

Topics in Natural Language Processing

Shay Cohen

Institute for Language, Cognition and Computation

University of Edinburgh

Lecture 7

Administrativa

I received suggested topics from most of you

- If you didn't send a topic yet, please send it as soon as possible
- Next thing: scheduling everybody and allocating brief paper responses
- I will try to allocate brief paper responses on the same topic you present (but different papers)
- It will not always work out

Last class

Semiring inference: CKY and the inside algorithms

$$\alpha(A, i, j) = \bigoplus_{\substack{k \\ A \rightarrow B, c}} \alpha(B, i, k) \otimes \alpha(c, k+1, j)$$

$i, j \in \mathbb{N}$, $A \in N$
 $A \rightarrow B \ c \in R$



Parsing as Weighted Logic Programming

$$\text{constit}(a, i, j) \oplus = \text{constit}(b, i, k) \otimes \text{constit}(c, k + 1, j) \otimes \text{rule}(a \rightarrow b c)$$

$$\text{constit}(a, i, i) \oplus = \text{rule}(a \rightarrow w_i)$$

Goal: $\text{constit}(S, 0, n)$

Example of a Weighted Logic Programme

We are given a sequence w_1, \dots, w_n of some symbols.

$$\text{prob}(b, i) \oplus = \text{prob}(a, i-1) \otimes \text{transition}(a \rightarrow b) \otimes \text{emission}(b, w_i)$$

$$\text{prob}(a, 1) \oplus = \text{start_state}(b) \otimes \text{emission}(a, w_1)$$

$$\alpha(B, i) = \sum_A \alpha(A, i-1) \times t(A \rightarrow B) \times \text{emission}(B \rightarrow w_i)$$

\uparrow state \uparrow point in the sentence

"Forward" algorithm

Viterbi

| | |
|---------------|----------------|
| $a \otimes b$ | $a \times b$ |
| $a \oplus b$ | $\max\{a, b\}$ |

Example of a Weighted Logic Programme

We are given a sequence w_1, \dots, w_n of some symbols.

$$\text{prob}(b, i) \oplus = \text{prob}(a, i - 1) \otimes \text{transition}(a \rightarrow b) \otimes \text{emission}(b, w_i)$$

$$\text{prob}(a, 1) \oplus = \text{start_state}(b) \otimes \text{emission}(a, w_1)$$

Hidden Markov models and the forward algorithm:

- emission are the emission probabilities
- transition are the transition probabilities
- start_state are the initial probabilities
- $\text{prob}(b, i)$ gives the probability $p(w_1, \dots, w_i, S_i = b)$

Weighted Logic Programming

Once we present the inference algorithm as a weighted logic program, there are general-purpose algorithms to solve it.

