## Natural Language Processing (CSE 490U): Sequence Models (II)

Noah Smith

© 2017

University of Washington nasmith@cs.washington.edu

January 30-February 3, 2017

1/63

## Mid-Quarter Review: Results

Thank you!

Going well:

- ► Content! Lectures, slides, readings.
- Office hours, homeworks, course structure.

イロン イロン イヨン イヨン 三日

2/63

Changes to make:

- Math (more visuals and examples).
- More structure in sections.
- Prerequisites.

#### Full Viterbi Procedure

Input:  $\boldsymbol{x}$ ,  $p(X_i \mid Y_i)$ ,  $p(Y_{i+1} \mid Y_i)$ 

Output:  $\hat{y}$ 

- 1. For  $i \in \langle 1, \ldots, \ell \rangle$ :
  - Solve for  $s_i(*)$  and  $b_i(*)$ .
    - Special base case for i = 1 to handle start state  $y_0$  (no max)
    - General recurrence for  $i \in \langle 2, \ldots, \ell 1 \rangle$
    - Special case for  $i = \ell$  to handle stopping probability
- 2.  $\hat{y}_{\ell} \leftarrow \operatorname*{argmax}_{y \in \mathcal{L}} s_{\ell}(y)$
- 3. For  $i \in \langle \ell, \ldots, 1 \rangle$ :
  - $\blacktriangleright \hat{y}_{i-1} \leftarrow b(y_i)$

	$x_1$	$x_2$	 $x_\ell$
y			
y'			
:			
$y^{last}$			

	$x_1$	$x_2$	 $x_\ell$
y	$s_1(y)$		
y'	$s_1(y')$		
÷			
$y^{last}$	$s_1(y^{last})$		

 $s_1(y) = p(x_1 \mid y) \cdot p(y \mid y_0)$ 

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ = 臣 = のへで

5/63

	$x_1$	$x_2$		$x_{\ell}$
y	$s_1(y)$	$s_2(y)$		
y'	$s_1(y')$	$s_2(y')$		
÷				
$y^{last}$	$s_1(y^{last})$	$s_2(y^{last})$		

$$s_i(y) = p(x_i \mid y) \cdot \max_{y' \in \mathcal{L}} p(y \mid y') \cdot \boxed{s_{i-1}(y')}$$

4 ロ ト 4 回 ト 4 直 ト 4 直 ト 道 の Q ()
6 / 63

	$x_1$	$x_2$	•••	$x_\ell$
y	$s_1(y)$	$s_2(y)$		$s_\ell(y)$
y'	$s_1(y')$	$s_2(y')$		$s_\ell(y')$
÷				
$y^{last}$	$s_1(y^{last})$	$s_2(y^{last})$		$s_\ell(y^{last})$

$$s_{\ell}(y) = p(\bigcirc \mid y) \cdot p(x_{\ell} \mid y) \cdot \max_{y' \in \mathcal{L}} p(y \mid y') \cdot \boxed{s_{\ell-1}(y')}$$

・ロト (部) (言) (言) (言) (の) (の) (7/63)

#### Viterbi Asymptotics

Space:  $O(|\mathcal{L}|\ell)$ 

Runtime:  $O(|\mathcal{L}|^2 \ell)$ 

	$x_1$	$x_2$	 $x_{\ell}$
y			
y'			
:			
$y^{last}$			

 Instead of HMM parameters, we can "featurize" or "neuralize."

- Instead of HMM parameters, we can "featurize" or "neuralize."
- Viterbi instantiates an general algorithm called max-product variable elimination, for inference along a chain of variables with pairwise "links."

- Instead of HMM parameters, we can "featurize" or "neuralize."
- Viterbi instantiates an general algorithm called max-product variable elimination, for inference along a chain of variables with pairwise "links."
- ► Viterbi solves a special case of the "best path" problem.



