



人工知能特論II 第6回

二宮 崇

今日の講義の予定

- 確率的文法
 - 品詞解析
 - HMM
 - 構文解析
 - PCFG
- 教科書
 - 北研二(著) 辻井潤一(編) 言語と計算4 確率的言語モデル 東大出版会
 - C. D. Manning & Hinrich Schütze “FOUNDATIONS OF STATISTICAL NATURAL LANGUAGE PROCESSING” MIT Press, 1999
 - Christopher M. Bishop “PATTERN RECOGNITION AND MACHINE LEARNING” Springer, 2006



品詞解析

PART-OF-SPEECH (POS) TAGGING



品詞解析

- 品詞タガー

“I have a pen.”

トーカーナイザー

I have a pen .

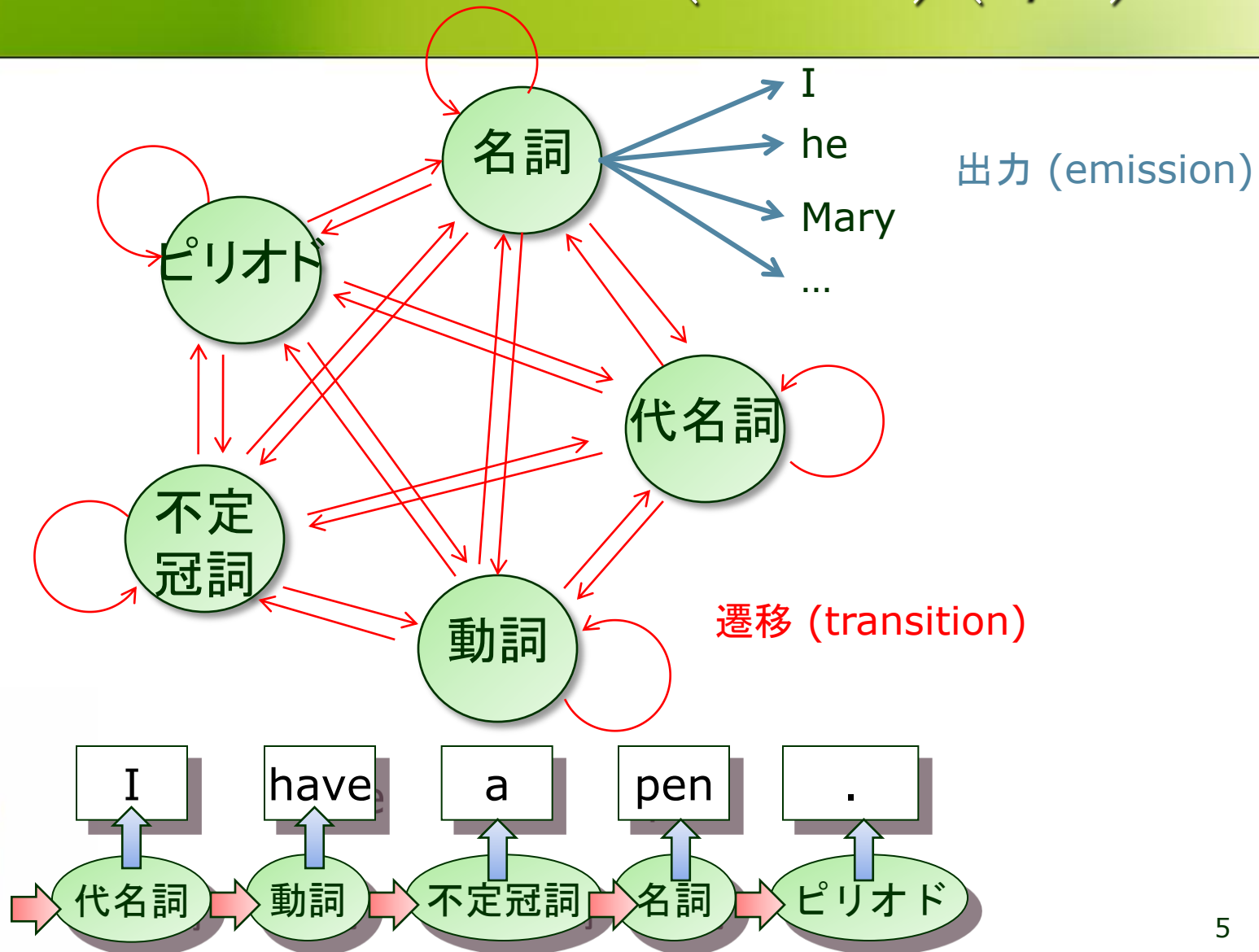
POSタガー

I have a pen .
代名詞 動詞 不定冠詞 名詞 ピリオド



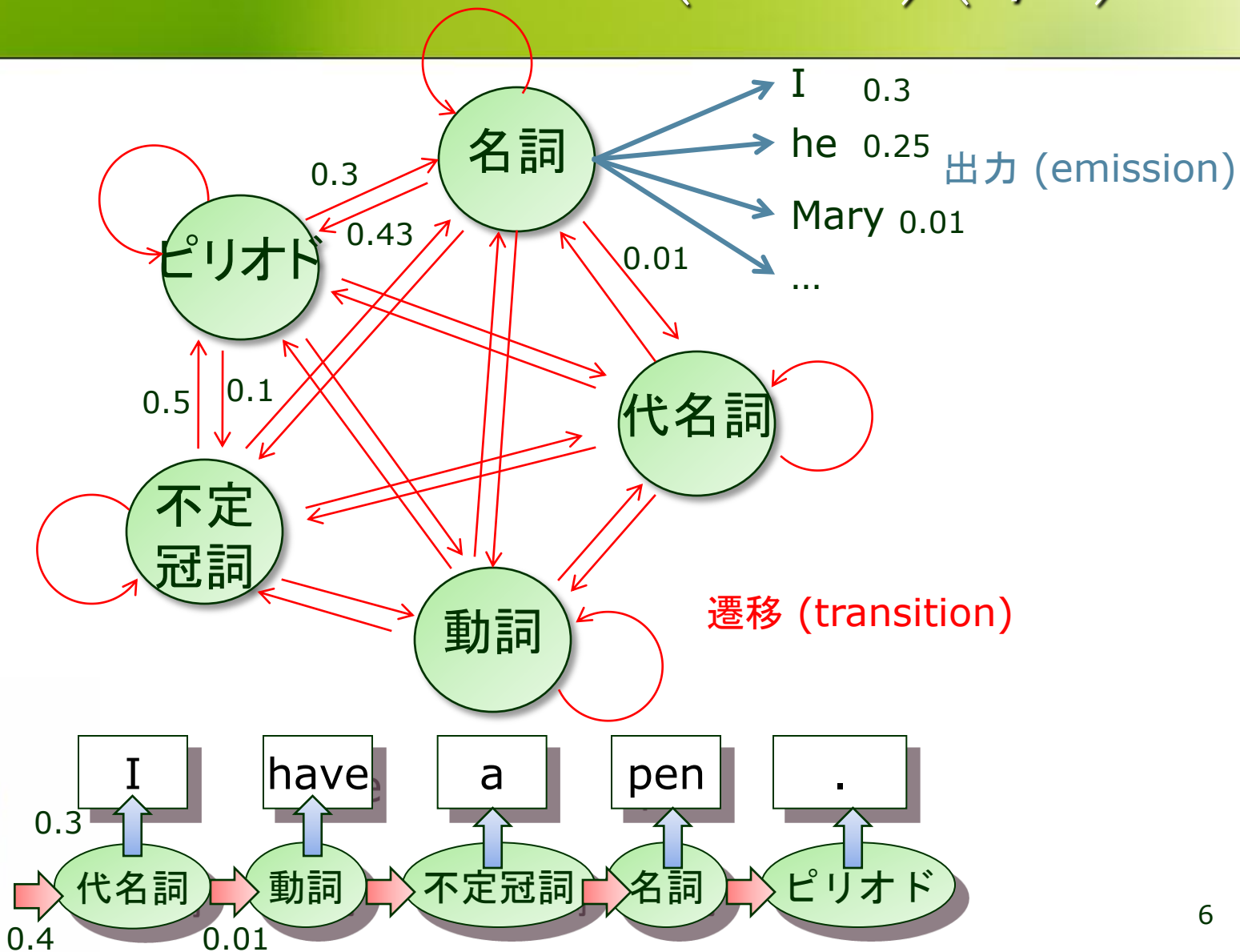
隠れマルコフモデル

