	93.87	95.00	92.87	94.40	96.50	54.10	100 %
term	60.33	45.73	52.27	48.34	55.03	74.10	95 %
ALL	77.83	78.20	81.00	78.68	81.77	60.09	97.86 %

Kleister-Charity dataset (Azure CV)

post_town	83.30	77.03	77.70	79.97	81.03	66.04	98 %
postcode	82.63	87.10	88.40	81.06	82.97	87.60	100 %
street_line	68.17	62.23	72.03	70.92	75.33	75.02	96 %
charity_name	72.40	75.93	78.03	78.82	79.10	67.00	99 %
charity_number	96.73	96.67	95.37	95.76	96.57	98.60	98 %
income	70.93	64.43	69.73	72.86	76.90	69.00	97 %
report_date	95.67	96.60	96.77	95.42	95.80	89.00	100 %
spending	61.67	67.30	68.60	71.20	74.33	73.00	92 %
ALL	80.10	78.33	81.50	80.74	82.97	78.16	97.45 %

Applica Kleister challenges

`https://github.com/applicaaai/kleister-charity`

`https://github.com/applicaaai/kleister-nda`