- The “agenda” algorithm
- Search algorithms
- Tabular algorithms

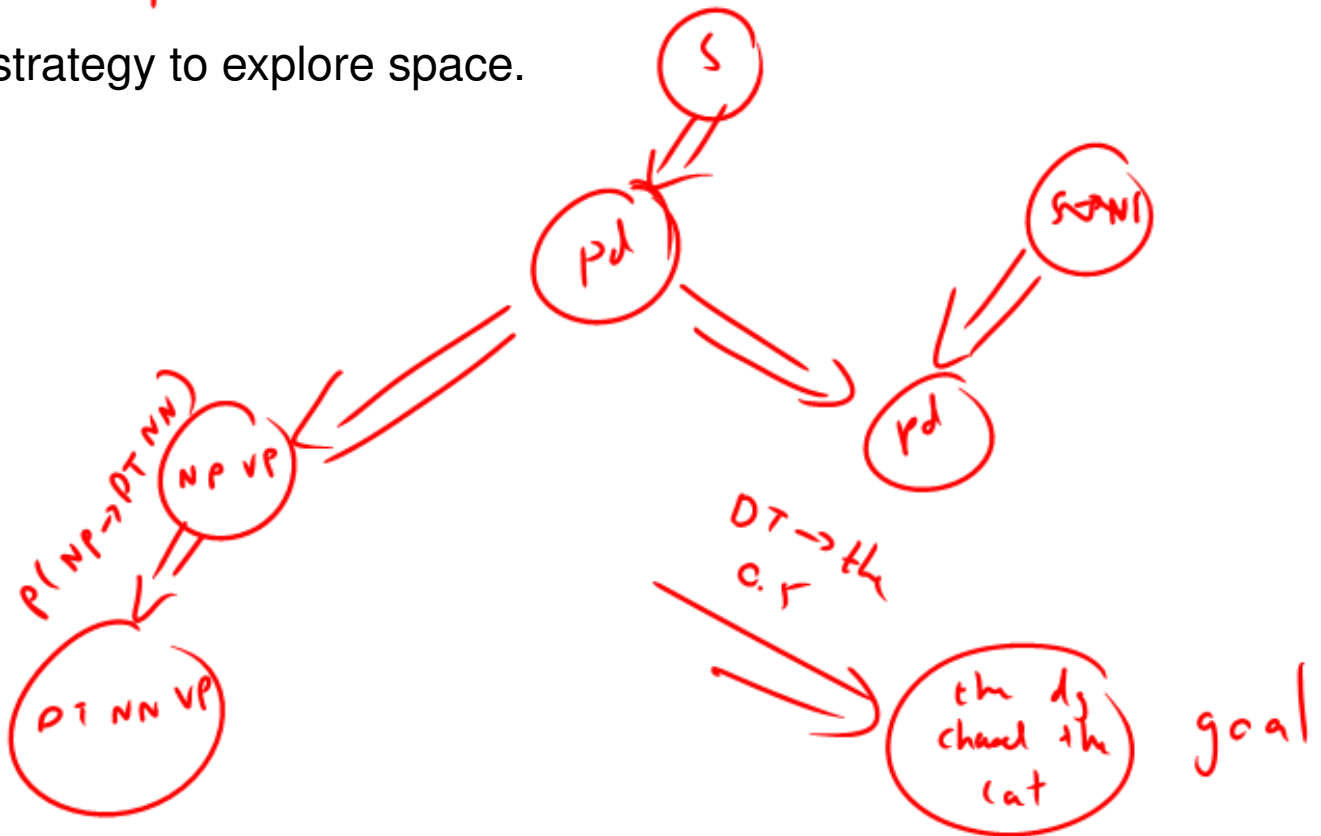
If that interests you, take a look at <http://www.dyna.org>.

Parsing Using Search

State space: *partial derivations*

a a A B c d E A

Need a strategy to explore space.

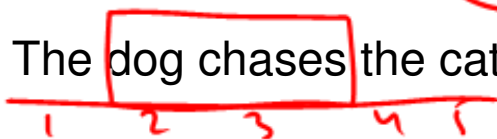


Inference Using Graphs and Hypergraphs

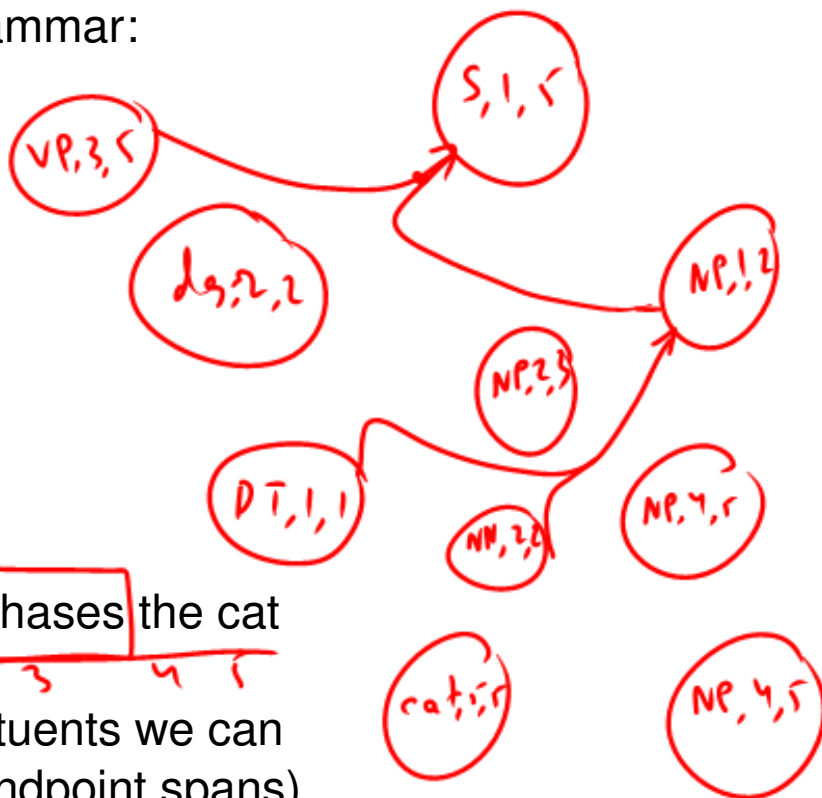
Consider the context-free grammar:

- S → NP VP
- NP → DT NN
- VP → VB NP
- VB → chases
- NN → dog | cat | chases
- DT → the
- NP → DT NN NN

We want to parse: The dog chases the cat



What are the potential constituents we can create? (nonterminals with endpoint spans)



Today's class

- Log-linear models and their estimation
- Regularisation

Estimation until now

- Count and normalise
- Corresponds to maximum likelihood estimate for multinomial models

A Type-based POS Tagging Model

- Want to model $p(\overset{y}{\text{tag}} \mid \overset{x}{\text{word}})$
- Can use a simple multinomial model, but...
- What about orthography and morphology?
- What about unseen words? ←

Problems ① Unseen words
with ~~tags~~ ② Not explicitly
count & normalize

linguistic structure

Linear Score for POS Tagging Model

$$\Omega = T \times V \quad T - \text{tag} \quad V - \text{vocabulary}$$

Linear score: $score(y|x) = \sum_{i=1}^d w_i g_i(x, y)$

$$g_{ik}(x, y) = I(x = \text{dog}, y = \text{VB})$$

$$I(y = \text{PN}, x = \text{capitalized})$$

Probability model:

$$p(x, y) = \frac{\exp(score(y|x))}{\sum_{x, y} \exp(score(y|x))}$$

$$I(y = \text{VB}, x \text{ ends in ion})$$

Estimation

We observe $(x_1, y_1), \dots, (x_n, y_n)$

What is the likelihood?

$$f(x_1, y_1, \dots, x_n, y_n | w) = \prod_{i=1}^n p(x_i, y_i | w) =$$
$$= \frac{\prod_{i=1}^n \exp\left(\sum_{j=1}^d w_j g_j(x_i, y_i)\right)}{z(w)}$$

$$z(w) = \sum_{x, y} \exp\left(\sum_{j=1}^d w_j g_j(x, y)\right)$$

Estimation

$$Z(w) = \sum_{x,y} \exp(w^T g(x,y))$$

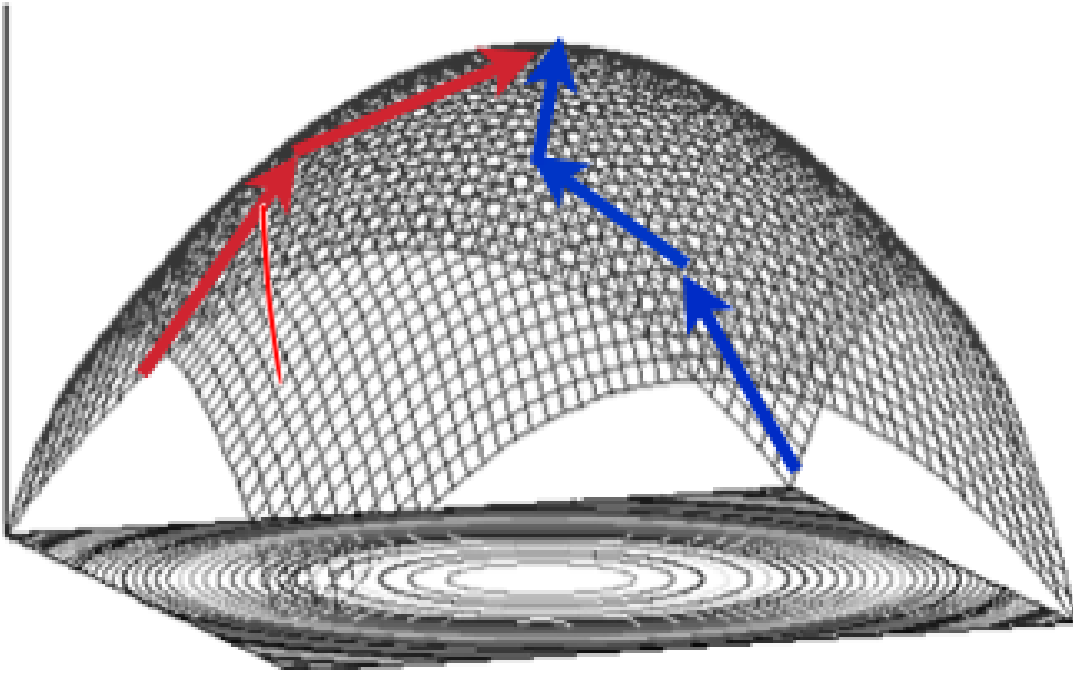
$$f(w|x_1, y_1, \dots, x_n, y_n) = \prod_{i=1}^n \frac{\exp(w^T g(x_i, y_i))}{Z(w)}$$

