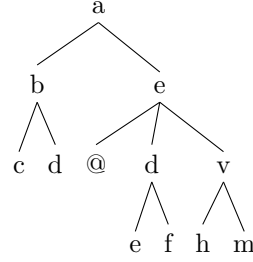


Expectation-Maximization in FSTA
The skeleton

I. Definitions

- Ranked alphabet: $\Sigma = \bigcup_k \Sigma_k$
 $\forall i \geq 0, \Sigma_i \subseteq \Sigma$
- $x[] \in T_\Sigma$ if $x \in \Sigma_0$
- $x[t_1, t_2, t_3 \dots t_k] \in T_\Sigma$ if $x \in \Sigma_k$ and $t_1, t_2, t_3 \dots t_k \in T_\Sigma$
- $\cdot \rightarrow$ context c: the context at the position of node n.



Context 1 = Nonroot (Node e) (context 2) ($[]$) ($([\boxed{d[e[], f[]]}, \boxed{v[h[], m[]]}])$)

\uparrow \uparrow
 t_d t_v

Context 2 = Nonroot (Node a) (context 3) ($([\boxed{b[c[], d[]]}])$)

\uparrow
 t_b

Context 3 = Root

II. Boolean Derivational Path

For any FSTA G with $G = \langle Q, \Sigma, I, \Delta \rangle$.

- Q : a finite set of states
- I : a finite set of initial state with $I \in Q$
- $\Delta: \bigcup_k Q^k \times \Sigma_k \times Q$

Define **under**: $T_\Sigma \times Q \rightarrow \{0, 1\}$

Given root node is at state q , whether this tree can be generated by this grammar G

- **under**($x[]$)(q) = $\Delta([], x, q)$ when $x \in \Sigma_0$
- **under**($x[t_1, t_2, t_3 \dots t_k]$)(q) = $\bigvee_{q_1 \in Q, q_2 \in Q \dots q_k \in Q} \text{under}(t_1, q_1) \wedge \text{under}(t_2, q_2) \wedge \dots \wedge \text{under}(t_k, q_k) \wedge \Delta([q_1, q_2 \dots q_k], x, q)$

Define **over**: $\text{context} \times Q \rightarrow \{0, 1\}$

Given a node in a tree is at state q , whether such node can have this context in grammar G

- **over**(a)(q) = $I(q)$ if $a = \text{root}$
- **over**(Nonroot (a) (context a') ($[l_1, l_2, l_3 \dots l_k]$) ($[r_1, r_2, r_3 \dots r_m]$))(q)
 $= \bigvee_{q_1, q_2 \dots q_k \in Q} \bigvee_{p_0 \in Q} \bigvee_{p_1, p_2 \dots p_m \in Q} \text{under}(l_1, q_1) \wedge \text{under}(l_2, q_2) \wedge \dots \wedge \text{under}(l_k, q_k)$
 $\wedge \text{under}(r_1, p_1) \wedge \text{under}(r_2, p_2) \wedge \dots \wedge \text{under}(r_m, p_m)$
 $\wedge \text{over}(\text{context } a', p_0)$
 $\wedge \Delta([q_1, q_2 \dots q_k, q, p_1, p_2 \dots p_m], a, p_0)$

III. Probabilistic Over and Under

- PFSTA

- $\langle Q, \Sigma, I, \Delta \rangle$
 - Q : a finite set of states
 - $I: Q \rightarrow [0, 1]$
 - $\Delta: \bigcup_k (Q \times \Sigma_k \times Q^k \rightarrow [0, 1])$
- *Normalized* iff
 - For all states q_0 , the total probability of leaving q_0 is 1

$$\sum_{\sigma \in \Sigma_k} \sum_{p^k} \Delta(q_0, \sigma, p^k) = 1$$

- $[I(q_1), I(q_2), \dots, I(q_n)]$ is a stochastic vector
- For any FSTA G with $G = \langle Q, \Sigma, I, \Delta \rangle$, define **under**: $T_\Sigma \times Q \rightarrow [0, 1]$
 Given the root node is at state q , the probability this tree can be generated by this grammar G
 - **under**($x[]$)(q) = $\Delta(q, x, [])$ when $x \in \Sigma_0$
 - **under**($x[t_1, t_2, t_3 \dots t_k]$)(q) = $\sum_{q_1 \in Q, q_2 \in Q, \dots, q_k \in Q} (\prod_{i=1}^k \text{under}(t_i, q_i)) \times \Delta(q, x, [q_1, q_2, \dots, q_k])$

