Natural Language Processing (CSE 517): Sequence Models

Noah Smith

© 2018

University of Washington nasmith@cs.washington.edu

May 2, 2018

Project

Include control characters in vocabulary, so $|\mathcal{V}|$ =136,755.

Extension on the dry run: Wednesday, May 9.

Mid-Quarter Review: Results

Thank you!

Going well:

- ► Lectures, examples, explanations of math, slides, engagement of the class, readings
- Unified framework, connections among concepts, up-to-date content, topic coverage

Changes to make:

- Posting slides before lecture
- ► Expectations on project

Sequence Models (Quick Review)

Models:

- ► Hidden Markov
- \blacktriangleright " $\phi(x,i,y,y')$ "
- Algorithm: Viterbi

Applications:

- part-of-speech tagging (Church, 1988)supersense tagging (Ciaramita and Altun, 2006)
- supersonse tagging (claratimes and vittari, 20
- ▶ named-entity recognition (Bikel et al., 1999)
- multiword expressions (Schneider and Smith, 2015)
- base noun phrase chunking (Sha and Pereira, 2003)

Learning:

► Supervised parameter estimation for HMMs

A problem with a long history: word-sense disambiguation.

A problem with a long history: word-sense disambiguation.

Classical approaches assumed you had a list of ambiguous words and their senses.

► E.g., from a dictionary

A problem with a long history: word-sense disambiguation.

Classical approaches assumed you had a list of ambiguous words and their senses.

► E.g., from a dictionary

Ciaramita and Johnson (2003) and Ciaramita and Altun (2006) used a lexicon called WordNet to define 41 semantic classes for words.

► WordNet (Fellbaum, 1998) is a fascinating resource in its own right! See http://wordnetweb.princeton.edu/perl/webwn to get an idea.

A problem with a long history: word-sense disambiguation.

Classical approaches assumed you had a list of ambiguous words and their senses.

► E.g., from a dictionary

Ciaramita and Johnson (2003) and Ciaramita and Altun (2006) used a lexicon called WordNet to define 41 semantic classes for words.

► WordNet (Fellbaum, 1998) is a fascinating resource in its own right! See http://wordnetweb.princeton.edu/perl/webwn to get an idea.

This represents a coarsening of the annotations in the Semcor corpus (Miller et al., 1993).

Example: box's Thirteen Synonym Sets, Eight Supersenses

- 1. box: a (usually rectangular) container; may have a lid. "he rummaged through a box of spare parts"
- 2. box/loge: private area in a theater or grandstand where a small group can watch the performance. "the royal box was empty"
- 3. box/boxful: the quantity contained in a box. "he gave her a box of chocolates"
- 4. corner/box: a predicament from which a skillful or graceful escape is impossible. "his lying got him into a tight corner"
- 5. box: a rectangular drawing. "the flowchart contained many boxes"
- 6. box/boxwood: evergreen shrubs or small trees
- 7. box: any one of several designated areas on a ball field where the batter or catcher or coaches are positioned. "the umpire warned the batter to stay in the batter's box"
- 8. box/box seat: the driver's seat on a coach. "an armed guard sat in the box with the driver"
- 9. box: separate partitioned area in a public place for a few people. "the sentry stayed in his box to avoid the cold"
- 10. box: a blow with the hand (usually on the ear). "I gave him a good box on the ear"
- 11. box/package: put into a box. "box the gift, please"
- 12. box: hit with the fist. "I'll box your ears!"
- 13. box: engage in a boxing match.

