

Feature Extraction for Pathology Reports Classification with Precise Negation Scope Detection

Olga Zamaraeva, Kristen Howell, & Adam Rhine
University of Washington
Department of Linguistics

DELPH-IN Summit in Oslo, Norway
August 7 2017

Overview

- Broad task: Classify pathology reports as lymph node positive or negative
 - Did the cancer spread to the lymph nodes?
- Core task: find concepts which could be important features (e.g. *malignancy*) in negated context
 - “*No malignancy was found in the lymph nodes.*”
 - “*Metastasis: not identified.*”

RQ 1: How good ERG is for negation detection?

- McKinlay et al (2012)
 - Used ERG to find negation scope
 - Evaluated on sentence level (intrinsic evaluation of negation scope detection)
 - Data: Biomedical literature
 - Events of interest are already identified in the data in previous stages.
 - Feature vectors are constructed for these events, and then the events are classified in terms of negation scope.
- Packard et al. (2012)
 - Negation scope in Sherlock Holmes stories.

RQ2: Is negation important (for classification)?

- **...and if so, can we classify better with ERG?**
- Maybe not:
 - Words that occur negated as well as not negated will likely not be selected as most informative
- Maybe yes:
 - Negated concepts may be good features themselves
 - (Also, machine learning classification is not the only scenario)

Dataset

- SEER program (<https://seer.cancer.gov/>)

	Total	LN	POS	NEG	UNK
Training	581	298	52	91	155
Test	293	137	26	44	67

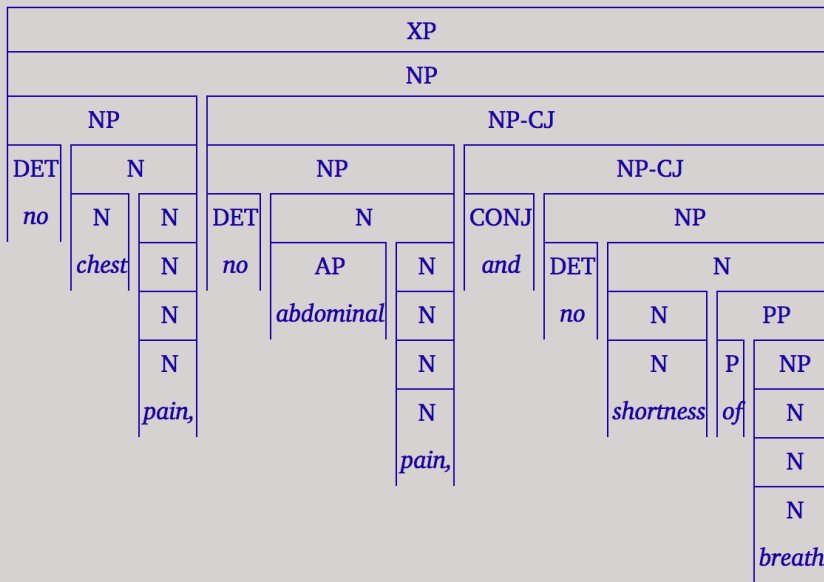
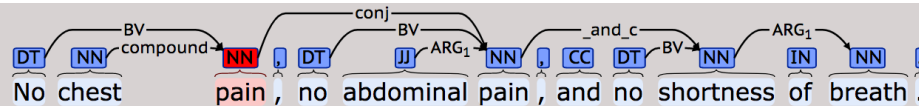
- Annotated on report level
 - Not annotated for negation scope
 - Makes our approach amenable only to extrinsic evaluation

(Extrinsic) Evaluation

- Compare to NegEx
 - Chapman et al. (2002)
 - Regular expressions-based
 - Very widely used
 - Easy to adapt to any English dataset
- Comparing rule-based to rule-based*
 - Independent of the dataset?
 - Not too many added heuristics?

