

Natural Language Processing

CSE 325/425



Sihong Xie

Lecture 19:

- Training of PCFG
- Lexicalization of PCFG

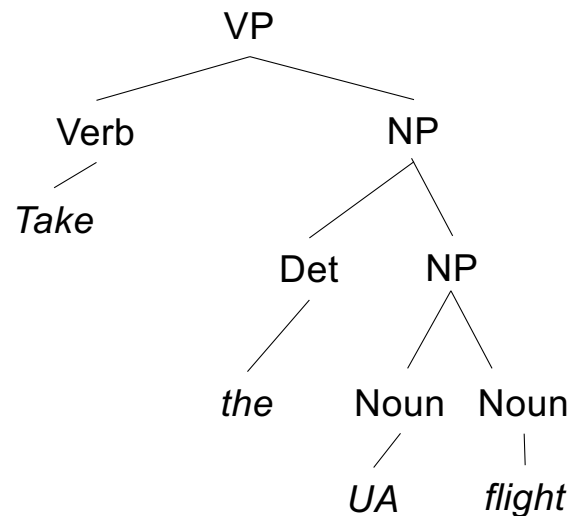
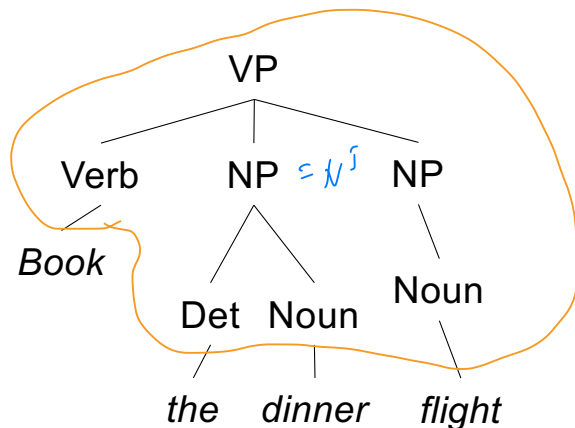
Training of PCFG

CYK with prob. on PCFG

$f = \text{Input} \rightarrow \text{Output}$

- MLE: given a training corpus with sentences and their parsing trees (e.g., Penn Treebank),
 - extract rules observed in the corpus.
 - probability of a rule: how often it appears / how often the LHS appears.

Example: with two training parsing trees



$Pr(NP \rightarrow \text{Noun})$

$$= \frac{\text{Count}(NP \rightarrow \text{Noun})}{\text{Count}(NP)}$$

$$= \frac{1}{4}$$

Training of PCFG

What if we have a large number of sentences without being parsed?

- Need to deal with the unknown trees as latent structures;
- Similar to the unknown POS-tag sequences in the learning of HMM.
- EM algorithm. $\beta_j(p, q)$ $\alpha_j(p, q)$ *outside*
 - E-step: run the inside and outside algorithms to find the inside and ~~output~~ probabilities. (Note: in HMM, this is the forward-backward algorithm).
 - M-step: estimate the rule probabilities based on the expectation of frequencies of occurrence using the inside/outside probabilities.

Training of PCFG

EM algorithm

- Focus on the M-step, given the inside and outside probabilities.

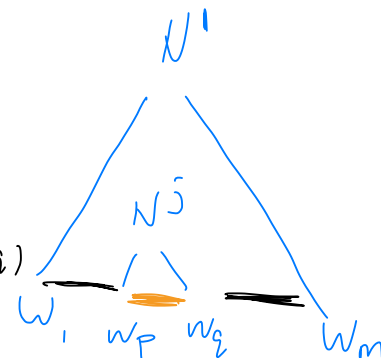
- Probability of a non-terminal N^j covering the range $[p, q]$:

$$1 \leq p \leq q \leq m$$

$$P(N^j \Rightarrow w_{pq} | N^1 \Rightarrow w_{1m}, G) = \frac{\alpha_j(p, q) \beta_j(p, q)}{\pi}$$

← given

$$\pi = \sum_j \alpha_j(p, q) \beta_j(p, q)$$



- Expected frequency of the non-terminal:

$$E(N^j \text{ is used in the derivation}) = \sum_{p=1}^m \sum_{q=p}^m \frac{\alpha_j(p, q) \beta_j(p, q)}{\pi}$$

$N^j \in N$ of PCFG extracted from some labeled corpus.

$$E[f(x)]$$

$$= \sum_{\pi} p_{\pi}(\pi) f(\pi) = \sum_{\pi} p_{\pi}(\pi)$$

$$(f(\pi) = 1)$$

Training of PCFG

EM algorithm

- Focus on the M-step, given the inside and outside probabilities.

dependent on previous est. of rule prob.

