



Leveraging Information Redundancy of Real-World Data Through Distant Supervision

Presentation: Ariel Cohen

Innovation and Data Unit, IT Department, AP-HP

Ariel Cohen – Innovation and Data Unit, IT Department, AP-HP Alexandrine Lanson – LIMICS – Sorbonne Université Emmanuelle Kempf - Department of Medical Oncology, Henri Mondor and Albert Chenevier Teaching Hospital, AP-HP Xavier Tannier - LIMICS – Sorbonne Université











Paris University Hospitals CDW & Information Redundancy

 (\mathbf{b})



HÔPITAUX DE PARIS

Distant Supervision Approach

- Use multiple data sources to build annotated datasets
- Faster than what can be produced by manual annotation
- Produce a **silver standard** dataset with **Noisy Labels**
- Distant Supervision

B

- Deep neural networks (DNNs) are susceptible to overfit on noisy labels
- Several efforts and methods have been developed to be able to learn from noisy labels



Objectives & Outline

- We propose a new text labeling method that leverages temporal information redundancy from external sources without using any other text information apart from dates themselves for the task of event extraction and classification
- Apply and test this method in order to reconstruct key events of cancer patients' pathway using clinical notes of a cohort of 380,000 patients. We train a classifier using different methods for noise label management, validate the end-to-end approach and compare it with a baseline classifier using an expert annotated corpus
 - Use this model for inference where information is not redundant or lacking



Methodology

- We consider a task of event extraction and classification, from an organizational data lake containing text and other sources of information about real-world events.
- Methodology
- Under the assumption that yellow type events are mentioned using a similar language distribution, we will leverage the text around the mention of the first event to train a model which allows us to identify the second yellow event from text.
- We use the time dimension as a pivot to align
- complementary information present in different data





Programmatic Annotation : Positive Examples

• General method for date alignment between an external source of structured data including event dates and a paired text corpus (e.g., linked by a person/client/user identifier).

Algorithm

- 1. Date recognition (NER) & Normalization
- 2. Select dates mentioned once in the document
- For each person join the dates from structured data with the ones of step 1
 > aligned dates

Person ID	Date	Medical Procedure
1	18/8/2015	Surgery
1	20/07/2019	Biopsy
2	3/10/2022	Biopsy

A male patient born on <u>5/1/1991</u> with a tumour in the upper pole of the right testicle found by chance during self-exploration, for which reason he attended a urology consultation on <u>5/10/2020</u> where a physical examination was performed,

[...] The patient was biopsied on 20/7/2019 at the hospital [...]

Linked by Person ID



Positive Examples

Programmatic Annotation : Negative Examples

2 methods for the selection **of alternative examples**

- 1. Random Selection (RS): A naive random selection among all non-matched dates in the previous step;
- 2. Proximity Selection (PS): A random selection weighted by a score of distance in text with respect to the matched dates.

Negative Examples

A male patient born on 4/1/1991 with a tumour in the upper pole of the right testicle found by chance during self-exploration, for which reason he attended a urology consultation on 5/10/2020 where a physical examination was performed, [...] The patient was biopsied on 20/7/2019 at the hospital [...]



Dataset Creation & Illustrative example

For each labelled date (positive and negative) we extract a **textual context** surrounding the date (snippet) and we mask the date to be classified

>> Silver dataset

>> Use of Noise Management techniques to learn from it

ID	Text	Silver label	Gold label
1	The patient was biopsied on	Biopsy	Biopsy
	<masked date=""> at the hospital</masked>		
2	The patient showed the results of the	Other	Biopsy
	biopsy done outside the hospital on		
	<masked date=""> []</masked>		
3	digestive endoscopy on <masked date="">.</masked>	Biopsy	Biopsy
	The pathological analysis indicates [] On		
	17.08.2019 the patient underwent abdomi-		
	nal surgery		
4	digestive endoscopy on February 15th,	Surgery	Surgery
	2019. The pathological analysis indicates		
	[] On <masked date=""> the patient under-</masked>		
	went abdominal surgery		
5	Mr. Smith came to consultation on	Other	Other
	<masked date=""></masked>		
6	Mrs. Dupont came to consultation on	Biopsy	Other
	<masked date=""></masked>		