*negation detection

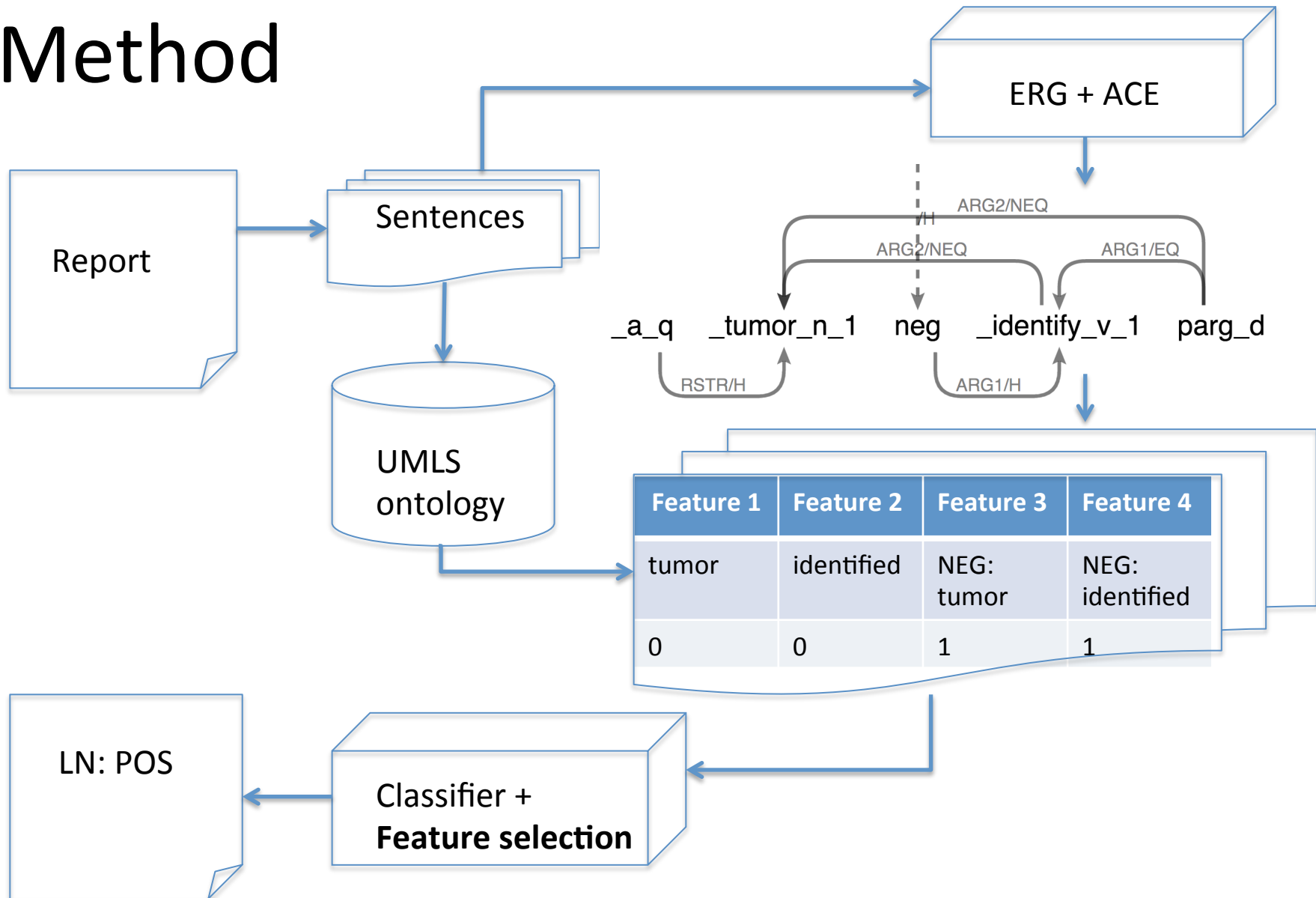
Example where NegEx fails



```
e3:
e3:unknown<0:61>[ARG x4]
_1:udef_q<0:61>[BV x4]
_2:_no_q<0:2>[BV x10]
e14:compound<3:14>[ARG1 x10, ARG2 x13]
_3:udef_q<3:8>[BV x13]
x13:_chest_n_1<3:8>[]
x10:_pain_n_of<9:14>[]
_4:udef_q<15:61>[BV x22]
x4:implicit_conj<15:61>[L-INDEX x10, R-INDEX x22]
_5:_no_q<15:17>[BV x27]
e30:_abdominal_a_1<18:27>[ARG1 x27]
x27:_pain_n_of<28:33>[]
x22:_and_c<34:37>[L-INDEX x27, R-INDEX x33]
_6:_no_q<38:40>[BV x33]
x33:_shortness_n_of<41:50>[ARG1 x38]
_7:udef_q<54:61>[BV x38]
x38:_breath_n_1<54:61>[]
```