- Instead of HMM parameters, we can "featurize" or "neuralize."
- Viterbi instantiates an general algorithm called max-product variable elimination, for inference along a chain of variables with pairwise "links."
- ► Viterbi solves a special case of the "best path" problem.
- ► Higher-order dependencies among *Y* are also possible.

$$s_i(y, y') = \max_{y'' \in \mathcal{L}} p(x_i \mid y) \cdot p(y \mid y', y'') \cdot s_{i-1}(y', y'')$$

### Applications of Sequence Models

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

### Parts of Speech

#### http://mentalfloss.com/article/65608/

master-particulars-grammar-pop-culture-primer



14/63

#### Parts of Speech

"Open classes": Nouns, verbs, adjectives, adverbs, numbers

イロン イロン イヨン イヨン 三日

15/63

- "Closed classes":
  - Modal verbs
  - Prepositions (on, to)
  - ▶ Particles (*off*, *up*)
  - Determiners (*the*, *some*)
  - Pronouns (she, they)
  - Conjunctions (and, or)

## Parts of Speech in English: Decisions

Granularity decisions regarding:

- verb tenses, participles
- plural/singular for verbs, nouns
- proper nouns
- comparative, superlative adjectives and adverbs

Some linguistic reasoning required:

- Existential there
- Infinitive marker to
- ► wh words (pronouns, adverbs, determiners, possessive whose) Interactions with tokenization:
  - Punctuation
  - Compounds (Mark'll, someone's, gonna)

Penn Treebank: 45 tags,  $\sim$ 40 pages of guidelines (Marcus et al., 1993)

## Parts of Speech in English: Decisions

Granularity decisions regarding:

- verb tenses, participles
- plural/singular for verbs, nouns
- proper nouns
- comparative, superlative adjectives and adverbs

Some linguistic reasoning required:

- Existential there
- Infinitive marker to
- ► wh words (pronouns, adverbs, determiners, possessive whose) Interactions with tokenization:
  - Punctuation
  - ► Compounds (*Mark'll, someone's, gonna*)
- ▶ Social media: hashtag, at-mention, discourse marker (*RT*), URL, emoticon, abbreviations, interjections, acronyms
   Penn Treebank: 45 tags, ~40 pages of guidelines (Marcus et al., 1993)

TweetNLP: 20 tags, 7 pages of guidelines (Gimpel et al. 2011)

#### Example: Part-of-Speech Tagging

ikr smh he asked fir yo last name

so he can add u on fb lololol

## Example: Part-of-Speech Tagging



## Example: Part-of-Speech Tagging



## Why POS?

- ► Text-to-speech: record, lead, protest
- ▶ Lemmatization:  $saw/V \rightarrow see$ ;  $saw/N \rightarrow saw$
- Quick-and-dirty multiword expressions: (Adjective | Noun)\* Noun (Justeson and Katz, 1995)
- Preprocessing for harder disambiguation problems:
  - ► The Georgia branch had taken **on** loan commitments ....
  - ► The average of interbank offered rates plummeted ....

Define a map  $\mathcal{V} \to \mathcal{L}$ .

Define a map  $\mathcal{V} \to \mathcal{L}$ .

How to pick the single POS for each word? E.g., raises, Fed, ....

Define a map  $\mathcal{V} \to \mathcal{L}$ .

How to pick the single POS for each word? E.g., raises, Fed, ....

Penn Treebank: most frequent tag rule gives 90.3%, 93.7% if you're clever about handling unknown words.

Define a map  $\mathcal{V} \to \mathcal{L}$ .

How to pick the single POS for each word? E.g., raises, Fed, ...

Penn Treebank: most frequent tag rule gives 90.3%, 93.7% if you're clever about handling unknown words.

All datasets have some errors; estimated upper bound for Penn Treebank is 98%.