Dataset Creation

B





<SEP>

Application use case

- The CDW of the Greater Paris University Hospitals (Assistance Publique-Hôpitaux de Paris — AP-HP) contains the EHR of 380,000 patients with cancer and around 35,000 new cancer cases each year
- We address the problem of cancer patient reconstruction and focus on diagnosis date (biopsy) because of its importance in oncology research (e.g., implication in treatment effectiveness studies, survival analysis, epidemiology, etc.).



ASSISTANCE DE PARIS

Experiment Setting

- Test and Development corpus (101 & 60 documents respectively).
- This corresponds to 1,474 date entities annotated in a binary way for the test set and 680 for the development set (4.5% and 10% correspond to biopsies respectively).
- Apply the **programmatic annotation method** on 4 different settings (data sources and negative examples selection)
- For each produced dataset we test different training configurations:
 - Regular training with CE loss
 - NCE-RCE loss
 - O2U strategy with CE loss
 - O2U strategy with NCE-RCE loss
 - LRT strategy
 - Multiple training with different number of examples of the **development** set to measure **performance as a function of the size** of the **expert annotated** data
 - Compare the performance within patients with and without a biopsy done at hospital

Experiment setting

Results

• 2 training datasets of 10,850 samples (50% biopsies) and 2 of 16,200 samples (33% biopsies).

Method	PR-RS	PR-PS	PR-SP-RS	PR-SP-PS
CE	0.54 (0.53 - 0.54)	0.49 (0.44 - 0.54)	0.46 (0.45 - 0.48)	0.44 (0.42 - 0.45)
NCE-RCE	0.55 (0.53 - 0.60)	0.58 (0.55 - 0.60)	0.60 (0.58 - 0.70)	0.69 (0.40 - 0.74)
O2U - NCE-RCE	0.51 (0.46 - 0.58)	0.56 (0.50 - 0.61)	0.68 (0.67 - 0.71)	0.64 (0.63 - 0.71)
O2U - CE	0.58 (0.50 - 0.59)	0.56 (0.54 - 0.59)	0.67 (0.60 - 0.72)	0.68 (0.63 - 0.71)
LRT	0.56 (0.47 - 0.62)	0.58 (0.56 - 0.62)	0.65 (0.57 - 0.70)	0.70 (0.65 - 0.70)

Results

Table 2: Median F1-score ([min-max] over 5 iterations on the test set) comparison between methods using different training datasets. These are: i. PR-RS: Pathology Reports and Random Selection; ii. PR-PS: Pathology Reports and Proximity Selection; iii. PR-SP-RS: Pathology Reports, Surgery Procedures and Random Selection; iv. PR-SP-PS: Pathology Reports, Surgery Procedures and Proximity Selection.

Class	Precision	Recall	Support
Biopsy	0.91	0.62	34
Other	0.96	0.99	280

Table 3: Performance of the programmatic annotation for the PR-RS dataset, evaluated on development set.



Results: a baseline with an expert annotated dataset

- Model performance increases in function of the number of examples used for training using an expert annotated dataset
- The distant supervised approaches perform comparably to the results obtained with 450 clean labeled entities
 - Performance increase is slower when using more than 20 documents.



