Our Quest for Interpretable Natural Language Processing

Mihai Surdeanu

April 2020





Recent ML developments: a deal with the devil



Interpretability

- Interpretability is an overloaded term in machine learning (ML)
- But we can classify it roughly in two classes

1. Post hoc interpretability (Ribeiro et al., 2016)



Figure 3: Toy example to present intuition for LIME. The black-box model's complex decision function f (unknown to LIME) is represented by the blue/pink background, which cannot be approximated well by a linear model. The bold red cross is the instance being explained. LIME samples instances, gets predictions using f, and weighs them by the proximity to the instance being explained (represented here by size). The dashed line is the learned explanation that is locally (but not globally) faithful.

2. Global interpretability: Converting the statistical model into something interpretable (aka "knowledge distillation")



(Craven and Shavlik, 1996) (Hinton et al., 2015)

Comparison

- Post hoc interpretability
 - May provide only an approximate explanation (Netflix lies to you ^(C))
 - You can't fix a problem with the original model when identified
- Global interpretability
 - May lose some performance in the conversion
 - Allows corrections to the model when problems are discovered

Comparison

- Post hoc interpretability
 - May provide only an approximate explanation (Netflix lies to you ^(C))
 - You can't fix a problem with the original model when identified
- Global interpretability
 - May lose some performance in the conversion
 - Allows corrections to the model when problems discovered

Comparison

- Post hoc interpretability_
 - May provide of (Netflix lies to
 - You can't fix a probly the original model when identified
- Global interpretability
 - May lose some performance in the conversion
 - Allows corrections to the model when problems discovered

nn

Globally interpretable models for natural language processing

- 1. Humans can understand the model
- 2. Humans can change the model

Academia vs. industry

Implementa@ons*of*En@ty*Extrac@on*



Why do we care about global interpretability?

- Most of today's meaningful projects are interdisciplinary, e.g., ML + medicine
- We can't expect an expert in another domain to understand (and fix!) our statistical classifiers
- But we need to iterate quickly...



USE CASE: MACHINE READING FOR CANCER RESEARCH

Why cancer?



\$200 billion have been invested in cancer research since (R. Barzilay, NAACL 2016)₄

Why cancer?

Cancer Death Rates* Among Women, US, 1930 – 2005



*Age-adjusted to the 2000 US standard population. Source: US Mortality Data 1960-2005, US Mortality Volumes 1930-1959, National Center for Health Statistics, Centers for Disease Control and Prevention, 2008.

Why cancer?

Cancer Death Rates* Among Men, US, 1930 – 2009



*Age-adjusted to the 2000 US standard population. Source: US Mortality Data 1960-2009, US Mortality Volumes 1930-1959, National Center for Health Statistics, Centers for Disease Control and Prevention.

Why the slow progress?



Publications indexed by PubMed each year since 1995

Why the slow progress?



90% are never cited!¹

¹ http://www.smithsonianmag.com/smart-news/half-academic-studies-are-never-read-more-three-people-180950222/?no-ist 18



"A knotty puzzle may hold a scientist up for a century, when it may be that a colleague has the solution already and is not even aware of the puzzle that it might solve."

– Isaac Asimov, The Robots of Dawn

We need machine reading

• If humans can't process this much information, then machines must help!

Machine reading for biomedical literature



















Few grammar rules

Туре	Syntax	Surface	Total	
Entities	0	15		15
Generic entities	0	2		2
Modifications	0	6		6
Mutants	0	9		9
Total entities	0	32		32
Simple event				26
Binding Inese rule	Later I will discuss our work on using learn such rules using			37
Hydrolysis Later I v				10
Translocation Using				12
Positive regu	representation learning.			20
Negative regu				17
Total events	95	27		122
Total	95	59		154

How well does machine reading work for a complete reading task?



But is machine reading *actually* useful?

