
Recursive neural networks for semantic interpretation

Sam Bowman

Department of Linguistics and NLP Group
Stanford University

*with help from Chris Manning, Chris Potts, Richard Socher,
Jeffrey Pennington, J.T. Chipman*

Recent progress on deep learning

Neural network models are starting to seem pretty good at capturing aspects of meaning.

From Stanford NLP alone:

- Sentiment (EMNLP '11, EMNLP '12, EMNLP '13)
 - Paraphrase detection (NIPS '11)
 - Knowledge base completion (NIPS '13, ICLR '13)
 - Word–word translation (EMNLP '13)
 - Parse evaluation (NIPS '10, NAACL '12, ACL '13)
 - Image labelling (ICLR '13)
-

Recent progress on deep learning

Wired, Jan 2014:

Where will this next generation of researchers take the deep learning movement? The big potential lies in deciphering the words we post to the web — the status updates and the tweets and instant messages and the comments — and there's enough of that to keep companies like Facebook, Google, and Yahoo busy for an awfully long time.

Today

Can these techniques learn models for general purpose NLU?

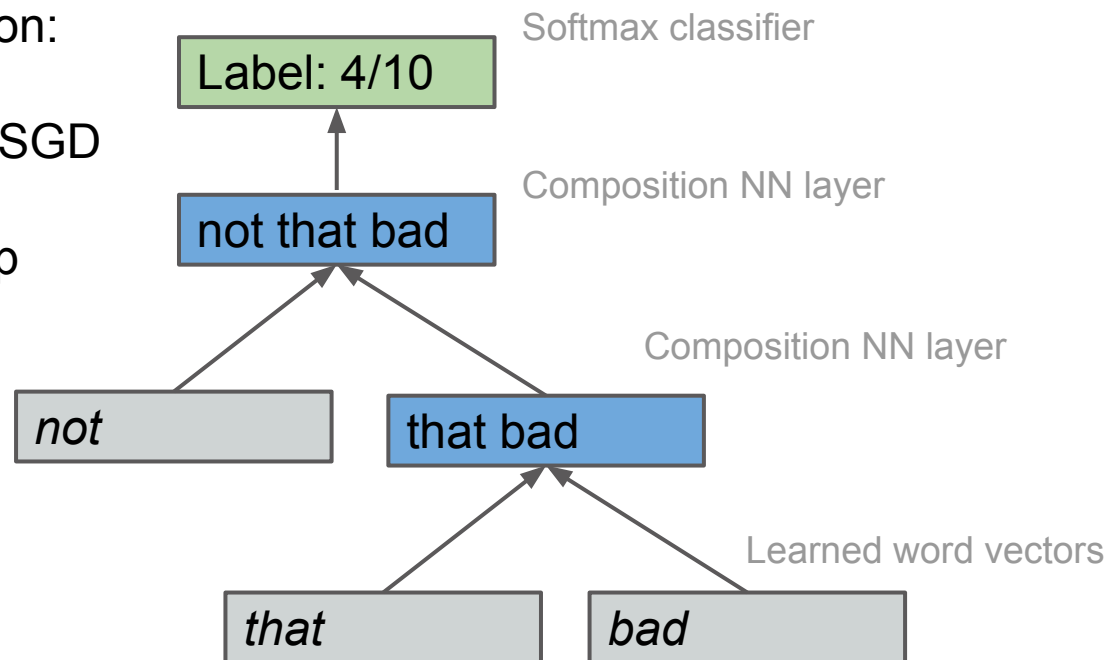
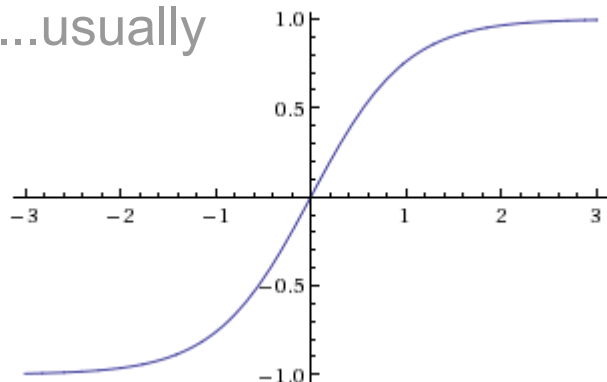
- **Survey: Deep learning models for NLU**
 - Experiment: Can RNTNs learn to reason with quantifiers (in an ideal world)?
 - Experiment: Can RNTNs learn the natural logic *join* operator?
 - Experiment: How do these models do on a challenge dataset?
-

Recursive neural networks for text

- Words and constituents are ~50 dimensional vectors.
- RNN composition function:
 $y = f(Mx + b)$
- Optimize with AdaGrad SGD or L-BFGS
- Gradients from backprop (through structure)