Example: box's Thirteen Synonym Sets, Eight Supersenses

- 1. box: a (usually rectangular) container; may have a lid. "he rummaged through a box of spare parts" \leadsto N.ARTIFACT
- 2. box/loge: private area in a theater or grandstand where a small group can watch the performance. "the royal box was empty" ->> N.ARTIFACT
- 3. box/boxful: the quantity contained in a box. "he gave her a box of chocolates" --> N.QUANTITY
- 4. corner/box: a predicament from which a skillful or graceful escape is impossible. "his lying got him into a tight corner" ->> N.STATE
- 5. box: a rectangular drawing. "the flowchart contained many boxes" ->> N.SHAPE
- 6. box/boxwood: evergreen shrubs or small trees → N.PLANT
- 7. box: any one of several designated areas on a ball field where the batter or catcher or coaches are positioned. "the umpire warned the batter to stay in the batter's box" NARTIFACT
- 8. box/box seat: the driver's seat on a coach. "an armed guard sat in the box with the driver" \leadsto N.ARTIFACT
- 9. box: separate partitioned area in a public place for a few people. "the sentry stayed in his box to avoid the cold" → N.ARTIFACT
- 10. box: a blow with the hand (usually on the ear). "I gave him a good box on the ear" \sim N.ACT
- 11. box/package: put into a box. "box the gift, please" \sim V.CONTACT
- 12. box: hit with the fist. "I'll box your ears!" → V.CONTACT
- 13. box: engage in a boxing match. ~ V.COMPETITION

Supersense Tagging Example

```
Clara Harris , one of the guests in the N.PERSON N.PERSON
```

```
box , stood up and demanded N.ARTIFACT V.MOTION V.COMMUNICATION
```

```
water ...
N.SUBSTANCE
```

Ciaramita and Altun's Approach

Features at each position in the sentence:

- word
- "first sense" from WordNet (also conjoined with word)
- ► POS, coarse POS
- ► shape (case, punctuation symbols, etc.)
- previous label

All of these fit into " $\phi(x, i, y, y')$."

Supervised Training of Sequence Models (Discriminative)

Given: annotated sequences $\langle \langle \boldsymbol{x}_1, \boldsymbol{y}_1, \rangle, \dots, \langle \boldsymbol{x}_n, \boldsymbol{y}_n \rangle \rangle$

Assume:

$$\operatorname{predict}(\boldsymbol{x}) = \underset{\boldsymbol{y} \in \mathcal{L}^{\ell+1}}{\operatorname{argmax}} \sum_{i=1}^{\ell+1} \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, i, y_i, y_{i-1})$$
$$= \underset{\boldsymbol{y} \in \mathcal{L}^{\ell+1}}{\operatorname{argmax}} \mathbf{w} \cdot \sum_{i=1}^{\ell+1} \boldsymbol{\phi}(\boldsymbol{x}, i, y_i, y_{i-1})$$
$$= \underset{\boldsymbol{y} \in \mathcal{L}^{\ell+1}}{\operatorname{argmax}} \mathbf{w} \cdot \boldsymbol{\Phi}(\boldsymbol{x}, \boldsymbol{y})$$

Estimate: w

Perceptron

Perceptron algorithm for classification:

- ▶ For $t \in \{1, ..., T\}$:
 - ▶ Pick i_t uniformly at random from $\{1, ..., n\}$.

 - $\qquad \qquad \mathbf{w} \leftarrow \mathbf{w} \alpha \left(\phi(\boldsymbol{x}_{i_t}, \hat{\ell}_{i_t}) \phi(\boldsymbol{x}_{i_t}, \ell_{i_t}) \right)$

Structured Perceptron

Collins (2002)

Perceptron algorithm for classification structured prediction:

- ▶ For $t \in \{1, ..., T\}$:
 - ▶ Pick i_t uniformly at random from $\{1, ..., n\}$.
 - $\qquad \qquad \hat{\boldsymbol{y}}_{i_t} \leftarrow \operatorname*{argmax}_{\boldsymbol{y} \in \mathcal{L}^{\ell+1}} \mathbf{w} \cdot \boldsymbol{\Phi}(\boldsymbol{x}_{i_t}, \boldsymbol{y})$
 - $\qquad \qquad \mathbf{w} \leftarrow \mathbf{w} \alpha \left(\mathbf{\Phi}(\boldsymbol{x}_{i_t}, \boldsymbol{\hat{y}}_{i_t}) \mathbf{\Phi}(\boldsymbol{x}_{i_t}, \boldsymbol{y}_{i_t}) \right)$