Mutual exclusivity intuition



The Mutex algorithm (by Emek Demir and colleagues)

Mutex insight: If a tumor "wants" to disable a mechanism, it will mutate something upstream, but it generally won't "pay" for two mutations that do the same thing. So mutually exclusive mutations plus a good model can tell us which mechanisms the tumor disables.



How Mutex works

- Mutex does a graph search on the signal g network to find subgraphs of genes that
 - are altered in mutually exclusive manner, and
 - have a common downstream signaling target.

Machine reading contributes here!

Patient data




Brain (GBM)

Machine reading suggests novel hypotheses that are missed by the authors of the individual publications.

LEARNING INTERPRETABLE MODELS

 It took us 1 – 2 person months to build the grammar of 150+ rules for the biomedical domain

- In many cases, one does not have this time
 - Scenario 1: But training data exists
 - Scenario 2: And training data does not exist either...

Motivation

 It took us 1 – 2 person months to build the grammar of 154 rules for the biomedical domain

- In many cases, one does not have this time
 - Scenario 1: But training data exists (2 papers)
 - Scenario 2: And training data does not exist either...

SnapToGrid: From Statistical to Interpretable Models for Biomedical Information Extraction

Marco A. Valenzuela-Escárcega, Gus Hahn-Powell, Dane Bell, Mihai Surdeanu University of Arizona Tucson, AZ 85721, USA {marcov, hahnpowell, dane, msurdeanu}@email.arizona.edu

A three-step process

- 1. Train a statistical classifier
 - a) Aggressive feature selection using regularization
- 2. Convert model to rules
 - a) Convert features to rules
 - b) "Snap to grid": throw away most statistics by discretizing feature weights
- 3. Model editing

BioNLP 2009 Shared Task



BioNLP 2009 Shared Task



BioNLP 2009 Shared Task



Step 1: train statistical classifiers

• We use logistic regression for both classifiers

- Feature selection through regularization
 - L1 regularization: aggressive feature selection; slightly lower performance

Machine learning features

Triggers

- Token
 - Word, lemma, gazetteer
- Surface
 - Token features for a window around token of interest
 - Bigrams
- Syntax
 - Dependency paths up to depth 2
 - Token features for token at the end of each path
- Bag-of-words and entity count
 - For whole sentence
 - For window surrounding token of interest

Relations

- Path
 - Shortest dependency path between nodes of interest
 - With and without lexicalization
 - Path length
- Surface
 - Words surrounding and between nodes of interest
- Consistency
 - Soft constraints on edges between trigger and arguments (e.g., only regulations have causes)
- N-gram
 - Dependency, token, dependency
 - Token, dependency, token



"Passive subject of a phosphorylation trigger that is a protein"

- Features are just patterns
- We simply rewrote them using Odin syntax



type : dependency
 label: Phosphorylation
 pattern: |
 trigger:Phosphorylation
 theme:Protein = >nsubjpass

- We know have a decision list classifier. But:
- Feature weights are unbounded continuous values
 - Useful for resolving conflicts
 - But nearly impossible to understand/modify
- Rules still need to vote
 - We normalize and discretize weights ("votes") using Scott's rule (used for the generation of bins in histograms)



- Binding

 PROTEIN recruits PROTEIN
- Localization
 - PROTEIN is recruited to the cytoplasm

```
# vote: +2
- name: Binding_1
label: Binding
type: token
action: countMentions
pattern: |
```

```
[lemma=recruit & tag=/^(V|N|J)/]
```

```
# vote: +1
- name: Localization_1
label: Localization
type: token
action: countMentions
pattern: |
[lemma=recruit & tag=/^(V|N|J)/]
```

Step 3: Model editing

 Two linguists were given the task of improving the generated rules

- Constraints:
 - Only have access to the model. No peeking at the training data
 - Approximately one hour to work on the task

Expert recommendations

- Generalize syntactic patterns, e.g., participants in events may be heads or modifiers of noun phrases
 - E.g.: "K-Ras" or "the K-Ras protein"
- Eliminate trigger rules that were not sufficiently discriminative
- Make rules robust to common parsing mistakes
- (22 total specific recommendations)









Summary

- Converting statistical models into a deterministic decision list classifier does impact performance negatively
- But keeping the human in the loop allows us to recover almost all the lost performance. And we end up with an interpretable model!