$$f(x) = \tanh(x)$$

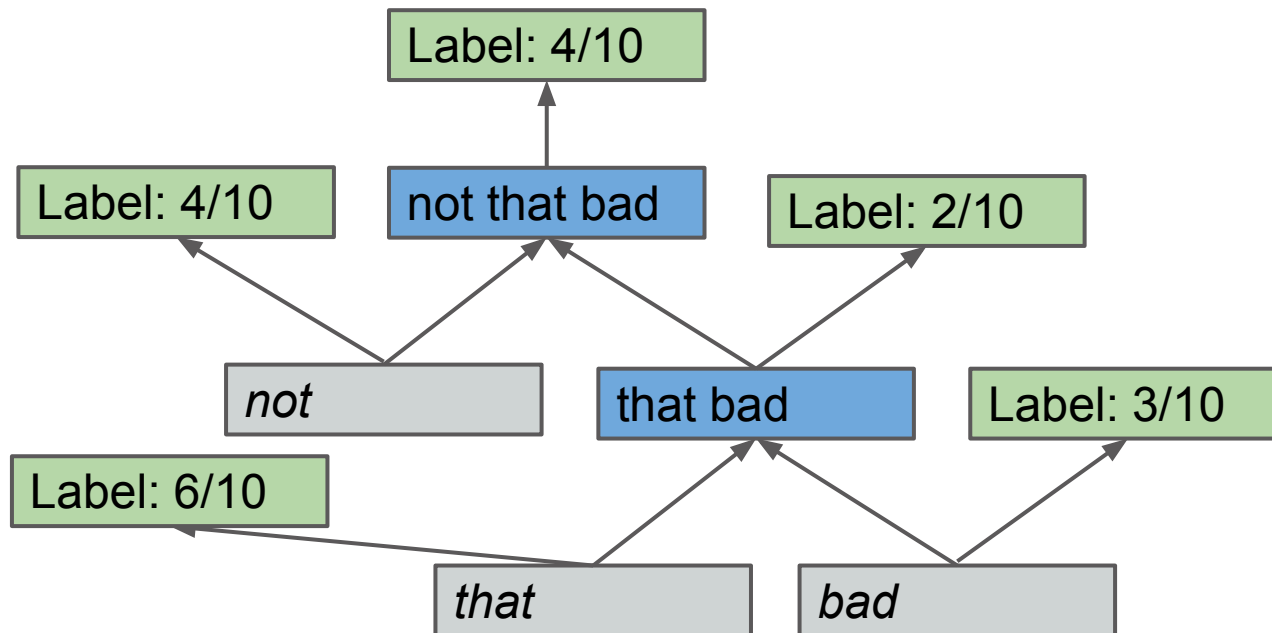
...usually



Recursive neural networks for text

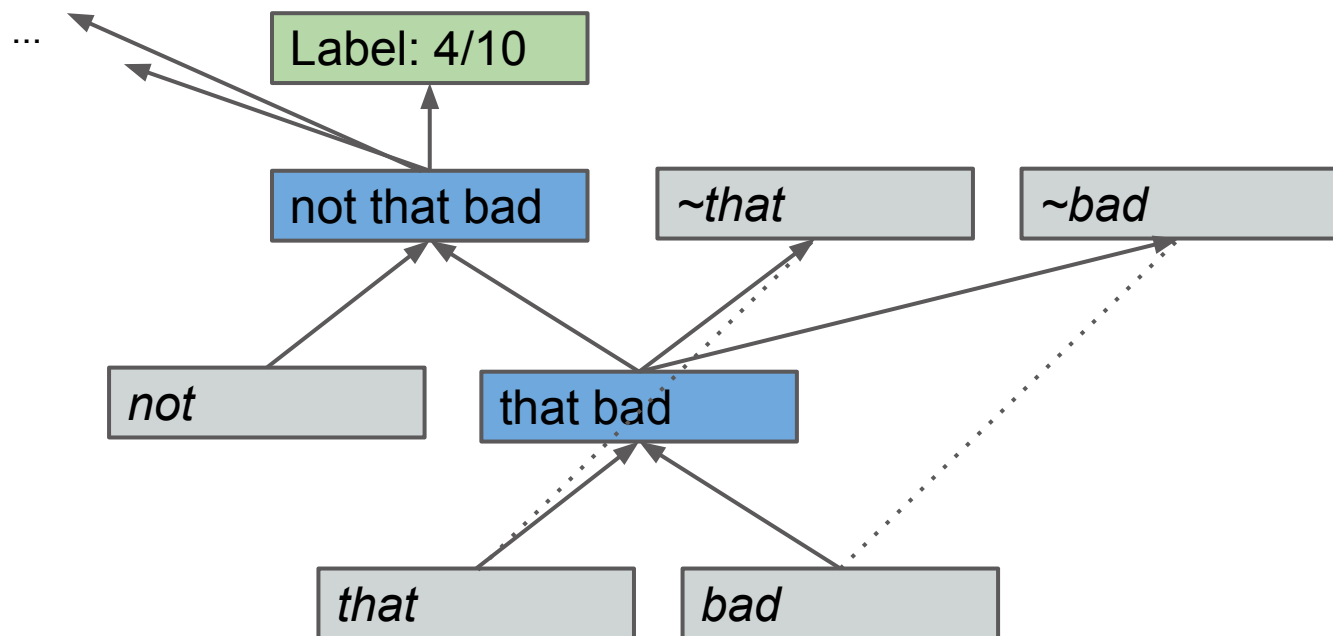
Supervision for everyone!