Hidden Markov Model (HMM) (1/3)



隠れマルコフモデル

Hidden Markov Model (HMM) (2/3)



隠れマルコフモデル

Hidden Markov Model (HMM) (3/3)

- Q : 状態の有限集合
- Σ : 出力記号の有限集合
- π_q : 文頭が状態 q になる確率
- $a_{q,r}$: 状態 q から状態 r への遷移確率
 - $\sum_{r \in Q} a_{q,r} = 1$
- $b_{q,o}$: 状態 q における記号 o の出力確率
 - $\sum_{o \in \Sigma} b_{q,o} = 1$



状態記号列に対する確率の例

π

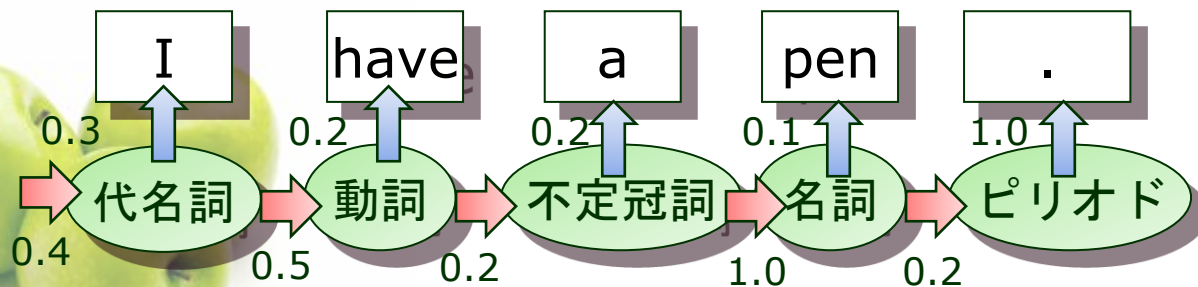
名詞	代名詞	動詞	不定冠詞	ピリオド
0.2	0.4	0.1	0.3	0.0

a

遷移元 \ 遷移先	名詞	代名詞	動詞	不定冠詞	ピリオド
名詞	0.5	0.0	0.3	0.0	0.2
代名詞	0.1	0.0	0.5	0.2	0.2
動詞	0.2	0.2	0.2	0.2	0.2
不定冠詞	1.0	0.0	0.0	0.0	0.0
ピリオド	0.5	0.0	0.5	0.0	0.0