Fine. But how do you do this with neural networks that do not have explicit features?

Exploring Interpretability in Event Extraction: Multitask Learning of a Neural Event Classifier and an Explanation Decoder

Zheng Tang *, Gus Hahn-Powell [†], Mihai Surdeanu *

* Department of Computer Science [†] Department of Linguistics University of Arizona {zhengtang, hahnpowell, msurdeanu}@email.arizona.edu Intuition: jointly training for prediction and interpretability!

- 1. Use neural classifiers instead of the LRs in the previous paper ☺
- 2. Decode rules from natural language texts (reusing ideas from machine translation)
 - Source language original texts
 - "Target language" grammar rule that matches
- 3. Train them jointly

Encoder-decoder neural architectures for machine translation



Figure 1: Our model reads an input sentence "ABC" and produces "WXYZ" as the output sentence. The model stops making predictions after outputting the end-of-sentence token. Note that the LSTM reads the input sentence in reverse, because doing so introduces many short term dependencies in the data that make the optimization problem much easier.

I eat cake

Our "source" and "target languages"

(target)

Phosphorylation_syntax_1a_noun:

trigger = [lemma=/phosphorylation/ !word=/(?i)^(de|auto)/]

theme:BioChemicalEntity = prep_of appos? /nn|conj_(and|or|nor)|cc/{,2}



(source)









 $context_1^E$

A

]lstm∣**≮**







Task 1 : Event Extraction

 $context_0^E$

LSTM

 $context_2^E$

LSTM

 $context_3^E$

LSTM

 $context_{4}^{E}$

LSTM



Task 1: Event Extraction

70




Data

- For event classification: BioNLP 2013 (similar to the data used in the previous paper)
- For rule decoding: pairs of (sentence, rule) extracted by our machine reading system with manually written rules
 - Some come from the BioNLP dataset and are aligned with the gold annotations in BioNLP 2013
 - Approximately 15K pairs come from other publications ("silver" data)

Results for event classification

We extract 3 types of events

	Phosphorylation (P)			Localization (L)			Gene Expression (GE)		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F 1
Rule baseline	92.68	48.12	63.35	66.13	44.44	53.16	51.08	69.79	58.98
T1	87.78	49.38	63.20	100.00	4.04	7.77	89.32	64.30	74.77
T1 + Silver	62.75	82.50	71.28	54.55	34.34	42.15	68.43	74.31	71.25
T1 + Silver + T2	84.38	68.75	75.77	76.60	39.39	52.03	69.92	71.24	70.58

- Rule baseline system with manually written rules
- T1 supervised neural event classifier
- T1 + Silver semi-supervised neural event classifier
- T1 + Silver + T2 semi-supervised neural event classifier trained jointly with the rule decoder

Results for event classification

	Phosphorylation (P)			Localization (L)			Gene Expression (GE)		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Rule baseline	92.68	48.12	63.35	66.13	44.44	53.16	51.08	69.79	58.98
T1	87.78	49.38	63.20	100.00	4.04	7.77	89.32	64.30	74.77
T1 + Silver	62.75	82.50	71.28	54.55	34.34	42.15	68.43	74.31	71.25
T1 + Silver + T2	84.38	68.75	75.77	76.60	39.39	52.03	69.92	71.24	70.58
	Semi-supervised learning helps!								

Results for event classification

	Phosphorylation (P)			Localization (L)			Gene Expression (GE)		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Rule baseline	92.68	48.12	63.35	66.13	44.44	53.16	51.08	69.79	58.98
T1	87.78	49.38	63.20	100.00	4.04	7.77	89.32	64.30	74.77
T1 + Silver	62.75	82.50	71.28	54.55	34.34	42.15	68.43	74.31	71.25
T1 + Silver + T2	84.38	68.75	75.77	76.60	39.39	52.03	69.92	71.24	70.58
Jointly training for prediction and interpretability helps prediction!									