- ~10k sentences
- ~200k sentiment labels from mechanical Turk



Recursive neural networks for text

- Recursive autoencoder
- Two objectives: Classification and reconstruction

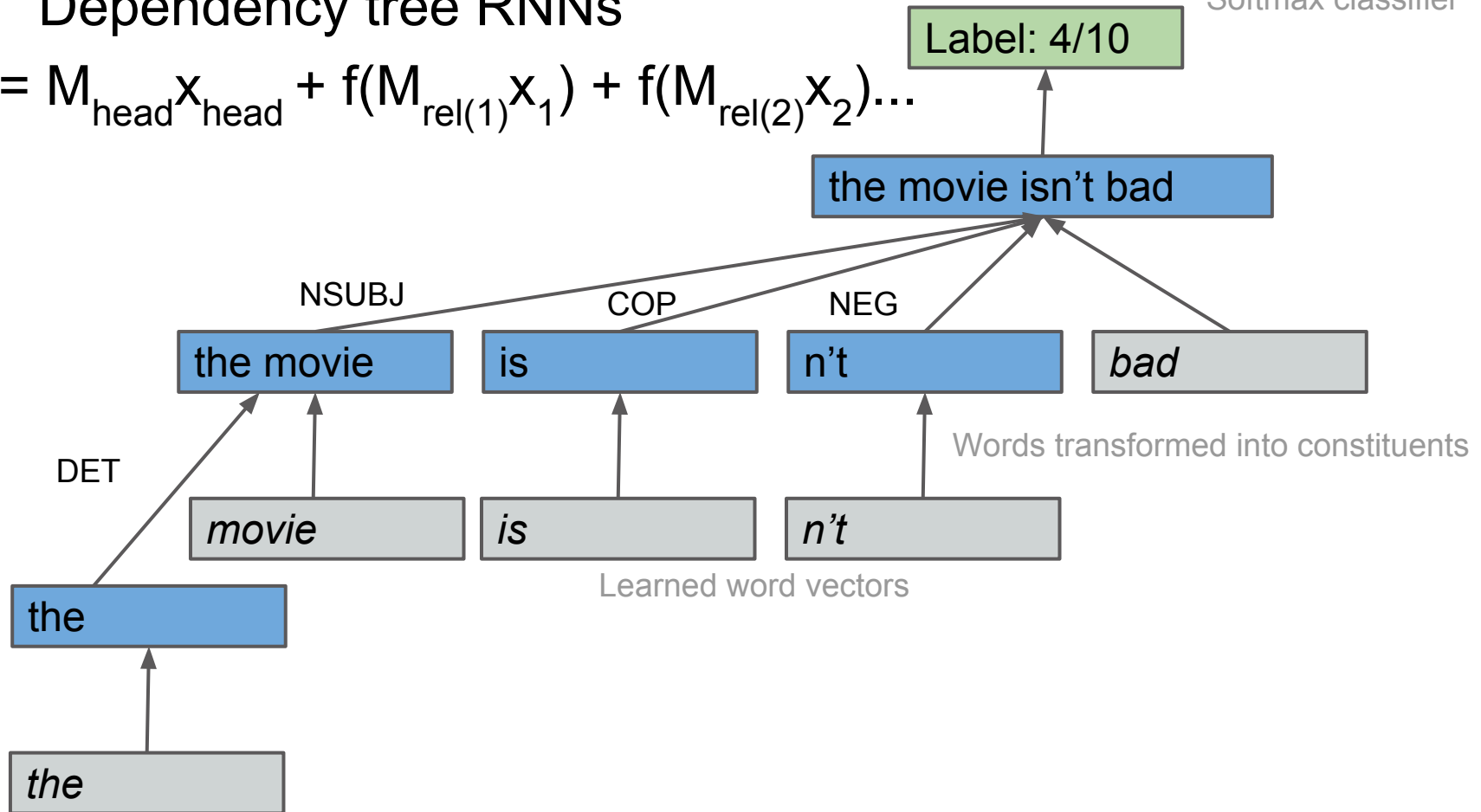


Recursive neural networks for text

- Dependency tree RNNs

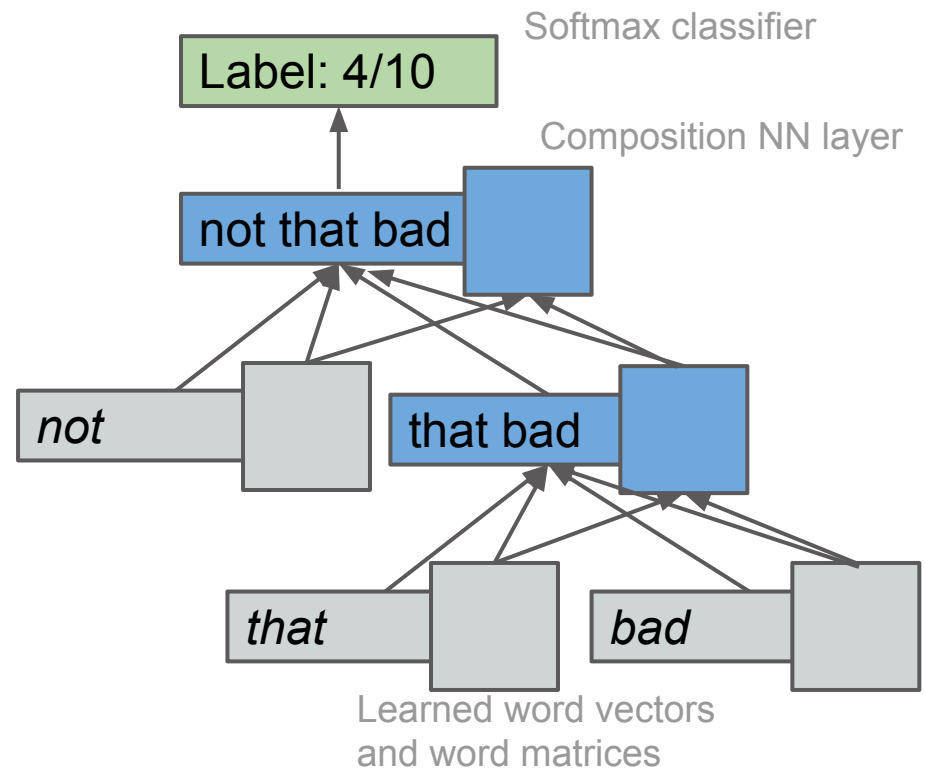
$$y = M_{\text{head}} x_{\text{head}} + f(M_{\text{rel}(1)} x_1) + f(M_{\text{rel}(2)} x_2) \dots$$

Softmax classifier



Recursive neural networks for text

- Matrix-vector RNN
composition functions:
 $y = f(M_v[Ba; Ab])$
 $Y = M_m[A; B]$

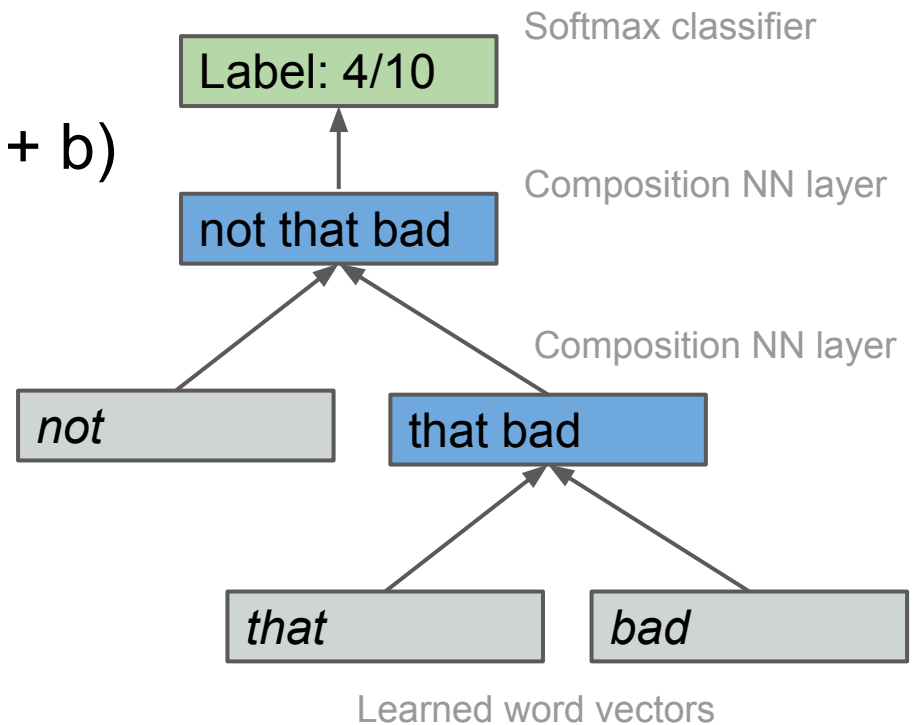


Recursive neural networks for text

- Recursive neural tensor network composition

function:

$$y = f(x_1 M^{[1 \dots N]} x_2 + Mx + b)$$



Recursive neural networks for text

And more:

- Convolutional RNNs (Kalchbrenner, Grefenstette, and Blunsom 2014)
- Bilingual objectives (Hermann and Blunsom 2014)

...

And this isn't even considering model structures for language modeling or speech recognition...