b

状態 \ 出力	a	...	have	...	I	...	pen
名詞	0.01	...	0.0	...	0.01	...	0.1	...	0.0	...
代名詞	0.0	...	0.0	...	0.3	...	0.0	...	0.0	...
動詞	0.0	...	0.2	...	0.0	...	0.0	...	0.0	...
不定冠詞	0.2	...	0.0	...	0.0	...	0.0	...	0.0	...
ピリオド	0.0	...	0.0	...	0.0	...	0.0	...	1.0	...



$$0.4 * 0.3 * 0.5 * 0.2 * 0.2 * 0.2 * 1.0 * 0.1 * 0.2 * 1.0 = 0.0000096$$

状態記号列の確率と 生成確率

- 状態と記号の列が与えられた時

状態記号列: $q_1 o_1 q_2 o_2 \cdots q_T o_T$

$$\begin{aligned} p(q_1 o_1 q_2 o_2 \cdots q_T o_T) &= \pi_{q_1} b_{q_1, o_1} a_{q_1, q_2} b_{q_2, o_2} \cdots a_{q_{T-1}, q_T} b_{q_T, o_T} \\ &= \pi_{q_1} a_{q_1, q_2} \cdots a_{q_{T-1}, q_T} b_{q_1, o_1} b_{q_2, o_2} \cdots b_{q_T, o_T} \\ &= \pi_{q_1} \prod_{t=2}^T a_{q_{t-1}, q_t} \prod_{t=1}^T b_{q_t, o_t} \end{aligned}$$

- 記号列のみが与えられた時

記号列: $o_1 o_2 \cdots o_T$

$$p(o_1 o_2 \cdots o_T) = \sum_{q_1 \in Q, q_2 \in Q, \cdots, q_T \in Q} p(q_1 o_1 q_2 o_2 \cdots q_T o_T) \quad (\text{生成確率})$$



品詞解析

- 品詞解析 (入力: $o_1o_2\dots o_T$)

$$\tilde{q}_1\tilde{q}_2\cdots\tilde{q}_T = \arg \max_{q_1 \in Q, q_2 \in Q, \dots, q_T \in Q} p(q_1o_1q_2o_2\cdots q_To_T)$$



HMMの教師無し学習

Unsupervised Learning of HMMs

- パラメータ推定

- 訓練データ (入力)

訓練データ: $x_1 x_2 \cdots x_n$

x_i : 文を表す記号列(単語列)。 $x_i = o_{i1} o_{i2} o_{i3} \cdots o_{iT_i}$ とする。

T_i : x_i の記号列長

- パラメータ (出力)

$$\begin{aligned}\pi, a, b &= \arg \max_{\pi, a, b} \prod_{i=1}^n p(x_i) \\ &= \arg \max_{\pi, a, b} \prod_{i=1}^n p(o_{i1} o_{i2} \cdots o_{iT_i}) \\ &= \arg \max_{\pi, a, b} \prod_{i=1}^n \sum_{q_1 \in Q, q_2 \in Q, \dots, q_{T_i} \in Q} p(q_1 o_{i1} q_2 o_{i2} \cdots q_{T_i} o_{iT_i})\end{aligned}$$



HMMの教師付学習

Supervised Learning of HMMs

- パラメータ推定

- 訓練データ (入力)

訓練データ: $x_1 x_2 \cdots x_n, y_1 y_2 \cdots y_n$

x_i : 文を表す記号列(単語列)。 $x_i = o_{i1} o_{i2} o_{i3} \cdots o_{iT_i}$ とする。

T_i : x_i の記号列長

y_i : x_i に対応する正解状態列(正解品詞列)。 $y_i = q_{i1} q_{i2} q_{i3} \cdots q_{iT_i}$ とする。

- パラメータ (出力)

$$\begin{aligned} \pi, a, b &= \arg \max_{\pi, a, b} \prod_{i=1}^n p(x_i, y_i) \\ &= \arg \max_{\pi, a, b} \prod_{i=1}^n p(q_{i1} o_{i1} q_{i2} o_{i2} \cdots q_{iT_i} o_{iT_i}) \end{aligned}$$



確率的文脈自由文法 (PCFG)

PROBABILISTIC CONTEXT FREE GRAMMAR (PCFG)