Results for rule decoding

	BLEU	Exact Matches	Non-exact Explainable Matches
Р	93.80	86.11	2/15
L	83.78	84.33	1/9
GE	78.99	76.45	10/43

Rule error analysis

Hand-written Rule	Decoded Rule
trigger = [lemma = /phosphorylate/ & ! word = /(?i)^(de auto)/ & tag = /^(V JJ)/ & ! mention = ModificationTrigger]	trigger = [lemma = /phosphorylate/ & ! word = /(?i)^(de auto)/ & tag = /^(V JJ)/ & ! = ModificationTrigger]
theme : BioChemicalEntity = > nsubjpass prep_of ? /conj_(and or nor) nn cc/ { , 2 }	theme : BioChemicalEntity = > nsubjpass prep_of ? /conj_(and or nor) nn cc/ { , 2 }
trigger = [lemma = /detect localiz locat releas secret translocat/ & ! word = /(?i)^de/]	trigger = [lemma = /detect localiz locat releas secret translocat/ & ! word = /(?i)^de/ & ! outgoing = /prep_(by of)/]
theme : BioChemicalEntity = prep_of ? appos ? /conj_(and or nor) cc nn/ { , 2 }	theme : BioChemicalEntity = < /conj_(and or nor)/ ? /conj_(and or nor) cc nn prep_of/ { , 2 }
<pre>trigger = [lemma = / phosphorylation / & ! word = / (? i) ^ (de auto) / & ! outgoing = / prep_(by of) /]</pre>	trigger = [lemma = / phosphorylate / & ! word = / (? i) ^ (de auto) / & tag = / ^ (V JJ) / & ! mention = ModificationTrigger]
theme : BioChemicalEntity = < / conj_ (and or nor) / ? / conj_ (and or nor) cc nn prep_of / { 2 } site : Site ? = nn < dobj ? / prep_ (at on) / num ?	cause : BioChemicalEntity ? = < xcomp ? (nsubj agent < vmod) / appos nn conj_ (and or nor) cc / { 2 }
	theme : BioChemicalEntity = (dobj xcomp) / conj_ (and or nor) dep cc nn prep_of / { 2 } (>> [word = by]) { 2 } site : Site ? = dobj ? / prep_ (at on) nn conj_ (and or nor) cc / { 2 }

Green – missed text in the decoded rule **Red** – hallucinated text

Summary

- A neural approach that jointly trains for prediction and interpretability
- The joint training improves prediction!

Motivation

 It took us 1 – 2 person months to build the grammar of 154 rules for the biomedical domain

- In many cases, one does not have this time
 - Scenario 1: But training data exists
 - Scenario 2: And training data does not exist either...

Lightly-supervised Representation Learning with Global Interpretability

Andrew Zupon, Maria Alexeeva, Marco A. Valenzuela-Escárcega, Ajay Nagesh, and Mihai Surdeanu

University of Arizona, Tucson, AZ, USA

{zupon, alexeeva, marcov, ajaynagesh, msurdeanu}@email.arizona.edu

Architecture and walkthrough example



Traditional rule bootstrapping algorithm for entity classification



Learning to read

We want to combine:

- The advantages of representation learning, aka "word embeddings"
 - Neural network language models handle unsupervised data well
- The interpretability of our current approach
 - Produce patterns in the end
- Keep the human in the loop, but minimally
 - Just a few examples

Starting point: distributional hypothesis

- Distributional hypothesis:
 - By looking at a word's context, one can infer its meaning (Harris, 1954)
 - You shall know a word by the company it keeps (Firth, 1957)

Example

- tasty X
- X with butter
- X and coffee
- greasy X

Example

- tasty X
- X with butter
- X and coffee
- greasy X



Starting point: word2vec, skip-gram

You shall know a word by the company it keeps. - Firth, 1957

You shall know the company by the word it keeps.