Figure 3: Median F1-score ([min-max] over 5 iterations) infunction of number of training examples using an expertlabeled dataset.ASSISTANCE IN HOPITAUX

Results

Results: bias

- Model performance for two patient groups: those who underwent at least one hospital biopsy, and the others.
- Higher recall rates among the first group

	w/ biopsy		w/o biopsy	
Method	Prec.	Rec.	Prec.	Rec.
CE	0.38	0.90	0.27	0.53
NCE-RCE	0.61	0.90	0.59	0.69
O2U - NCE-RCE	0.54	0.93	0.57	0.69
O2U - CE	0.56	0.93	0.60	0.69
LRT	0.64	0.90	0.59	0.69
EL dataset - CE	0.69	0.93	0.62	0.75

Table 4: Median precision and recall (over 5 iterations) for patients with and without a biopsy done at the hospital, evaluated on test set. Model trained on PR-SP-PS and EL datasets.



Results

Discussion & Conclusions

- Our domain agnostic method helps to minimize the necessary annotation effort in a context when experts' available time is scarce, such as healthcare; however, we reduce expert annotation effort in exchange for an industry knowledge investment.
- All models, including the expert labeled one, **underperform** when evaluated on **patients without biopsy procedures** done at the hospital.
- The absence of any NER procedure for the concepts to identify allows us to learn complex patterns from text
- The method is language- and vocabulary-independent, and therefore directly applicable to other CDWs, to retrieve other important dates or in another dataset that includes event dates and a paired text corpus.
- Successful results comparable with models fitted on an expert annotated corpus
- The benefit and need of using **noise management** methods is also demonstrated. We **improve performance up to 59%**, all approaches are distinctly better than a classic training using the CE loss
- It has been experimentally shown that the use of **multiple sources** combined together reach **better results**
- Other techniques of noise management and semi-supervised approaches, as well as a better bias understanding, are still to be explored.









Ó

Noise Management

- 1. <u>Robust Loss Function</u>: Normalized Cross Entropy Reverse Cross Entropy (NCE-RCE) loss, an Active-Passive loss introduced by Ma et al. (2020), with theoretical robustness properties
- 2. <u>Sample Selection</u>: O2U-Net approach (Huang et al.,2019). Multiple rounds changing between Overfitting and Underfitting using a cyclical learning rate. Sample losses are ranked and the top k% is then removed from the dataset. Then a final step of training is done on cleaned data.
- **Noisy Labels** 3. Label Refurbishment: a classifier trained on noisy data has low confidence in the label of a sample, that label is likely to be false Zheng et al. (2020). The method applies a Likelihood Ratio Test (LRT) on noisy classifier predictions to check label purity, the likelihood ratio is compared with a predetermined threshold δ ; then, it corrects wrong labels for future training



Clinical Data Warehouse (CDW)



npj Digital Medicine

ARTICLE OPE

International electronic health record-derived COVID-19 clinical course profiles: the 4CE consortium

Gabriel A. Brat (14), Griffin M. Weber^{1,43}, Nils Gehlenborg¹, Paul Avillach (16), Nathan P. Palmer¹, Luca Chiovato^{2,3}, James Cimin Lemuel R. Waitman⁵, Glibert S. Omenn (17), Alberto Malovini⁷, Jason H. Moora^{2,8}, Brett K. Beaulieu-Jones (10), Valentina Tibollo²,

Journal of the American Medical Informatics Association, 2024, 1–11 https://doi.org/10.1093/jamia/ocae069 Research and Applications

Research and Applications

Collaborative and privacy-enhancing workflows on a clinical data warehouse: an example developing natural language processing pipelines to detect medical conditions

Thomas Petit-Jean (), MSc^{1,*}, Christel Gérardin (), MD, PhD^{1,2}, Emmanuelle Berthelot (), MD³, Gilles Chatellier (), MD^{1,4}, Marie Frank (), MD⁵, Xavier Tannier (), PhD⁶, Emmanuelle Kempf (), MD, PhD^{5,7}, Romain Bey (), PhD¹

