Dogle Learning on the Web

Fernando Pereira

Thanks to

Koby Crammer, Kuzman Ganchev, João Graça, Yi Liu, Marius Pasça, Franz Och, Joseph Reisinger, Deepak Ravichandran, Stefan Riezler, Partha Talukdar, Ben Taskar, Alexander Vasserman, and other colleagues at Google and University of Pennsylvania. Penn work funded by NSF



Web Inference and Learning

- Building and using the Web requires algorithms that can draw reliable inferences from the wealth of evidence implicit in Web data
 - How to interpret this term in this context?
 - Does this sentence answer that question?
 - Will this user click on that ad?
- Learning: create concise representations to support good inferences



Beyond Supervised Learning

- Explicit annotation is costly and misleading
 - Expert annotation: difficult to agree on annotation criteria, repeated failures to achieve self-sustainability:
 - User annotation: what's its immediate value to the user, attractive to spammers
- Learning from what users do: we constantly seek and organize information
 - Examples: machine translation, query expansion
- Lots of (mostly) unlabeled data



Learning from Parallel Data



Statistical Machine Translation

- Inference: is this phrase a good translation of that phrase in this context?
- Indirect evidence: translation pairs, monolingual text
- Memory-based (non-parametric) learning:
 - Bilingual: phrase translation table
 - Monolingual: language model



Statistical Machine Translation

Machine learning from human communication

• Parallel texts:

SEHR GEEHRTER GAST! KUNST, KULTUR UND KOMFORT IM HERZEN BERLIN.

DEAR GUESTS, ART, CULTURE AND LUXURY IN THE HEART OF BERLIN.

DIE ÖRTLICHE NETZSPANNUNG BETRÄGT 220/240 VOLT BEI 50 HERTZ. THE LOCAL VOLTAGE IS 220/240 VOLTS 50 HZ.

ΕN

Building a Translator

• Align (cf. genomic alignments...)



- Extract phrase pairs from alignments
- Model the statistics of the target language
- To translate a source text:
 - search over phrase-to-phrase translations
 - filter with model of target language 🥥



Meaning from Web Search

- Term mismatches between queries and documents in information retrieval
 - Query terms and relevant document terms differ
- Query expansion: add search terms known to occur in relevant documents
 - Increase recall
 - Decrease query term ambiguity



Query Expansion by Translation

- Training
 - Align queries and clicked snippets
 - Create translation tables
 - Train *n*-gram language model from query logs
- Use
 - "Translate" many queries
 - Extract and store table of "translated" terms in context



From query to snippets





Ambiguity Resolution

• 5-best phrase-level translations

(herbs,herbs)(for,for)(chronic,chronic)(heartburn, heartburn) (herbs,herb)(for,for)(chronic,chronic)(heartburn, heartburn) (herbs,remedies)(for,for)(chronic,chronic)(heartburn, heartburn) (herbs,medicine)(for,for)(chronic,chronic)(heartburn, heartburn) (herbs,supplements)(for,for)(chronic,chronic)(heartburn, heartburn)

(herbs,herbs)(for,for)(mexican,mexican)(cooking,cooking) (herbs,herbs)(for,for))(cooking,cooking)(mexican,mexican (herbs,herbs)(for,for)(mexican,mexican)(cooking,food) (mexican,mexican)(herbs,herbs)(for,for)(cooking,cooking) (herbs,spices)(for,for)(mexican,mexican)(cooking,cooking)



Beyond Parallel Corpora



Reading the Web



- Elementary semantic inference
- First step: what are the possible classes for each instance?

Someone Told Us

• Text patterns (Hearst 92, Van Durme & Pasça 08)



Informative Co-occurrences

- WebTables (Cafarella et al. 08)
 - 154M HTML tables from Web pages
 - Cluster instances in table columns





Combining Information Sources

- Bootstrapping:
 - Seed set of (instance, class) pairs
 - Compute instance similarity from additional sources
 - Use similarity to infer new (instance, class) pairs
- Approach: label propagation in a graph
 - instance nodes
 - cluster/class nodes
 - bipartite structure





















































Iteration 2

924k (class, instance) pairs extracted from 100M web documents

A8 Propagation WebTables

74M (class, instance) pairs extracted from WebTables

$$MRR = \frac{1}{|\text{test-set}|} \sum_{v \in \text{test-set}} \frac{1}{\text{rank}_v(\text{class}(v))}$$

• Instances found solely by label propagation:

Class	Precision at 100	
	$\underbrace{\mathbf{A8} \text{ extractions}}_{\mathbf{A8}}$	
Book Publishers	UNIVERSITY OF PENNSYLVANIA 87.36	
Federal Agencies	29.89	
NFL Players	94.95	
Scientific Journals	90.82	
Mammal Species	84.27	

Scientific Journals	Journal of Physics, Nature, Structural and Molecular
	Biology, Sciences Sociales et santé, Kidney and Blood
	Pressure Research, American Journal of Physiology–
	Cell Physiology
NFL Players	Tony Gonzales, Thabiti Davis, Taylor Stubblefield,
	Ron Dixon, Rodney Hannah
Book Publishers	Small Night Shade Books, House of Anansi Press,
	Highwater Books, Distributed Art Publishers, Copper
	Canyon Press

• Label propagation by-product:

Seed Class	Non-Seed Class Labels Discovered
Book Publishers	small presses, journal publishers, educational pub-
	lishers, academic publishers, commercial publishers
NFL Players	sports figures, football greats, football players, backs,
	quarterbacks
Scientific Journals	prestigious journals, peer-reviewed journals, refereed
	journals, scholarly journals, academic journals