#### Supervised Training of Hidden Markov Models

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

$$p(\boldsymbol{x}, \boldsymbol{y}) = \prod_{i=1}^{\ell+1} \theta_{x_i|y_i} \cdot \gamma_{y_i|y_{i-1}}$$

Parameters: for each state/label  $y \in \mathcal{L}$ :

- ▶  $\theta_{*|y}$  is the "emission" distribution, estimating  $p(x \mid y)$  for each  $x \in \mathcal{V}$
- ▶  $\gamma_{*|y}$  is called the "transition" distribution, estimating  $p(y' \mid y)$  for each  $y' \in \mathcal{L}$

#### Supervised Training of Hidden Markov Models

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

$$p(\boldsymbol{x}, \boldsymbol{y}) = \prod_{i=1}^{\ell+1} \theta_{x_i|y_i} \cdot \gamma_{y_i|y_{i-1}}$$

Parameters: for each state/label  $y \in \mathcal{L}$ :

- ▶  $\theta_{*|y}$  is the "emission" distribution, estimating  $p(x \mid y)$  for each  $x \in \mathcal{V}$
- ▶  $\gamma_{*|y}$  is called the "transition" distribution, estimating  $p(y' \mid y)$  for each  $y' \in \mathcal{L}$

Maximum likelihood estimate: count and normalize!

# TnT, a trigram HMM tagger with smoothing: 96.7% (Brants, 2000)

#### Back to POS

TnT, a trigram HMM tagger with smoothing: 96.7% (Brants, 2000)

State of the art:  $\sim$ 97.5% (Toutanova et al., 2003); uses a feature-based model with:

- capitalization features
- spelling features
- name lists ("gazetteers")
- context words
- hand-crafted patterns

#### Back to POS

TnT, a trigram HMM tagger with smoothing: 96.7% (Brants, 2000)

State of the art:  ${\sim}97.5\%$  (Toutanova et al., 2003); uses a feature-based model with:

- capitalization features
- spelling features
- name lists ("gazetteers")
- context words
- hand-crafted patterns

There might be very recent improvements to this.

Parts of speech are a minimal syntactic representation.

Sequence labeling can get you a lightweight *semantic* representation, too.

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
- 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"  $\rightsquigarrow$  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"  $\rightsquigarrow$  N.ARTIFACT
- 3. box/boxful: the quantity contained in a box. "he gave her a box of chocolates"  $\rightsquigarrow N.QUANTITY$
- 4. corner/box: a predicament from which a skillful or graceful escape is impossible. "his lying got him into a tight corner"  $\rightsquigarrow$  N.STATE
- 5. box: a rectangular drawing. "the flowchart contained many boxes"  $\rightarrow$  N.SHAPE
- 6. box/boxwood: evergreen shrubs or small trees  $\rightsquigarrow$  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"  $\rightsquigarrow$  N.ARTIFACT
- 8. box/box seat: the driver's seat on a coach. "an armed guard sat in the box with the driver"  $\rightsquigarrow$  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"  $\rightsquigarrow$  N.ARTIFACT
- 10. box: a blow with the hand (usually on the ear). "I gave him a good box on the ear"  $\rightsquigarrow$   $\rm N.ACT$
- 11. box/package: put into a box. "box the gift, please" ~> V.CONTACT
- 12. box: hit with the fist. "I'll box your ears!"  $\rightsquigarrow$  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

イロト 不得 とくき とくき とうき

38 / 63

water N.SUBSTANCE

.

#### Ciaramita and Altun's Approach

Features at each position in the sentence:

- word
- "first sense" from WordNet (also conjoined with word)

39 / 63

- POS, coarse POS
- shape (case, punctuation symbols, etc.)
- previous label

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

#### Featurizing HMMs

Log-probability score of y (given x) decomposes into a sum of local scores:

$$\operatorname{score}(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i=1}^{\ell+1} \underbrace{(\log p(x_i \mid y_i) + \log p(y_i \mid y_{i+1}))}_{(1)}$$
(1)

Featurized HMM:

$$\operatorname{score}(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i=1}^{\ell+1} \underbrace{(\mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, i, y_i, y_{i-1}))}_{i=1} (2)$$
$$= \mathbf{w} \cdot \sum_{i=1}^{\ell+1} \boldsymbol{\phi}(\boldsymbol{x}, i, y_i, y_{i-1})$$
$$\operatorname{global features, } \boldsymbol{\Phi}(\boldsymbol{x}, \boldsymbol{y})$$
(3)

40 / 63

#### What Changes?

Algorithmically, not much!