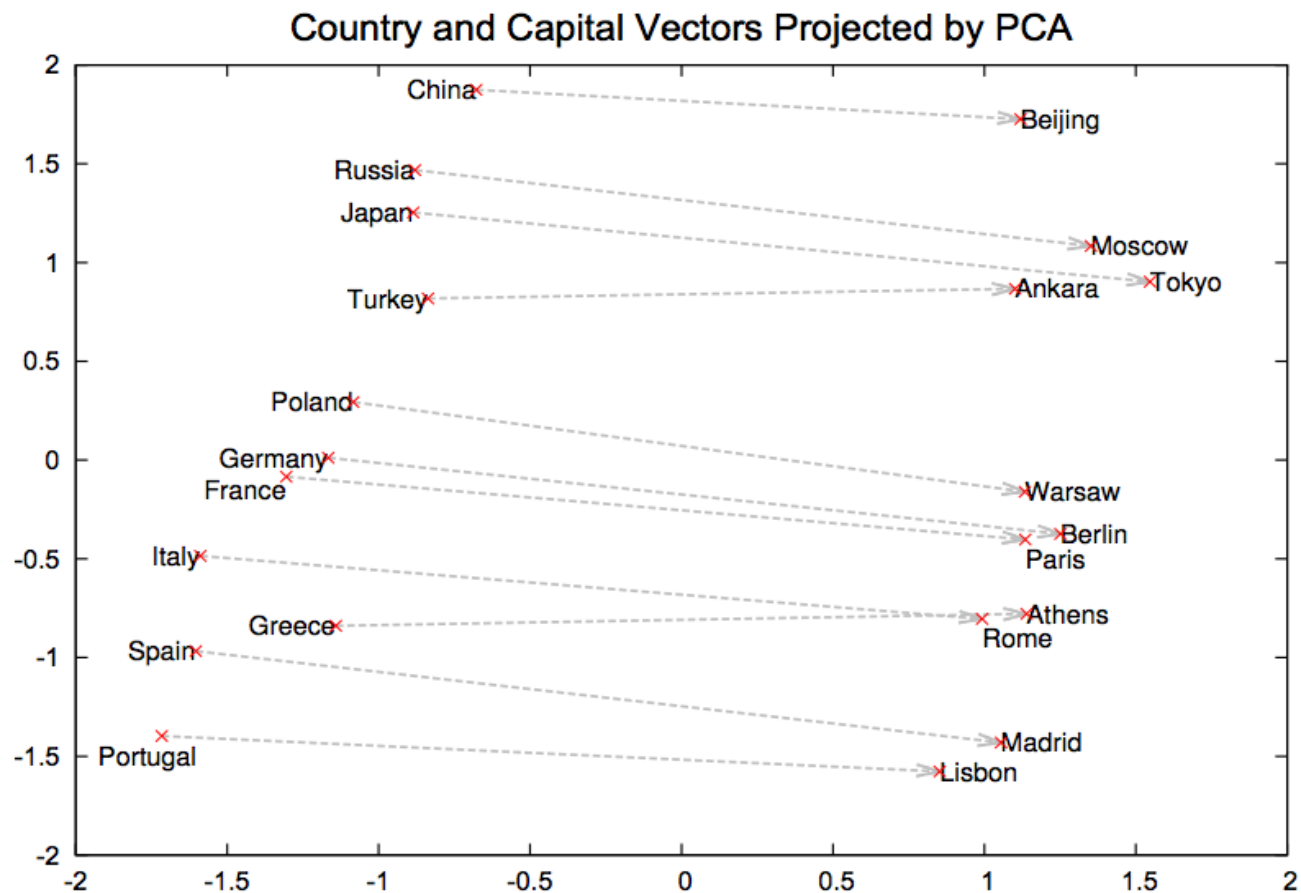
Today

Can these techniques learn models for general purpose NLU?

- Survey: Deep learning models for NLU
 - **Experiment: Can RNTNs learn to reason with quantifiers (in an ideal world)?**
 - Experiment: Can RNTNs learn the natural logic *join* operator?
 - Experiment: How do these models do on a challenge dataset?
-

The problem

Mikolov et al. 2013, NIPS



The problem

The Mikolov et al. result:

- Paris - France + Spain = Madrid
 - Paris - France + USA = ?
 - most - some + all = ?
 - not = ?
-

The problem

- Relatively little work to date on the expressive power of this kind of model.
 - The goal of the project:
 - Can the representation learning systems used in practice capture every aspect of meaning that formal semantics says language users need?
 - This talk:
 - Can RNNs learn to accurately reason with quantification and monotonicity?
-





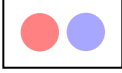


Strict unambiguous NLI

- Hard to test on `world` \leftrightarrow *sentence*. (Why?)
 - What about *sentence* \leftrightarrow *sentence*?
 - Natural language inference (NLI):
 - Doing logical inference where the logical formulae are represented using natural language.
 - (as formalized for NLP here by MacCartney, '09)
 - Framed as classification task:
 - All dogs bark and Fido is a dog. \square Fido barks.
 - No dog barks. \equiv All dogs don't bark.
 - No dog barks. $?$ Some dog barks.
-

Strict unambiguous NLI

- MacCartney's seven possible relations between phrases/sentences:

Slide from Bill MacCartney

	$x \equiv y$	equivalence	<i>couch</i> \equiv <i>sofa</i>
	$x \sqsubset y$	forward entailment (strict)	<i>crow</i> \sqsubset <i>bird</i>
	$x \supset y$	reverse entailment (strict)	<i>European</i> \supset <i>French</i>
	$x \wedge y$	negation (exhaustive exclusion)	<i>human</i> \wedge <i>nonhuman</i>
	$x \mid y$	alternation (non-exhaustive exclusion)	<i>cat</i> \mid <i>dog</i>
	$x \smile y$	cover (exhaustive non-exclusion)	<i>animal</i> \smile <i>nonhuman</i>
	$x \# y$	independence	<i>hungry</i> $\#$ <i>hippo</i>

Monotonicity (a quick reminder)

- A way of using lexical knowledge to reason about sentences.
 - Given: black dogs \sqsubset dogs, dogs \sqsubset animals
 - Upward monotone:
 - *some dogs bark \sqsubset some animals bark*
 - Downward monotone:
 - *all dogs bark \sqsubset all black dogs bark*
 - *Non-monotone:*
 - *most dogs bark $\#$ most animals bark*
 - *most dogs bark $\#$ most black dogs bark*
-

Strict unambiguous NLI

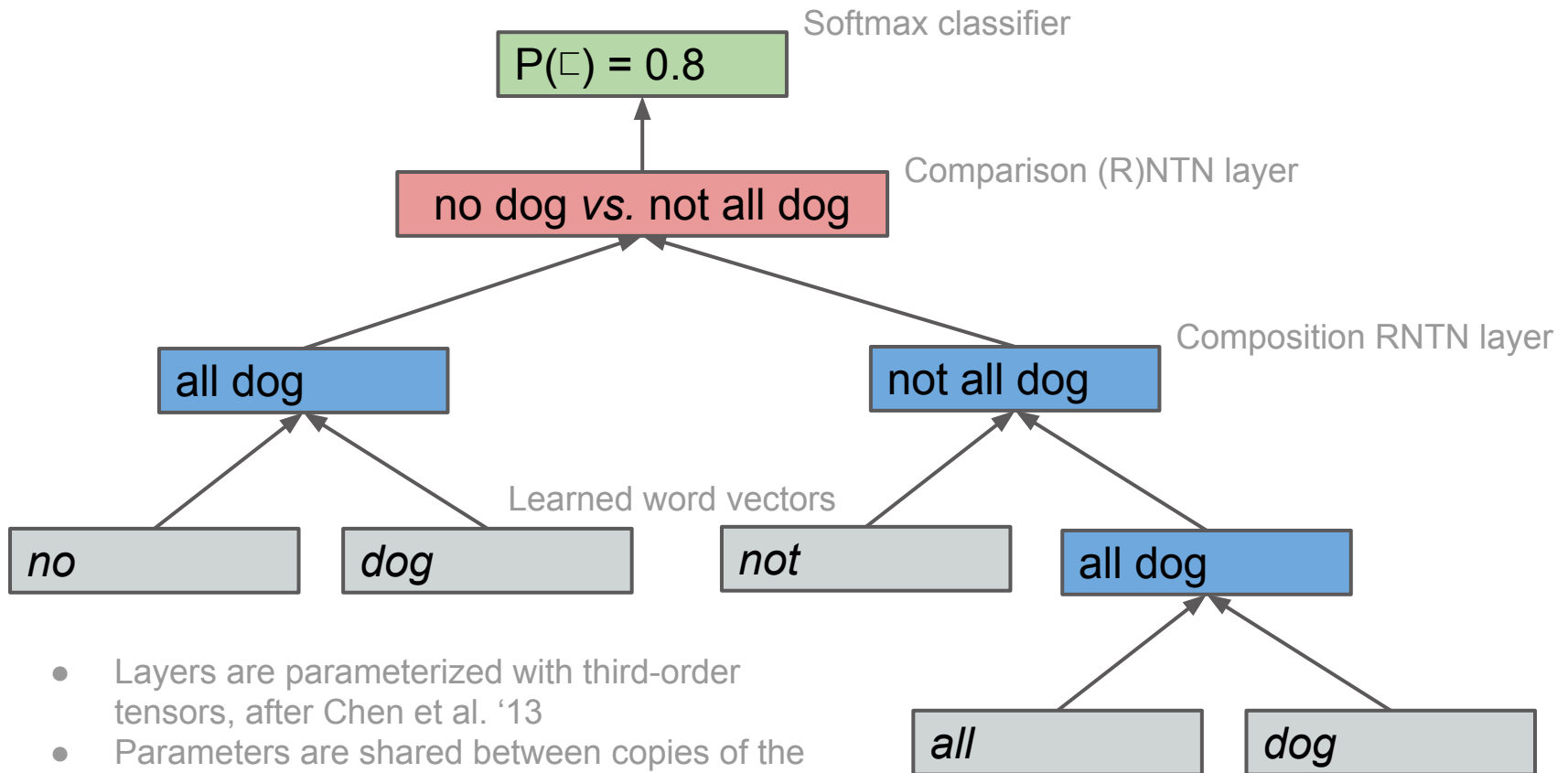
Strip away everything *e/se* that makes natural language hard:

- Small, unambiguous vocabulary
 - No morphology (no tense, no plurals, no agreement..)
 - No pronouns/references to context
 - Unlabeled constituency parses are given in data
-

The setup

- Small (~50 word) vocabulary
 - Three basic types:
 - Quantifiers: *some, all, no, most, two, three*
 - Predicates: *dog, cat, animal, live, European, ...*
 - Negation: *not*
 - Handmade dataset, 12k sentence pairs, grouped into templates.
 - All sentences of the form *QPP*, with optional negation on each predicate:
 - ((some x) bark) # ((some x) (not bark))*
 - ((some dog) bark) # ((some dog) (not bark))*
 - ((most (not dog)) European) ⊃ ((most (not dog)) French)*
-

The model: an RNTN for NLI



- Layers are parameterized with third-order tensors, after Chen et al. '13
 - Parameters are shared between copies of the composition layer
 - Input word vectors are initialized randomly and learned.
-

Five experiments

- *All-in*: train and test on all data. \Rightarrow 100%
- *All-split*: train on 85% of each pattern, test on rest.
 \Rightarrow 100%

(most dog) bark | (no dog) alive

(all cat) French \sqsupset (some cat) European

(most dog) French | (no dog) European

Five experiments

- *One-set-out*: hold out one pattern for testing only, split remaining data 85/15.
 - (most x) European | (no x) European
 - *One-subclass-out*: hold out one set of patterns for testing only, split remaining data 85/15.
 - (most x) y | (no x) y
 - *One-pair-out*: hold out one every pattern with a given pair of quantifiers for testing only, split rest.
 - (most (not x)) y # (no x) $z \dots$
-

Pilot results

Target dataset	Data evaluated	SET-OUT	SUBCL.-OUT	PAIR-OUT
<i>(most x) bark</i> <i>(no x) bark</i>	target dataset only	100%	100%	93.6%
	all held out datasets	(100%)	36.8%	78.8%
	all test data	99.8%	95.9%	93.8%
<i>(two x) bark</i> # <i>(all x) bark</i>	target dataset only	0%	100%	94.7%
	all held out datasets	(0%)	100%	62.7%
	all test data	97.5%	99.3%	93.0%
<i>(some x) bark</i> ^ <i>(no x) bark</i>	target dataset only	0%	0%	0%
	all held out datasets	(0%)	0%	25.2%
	all test data	97.7%	94.0%	85.5%

MacCartney's join:

$(most\ x)\ y \sqsubset (some\ x)\ y, (some\ x)\ y \wedge (no\ x)\ y$

$\models (most\ x)\ y \mid (no\ x)\ y$

$(some\ x)\ y \sqsupset (most\ x)\ y, (most\ x)\ y \mid (no\ x)\ y$

$\models (some\ x)\ y \{ \sqsupset \wedge \mid \# - \} (no\ x)\ y$

Today

Can these techniques learn models for general purpose NLU?

- Survey: Deep learning models for NLU
 - Experiment: Can RNTNs learn to reason with quantifiers (in an ideal world)?
 - **Experiment: Can RNTNs learn the natural logic *join* operator?**
 - Experiment: How do these models do on a challenge dataset?
-

Extra experiments: MacC's Join

\bowtie	\equiv	\sqsubset	\sqsupset	\wedge	\mid	\cup	$\#$
\equiv	\equiv	\sqsubset	\sqsupset	\wedge	\mid	\cup	$\#$
\sqsubset	\sqsubset	\sqsubset	$\equiv \sqsubset \sqsupset \mid \#$	\mid	\mid	$\sqsubset \wedge \mid \cup \#$	$\sqsubset \mid \#$
\sqsupset	\sqsupset	$\equiv \sqsubset \sqsupset \cup \#$	\sqsupset	\cup	$\sqsupset \wedge \mid \cup \#$	\cup	$\sqsupset \cup \#$
\wedge	\wedge	\cup	\mid	\equiv	\sqsupset	\sqsubset	$\#$
\mid	\mid	$\sqsubset \wedge \mid \cup \#$	\mid	\sqsubset	$\equiv \sqsubset \sqsupset \mid \#$	\sqsubset	$\sqsubset \mid \#$
\cup	\cup	\cup	$\sqsupset \wedge \mid \cup \#$	\sqsupset	\sqsupset	$\equiv \sqsubset \sqsupset \cup \#$	$\sqsupset \cup \#$
$\#$	$\#$	$\sqsubset \cup \#$	$\sqsupset \mid \#$	$\#$	$\sqsupset \mid \#$	$\sqsubset \cup \#$	$\equiv \sqsubset \sqsupset \wedge \mid \cup \#$

MacCartney's join table: $aRb \ \& \ bR'c \Rightarrow a\{join(R,R')\}c$

Cells that contain $\#$ represent uncertain results and can be approximated by just $\#$.

Extra experiments: Lattices with join

EXTRACTED RELATIONS:

$b \equiv b$

$b \cup c$

$b \cup d$

$b \supseteq e$

$c \cup d$

$c \supseteq e$

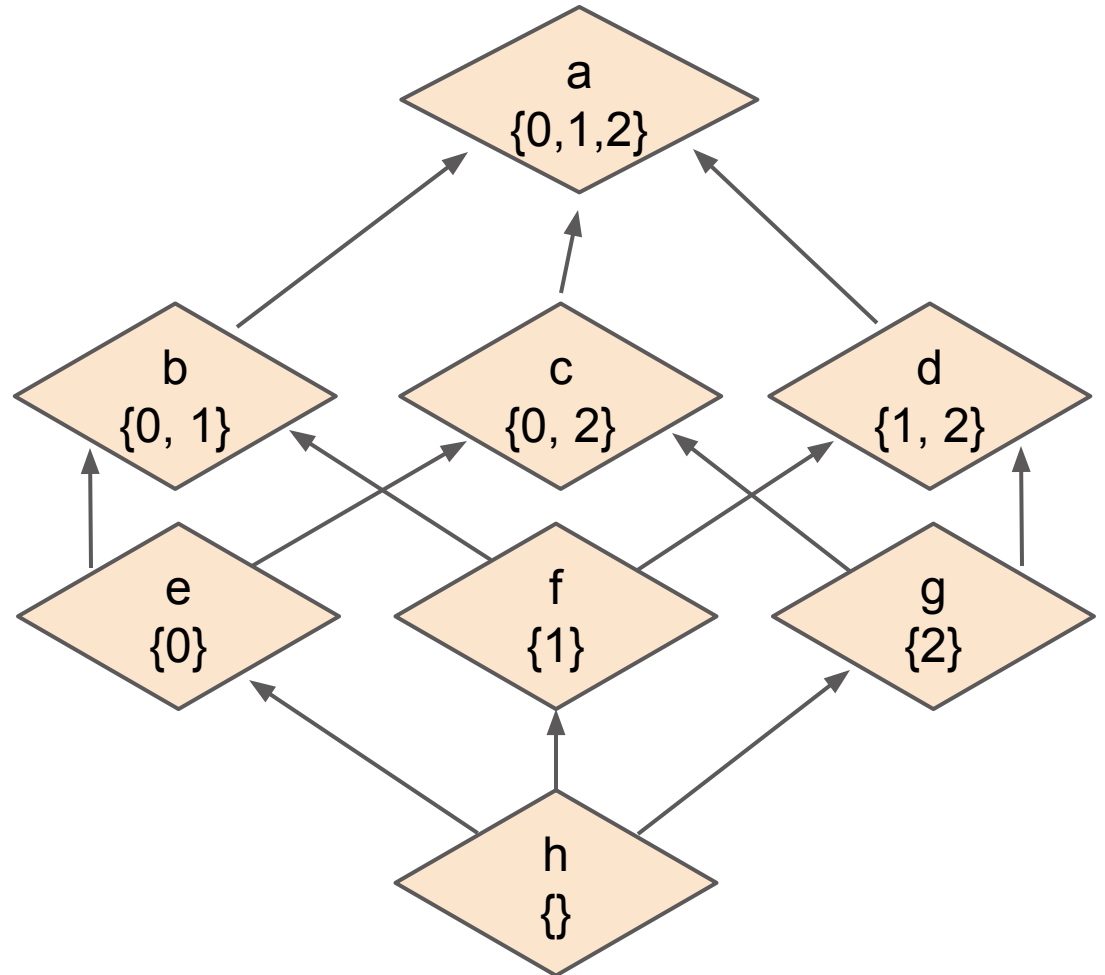
$c \wedge f$

$c \supseteq g$

$e \sqsubseteq b$

$e \sqsubseteq c$

...



Extra experiments: Lattices with join

EXTRACTED RELATIONS:

$b \equiv b$

$b \cup c$

$b \cup d$

$b \supseteq e$

$c \cup d$

$c \supseteq e$

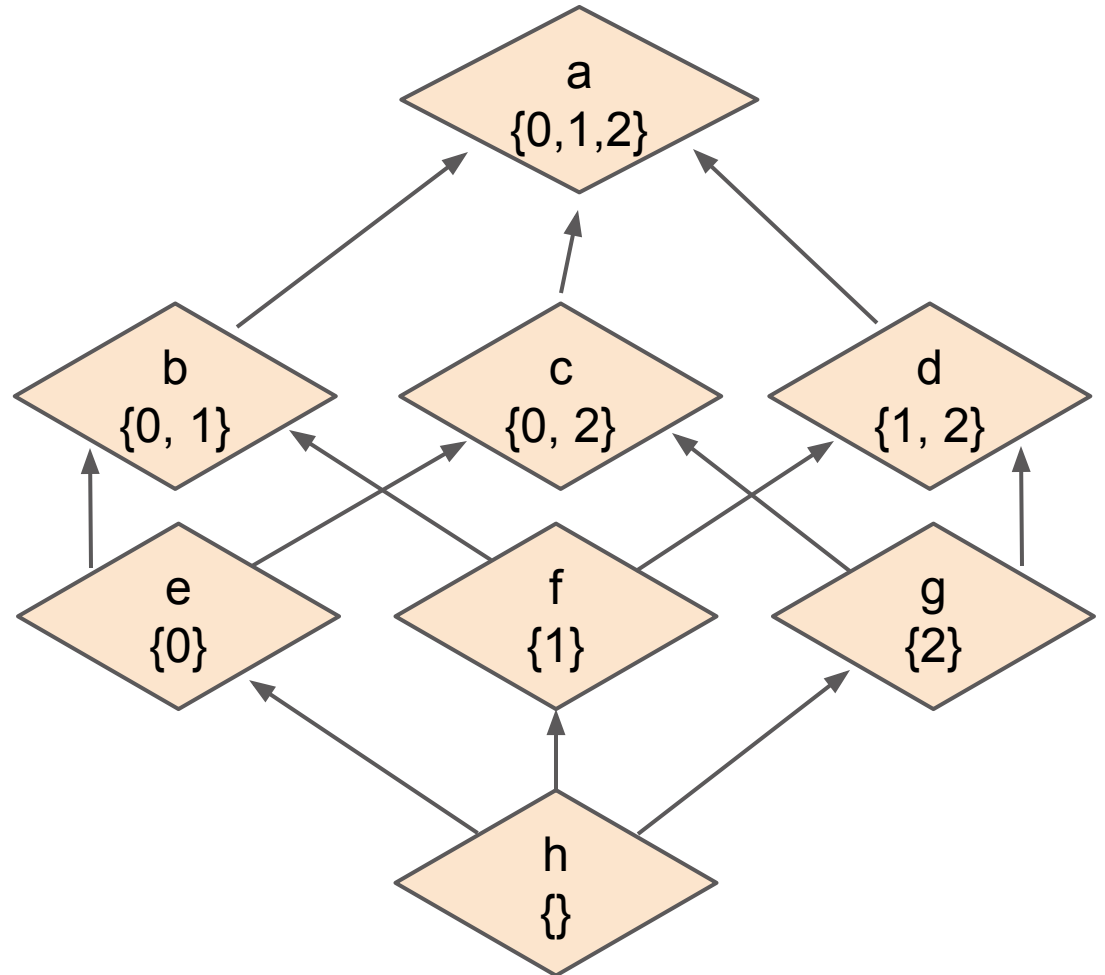
$c \wedge f$

$c \supseteq g$

$e \sqsubseteq b$

$e \sqsubseteq c$

...



Extra experiments: Lattices with join

TRAIN:

$b \cup c$

$c \cup d$

$c \sqsupset e$

$c \sqsupset g$

$e \sqsubset c$

...