CFG: 構文木

簡単なCFGの例

S → SUBJ VP1

S → SUBJ V

SUBJ → NP が

VP1 → OBJ1 V

OBJ1 → NP を

NP → S NP

V → 送った

V → 読んだ

NP → 香織

NP → 恵

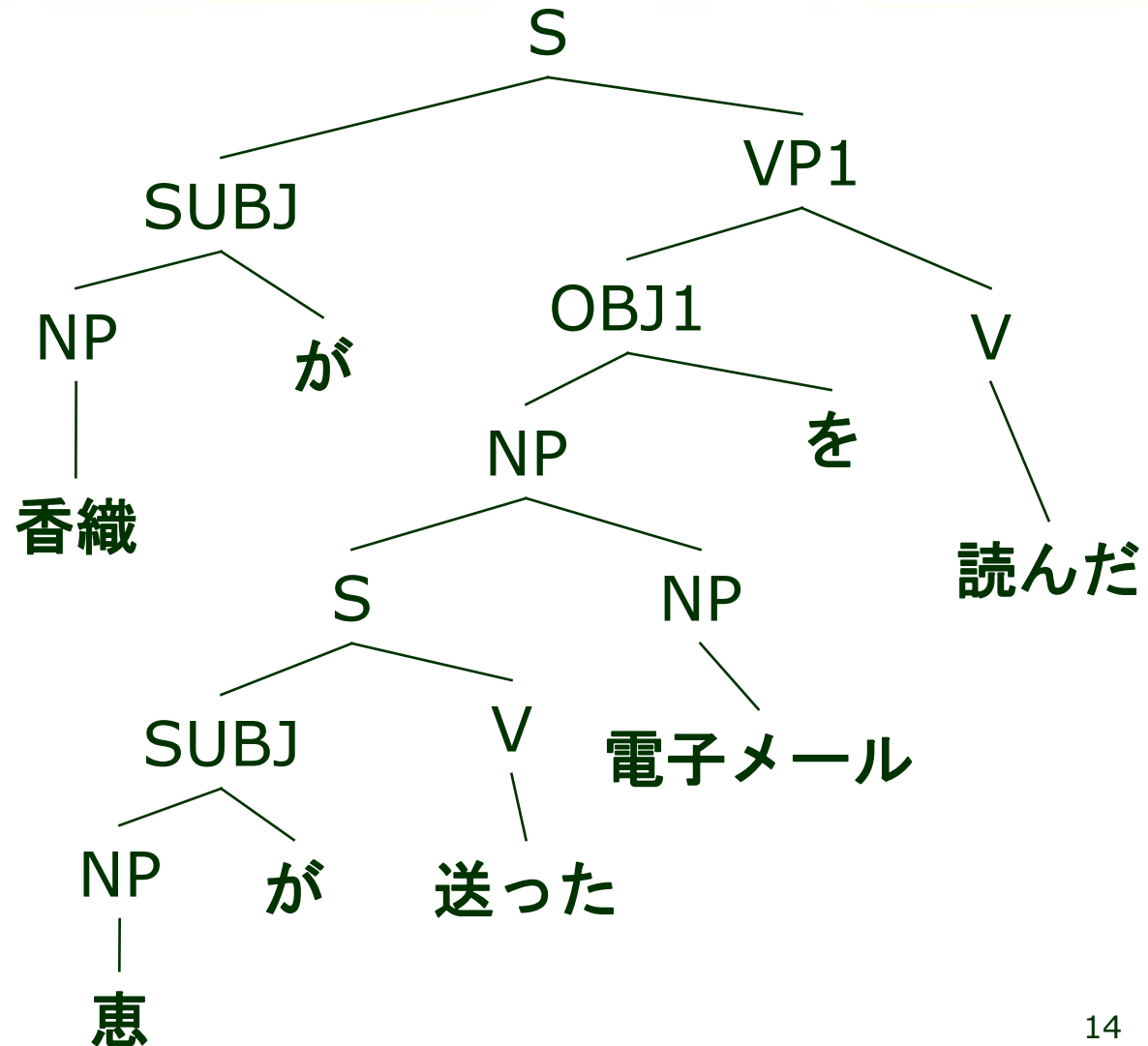
NP → 電子メール

NP → プレゼント

NP → 香織 NP1

NP → 恵 NP1

NP1 → と NP



確率的CFG (PCFG)

- CFGの書換規則の適用確率をパラメータ化した文法
- 構文木の確率は、適用された書換規則のパラメータの積
- 各パラメータは $0.0 \leq \theta \leq 1.0$

