User-Centric Ontology Population

KENNETH CLARKSON, **ANNA LISA GENTILE**, DANIEL GRUHL, PETAR RISTOSKI, JOSEPH TERDIMAN, STEVE WELCH

IBM RESEARCH

Motivation

"The Semantic Web provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries" - Tim Berners-Lee 2001

17 years later still vast amount of valuable unstructured and semi-structured data is published on the Web

Goal: automatically extract semantic data from text



Input Text Corpus

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis pharetra, erat ac ante, tristique ultricies erat eros nec turpis. Maecenas eu libero vel mi consequat sit amet, consectetur adipiscing elit. Quisque dapibus aliquam quam, quis laoree sollicitudin dictum. Nam quam diam, malesuada at consectetur sit amet, varius quam. Suspendisse consectetur iaculis ornare. Quisque augue ligula, lacinia eg pharetra ipsum consectetur orci fringilla tempor. Donec at velit id eros suscipit ar congue. Cras ornare ligula quis lacus dignissim fermentum rhoncus nec mai tristique. Proin lacus felis, auctor ac sollicitudin a, elementum in ante. Maecena eleifend nec. Donec nisi neque, auctor in vestibulum eget, condimentum at erat. C per conubia nostra, per inceptos himenaeos.

Maecenas aliquam adipiscing metus ut ultrices. Praesent quam velit, dictum nor est a leo aliquet pretium. Suspendisse odio purus, iaculis at venenatis id, semper consequat. Suspendisse vitae elit ipsum, non vestibulum ligula. Aenean interd pellentesque mi, in convallis quam erat eget nisl. Nam ante magna, interdum vel e portitor urna ac facilisis. Vivamus malesuada sapien at tortor viverra adipi pellentesque quis, viverra nec tellus. Sed accumsan purus at dolor mollis blandi aliquam eget arcu. Nulla facilisi.

Target Ontology





Entity Detection

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis pharetra, erat ac at eros nec turpis. Maecenas eu libero vel mi conseguat ante, tristique u sit amet, consected adipiscing elit. Quisque dapibus aliquam quam, quis laoree sollicitudin dictum. Nam quam diam, malesuada at consectetur sit amet, varius quam. Suspendisse consectetur iaculis ornare. Quisque aug lacinia eq pharetra ipsum consectetur orci fringilla tempor. Donec at velit and auscipit an congu ornare ligula quis lacus dignissim fermentum rhoncus nec mai tristique. acus felis, au sollicitudin a, elementum i */aecena eleifend nec. Donec nisi neguc, sacco, in vestibulum eget, condir t erat. C per conubia nostra, per inceptos himenaeos.

Maecenas aliquam adipiscing metus ut ultrices. Praesent auam velit, dictum nor est a leo aliquet pretium. Suspendisse odio purus, iaculi consequat. Suspendisse vitae elit ipsum, non vestibulum ligula. Aenean interd pellentesoue mi. in convallis quam erat eget nisl. Nam ante magna, interdum vel e portitor facilisis. Vivamu pellentesque quis, viverra nec tellu aliquam eget arcu. Nulla facilisi.

Target Ontology





User-defined Entity Grouping



Target Ontology





Ontology Alignment Target Ontology



User concepts

Back and neck pain and tightness, After visiting the chiropractor and message therapist with no relief, I miraculously made the connection that it was the Ambien and not stress that was causing my pain I have been off the med. for less than a week and feel great. The back painis completely gone.

No side effects. I have been on it long term, although I do stop taking it occasionally. It has been a blessing for me. Always wake up feeling refreshed.

Next day fatigue, bloating, stomach pain, memory loss, racing heart, no sleep without it!

Memory loss, confusion, sallucinations, headache. Thought I was asleep, actually made multiple phone calls and ended up driving on freeway. Remember bits a pieces of night such as the feeling as though I wasn't behind wheel of my car but knew I was on the freeway. White lines were lifting up and going multiple directions, seeing double, vivid colors, dog moving on my ph picture and only GOD knows what else. I cry now every time I think about danger I put myself and others in. So now I have anxiety and panic when I sleep or take meds- which my doc gv me clonazepamon Day 2. Ambien should be removed off market and def not prescribed to anyone!





Ontology Population







Knowledge Discovery Process (Fayyad et al. 1996)





How good are machines?

~80% accuracy

~85% accuracy

~90% accuracy



Is 80% enough?