Viterbi recurrence before and after:

$$s_{1}(y) = p(x_{1} \mid y) \cdot p(y \mid y_{0})$$
  

$$s_{i}(y) = p(x_{i} \mid y) \cdot \max_{y' \in \mathcal{L}} p(y \mid y') \cdot \boxed{s_{i-1}(y')}$$
  

$$s_{\ell}(y) = p(\bigcirc \mid y) \cdot p(x_{\ell} \mid y) \cdot \max_{y' \in \mathcal{L}} p(y \mid y') \cdot \boxed{s_{\ell-1}(y')}$$

Now:

$$s_{1}(y) = \exp \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, 1, y, y_{0})$$
  

$$s_{i}(y) = \max_{y' \in \mathcal{L}} \exp \left[\mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, i, y, y')\right] \cdot \boxed{s_{i-1}(y')}$$
  

$$s_{\ell}(y) = \max_{y' \in \mathcal{L}} \exp \left[\mathbf{w} \cdot \left(\boldsymbol{\phi}(\boldsymbol{x}, \ell, y, y') + \boldsymbol{\phi}(\boldsymbol{x}, \ell+1, \bigcirc, y)\right)\right] \cdot \boxed{s_{\ell-1}(y')}$$

#### Supervised Training of Sequence Models (Discriminative)

Given: annotated sequences 
$$\langle\langle m{x}_1,m{y}_1,
angle,\ldots,\langlem{x}_n,m{y}_n
angle
angle$$

Assume:

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

42 / 63

Estimate: w

#### Perceptron

Perceptron algorithm for classification:

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

$$\begin{array}{l} \bullet \quad \hat{\ell}_{i_t} \leftarrow \operatorname*{argmax}_{\ell \in \mathcal{L}} \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}_{i_t}, \ell) \\ \bullet \quad \mathbf{w} \leftarrow \mathbf{w} - \alpha \left( \boldsymbol{\phi}(\boldsymbol{x}_{i_t}, \hat{\ell}_{i_t}) - \boldsymbol{\phi}(\boldsymbol{x}_{i_t}, \ell_{i_t}) \right) \end{array}$$

#### Structured Perceptron

Collins (2002)

Perceptron algorithm for classification structured prediction:

► For t ∈ {1,...,T}:
Pick i<sub>t</sub> uniformly at random from {1,...,n}.
$$\hat{y}_{i_t} \leftarrow \underset{y \in \mathcal{L}^{\ell+1}}{\operatorname{argmax}} \mathbf{w} \cdot \Phi(x_{i_t}, y)$$
►  $\mathbf{w} \leftarrow \mathbf{w} - \alpha \left( \Phi(x_{i_t}, \hat{y}_{i_t}) - \Phi(x_{i_t}, y_{i_t}) \right)$ 

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

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

(ロ)、<</li>
 (目)、
 (目)、
 (日)、
 (日)、

#### Back to Supersenses

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

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

#### Segmentations

Segmentation:

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

$$\left\langle \begin{array}{c} \boldsymbol{x}_{1:\ell_1}, \\ \boldsymbol{x}_{(1+\ell_1):(\ell_1+\ell_2)}, \\ \boldsymbol{x}_{(1+\ell_1+\ell_2):(\ell_1+\ell_2+\ell_3)}, \dots, \\ \boldsymbol{x}_{(1+\sum_{i=1}^{m-1}\ell_i):\sum_{i=1}^{m}\ell_i} \end{array} \right\rangle$$
(4)

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

Application: word segmentation for writing systems without whitespace.