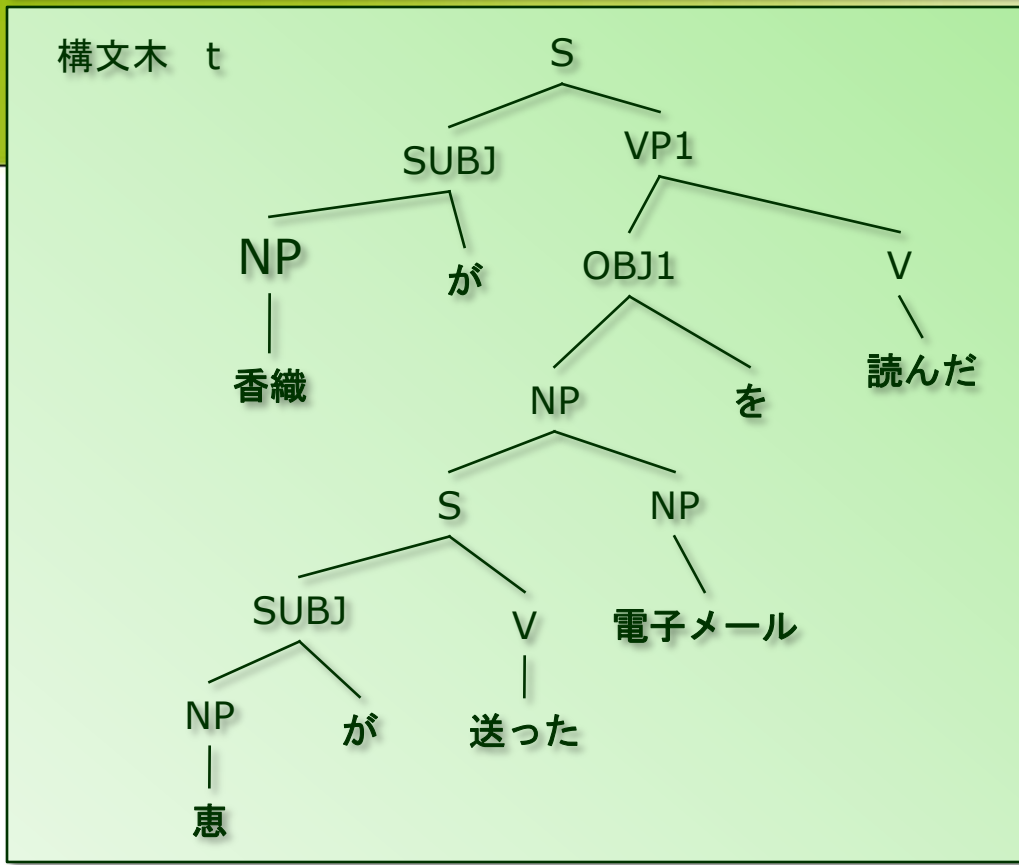
簡単なCFGの例	パラメータ
S → SUBJ VP1	$\theta_{S \rightarrow \text{SUBJ VP1}}$
S → SUBJ V	$\theta_{S \rightarrow \text{SUBJ V}}$
SUBJ → NP が	$\theta_{\text{SUBJ} \rightarrow \text{NP が}}$
VP1 → OBJ1 V	$\theta_{\text{VP1} \rightarrow \text{OBJ1 V}}$
OBJ1 → NP を	$\theta_{\text{OBJ1} \rightarrow \text{NP を}}$
NP → S NP	$\theta_{\text{NP} \rightarrow \text{S NP}}$
V → 送った	$\theta_{V \rightarrow \text{送った}}$
V → 読んだ	$\theta_{V \rightarrow \text{読んだ}}$
NP → 香織	$\theta_{\text{NP} \rightarrow \text{香織}}$
NP → 恵	$\theta_{\text{NP} \rightarrow \text{恵}}$
NP → 電子メール	$\theta_{\text{NP} \rightarrow \text{電子メール}}$
NP → プレゼント	$\theta_{\text{NP} \rightarrow \text{プレゼント}}$
NP → 香織 NP1	$\theta_{\text{NP} \rightarrow \text{香織 NP1}}$
NP → 恵 NP1	$\theta_{\text{NP} \rightarrow \text{恵 NP1}}$
NP1 → と NP	$\theta_{\text{NP1} \rightarrow \text{と NP}}$



構文木の確率

文 x = "香織が恵が送った電子メールを読んだ"

簡単なCFGの例	パラメータ
S → SUBJ VP1	$\theta_{S \rightarrow \text{SUBJ VP1}}$
S → SUBJ V	$\theta_{S \rightarrow \text{SUBJ V}}$
SUBJ → NP が	$\theta_{\text{SUBJ} \rightarrow \text{NP が}}$
VP1 → OBJ1 V	$\theta_{\text{VP1} \rightarrow \text{OBJ1 V}}$
OBJ1 → NP を	$\theta_{\text{OBJ1} \rightarrow \text{NP を}}$
NP → S NP	$\theta_{\text{NP} \rightarrow \text{S NP}}$
V → 送った	$\theta_{V \rightarrow \text{送った}}$
V → 読んだ	$\theta_{V \rightarrow \text{読んだ}}$
NP → 香織	$\theta_{\text{NP} \rightarrow \text{香織}}$
NP → 恵	$\theta_{\text{NP} \rightarrow \text{恵}}$
NP → 電子メール	$\theta_{\text{NP} \rightarrow \text{電子メール}}$
NP → プレゼント	$\theta_{\text{NP} \rightarrow \text{プレゼント}}$
NP → 香織 NP1	$\theta_{\text{NP} \rightarrow \text{香織 NP1}}$
NP → 恵 NP1	$\theta_{\text{NP} \rightarrow \text{恵 NP1}}$
NP1 → と NP	$\theta_{\text{NP1} \rightarrow \text{と NP}}$



$$\begin{aligned}
 P(t) = & \theta_{S \rightarrow \text{SUBJ VP1}} \times \theta_{\text{SUBJ} \rightarrow \text{NP が}} \times \theta_{\text{NP} \rightarrow \text{香織}} \times \\
 & \theta_{\text{VP1} \rightarrow \text{OBJ1 V}} \times \theta_{\text{OBJ1} \rightarrow \text{NP を}} \times \theta_{\text{NP} \rightarrow \text{S NP}} \times \\
 & \theta_{S \rightarrow \text{SUBJ V}} \times \theta_{\text{SUBJ} \rightarrow \text{NP が}} \times \theta_{\text{NP} \rightarrow \text{恵}} \times \\
 & \theta_{V \rightarrow \text{送った}} \times \theta_{\text{NP} \rightarrow \text{電子メール}} \times \theta_{V \rightarrow \text{読んだ}}
 \end{aligned}$$