This can be viewed as stochastic subgradient descent on the structured hinge loss:

$$\sum_{i=1}^n \underbrace{\max_{oldsymbol{y} \in \mathcal{L}^{\ell_i+1}} \mathbf{w} \cdot oldsymbol{\Phi}(oldsymbol{x}_i, oldsymbol{y})}_{ ext{fear}} - \underbrace{\mathbf{w} \cdot oldsymbol{\Phi}(oldsymbol{x}_i, oldsymbol{y}_i)}_{ ext{hope}}$$

Back to Supersenses

```
Clara
       Harris
                , one of the
                                guests
                                              the
      N.PERSON
                                N.PERSON
                                           demanded
     box
                stood
                          up
                                  and
  N. ARTIFACT
                       V.MOTION
                                       V.COMMUNICATION
     water
  N.SUBSTANCE
```

Shouldn't Clara Harris and stood up be respectively "grouped"?

Segmentations

Segmentation:

- ▶ Input: $\boldsymbol{x} = \langle x_1, x_2, \dots, x_\ell \rangle$

where $\ell = \sum_{i=1}^{m} \ell_i$.

Application: word segmentation for writing systems without whitespace.

Segmentations

Segmentation:

- ▶ Input: $\boldsymbol{x} = \langle x_1, x_2, \dots, x_\ell \rangle$

where $\ell = \sum_{i=1}^{m} \ell_i$.

Application: word segmentation for writing systems without whitespace.

With arbitrarily long segments, this does not look like a job for $\phi(x, i, y, y')$!

Segmentation as Sequence Labeling

Ramshaw and Marcus (1995)

Two labels: B ("beginning of new segment"), I ("inside segment")

$$\blacktriangleright \ \ell_1=4, \ell_2=3, \ell_3=1, \ell_4=2 \longrightarrow \langle \mathsf{B}, \mathsf{I}, \mathsf{I}, \mathsf{I}, \mathsf{B}, \mathsf{I}, \mathsf{I}, \mathsf{B}, \mathsf{B}, \mathsf{I} \rangle$$

Three labels: B, I, O ("outside segment")

Five labels: B, I, O, E ("end of segment"), S ("singleton")

Segmentation as Sequence Labeling

Ramshaw and Marcus (1995)

Two labels: B ("beginning of new segment"), I ("inside segment")

$$\blacktriangleright \ \ell_1=4, \ell_2=3, \ell_3=1, \ell_4=2 \longrightarrow \langle \mathsf{B}, \mathsf{I}, \mathsf{I}, \mathsf{I}, \mathsf{B}, \mathsf{I}, \mathsf{I}, \mathsf{B}, \mathsf{B}, \mathsf{I} \rangle$$

Three labels: B, I, O ("outside segment")

Five labels: B, I, O, E ("end of segment"), S ("singleton")

Bonus: combine these with a label to get labeled segmentation!

Named Entity Recognition as Segmentation and Labeling

An older and narrower subset of supersenses used in information extraction:

- person,
- location,
- organization,
- ▶ geopolitical entity,
- ...and perhaps domain-specific additions.

Named Entity Recognition

With $\underline{\text{Commander Chris Ferguson}}$ at the helm , $\underline{\text{person}}$

←□ → ←□ → ← □ → ← □ → −

Named Entity Recognition





Named Entity Recognition: Evaluation

```
      10
      11
      12
      13
      14
      15
      16
      17
      18
      19

      rescue Britons stranded by Eyjafjallajökull 's volcanic ash cloud .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      .
      <td
```

Segmentation Evaluation

Typically: precision, recall, and F_1 .

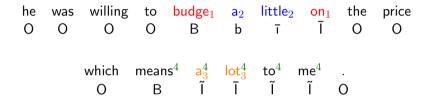
Multiword Expressions

Schneider et al. (2014b)