Define **over**: $\text{context} \times Q \rightarrow [0, 1]$

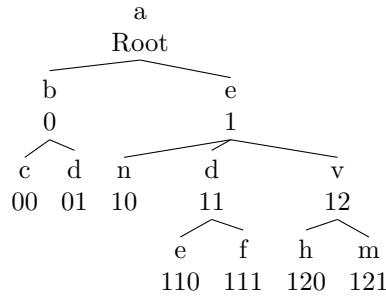
Given a node in a tree is at state q , the probability of such node can have a specific context in grammar G

- **over**(a)(q) = $I(q)$ if $a = \text{root}$
- **over**(Nonroot (a) (context a') ($[l_1, l_2, l_3 \dots l_k]$) ($[r_1, r_2, r_3 \dots r_m]$))(q)
 $= \sum_{q_1, q_2, \dots, q_k \in Q} \sum_{p_0 \in Q} \sum_{p_1, p_2, \dots, p_m \in Q} (\prod_{i=1}^k \text{under}(l_i, q_i)) \times (\prod_{j=1}^m \text{under}(r_j, p_j))$
 $\times \Delta(p_0, a, [q_1, q_2, \dots, q_k, q, p_1, p_2, \dots, p_m]) \times \text{over}(a', p_0)$

IV. Necessary functions

- **NodeIndex** $a = \text{Root} \mid \text{Non-Root Int}$

Extending Huffman encoding: n-ary branching tree: $\{0, 1, 2, 3, 4 \dots n\}$



- **getContext** :: **Tree** → **NodeIndex** → **Context**
- **getSubtree**::**Tree** → **NodeIndex** → [**Tree**]
- **getLabel** :: **Tree** → **NodeIndex** → **Node**
- **isLabel** :: **Tree** → (**Node**, **NodeIndex**) → **Int**
- **indexList**:: **Tree** → [(**NodeIndex**, **Node**)]

V. EM in FSTA

- Problem Setting:

- The complete data: given a set of annotated trees **with the states information** of each nodes, estimate the distribution of Δ and I to maximize the likelihood of the observed data
- The incomplete data: given a set of annotated trees, estimate the probability distribution of Δ and I with the tree-set.

- Formula

- The MLE in PFSTA (Following Chi & Geman 1997)

Given the set of annotated trees $\{t_1, t_2, t_3 \dots t_n\}$ along with states information, the likelihood is

$$\begin{aligned}
 L &= L(p; t_1, t_2, t_3 \dots t_n) \\
 &= \prod_{i=1}^n \prod_{trans \in Q \times \Sigma_k \times Q^k} p(trans)^{freq(trans; w_i)} \\
 \log(L) &= \sum_r \sum_{i=1}^n freq(trans; w_i) \log(p(trans))
 \end{aligned}$$

$$\begin{aligned}
 \forall q \in Q, I(q) &= \frac{count(Q_0=q)}{\sum_{q' \in Q} count(Q_0=q')} \\
 \Delta(q_1, x, [p_1, p_2, p_3 \dots p_k]) &= \frac{count(q_1, x, [p_1, p_2, p_3 \dots p_k])}{\sum_{s_1, s_2 \dots s_k \in Q} \sum_{x' \in \Sigma_k} count(q_1, x', [s_1, s_2, s_3 \dots s_k])}
 \end{aligned}$$

- Expected count

$$\begin{aligned}
 \forall q \in Q, \text{ecount}(Q_0 = q \mid T) &= \Pr(Q_0 = q \mid T) = \frac{I(q) \times \text{under}(T)(q)}{Pr(T)} \\
 &= \sum_{i \in \{Index\}} \Pr(\text{isLabel}(i, \mathbf{x}), Q_i = q, Q_{i+1} = [p_1, \dots, p_k] \mid T) \\
 &= \sum_{i \in \{Index\}} \frac{\text{isLabel}(i, \mathbf{x}) \times \text{over}(\text{getContext}(T, i), q) \times \Delta(q, x, [p_1, p_2, p_3 \dots p_k]) \times \prod_{j \in \text{getSubtree}(T, i)} \text{under}(j, p_j)}{Pr(T)}
 \end{aligned}$$

- Algorithm: Initialize the parameters of Δ and I

while not converge (\leftarrow need threshold)

Calculate the expected counts of each rule under such parametrization

Update the parameters using the expected counts and MLE