書換規則の制約

- Sからの生成
 - S→SUBJ VP1を使う場合
 - S→SUBJ Vを使う場合
- Sからの生成は2通りしかない
ので、
 - $\theta_{S \rightarrow \text{SUBJ VP1}} + \theta_{S \rightarrow \text{SUBJ V}} = 1.0$
- 同様にNPも
 - $\theta_{\text{NP} \rightarrow \text{S NP}} + \theta_{\text{NP} \rightarrow \text{香織}} + \theta_{\text{NP} \rightarrow \text{恵}} + \theta_{\text{NP} \rightarrow \text{電子メール}} + \theta_{\text{NP} \rightarrow \text{プレゼント}} + \theta_{\text{NP} \rightarrow \text{香織 NP1}} + \theta_{\text{NP} \rightarrow \text{恵 NP1}} = 1.0$

簡単なCFGの例	パラメータ
S → SUBJ VP1	$\theta_S \rightarrow \text{SUBJ VP1}$
S → SUBJ V	$\theta_S \rightarrow \text{SUBJ V}$
SUBJ → NP が	$\theta_{\text{SUBJ}} \rightarrow \text{NP が}$
VP1 → OBJ1 V	$\theta_{\text{VP1}} \rightarrow \text{OBJ1 V}$
OBJ1 → NP を	$\theta_{\text{OBJ1}} \rightarrow \text{NP を}$
NP → S NP	$\theta_{\text{NP}} \rightarrow \text{S NP}$
V → 送った	$\theta_V \rightarrow \text{送った}$
V → 読んだ	$\theta_V \rightarrow \text{読んだ}$
NP → 香織	$\theta_{\text{NP}} \rightarrow \text{香織}$
NP → 恵	$\theta_{\text{NP}} \rightarrow \text{恵}$
NP → 電子メール	$\theta_{\text{NP}} \rightarrow \text{電子メール}$
NP → プレゼント	$\theta_{\text{NP}} \rightarrow \text{プレゼント}$
NP → 香織 NP1	$\theta_{\text{NP}} \rightarrow \text{香織 NP1}$
NP → 恵 NP1	$\theta_{\text{NP}} \rightarrow \text{恵 NP1}$
NP1 → と NP	$\theta_{\text{NP1}} \rightarrow \text{と NP}$



書換規則の制約

- CFG書換規則を $A \rightarrow \alpha$ と表したとき、(Aは非終端記号、 α は非終端記号列)すべての非終端記号Aに対し、

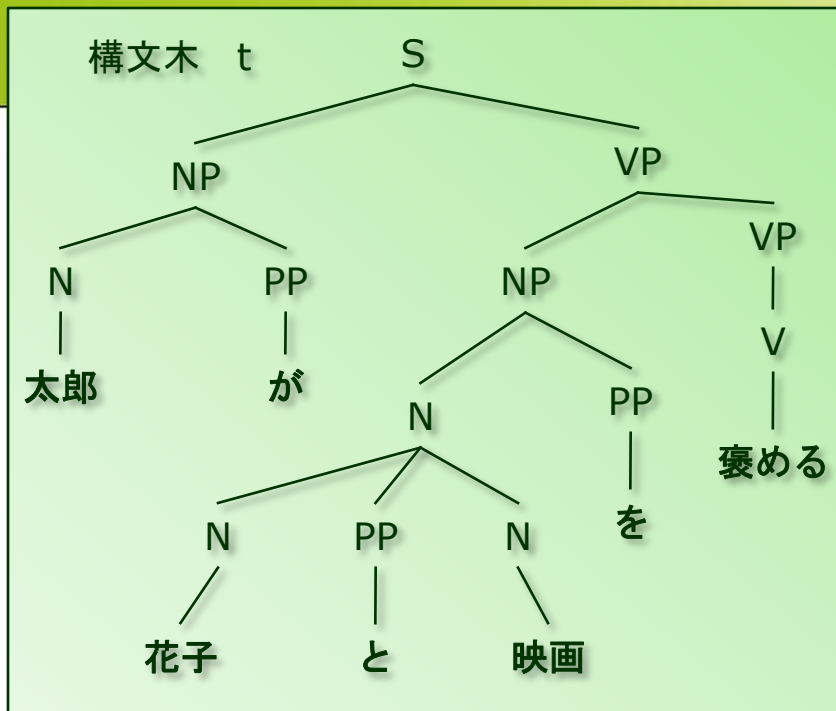
$$\sum_{\alpha} \theta_{A \rightarrow \alpha} = 1.0$$

