

Lecture 7.5.

HMM's to SCFG's

Note Title

1/26/2010

$$SCFG = (\mathcal{V}, S, s, \mathcal{R})$$

Chomsky Normal Form

- cf. $\left\{ \begin{array}{l} \mathcal{V} \sim \text{finite set of terminals } \mathcal{V}_0 \text{ (Sam, thinks, snow)} \\ S \sim \text{finite set of non-terminals } S_0 = \{S, NP, VP, V\} \\ \mathcal{R} - \text{finite set of production, } A \rightarrow BC \in S_0, \text{ or } A \rightarrow X \in S_0 \\ s \in S \text{ is the start symbol (} S_0 = S \text{)} \end{array} \right.$

+ Production probabilities $P(A \rightarrow \beta) = P(\beta|A)$
for each $A \rightarrow \beta \in \mathcal{R}$

generate $r_1 \dots r_n$ is sequence of productions used to generate a tree Ψ , then

$$P_G(\Psi) = \prod_{r \in \mathcal{R}} P(r)^{fr(\Psi)}$$

$\sum_{\Psi} P_G(\Psi) = 1$ $fr(\Psi) = \text{no. times } r \text{ is used.}$

What do we want to compute?

1. What is the prob. $P_G(\omega)$ of string ω

$$P_G(\omega) = P_G(s \Rightarrow^* \omega) = \sum_{\Psi \in \Psi_G(\omega)} P_G(\Psi)$$

2. What is the most probable parse $\Psi(\omega)$ of string ω ?

$$\hat{\Psi}(\omega) = \arg \max_{\Psi \in \Psi_G(\omega)} P_G(\Psi)$$

$\Psi_G(\omega)$ - set of parse trees for ω generated by G
 $P_G(\Psi)$ - prob of tree Ψ w.r.t. grammar G .

SCFG "inside" algorithm

generalization of DP

Goal: compute $P(\omega) = \sum_{\Psi \in \Psi_G(\omega)} P(\Psi) = P(s \Rightarrow^* \omega)$

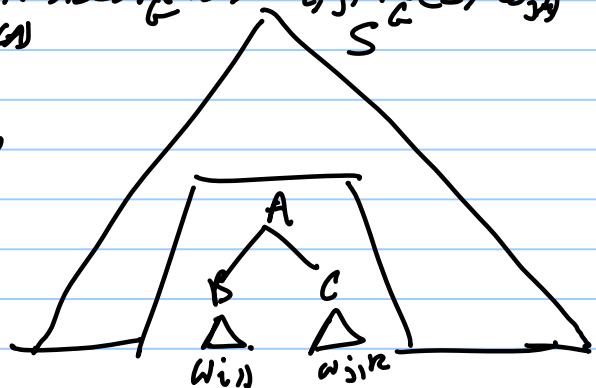
Data structure: table $P(A \Rightarrow^* \omega_{ij})$ for $A \in S, 0 \leq i < j \leq n$

Base case: $P(A \Rightarrow^* \omega_{i,i+1}) = P(A \rightarrow \omega_i), i = 0, \dots, n-1$

Return: $P_G(A \Rightarrow^* \omega_{ij}) = \sum_{k=i+1}^j \sum_{A \rightarrow BC \in \mathcal{R}} P_G(A \rightarrow BC) P_G(B \Rightarrow^* \omega_{i,k}) P_G(C \Rightarrow^* \omega_{k,j})$

Return: $P(s \Rightarrow^* \omega_{0,n})$

$P_G(A \Rightarrow^* \omega_{i,j})$ is called the "inside probability"



(2) Learning: Dataset of Examples

Assume the data is visible.

Count how frequently a rule is applied

Rule: $S \rightarrow NP VP$ Count 3 Frequency 1
 $NP \rightarrow Rice$ Count 2 " 2/3

This is also the MLE for the SCFG
 i.e. $P(\text{data} | \text{parameters } \theta)$

MLE is
 consistent
 asymptotically
 optimal
 (lowest variance)

But, this requires that the data is 'visible' or has been labeled. (like first HMM case)

Unsupervised Training EM (again)

$Z = (x, y)$ x - observed data,
 y - hidden variables (data)
 θ - parameters

EM algorithm: given visible data x :

1. Guess initial value θ_0 of parameters (rule prob)

2. Repeat:

E-step: For all $y_1 \dots y_n \in Y$ generate pseudo-data $(x, y_1) \dots (x, y_n) \sim (x, y)$ 'frequency' $P_{\theta_i}(y|x)$

M-step: set θ_{i+1} to MLE from the pseudo-data

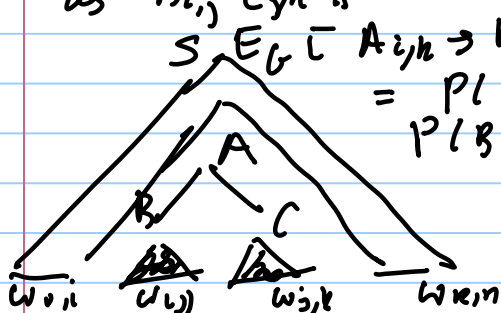
Sometimes can perform maximization directly from sufficient statistics (expected products / frequencies)