- Probability of a rule deriving words in the range $[p, q]$:

$$P(N^j \rightarrow N^r N^s \xRightarrow{*} w_{pq} | N^1 \xRightarrow{*} w_{1m}, G)$$

$$= \frac{\sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\pi}$$

Expected frequency of a rule:

$$E(N^j \rightarrow N^r N^s, N^j \text{ used})$$

$$= \frac{\sum_{p=1}^{m-1} \sum_{q=p+1}^m \sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\pi}$$

$$1 \leq p \leq q \leq m$$

- Estimation of the conditional probability of RHS given the LHS:

$$\hat{P}(N^j \rightarrow N^r N^s) = \frac{\sum_{p=1}^{m-1} \sum_{q=p+1}^m \sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\sum_{p=1}^m \sum_{q=p}^m \alpha_j(p, q) \beta_j(p, q)}$$

① from prev slide

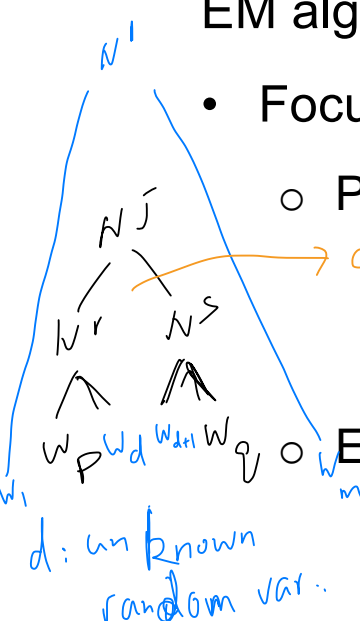
② " π " is canceled out.

$$\sum_p \sum_q P(N^j \rightarrow N^r N^s)$$

$$\xrightarrow{*} w_{pq}$$

$$| N^1 \xRightarrow{*} w_{1m}$$

$$G)$$



EM algorithm

EM algorithm

- Focus on the M-step, given the inside and outside probabilities.
 - Probability of the rule $N^j \rightarrow w^k$ deriving the word w^k somewhere in the sentence:

$$\begin{aligned} P(N^j \rightarrow w^k | N^1 \xRightarrow{*} w_{1m}, G) &= \frac{\sum_{h=1}^m \alpha_j(h, h) P(N^j \rightarrow w_h, w_h = w^k)}{\pi} \\ &= \frac{\sum_{h=1}^m \alpha_j(h, h) P(w_h = w^k) \beta_j(h, h)}{\pi} \end{aligned}$$

- Conditional probability of w^k given N^j :

$$\hat{P}(N^j \rightarrow w^k) = \frac{\sum_{h=1}^m \alpha_j(h, h) P(w_h = w^k) \beta_j(h, h)}{\sum_{p=1}^m \sum_{q=p}^m \alpha_j(p, q) \beta_j(p, q)}$$

↖ from slide 4 and "Ti" is
canceled out

Questioning the PCFG assumptions

- Properties of PCFG:

- Place invariance $\forall k \ P(N_{k(k+c)}^j \rightarrow \zeta)$ is the same
- Context-free $P(N_{kl}^j \rightarrow \zeta | \text{anything outside } k \text{ through } l) = P(N_{kl}^j \rightarrow \zeta)$
- Ancestor-free $P(N_{kl}^j \rightarrow \zeta | \text{any ancestor nodes outside } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$

- The probability of applying a rule is not independent of place

subject

- (a) **She's** able to take her baby to work with her. $\text{Pr}(\text{subject: NP} \rightarrow \text{Pronoun}) = 91\%$
- (b) Uh, **my wife** worked until we had a family. $\text{Pr}(\text{subject: NP} \rightarrow \text{Det Noun}) = 9\%$

object

- (a) Some laws absolutely prohibit **it**. $\text{Pr}(\text{object: NP} \rightarrow \text{Pronoun}) = 34\%$
- (b) All the people signed **confessions**. $\text{Pr}(\text{object: NP} \rightarrow \text{Det Noun}) = 66\%$

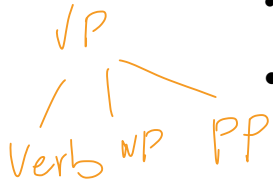
Questioning the assumptions

- Properties of PCFG:

- Place invariance $\forall k \ P(N_{k(k+c)}^j \rightarrow \zeta)$ is the same
- Context-free $P(N_{kl}^j \rightarrow \zeta | \text{anything outside } k \text{ through } l) = P(N_{kl}^j \rightarrow \zeta)$
- Ancestor-free $P(N_{kl}^j \rightarrow \zeta | \text{any ancestor nodes outside } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$

- The probability of applying a rule is not independent of context (words)

- “Moscow sent more than 100,000 soldiers into Afghanistan”
- should (“into Afghanistan”) be attached to “sent” (VP attachment) or “more than 100,000 soldiers” (NP attachment)?
- $\text{Pr}(\text{NP attachment}) > \text{Pr}(\text{VP attachment})$ according to training corpora,
- but in this sentence, it should be a VP attachment.