とする。



構文木の確率

CFG G	パラメータ	値
S → NP VP	$\theta_{S \rightarrow NP VP}$	1.0
NP → N PP	$\theta_{NP \rightarrow N PP}$	1.0
N → N PP N	$\theta_{N \rightarrow N PP N}$	0.1
VP → NP VP	$\theta_{VP \rightarrow NP VP}$	0.3
VP → V	$\theta_{VP \rightarrow V}$	0.7
PP → が	$\theta_{PP \rightarrow が}$	0.5
PP → を	$\theta_{PP \rightarrow を}$	0.3
PP → と	$\theta_{PP \rightarrow と}$	0.2
N → 太郎	$\theta_{N \rightarrow 太郎}$	0.3
N → 花子	$\theta_{N \rightarrow 花子}$	0.2
N → 映画	$\theta_{N \rightarrow 映画}$	0.4
V → 褒める	$\theta_{V \rightarrow 褒める}$	0.3
V → 見る	$\theta_{V \rightarrow 見る}$	0.7



$$\begin{aligned}
 p(t) &= \theta_{S \rightarrow NP VP} \times \theta_{NP \rightarrow N PP} \times \theta_{N \rightarrow 太郎} \times \\
 &\quad \theta_{PP \rightarrow が} \times \theta_{VP \rightarrow NP VP} \times \theta_{NP \rightarrow N PP} \times \\
 &\quad \theta_{N \rightarrow N PP N} \times \theta_{N \rightarrow 花子} \times \theta_{PP \rightarrow と} \times \\
 &\quad \theta_{N \rightarrow 映画} \times \theta_{PP \rightarrow を} \times \theta_{VP \rightarrow V} \times \theta_{V \rightarrow 褒める} \\
 &= 1.0 \times 1.0 \times 0.3 \times 0.5 \times 0.3 \times 1.0 \times \\
 &\quad 0.1 \times 0.2 \times 0.2 \times 0.4 \times 0.3 \times 0.7 \times 0.3 \\
 &= 0.000004536
 \end{aligned}$$

構文木 t の確率

- $C(r; t)$: CFG $\langle V_N, V_T, P, \sigma \rangle$ の書換規則 $r \in P$ が構文木 t 中で使われた回数

$$p(t) = \prod_{r \in P} \theta_r^{C(r; t)}$$



文の生成確率

- ある文 x に対し、 x を導出する全ての構文木集合を $T(x)$ としたとき、

$$p(x) = \sum_{t \in T(x)} p(t) \quad 0.0 \leq p(x) \leq 1.0$$

- PCFGによって生成されうる全ての構文木集合を T としたとき、

$$\sum_{t \in T} p(t) \approx 1.0$$



構文解析

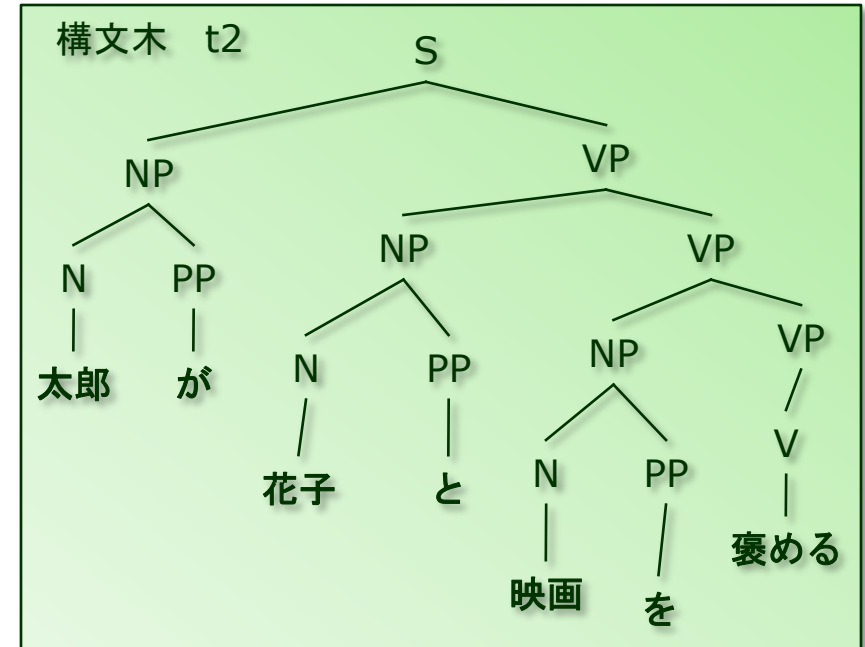
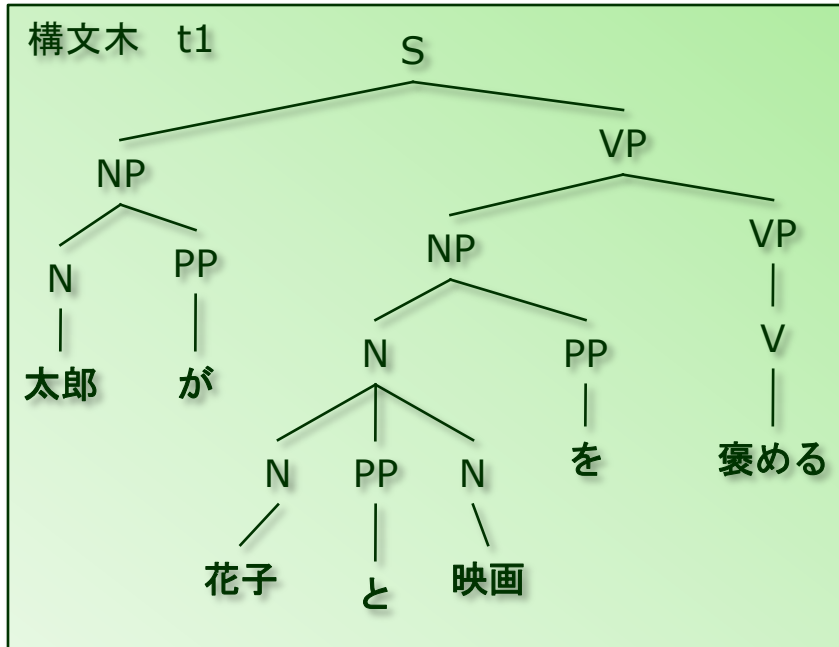
- ある文 x に対し、CFG $\langle V_N, V_T, P, \sigma \rangle$ を用いて x を導出できる全ての構文木集合を $T(x)$ としたとき、