Semantic Constraints for Better Classes

Semantic Constraints for Better Classes

Semantic Constraints for Better Classes

Experiments with Public Sources

- Make Talukdar *et al.* results easy to reproduce and extend
 - Freebase: multiple-sourced relational tables
 - Pantel et al. 09 gold-standard hypernyms
 - TextRunner (U. of Washington): hypernyms from open-domain extraction
 - YAGO (Suchanek *et al.*, 07): entityattribute knowledge base curated from Wikipedia and Wordnet

 Qualitative evidence that attribute nodes propagate the right information

YAGO	Top-2 WordNet Classes Assigned by MAD
Attribute	(example instances for each class are shown in brackets)
has_currency	wordnet_country_108544813 (Burma, Afghanistan)
	wordnet_region_108630039 (Aosta Valley, Southern Flinders Ranges)
works_at	wordnet_scientist_110560637 (Aage Niels Bohr, Adi Shamir)
	wordnet_person_100007846 (Catherine Cornelius, Jamie White)
has_capital	wordnet_state_108654360 (Agusan del Norte, Bali)
	wordnet_region_108630039 (Aosta Valley, Southern Flinders Ranges)
born_in	wordnet_boxer_109870208 (George Chuvalo, Fernando Montiel)
	wordnet_chancellor_109906986 (Godon Brown, Bill Bryson)
has_isbn	wordnet_book_106410904 (Past Imperfect, Berlin Diary)
	wordnet_magazine_106595351 (Railway Age, Investors Chronicle)

Propagation Objective

• MAD [Talukdar & Crammer 09] (simplified)

$$\arg\min_{\hat{\boldsymbol{Y}}} \sum_{l=1}^{m+1} \left[\| \boldsymbol{S} \hat{\boldsymbol{Y}}_{l} - \boldsymbol{S} \boldsymbol{Y}_{l} \|^{2} + \mu_{1} \sum_{u,v} \boldsymbol{M}_{uv} (\hat{\boldsymbol{Y}}_{ul} - \hat{\boldsymbol{Y}}_{vl})^{2} + \mu_{2} \| \hat{\boldsymbol{Y}}_{l} - \boldsymbol{R}_{l} \|^{2} \right]$$

- m labels, +1 dummy label
- $M = W^{\top} + W$ is the symmetrized weight matrix
- $\hat{\boldsymbol{Y}}_{vl}$: weight of label l on node v
- \mathbf{Y}_{vl} : seed weight for label l on node v
- S: diagonal matrix, nonzero for seed nodes
- \mathbf{R}_{vl} : regularization target for label l on node v

A Propagation Algorithm

Inputs $\boldsymbol{Y}, \boldsymbol{R} : |V| \times (|L|+1), \, \boldsymbol{W} : |V| \times |V|, \, \boldsymbol{S} : |V| \times |V|$ diagonal $\hat{\boldsymbol{Y}} \leftarrow \boldsymbol{Y}$ $\boldsymbol{M} = \boldsymbol{W} + \boldsymbol{W}^{\top}$ $Z_v \leftarrow \boldsymbol{S}_{vv} + \mu_1 \sum_{u \neq v} \boldsymbol{M}_{vu} + \mu_2 \quad \forall v \in V$ repeat for all $v \in V$ do $\hat{\boldsymbol{Y}}_v \leftarrow \frac{1}{Z_v} \left((\boldsymbol{S}\boldsymbol{Y})_v + \mu_1 \boldsymbol{M}_v. \hat{\boldsymbol{Y}} + \mu_2 \boldsymbol{R}_v \right)$ end for

until convergence

- Some details of the construction of the input matrices omitted for simplicity
- Converges under reasonable assumptions
- Many variants, alternative objectives (Subramanya and Bilmes 2008)

Good News

- Label propagation can combine multiple information sources effectively
- Useful coverage of class-instance relations, much bigger than in previous work
- Embarrassingly parallel algorithms
- Graph representation can encode a variety of useful constraints

Limitations

- Can't express "few classes per instance"
- How to classify instances in context?
 - Whistler paintings vs Whistler skiers
- Propagation is additive, averaging
 - Algorithmically nice (convexity, convergence)
 - But it can't "push back" to express incompatibilities between classes
 - artist ⊃ painter ⊄ ski-resort

Few Classes Per Instance

- If a classification is to be informative, it must have limited ambiguity
- Posterior regularization: constrain the ambiguity of final labeling
- POS tagging pilot (Graça et al. 09): motivated by this work, but easier to test
- Penalize the ℓ_1/ℓ_∞ norm of the posterior distribution of classes given instances

Contextual Classification

- Graph-based approach:
 - One node per mention (lots of nodes!)
 - Link mention nodes that have similar contexts
- Language modeling approach:
 - Class labels are also terms
 - Model probability of class given context
 - Find most likely class for each mention given its context

Back to Supervised Learning

- Where do edge weights come from?
- These experiments: heuristic scoring functions drawn from language modeling and information retrieval
- Can we learn edge scoring functions?
 - Minkov & Cohen: learning from random walks
 - McCallum's group: learning factor potentials by M-H sampling
 - Learn from user feedback (Talukdar et al., SIGMOD 2010)

Summary

- First steps in inferring broad-coverage semantic relationships from the actions of Web users:
 - What they write
 - How they interact with search results
- Multiple correlated sources provide a wealth of indirect supervision
- (Some) graph-based algorithms scale up
- Related work:
 - Wang & Cohen's SEAL and follow-ups
 - Weikum lab's SOFIE

Current Work

- Web-scale contextual semantic annotation
 - is-a, co-reference
- Combining multiple relationships
 - The distinct senses of "Whistler" belong to disjoint co-reference classes
- Probabilistic interpretations
 - Relational factor graphs
 - Non-parametric Bayesian models