Questioning the assumptions

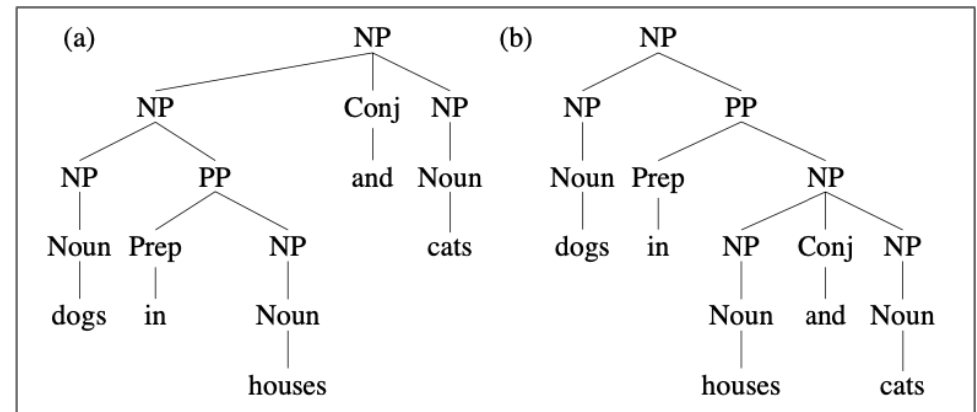
- Properties of PCFG:

- Place invariance $\forall k \ P(N_{k(k+c)}^j \rightarrow \zeta)$ is the same
- Context-free $P(N_{kl}^j \rightarrow \zeta | \text{anything outside } k \text{ through } l) = P(N_{kl}^j \rightarrow \zeta)$
- Ancestor-free $P(N_{kl}^j \rightarrow \zeta | \text{any ancestor nodes outside } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$

- The probability of applying a rule is not independent of context (words)

- Similar problem with conjunctions

- "dogs in houses and cats"

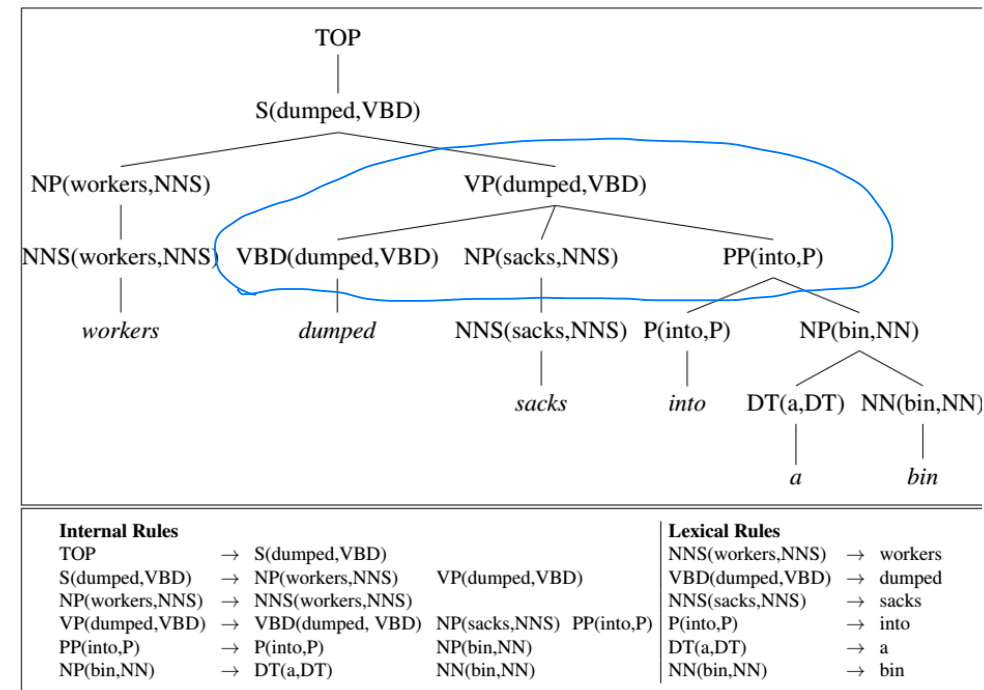


$$\begin{aligned} & \text{h/w 5 : } \Pr(\text{Tree 1}) \\ & = \Pr(\text{Tree 2}) \quad \text{PCFG} \end{aligned}$$

Probabilistic lexicalized CFG

$VP \rightarrow VBD \ NP \ PP$
 $VP \rightarrow VBD \ NP$

- Make the probability of a rule specific to the context (lexicons or words).
 - so that a parsing rule is selected depending on the contexts.
- A non-terminals in a rule adds:
 - head-word: the semantic center of the constituency;
 - head-POS-tag: how the head word is used.
- Copy the original ~~CFG~~ multiple times for all possible lexicalizations.