$$\begin{aligned}\tilde{t} &= \arg \max_{t \in T(x)} p(t) \\ &= \arg \max_{t \in T(x)} \prod_{r \in P} \theta_r^{C(r;t)}\end{aligned}$$



構文解析の例

- “太郎が花子と映画を褒める”に対する構文木は次の二種類しかない



$$\begin{aligned}
 p(t1) &= \theta_{S \rightarrow NP VP} \times \theta_{NP \rightarrow N PP}^2 \times \theta_{N \rightarrow \text{太郎}} \times \\
 &\quad \theta_{PP \rightarrow \text{が}} \times \theta_{VP \rightarrow NP VP} \times \theta_{N \rightarrow N PP N} \times \\
 &\quad \theta_{N \rightarrow \text{花子}} \times \theta_{PP \rightarrow \text{と}} \times \theta_{N \rightarrow \text{映画}} \times \\
 &\quad \theta_{PP \rightarrow \text{を}} \times \theta_{VP \rightarrow V} \times \theta_{V \rightarrow \text{褒める}} \\
 &= 1.0 \times 1.0^2 \times 0.3 \times 0.5 \times 0.3 \times 0.1 \times 0.2 \times 0.2 \times \\
 &\quad 0.4 \times 0.3 \times 0.7 \times 0.3 \\
 &= 0.000004536
 \end{aligned}$$

$$\begin{aligned}
 p(t2) &= \theta_{S \rightarrow NP VP} \times \theta_{NP \rightarrow N PP}^3 \times \theta_{N \rightarrow \text{太郎}} \times \\
 &\quad \theta_{PP \rightarrow \text{が}} \times \theta_{VP \rightarrow NP VP}^2 \times \theta_{N \rightarrow \text{花子}} \times \\
 &\quad \theta_{PP \rightarrow \text{と}} \times \theta_{N \rightarrow \text{映画}} \times \theta_{PP \rightarrow \text{を}} \times \\
 &\quad \theta_{VP \rightarrow V} \times \theta_{V \rightarrow \text{褒める}} \\
 &= 1.0 \times 1.0^3 \times 0.3 \times 0.5 \times 0.3^2 \times 0.2 \times 0.2 \times 0.4 \times \\
 &\quad 0.3 \times 0.7 \times 0.3 \\
 &= 0.000013608
 \end{aligned}$$

PCFGの教師無し学習

Unsupervised Learning of PCFG

- パラメータ推定

- 訓練データ (入力)

訓練データ: $x_1 x_2 \cdots x_n$

x_i : 文

- パラメータ (出力)

$$\begin{aligned}\theta &= \arg \max_{\theta} \prod_{i=1}^n p(x_i) \\ &= \arg \max_{\theta} \prod_{i=1}^n \sum_{t \in T(x_i)} p(t) \\ &= \arg \max_{\theta} \prod_{i=1}^n \sum_{t \in T(x_i)} \prod_{r \in P} \theta_r^{C(r;t)}\end{aligned}$$



PCFGの教師付パラメータ推定

Supervised Learning of PCFG

- 教師付学習

- 訓練データ (入力)

訓練データ: $x_1 x_2 \cdots x_n, y_1 y_2 \cdots y_n$

x_i : 文

y_i : x_i に対する正解構文木

- パラメータ (出力)

$$\begin{aligned}\theta &= \arg \max_{\theta} \prod_{i=1}^n p(x_i, y_i) \\ &= \arg \max_{\theta} \prod_{i=1}^n p(y_i) \\ &= \arg \max_{\theta} \prod_{i=1}^n \prod_{r \in P} \theta_r^{C(r; y_i)}\end{aligned}$$



まとめ

- 品詞解析
 - HMM
- 構文解析
 - PCFG
- 資料

<http://aiweb.cs.ehime-u.ac.jp/~ninomiya/ai2/>