#### Segmentations

Segmentation:

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

$$\left< \begin{array}{c} \boldsymbol{x}_{1:\ell_1}, \\ \boldsymbol{x}_{(1+\ell_1):(\ell_1+\ell_2)}, \\ \boldsymbol{x}_{(1+\ell_1+\ell_2):(\ell_1+\ell_2+\ell_3)}, \dots, \\ \boldsymbol{x}_{(1+\sum_{i=1}^{m-1}\ell_i):\sum_{i=1}^{m}\ell_i} \end{array} \right>$$
(4)

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  $\phi(\pmb{x},i,y,y')!$ 

## Segmentation as Sequence Labeling

Ramshaw and Marcus (1995)

Two labels: B ("beginning of new segment"), I ("inside segment")  $\blacktriangleright \ \ell_1 = 4, \ell_2 = 3, \ell_3 = 1, \ell_4 = 2 \longrightarrow \langle B, I, I, B, B, I, B, B, I \rangle$ 

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")  $\blacktriangleright \ \ell_1 = 4, \ell_2 = 3, \ell_3 = 1, \ell_4 = 2 \longrightarrow \langle B, I, I, B, B, I, B, B, I \rangle$ 

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:

50/63

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

## Named Entity Recognition

With Commander Chris Ferguson at the helm ,

person

Atlantistouched down at Kennedy Space Center .spacecraftlocation

## Named Entity Recognition



Atlantis<br/>spacecrafttouched down at<br/>locationKennedy Space Center<br/>location.BOOBIIO

イロト 不得下 イヨト イヨト 二日

52 / 63

### Named Entity Recognition: Evaluation

rescue Britons stranded by Eyjafjallajökull 's volcanic ash cloud . 0 B 0 В 0 0 0 0 B В 0 0  $\mathbf{O}$ 

#### 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 <thing<sub>1</sub>> is another man's <thing<sub>2</sub>>.

# 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]<sub>NP</sub> reckons [the current account deficit]<sub>NP</sub> will narrow to
[only $ 1.8 billion]<sub>NP</sub> in [September]<sub>NP</sub>
```

(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.
- ► Runtime depends on the size of *L* and the number of consecutive labels that features can depend on.

#### To-Do List

#### Read: Jurafsky and Martin (2016b,a)

#### 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. URL http://people.csail.mit.edu/mcollins/6864/slides/bikel.pdf.

Thorsten Brants. TnT - a statistical part-of-speech tagger. In Proc. of ANLP, 2000.

- 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.
- Kevin Gimpel, Nathan Schneider, Brendan O'Connor, Dipanjan Das, Daniel Mills, Jacob Eisenstein, Michael Heilman, Dani Yogatama, Jeffrey Flanigan, and Noah A. Smith. Part-of-speech tagging for Twitter: Annotation, features, and experiments. In Proc. of ACL, 2011.
- Daniel Jurafsky and James H. Martin. Information extraction (draft chapter), 2016a. URL https://web.stanford.edu/~jurafsky/slp3/21.pdf.

#### References II

- Daniel Jurafsky and James H. Martin. Part-of-speech tagging (draft chapter), 2016b. URL https://web.stanford.edu/~jurafsky/slp3/10.pdf.
- John S. Justeson and Slava M. Katz. Technical terminology: Some linguistic properties and an algorithm for identification in text. *Natural Language Engineering*, 1:9–27, 1995.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of English: the Penn treebank. *Computational Linguistics*, 19(2): 313–330, 1993.
- 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.
- 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.

#### References III

- Fei Sha and Fernando Pereira. Shallow parsing with conditional random fields. In *Proc. of NAACL*, 2003.
- Kristina Toutanova, Dan Klein, Christopher D. Manning, and Yoram Singer. Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proc. of NAACL*, 2003.