Probabilistic lexicalized CFG

- Problem with estimating the probabilities of the lexicalized rules
 - $VP(dumped, VBD) \rightarrow VBD(dumped, VBD) NP(sacks, NNS) PP(into, P)$
 - $\Pr(VP(dumped, VBD) \rightarrow VBD(dumped, VBD) NP(sacks, NNS) PP(into, P))$
 $= \frac{\text{Count of } (VP(dumped, VBD) \rightarrow VBD(dumped, VBD) NP(sacks, NNS) PP(into, P))}{\text{Count of } (VP(dumped, VBD))}$
 - How often, in the training corpus, can you see the expansion $VP \rightarrow VBD NP PP$, and the head word is “*dumped*”, which is POS-tagged as VBD, and the head word of NP is “*sacks*”, which is POS-tagged as NNS, and the head word of PP is “*into*”, which is POS-tagged as Preposition?
 - Data sparsity \Rightarrow need to make some independence assumption.
 - Same spirit as how HMM decomposes the joint probability of a long sequence.

$$\Pr(q_1 \rightarrow q_2 \rightarrow q_3 \rightarrow \dots \rightarrow q_T) = \prod_{t=1}^{T-1} \Pr(q_{t+1} | q_t)$$

$q_1 \perp\!\!\!\perp q_3 \mid q_2$

$\underbrace{\text{STOP}}_{\text{Independent events}} \leftarrow H \rightarrow \underbrace{\text{NP(sacks, NNS) PP(into, P)}}_{\text{Independent events}} \text{STOP}$

Collins Parser

- Specifies a particular way to calculate the probability of a rule

- $VP(\text{dumped}, VBD) \rightarrow VBD(\text{dumped}, VBD) \text{ NP(sacks, NNS) PP(into, P)}$
 $\Rightarrow VP(\text{dumped}, VBD) \rightarrow \text{STOP } \boxed{VBD(\text{dumped}, VBD)} \text{ NP(sacks, NNS) PP(into, P) STOP}$
- $\Pr(VP(\text{dumped}, VBD) \rightarrow \text{STOP } VBD(\text{dumped}, VBD) \text{ NP(sacks, NNS) PP(into, P)}) \text{ STOP}$

$$= \Pr_H(VBD | VP, \text{dumped}) \times$$

$$\Pr_L(\text{STOP} | VP, VBD, \text{dumped}) \times$$

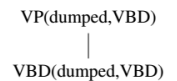
$$\Pr_R(\text{NP(sacks, NNS)} | VP, VBD, \text{dumped}) \times$$

$$\Pr_R(\text{PP(into, P)} | VP, VBD, \text{dumped}) \times$$

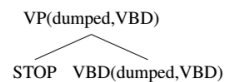
$$\Pr_R(\text{STOP} | VP, VBD, \text{dumped})$$

- Each of the above probability can be estimated using MLE
 - Less subject to data sparsity
 - since the joint events are more likely.

1. Generate the head $VBD(\text{dumped}, VBD)$ with probability
 $P(H | LHS) = P(VBD(\text{dumped}, VBD) | VP(\text{dumped}, VBD))$



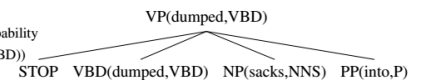
2. Generate the left dependent (which is STOP, since there isn't one) with probability
 $P(\text{STOP} | VP(\text{dumped}, VBD) VBD(\text{dumped}, VBD))$



3. Generate right dependent NP(sacks, NNS) with probability
 $P_r(\text{NP(sacks, NNS)} | VP(\text{dumped}, VBD), VBD(\text{dumped}, VBD))$



4. Generate the right dependent PP(into, P) with probability
 $P_r(\text{PP(into, P)} | VP(\text{dumped}, VBD), VBD(\text{dumped}, VBD))$



5. Generate the right dependent STOP with probability
 $P_r(\text{STOP} | VP(\text{dumped}, VBD), VBD(\text{dumped}, VBD))$