What is the log-likelihood?

$$\begin{aligned} L(w|x_1, y_1, \dots, x_n, y_n) &= \frac{1}{n} \sum_{i=1}^n \log \left(\frac{\exp(w^T g(x_i, y_i))}{Z(w)} \right) = \\ &= \frac{1}{n} \sum_{i=1}^n (w^T g(x_i, y_i) - \log Z(w)) = \\ &= \left[\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d w_j g_j(x_i, y_i) \right] - \log Z(w) \end{aligned}$$

Maximising the log-likelihood

$$w^* = \arg \max_w L(w|x_1, y_1, \dots, x_n, y_n)$$



Maximising the log-likelihood

Many of the maximisation algorithms are a variant of the update:

$$w^{(t+1)} \leftarrow w^{(t)} + \mu v$$

where $v \in \mathbb{R}^d$ and $v_i = \frac{\partial L}{\partial w_i} (w^{(t)})$.

Estimation

$$f(x) = c \Rightarrow f'(x) = 0$$

$$f(x) = ax \Rightarrow f'(x) = a$$

What is the average log-likelihood?

$$L(w|x_1, y_1, \dots, x_n, y_n) = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^d w_j g_j(x_i, y_i) - \log Z(w) \right)$$

$$\begin{aligned} \log f(x) &= \\ &= \frac{1}{f(x)} \cdot f'(x) \end{aligned}$$

What is the derivative?

$$\frac{\partial L}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n 0 + 0 + 0 + \dots + g_j(x_i, y_i) + 0 \dots + 0$$

$$- \frac{1}{Z(w)} \cdot \frac{\partial Z}{\partial w_j} = \left[\frac{1}{n} \sum_{i=1}^n g_j(x_i, y_i) \right] - \frac{1}{Z(w)} \cdot \frac{\partial Z}{\partial w_j}$$

chain rule

Derivative of $Z(w)$

$$\begin{aligned} [e^{f(x)}]' &= \\ &= f'(x) \cdot e^{f(x)} \end{aligned}$$

$$Z(w) = \sum_{x,y} \exp \left(\underbrace{\sum_{j=1}^d w_j g_j(x,y)} \right)$$

$$\frac{\partial Z}{\partial w_j}(w) = \sum_{x,y} \exp \left(\sum_{j=1}^d w_j g_j(x,y) \right) \cdot g_j(x,y)$$

Gradient of average log-likelihood

$$\frac{\partial L}{\partial w_j} = \left(\frac{1}{n} \sum_{i=1}^n g_j(x_i, y_i) \right) - \sum_{x,y} \frac{\exp(\sum_{k=1}^d w_k g_k(x, y))}{Z(w)} g_j(x, y)$$

$\frac{1}{n} \sum_{i=1}^n g_j(x_i, y_i) = E_{\tilde{p}}[g_j]$ where $\tilde{p}(x, y) = \frac{\text{count}(x, y)}{n}$

$\sum_{x,y} \frac{\exp(\sum_{k=1}^d w_k g_k(x, y))}{Z(w)} g_j(x, y) = \sum_{x,y} \tilde{p}(x, y) g_j(x, y)$

$\sum_{x,y} \frac{\exp(\sum_{k=1}^d w_k g_k(x, y))}{Z(w)} g_j(x, y) = \sum_{x,y} p(x, y | w) g_j(x, y) = E_{p(\cdot|w)}[g_j(x, y)]$

$p(x, y | w)$

Gradient of average log-likelihood

$$\frac{\partial L}{\partial w_j} = \left(\frac{1}{n} \sum_{i=1}^n g_j(x_i, y_i) \right) - \sum_{x,y} \frac{\exp(\sum_{k=1}^d w_k g_k(x, y))}{Z(w)} g_j(x, y)$$

$$\frac{1}{n} \sum_{i=1}^n g_j(x_i, y_i) =$$

$$\sum_{x,y} \frac{\exp(\sum_{k=1}^d w_k g_k(x, y))}{Z(w)} g_j(x, y) =$$

Therefore, the gradient is the difference between empirical expectations and expectations under the model