TEST:

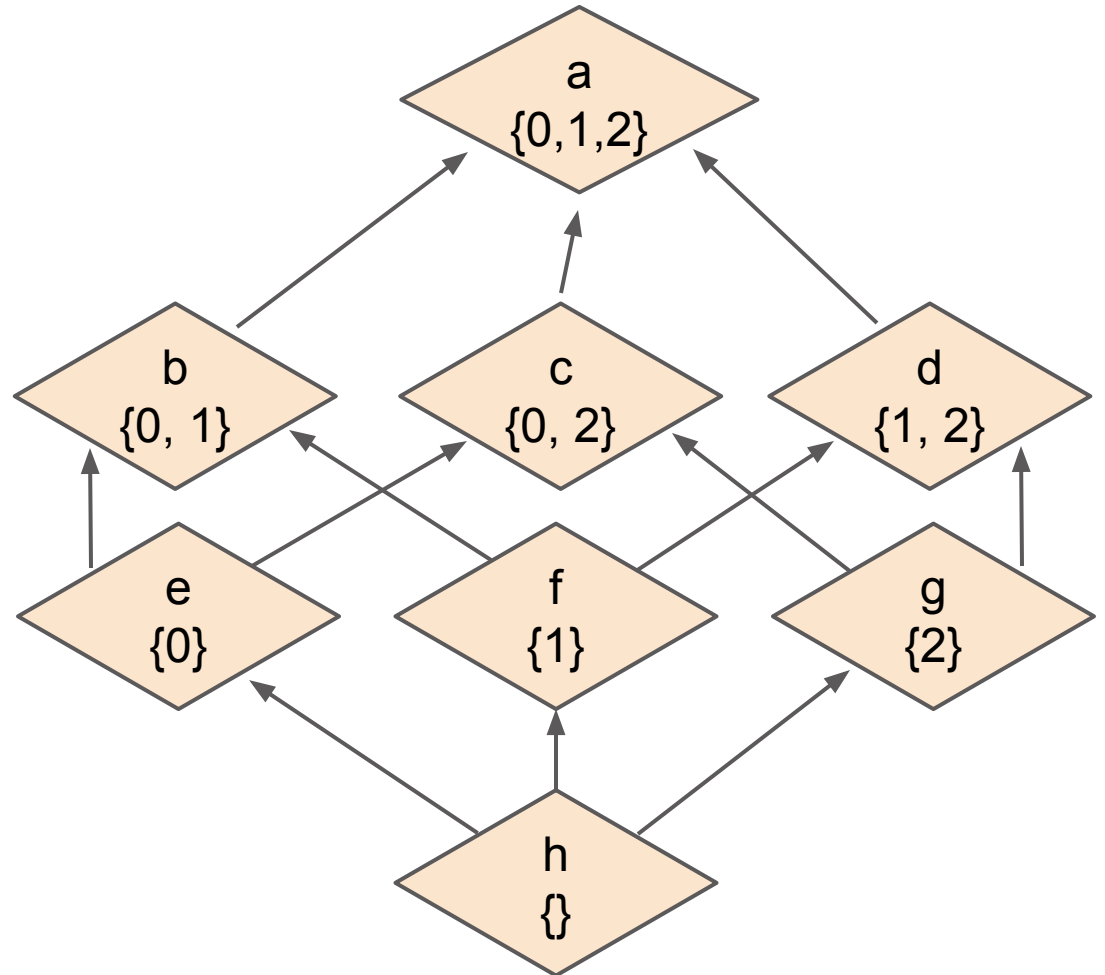
$b \equiv b$

$b \cup d$

$b \sqsupset e$

$c \wedge f$

$e \sqsubset b$



Extra experiments: Lattices with join

TRAIN:

$b \cup c$

$c \cup d$

$c \sqsupset e$

$c \sqsupset g$

$e \sqsubset c$

...

TEST:

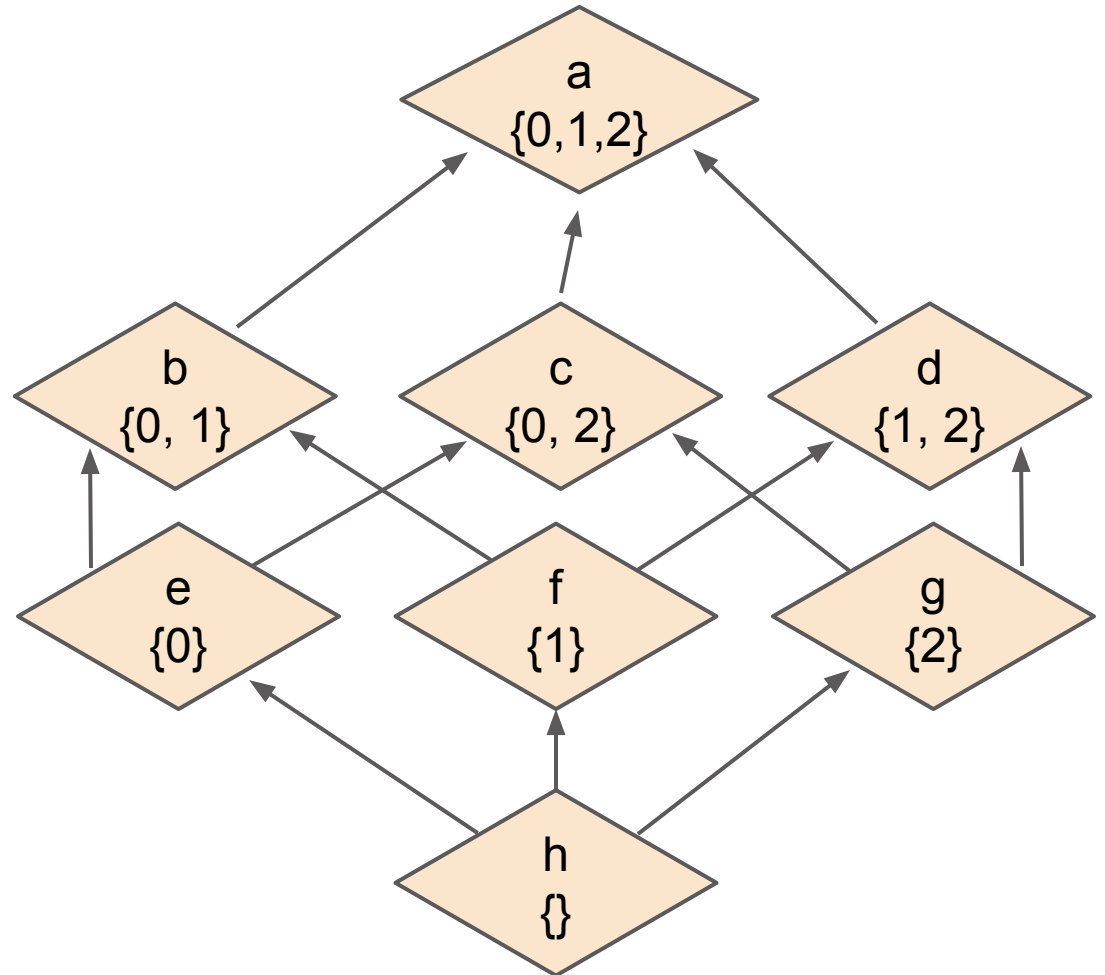
$b \equiv b$

$b \cup d$

~~$b \sqsupset e$~~

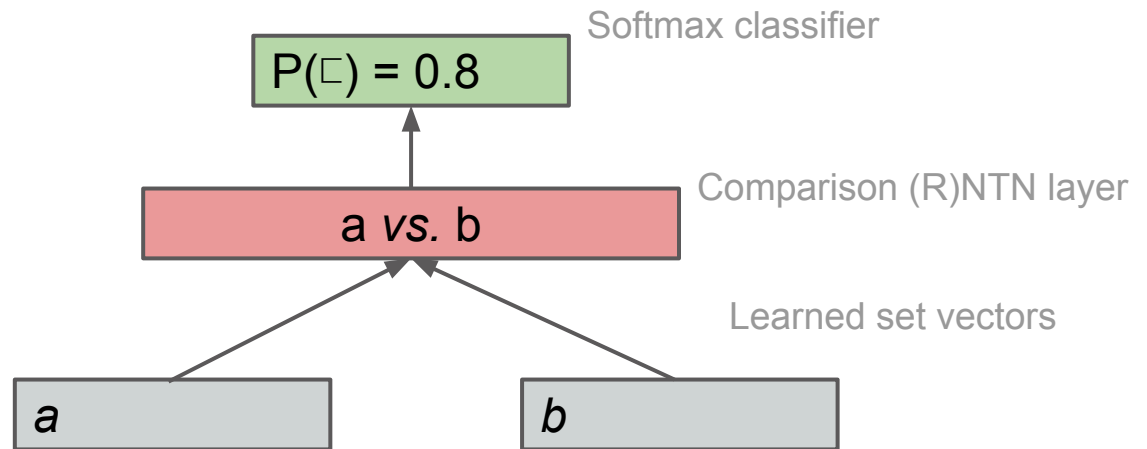
$c \wedge f$

$e \sqsubset b$



Extra experiments: Lattices with join

- Same model as in the monotonicity experiments above, but no composition/internal structure in the sentences.
- Lattice with 50 sets/nodes, 50% of data held out for testing.
⇒ 100% accuracy



Today

Can these techniques learn models for general purpose NLU?

- Survey: Deep learning models for NLU
 - Experiment: Can RNTNs learn to reason with quantifiers (in an ideal world)?
 - Experiment: Can RNTNs learn the natural logic *join* operator?
 - **Experiment: How do these models do on a challenge dataset?**
-

SemEval SICK

- NLP challenge dataset:
 - 10,000 sentence pairs labeled:
 - {forward entailment, contradiction, neutral}
 - “Sentences involving compositional knowledge” challenge:
 - No idioms, no named entities, no anaphora, tense doesn’t matter.
 - Requires general knowledge about word meaning and hypernymy, but no factoid knowledge.
-

SemEval SICK data

CONTRADICTION:

The woman in a red costume is leaning against a brick wall and playing an instrument.

The woman in a red costume is not leaning against a brick wall and is not playing an instrument.

NEUTRAL:

The player is dunking the basketball into the net and a crowd is in background.

A man with a jersey is dunking the ball at a basketball game.

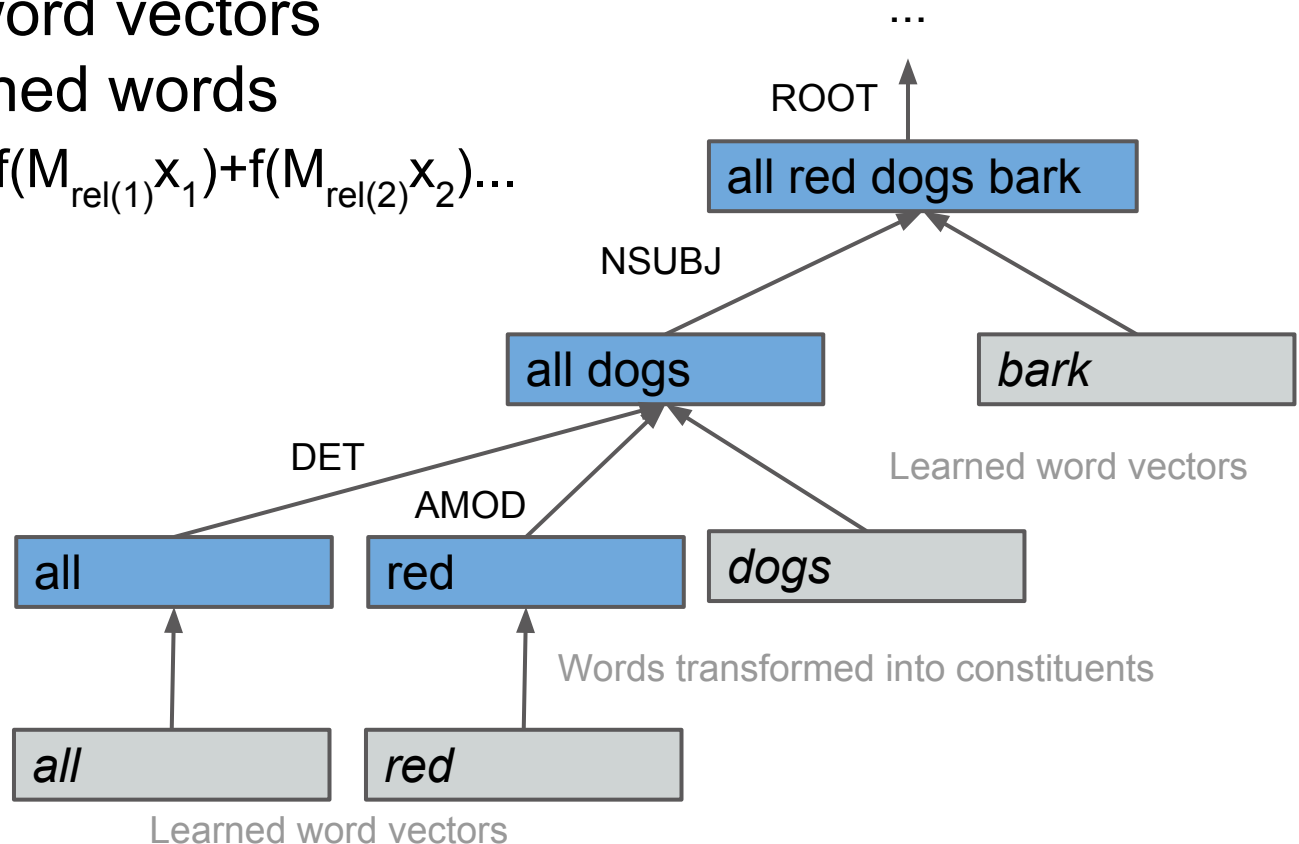
ENTAILMENT:

Four kids are doing backbends in the park

Four children are doing backbends in the park

SemEval SICK model

- Dependency tree RNNs
- Pretrained word vectors
- Partially-trained words
- $y = M_{\text{head}}x_{\text{head}} + f(M_{\text{rel}(1)}x_1) + f(M_{\text{rel}(2)}x_2) + \dots$



Results so far... eh?

- String inclusion baseline: 55.2%
- Most frequent class (Neutral): 56.4%
- Best dependency tree RNN: 74.5%
- Best SemEval result (Uillinois): 84.6%

But!

- No alignment or word sense disambiguation
-

Deep learning logistics

- There isn't any library yet that can do everything you'll need well.
 - But! Research code is available in MATLAB and Java
- Training monotonicity and SICK models: 4-18 hrs
- Lots of knobs to twiddle:
 - Stochastic optimization (AdaGrad/SGD) v. batch (L-BFGS)
 - Number of layers, dimensionality, L1 v. L2
 - Type of nonlinearity
 - Train/test split
 - DepTree RNNs: diagonal v. square matrices
- ...

Thanks!

Code is available for all three experiments.

sbowman@stanford.edu

Next steps

- Better formal characterizations of what it takes to learn to do inference
 - Better formal characterizations of the structures that can be learned
 - More types of network
 - More semantic phenomena
 - Test on natural language data
-