- Word2vec, skip-gram

Starting point: word2vec, skip-gram



P(city|Tucson) ++ P(place|Tucson) ++

... P(dog|Tucson) - -



Two important changes

- We will learn embeddings for *both named entities and patterns*
 - An entity's context is defined by the patterns that match it
- Added supervision in the objective function, to incorporate human-provided information
 - We have "seed" names in each category

Context as patterns

The city of Tucson is the place to vacation !



Objective function



Unsupervised, directly "inherited" from skip-gram

Light supervision from a few seed examples, iteratively expanded

EmBoot

SGD on the previous objective function



Entities



Promote entities closest to the seeds

Visualization of the training procedure

- On the CoNLL-2003 dataset
 - PER = purple
 - LOC = blue
 - ORG = green
 - MISC = red
- Human contribution: seed set with 10 entities in each category
 - PER: Clinton, Dole, ...
 - LOC: U.S., Germany, ...
 - ORG: Reuters, U.N., ...
 - MISC: Russian, German, ...






























Results, CoNLL dataset (4 classes)



Results, OntoNotes (11 classes)



Summary

- A lightly-supervised approach that jointly learns representations for entities and patterns that extract them
- State-of-the-art results for semi-supervised learning
- The rules can be edited by domain experts, and this leads to further improvements in performance (not shown)

Take-home message

- For large, inter-disciplinary projects we need to move beyond "black-box" methods to approaches that produce globally interpretable models
- We can produce such interpretable models using deep learning (best of both worlds?)

Bibliography

- 1. M.A. Valenzuela-Escarcega, G. Hahn-Powell, T. Hicks, and M. Surdeanu. A Domain-independent Rule-based Framework for Event Extraction. ACL-IJCNLP, 2015.
- 2. M. A. Valenzuela-Escarcega, G. Hahn-Powell, and M. Surdeanu. *Description of the Odin Event Extraction Framework and Rule Language*. <u>arXiv:1509.07513</u>, 2015.
- 3. D. Bell, G. Hahn-Powell, M. A. Valenzuela-Escarcega, and M. Surdeanu. *An Investigation of Coreference Phenomena in the Biomedical Domain*. LREC, 2016.
- 4. M. A. Valenzuela-Escarcega, G. Hahn-Powell, and M. Surdeanu. *Odin's Runes: A Rule Language for Information Extraction*. LREC, 2016.
- 5. G. Hahn-Powell, D. Bell, M. A. Valenzuela-Escarcega, and M. Surdeanu. *This before That: Causal Precedence in the Biomedical Domain*. BioNLP, 2016.
- 6. Large-scale Automated Machine Reading Discovers New Cancer Driving Mechanisms. Valenzuela-Escarcega, M. A.; Babur, O.; Hahn-Powell, G.; Bell, D.; Hicks, T.; Noriega-Atala, E.; Wang, X.; Surdeanu, M.; Demir, E.; and Morrison, C. T. Database: The Journal of Biological Databases and Curation. 2018.
- 7. Lightly Supervised Representation Learning with Global Interpretability. Zupon, A.; Alexeeva, M.; Valenzuela-Escarcega, M. A.; Nagesh, A.; and Surdeanu, M. In Proceedings of the 3rd Workshop on Structured Prediction for Natural Language Processing, 2019.
- 8. Z. Tang, M.A. Valenzuela-Escarcega, and M. Surdeanu. Exploring Interpretability in Event Extraction: Multitask Learning of a Neural Event Classifier and an Explanation Decoder. In preparation.

Many thanks to my collaborators!





















THANK YOU! QUESTIONS?

Conflict of interest disclosure

 M. Surdeanu discloses a financial interest in Lum.ai. This interest has been disclosed to the University of Arizona Institutional Review Committee and is being managed in accordance with its conflict of interest policies.