Innovation and Data Unit, IT Department, Assistance Publique-Hopitaux de Paris, Paris, 75012, France, ¹Institut Pierre-Louis
 d'Epidemiogie et de Sante Publique, INSERM, Sontonen Université, Prins, 78012, France, ¹Opartment of Cardiology, Hopital Bickitre,
 d'Epidemiogie et de Sante Publique, INSERM, Sontonen Université, Prins, 78012, France, ¹Department of Medical Informatics,
 diplante de Paris, Centre-Université de Paris, Paris, 7901, France, ²Department of Medical Informatics,
 diplante de Paris, Centre-Université de Sante Paris, Educationes, Paris, 7905, France, ¹Laborationes,
 diplante de Paris, Centre-Université de Sante Paris, Educationes, Paris, 7905, France, ¹Laborationes,
 diplante de Paris, Centre-Université de Sconsaissances por la e-Sante UniVCS, NISERM, Multicel Informatics,
 diplante de Paris, Centre, ¹Department of Medical Oncology, Henri Mondor and Albert Chenevier Teaching Hospital, Assistance
 Publique-Hopitaux de Paris, Centre, ¹Department of Medical Discourse,
 Henri Mondor and Albert Chenevier Teaching Hospital, Assistance
 Publique-Hopitaux de Paris, Centre, ¹Department of Medical Discourse,
 Henri Mondor and Albert Chenevier Teaching Hospital, Assistance
 Publique-Hopitaux de Paris, Centre, ¹Department of Medical Discourse,
 Hopitaux de Paris, Centre, ¹Department of Medical Discourse,
 Henri Mondor and Albert Chenevier Teaching Hospital, Assistance
 Publique Hopitaux de Paris, Centre,
 Department of Medical Discourse,
 Henri Mondor and Albert Chenevier Teaching Hospital, Assistance
 Hopitaux de Paris,
 Henri Scology,
 France,
 Teaching Hospitaux

Corresponding author: Thomas Petit-Jean, MSc, Innovation and Data Unit, Assistance publique-Hópitaux de Paris, 33 Boulevard de Picpus, Paris, 75012, rance (thomas.petitjean@aphpt)

ARTICLE OPEN

mental health

npj

() Check for up

Natural language processing of multi-hospital electronic health records for public health surveillance of suicidality

www.nature.com/npjmenta

Check for update

Romain Bey¹, Ariel Cohen^{1 ©}, Vincent Trebossen², Basile Dura¹, Pierre-Alexis Geoffroy^{34,5,6}, Charline Jean^{1,7,8}, Benjamin Landmanⁱ Thomas Petit-Jean¹, Gilles Chatellier^{1,9}, Kankoe Sallah¹⁰, Xavier Tannier¹¹, Aurelie Bourmaud^{9,12,13} and Richard Delorme^{2,14}

There is an urgent need to monitor the mental health of large populations, especially during crises such as the COVID-19 pandemic, to timely identify the most at-risk subgroups and to design targeted prevention campaigns. We therefore developed and validated surveillance indicators related to suicidality: the monthly number of hospitalisations caused by suicide attempts and the prevalence among them of five known risk factors. They were automatically computed analysing the electronic health records of fifteen university hospitals of the Paris area. France, using natural language processing algorithms based on artificial intelligence. We evaluated the relevance of these indicators conducting a retrospective cohort study. Considering 2911;920 records contained in a <u>common data warehouse, wet sets</u> of changes after the pandemic outbreak in the slope of the monthly number of suicide **Cancer Epidemiology**.

No changes in clinical presentation, treatment strategies and survival of pancreatic cancer cases during the SARS-COV-2 outbreak: A retrospective multicenter cohort study on real-world data

*EHR : Electronic Health Records

20

B

Other conclusions

- We confirm that it is possible to leverage information redundancy of an organizational data lake to build a programmatically annotated corpus and train ML models, minimizing the required expert time for the annotation task.
- Domain agnostic approach, interesting in settings with scarce experts' available time, huge amounts of data and an industry knowledge
- Structuring information from text **impacts downstream applications** as patient recruitment for clinical trials, treatment effectiveness studies, survival analysis or epidemiology studies
- Possible to **share** a method without sharing model weights or datasets



Conclusion