Use DP to calculate the expected probs:

$$E_{\theta} [I_{A \rightarrow B}(w)] = \sum_{0 \leq i < j < n} E_{\theta} [I_{A \rightarrow B}(w_{i,j})]$$

expected fraction of parses of w in which $A_{i,j}$ rewrites as $B_{i,j} C_{j,k}$ is

$$E_{\theta} [I_{A \rightarrow B}(w_{i,j})] = P(S \Rightarrow^* w_{i,j} A w_{j,k}) P(A \rightarrow B C) \\ = P(B \Rightarrow^* w_{i,j}) P(C \Rightarrow^* w_{j,k}) / P_r(w)$$



(3) Calculate $P_G(S \Rightarrow^* \omega_{0i} A \omega_{kn})$
 "outside probabilities"

Recursion from larger to smaller substrings in ω

Base case: $P(S \Rightarrow^* \omega_{0j} S \omega_{kn}) =$

Recursion: $P(S \Rightarrow^* \omega_{0j} C \omega_{kn}) =$

$$\sum_{i=1}^{j-1} \sum_{\substack{A, B \in S \\ A \rightarrow B C E R}} P(S \Rightarrow^* \omega_{0i} A \omega_{kn}) p(A \rightarrow B C) P(B \Rightarrow^* \omega_{ij})$$

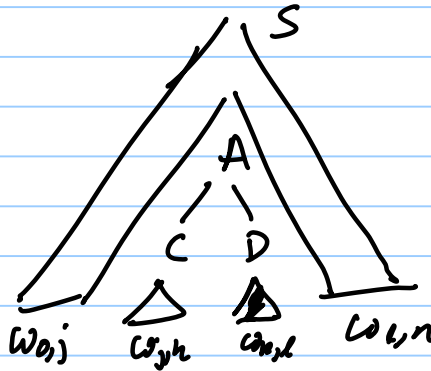
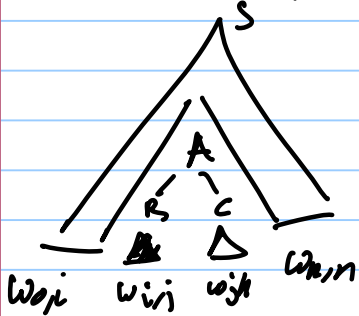
$$+ \sum_{l=k+1}^n \sum_{\substack{A, D \in S \\ A \rightarrow C D E R}} P(S \Rightarrow^* \omega_{0j} A \omega_{kn}) p(A \rightarrow C D) P(D \Rightarrow^* \omega_{kl})$$

Recursion in $P_G(S \Rightarrow^* \omega_{0i} A \omega_{kn})$

$P(S \Rightarrow^* \omega_{0j} C \omega_{kn}) =$

$$\sum_{i=0}^{j-1} \sum_{\substack{A, B \in S \\ A \rightarrow B C E R}} P(S \Rightarrow^* \omega_{0i} A \omega_{kn}) p(A \rightarrow B C) P(B \Rightarrow^* \omega_{ij})$$

$$+ \sum_{l=k+1}^n \sum_{\substack{A, D \in S \\ A \rightarrow C D E R}} P(S \Rightarrow^* \omega_{0j} A \omega_{kn}) p(A \rightarrow C D) P(D \Rightarrow^* \omega_{kl})$$



EM for SCFG:

Infer 'hidden structure' by ML of visible data.

1. guess initial rule probs.

2. repeat:

(a) parse a sample of sentences

(b) weight each parse by its cond. prob

(c) count rules used in each weighted parse, and estimate rule frequencies

- EM optimizes the likelihood of data D
- Each iteration is guaranteed not to decrease the likelihood of D . — but EM can get trapped in local minima
- The Inside-Outside algorithm can produce the expected counts without enumerating all parses

(4)

Important to realize that EM
follows from the Free Energy formula
 $-\log P(D|\theta)$

$$D = \{x_1, \dots, x_n\} \quad \leftarrow \text{visible data}$$
$$P(D|\theta) = \prod_{i=1}^n P(x_i|\theta)$$

$$P(D|\theta) = P(D, y|\theta) = P(D|y, \theta) P(y)$$

\uparrow
hidden

$$F(q, \theta) = \sum_y q(y) \log q(y) - \sum_y q(y) \log P(y, D|\theta)$$

E-step: compute: $q^{t+1}(y) = P(y|D, \theta^t)$
 $= \frac{P(y, D|\theta^t)}{P(D|\theta^t)}$

M-step: compute:
 $\theta^{t+1} = \underset{\theta}{\text{arg min}} \left(- \sum_y q^{t+1}(y) \log P(y, D|\theta) \right)$