1

Method



Feature Selection

Selection algorithm		with NegEx	with ERG
Variance Threshold (90%) (VT)	total	356	368
	negated	22	6
K-best (K=100)	total	100	100
	negated	17	7
Best percentile (10%)	total	488	408
	negated	88	13
Wrapper: AdaBoost (AB)	total	37	37
	negated	1	2
Wrapper: Random Forest (RF)	total	489	515
	negated	79	24
Wrapper: Logistic Regression (LR)	total	1257	1061
	negated	186	37

(Null? Incremental?) Results



↓ Classifiers

Feature selectors



FS→	VT		AB		RF		K=100		10%		LR	
	NE	ERG	NE	ERG	NE	ERG	NE	ERG	NE	ERG	NE	ERG
CL												
AB	0.74	0.76	0.71	0.73	0.69	0.74	0.69	0.64	0.69	0.72	0.77	0.69
DT	0.66	0.72	0.68	0.76	0.69	0.69	0.70	0.69	0.74	0.75	0.74	0.66
KNN	0.67	0.61	0.73	0.76	0.66	0.69	0.71	0.69	0.69	0.68	0.61	0.61
SVM	0.74	0.74	0.76	0.78	0.77	0.74	0.72	0.74	0.75	0.78	0.74	0.74
NB	0.64	0.66	0.42	0.39	0.69	0.69	0.34	0.36	0.64	0.68	0.55	0.59
NN	0.70	0.72	0.72	0.80	0.77	0.73	0.72	0.72	0.74	0.78	0.73	0.73
RF	0.73	0.75	0.76	0.80	0.78	0.77	0.74	0.74	0.76	0.77	0.75	0.76

F1 micro-average classification scores

Results on more constrained models

FS→	AB			RF			GB		
	BL	NE	ERG	BL	NE	ERG	BL	NE	ERG
AB	0.72	0.75	0.73	0.74	0.72	0.74	0.73	0.72	0.73
RF	0.77	0.76	0.74	0.78	0.75	0.75	0.72	0.71	0.73
GB	0.77	0.77	0.77	0.77	0.77	0.73	0.74	0.76	0.75
VT	0.77	0.77	0.76	0.77	0.75	0.75	0.74	0.75	0.74

Micro-average F1 scores

- No improvement from adding any negated features

Better(?) results

- Same dataset, filtered for “Laterlality Category”: Left or Right (lung, breast etc).
 - “Better” (bigger, more balanced) dataset,
 - but still small, reports are still different length etc.