- ▶ MW compounds: red tape, motion picture, daddy longlegs, Bayes net, hot air balloon, skinny dip, trash talk
- ▶ verb-particle: pick up, dry out, take over, cut short
- verb-preposition: refer to, depend on, look for, prevent from
- ▶ verb-noun(-preposition): pay attention (to), go bananas, lose it, break a leg, make the most of
- ▶ support verb: make decisions, take breaks, take pictures, have fun, perform surgery
- other phrasal verb: put up with, miss out (on), get rid of, look forward to, run amok, cry foul, add insult to injury, make off with
- ▶ PP modifier: above board, beyond the pale, under the weather, at all, from time to time, in the nick of time
- coordinated phrase: cut and dry, more or less, up and leave
- **conjunction/connective:** as well as, let alone, in spite of, on the face of it/on its face
- semi-fixed VP: smack <one>'s lips, pick up where <one> left off, go over <thing> with a fine-tooth(ed) comb, take <one>'s time, draw <oneself> up to <one>'s full height
- fixed phrase: easy as pie, scared to death, go to hell in a handbasket, bring home the bacon, leave of absence, sense of humor
- phatic: You're welcome. Me neither!
- ▶ proverb: Beggars can't be choosers. The early bird gets the worm. To each his own. One man's <thing₁> is another man's <thing₂>.

Sequence Labeling with Nesting

Schneider et al. (2014a)



Strong (subscript) vs. weak (superscript) MWEs.

One level of nesting, plus strong/weak distinction, can be handled with an eight-tag scheme.

Back to Syntax

Base noun phrase chunking:

[He]_{NP} reckons [the current account deficit]_{NP} will narrow to [only \$ 1.8 billion]_{NP} in [September]_{NP}

(What is a base noun phrase?)

"Chunking" used generically includes base verb and prepositional phrases, too.

Sequence labeling with BIO tags and features can be applied to this problem (Sha and Pereira, 2003).

Remarks

Sequence models are extremely useful:

- syntax: part-of-speech tags, base noun phrase chunking
- > semantics: supersense tags, named entity recognition, multiword expressions

All of these are called "shallow" methods (why?).

Remarks

Sequence models are extremely useful:

- syntax: part-of-speech tags, base noun phrase chunking
- ▶ semantics: supersense tags, named entity recognition, multiword expressions

All of these are called "shallow" methods (why?).

Issues to be aware of:

- Supervised data for these problems is not cheap.
- ▶ Performance always suffers when you test on a different style, genre, dialect, etc. than you trained on.
- ightharpoonup Runtime depends on the size of $\mathcal L$ and the number of consecutive labels that features can depend on.

References I

- Daniel M. Bikel, Richard Schwartz, and Ralph M. Weischedel. An algorithm that learns what's in a name. *Machine learning*, 34(1–3):211–231, 1999.
- Kenneth W. Church. A stochastic parts program and noun phrase parser for unrestricted text. In *Proc. of ANLP*, 1988.
- Massimiliano Ciaramita and Yasemin Altun. Broad-coverage sense disambiguation and information extraction with a supersense sequence tagger. In *Proc. of EMNLP*, 2006.
- Massimiliano Ciaramita and Mark Johnson. Supersense tagging of unknown nouns in WordNet. In *Proc. of EMNLP*, 2003.
- Michael Collins. Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms. In *Proc. of EMNLP*, 2002.
- Christiane Fellbaum, editor. WordNet: An Electronic Lexical Database. MIT Press, 1998.
- G. A. Miller, C. Leacock, T. Randee, and R. Bunker. A semantic concordance. In Proc. of HLT, 1993.
- Lance A Ramshaw and Mitchell P. Marcus. Text chunking using transformation-based learning, 1995. URL http://arxiv.org/pdf/cmp-lg/9505040.pdf.
- Nathan Schneider and Noah A. Smith. A corpus and model integrating multiword expressions and supersenses. In *Proc. of NAACL*, 2015.
- Nathan Schneider, Emily Danchik, Chris Dyer, and Noah A. Smith. Discriminative lexical semantic segmentation with gaps: Running the MWE gamut. *Transactions of the Association for Computational Linguistics*, 2:193–206, April 2014a.

References II

Nathan Schneider, Spencer Onuffer, Nora Kazour, Emily Danchik, Michael T. Mordowanec, Henrietta Conrad, and Noah A. Smith. Comprehensive annotation of multiword expressions in a social web corpus. In *Proc. of LREC*, 2014b.

Fei Sha and Fernando Pereira. Shallow parsing with conditional random fields. In Proc. of NAACL, 2003.