Introduce the human-in-the-loop





Introduce the human-in-the-loop

"Computers are incredibly fast, accurate, and stupid. Human beings are incredibly slow, inaccurate, and brilliant. Together they are powerful beyond imagination."

Einstein never said that





Proposed Solution

Given initial user conceptualization, the methodology supports:

- Finding candidate ontologies
- **Aligning** the user's conceptualization to a target ontologies
 - novel hierarchical classification approach
- Maintenance lifecycle
 - **build** (create new concepts)
 - change (splitting/merging concept)
 - grow (adding new instances to each concept)
 - from target ontologies
 - new facts extracted from unstructured data



Implementation





Aligning Input with a Target Ontology

Identify available ontologies

collective instance matching

Align user conceptualization using ML models

- Training data: instances of the target ontology
- Features : domain-specific word embeddings
- Classification strategies
 - Flat hierarchical classification
 - Top-down local classifier per parent node
 - Combine flat hierarchical with top-down local classifier per parent node
- Classifiers
 - SVM, Random Forests, Logistic Regression, Convolutional Neural Network



Flat Hierarchical Model

One model for each level of the hierarchy

- Simple
- High model complexity down the hierarchy



IBM

Top-Down Local Classifiers

One model for each parent in the hierarchy

- Simple
- Error propagation through levels



IBM

Combine Both Models

Combine flat hierarchical models with top-down local classifier

- flat model for level L-1
- local model for level L





Ontology Maintenance

Adding new instances

• Use existing models

Reassigning Instances

Leave-one-out validation

Generating new concepts

• If the class distribution is uniform then search for new concept

Merging concepts

• User's concepts aligned to the same target ontology concept should be merged

Concept splitting

- Use hierarchical clustering
- Refine until a criteria is met



Evaluation - Alignment

Task: label adverse drug events with preferred medical terms

Data:

- MedDRA ontology as a target ontology
- ADE groups extracted from "ask a patient blogs"

	User's conceptualization	MedDRA
#level1	17	27
#level2	62	304
#level3	106	1,444
#level4	169	20,935
#Instances	3,262	95,061



Evaluation metric: HITS@10 = proportion of correct mapping top 10 ranked suggestions

• Evaluate per each level of the hierarchy



Evaluation – Ontology Alignment

Baselines:

- String-based average-link matching
- Word embeddings
- LDA topic modeling

Evaluation metric: HITS@10

- proportion of correct mappings that appear in the top 10 ranked suggestions
- Evaluate per each level of the hierarchy



Results



IBM

10

Results: Level 4



IBM

Evaluation – Ontology Maintenance

Adding new instances – evaluate how precise the models can add new instances to the already aligned concepts

- Retrieved 298 **new** ADE from askapatient.com
- Measure HITS@k for each level of the hierarchy





Evaluation – Ontology Maintenance

Adding new concepts

- model's ability to suggest the user to add a new concept
- evaluation
 - 500 MedDRA instances that don't belong to the user's conceptualization (positive instances)
 - 500 instances that belong to the user's conceptualization

Results:

- Precision: 73.8%
- Recall: 84.6%
- F-score: 78.83%



Evaluation – Ontology Maintenance

Re-assigning Instances: evaluate the model's ability to reassign instances to other concepts.

The model identified 82 instances to be reassigned, from which 67 (81.7%) were accepted by the medical doctor

Examples:

- User errors: "stomach aches" was assigned to "Emotional disorder", which should be assigned to "Abdominal distension"
- Better matches: "sensitivity to light" was assigned to "Visual impairment", which was later reassigned to "Photophobia"



Further Use-Cases

- Maintain health and medical data
 - Adverse drug reactions
 - Drug brands
- Maintain e-shop product catalog and taxonomy
 - Map new features to an existing product catalog
 - Map new products in the product taxonomy
- Social media analysis
 - Identifying new trends
- Reviews analysis
 - Movies and actors



Conclusion

User-centric ontology population

Human-in-the-loop for each step

 Building, connecting and maintaining their conceptualization, using available ontologies

Novel hierarchical classification model

dynamically refined based on user interaction

The approach supports the user to achieve **nearly perfect performance**

The user has full control on their level of involvement in the process

 Trade-off between involvement/cost/time and performance/quality of results



User-Centric Ontology Population



Kenneth Clarkson, Anna Lisa Gentile, Daniel Gruhl, Petar Ristoski, Joseph Terdiman, Steve Welch