FS→	AB			RF			GB		
	CL	BL	NE	ERG	BL	NE	ERG	BL	NE
AB	0.87	0.86	0.86	0.83	0.77	0.87	0.84	0.86	0.85
RF	0.86	0.87	0.87	0.77	0.79	0.86	0.85	0.85	0.86
GB	0.87	0.87	0.87	0.80	0.78	0.84	0.87	0.87	0.87
Vt	0.87	0.87	0.86	0.81	0.79	0.86	0.86	0.86	0.86

Micro-average F1 scores

Issues

- Sentences are fragmented
 - Easy to make mistakes in tokenization
 - Difficult to parse in a meaningful way
- Dataset is small and unbalanced
 - Few ML algorithms are effective
 - **Very** hard to select features/not overfit
- All changes/tuning that we tried so far did not lead to much change in the results
 - Some negated features are among the selected ones but not many of them

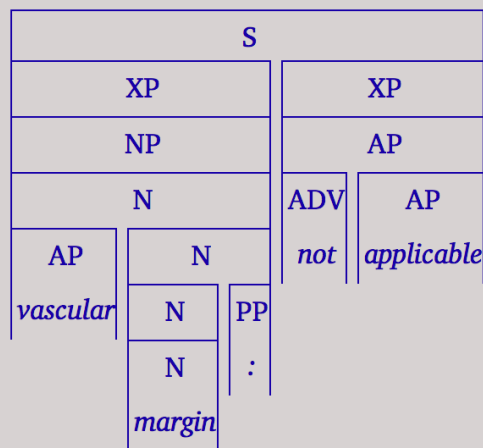
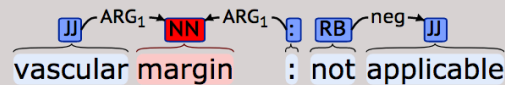
Parse Issues

- Parse coverage
 - 79% (22K/28K)
 - E.g. *Despite an FDA approved scoring guide that classifies 2+ immunohistochemistry results as equivocal, the literature suggests that up to 70% of these cases may be actually the fluorescence in situ hybridization (FISH) negative.*
 - E.g. *Necrosis of invasive component: not identified.*
 - fall back?
 - if fall back to NegEx, it adds more feature tokens than ERG (6K vs. 4K)
 - Compare only sentences for which have a parse?
 - Results still seem null/incremental
- Parse selection
 - Often one of the first parses is good but not the top parse.
 - E.g. *No **X or Y** were identified.*
 - Though, selecting widest scope did not help classification
 - Different top parse with fragmenting on/off
- MRS crawling is basic
 - Just looking at *neg* relation and its ARG1, ARG2 and *no* quantifier for nouns

Parse selection issues

- *Tumor: not identified*
 - Often times, *identified* is parsed as a post-head adjective (*not_c*) rather than a predicate (*neg*).

vascular margin: not applicable



e3:

e3:implicit_conj<0:31>[L-HNDL e7, L-INDEX e7, R-HNDL e6, R-INDEX e6]

e7:unknown<0:16>[ARG x9]

_1:udef_q<0:16>[BV x9]

e14:_vascular/jj_u_unknown<0:8>[ARG1 x9]

x9:_margin_n_1<9:15>[]

e15:_colon_p_namely<15:16>[ARG1 x9]

e6:unknown<17:31>[]

e19:neg<17:20>[ARG1 e21]

e21:_applicable_a_1<21:31>[]

Parse selection issues

- Passive voice with modals: top-ranked parse with fragments on:

	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td colspan="7" style="text-align: center;">XP</td></tr> <tr><td colspan="3" style="text-align: center;">NP</td><td colspan="4" style="text-align: center;">PP</td></tr> <tr><td colspan="3" style="text-align: center;">N</td><td colspan="4" style="text-align: center;">P</td></tr> <tr><td colspan="3" style="text-align: center;">AP</td><td colspan="4" style="text-align: center;">S</td></tr> <tr><td colspan="3" style="text-align: center;">regional</td><td colspan="4" style="text-align: center;">:</td></tr> <tr><td colspan="3" style="text-align: center;">N</td><td colspan="2" style="text-align: center;">NP</td><td colspan="2" style="text-align: center;">VP</td></tr> <tr><td colspan="3" style="text-align: center;">lymph</td><td colspan="2" style="text-align: center;">N</td><td colspan="2" style="text-align: center;">V</td></tr> <tr><td colspan="3" style="text-align: center;">nodes</td><td colspan="2" style="text-align: center;">nx</td><td colspan="2" style="text-align: center;">(cannot be</td></tr> <tr><td colspan="3"></td><td colspan="2"></td><td colspan="2" style="text-align: center;">VP</td></tr> <tr><td colspan="3"></td><td colspan="2"></td><td colspan="2" style="text-align: center;">V</td></tr> <tr><td colspan="3"></td><td colspan="2"></td><td colspan="2" style="text-align: center;">V</td></tr> <tr><td colspan="3"></td><td colspan="2"></td><td colspan="2" style="text-align: center;">V</td></tr> <tr><td colspan="3"></td><td colspan="2"></td><td colspan="2" style="text-align: center;">assessed).</td></tr> </table>	XP							NP			PP				N			P				AP			S				regional			:				N			NP		VP		lymph			N		V		nodes			nx		(cannot be							VP							V							V							V							assessed).		<p>e3:</p> <pre> e3:unknown<0:46>[ARG x5] _1:undef_q<0:20>[BV x5] e10:_regional_a_1<0:8>[ARG1 x5] e12:compound<9:20>[ARG1 x5, ARG2 x11] _2:undef_q<9:14>[BV x11] x11:_lymph/nn_u_unknown<9:14>[] x5:_node_n_1<15:20>[] e18:_colon_p_namely<20:21>[ARG1 e3, ARG2 e28] _3:undef_q<22:24>[BV x23] x23:_nx/fw_u_unknown<22:24>[] e28:neg<25:32>[ARG1 e31] e31:_can_v_modal<25:32>[ARG1 e33] e33:_assess_v_1<36:46>[ARG2 x23] e35:parg_d<36:46>[ARG1 e33, ARG2 x23] </pre>
XP																																																																																													
NP			PP																																																																																										
N			P																																																																																										
AP			S																																																																																										
regional			:																																																																																										
N			NP		VP																																																																																								
lymph			N		V																																																																																								
nodes			nx		(cannot be																																																																																								
					VP																																																																																								
					V																																																																																								
					V																																																																																								
					V																																																																																								
					assessed).																																																																																								

0

Parse selection issues

- Passive voice with modals: top-ranked parse with fragments off:

regional lymph nodes : nx (cannot be assessed) .

	S									
	NP					VP				
	N									
	N				PP		V	VP		
	N		N		P	N	(cannot	be	VP	
#	AP	N	N		:	N	V	V	V	
0	<i>regional</i>	<i>lymph</i>	<i>nodes</i>		<i>nx</i>	<i>assessed</i>	<i>assessed</i>	<i>assessed</i>	<i>assessed</i>	
<input type="checkbox"/>										

e23:

```

_1:udef_q<0:24>[BV x6]
e10:compound<0:20>[ARG1 x6, ARG2 x9]
_2:udef_q<0:14>[BV x9]
e15:_regional_a_1<0:8>[ARG1 x9]
x9:_lymph/nn_u_unknown<9:14>[]
x6:_node_n_1<15:20>[]
e16:_colon_p_namely<20:21>[ARG1 x6, ARG2 x17]
_3:udef_q<20:21>[BV x17]
x17:_nx/fw_u_unknown<22:24>[]
e23:neg<25:32>[ARG1 e3]
e3:_can_v_modal<25:32>[ARG1 e27]
e27:_assess_v_1<36:46>[ARG2 x6]
e29:parg_d<36:46>[ARG1 e27, ARG2 x6]

```

Conclusion

- ERG should be more useful than something like NegEx for finding negated items.
 - Without too many heuristics for crawling MRS?..
 - Even with heuristics, MRS is more general than surface strings
- How to achieve results that clearly show this?
 - So far, improving parse selection did not help.
 - Need to work more on tuning parameters for classifiers and (especially) feature selection