

# ベイズ確率文脈自由文法のための高速構文木サンプリング法 Split Position Slice Sampler

武井俊祐\* 牧野貴樹† 高木利久\*\*‡

Shunsuke Takei\* Takaki Makino† Toshihisa Takagi\*\*‡

**Abstract:** We propose a new tree sampling algorithm for Bayesian probabilistic context-free grammar (PCFG) called Split Position Slice Sampler. Split Position Slice Sampler is developed based on Beam Sampling method that is a fast MCMC sampling algorithm for Bayesian Hidden Markov Model, and adapted to Bayesian PCFG. This tree sampling method can be combined with Metropolis-Hastings sampler to constitute an efficient grammar sampling algorithm for Bayesian PCFG. Because this algorithm does not involve any approximation, more efficient inference is achieved without losing accuracy. We evaluate our approach by comparing with an existing method in a small artificial corpus.

**Keywords:** Bayesian inference, PCFG, MCMC

## 1 Introduction

本研究は、ベイズ拡張された確率文脈自由文法 (Bayesian PCFG) に対する効率的なサンプリングアルゴリズムを構築することで、高速かつ高精度な文法学習を可能にすることを目的とする。Bayesian PCFG モデルに対する精度の高いパラメータ推定手法としては、Johnson らの動的計画法を利用したマルコフ連鎖モンテカルロ法 (MCMC) による手法 [1] が知られているが、モデルが複雑になるにつれて増大する計算コストが問題となる。そこで我々は、ベイズ拡張された隠れマルコフモデル (Bayesian HMM) のための高速な MCMC 法として知られている Beam Sampling [2] を PCFG に応用することで Johnson らの手法を高速化する。Beam Sampling は、動的計画法と Slice Sampling [3] を利用することにより、Bayesian HMM の高速なパラメータ推定を可能にする手法であるが、これをそのままの形式でベイズ PCFG の枠組みへ応用しようとするとき、HMM において各時刻において割り当てられている補助変数をベイズ

PCFG で利用される内側確率表に直接対応付けることができずアルゴリズムが構築できない。本論文では内側確率表における終端位置と分割位置に補助変数を対応付ける形式として Split Position Slice Sampler という手法を新たに提案することで Beam Sampler の枠組みを拡張し、これを Johnson らの手法へ組み込むことによって高精度かつ高速な Bayesian PCFG モデルのパラメータ推定法を構築する。

### 1.1. Motivation

ベイズモデルとは、統計的機械学習分野において学習モデルにベイズ統計の手法を導入するもので、ベイズ化された学習モデルは、与えられたデータに対するモデルパラメータの推定値と不確実性の範囲を確率分布として表現する。結果としてベイズ統計における学習は、従来の最尤法による点推定に比べ汎化性能に優れ、過学習の問題も生じにくい。また、あらかじめ推定したい対象に何らかの知識がある場合、それを事前確率分布という形式でモデルに導入することでより良い推定が可能であるという利点もある。さらに、近年注目を集めているベイズモデルの拡張であるノンパラメトリックベイズモデルでは、無限次元のパラメータ空間を仮定し、データを表現する最適なモデルを推定することでモデル選択の問題を解決することができるため、様々な統計学習モ

\*東京大学大学院 新領域創成科学研究科 情報生命科学専攻, 〒277-8568 千葉県柏市柏の葉 5-1-5 総合研究棟 609, tel. 04-7136-3973,

e-mail: takei\_shunsuke@cb.k.u-tokyo.ac.jp, Department of Computational Biology, Graduate School of Frontier Science, The University of Tokyo, General Research Building 609, 5-1-5 Kashiwa-no-ha, Kashiwa, Chiba, 277-8568 JAPAN

†東京大学 総括プロジェクト機構

‡情報・システム研究機構ライフサイエンス統合データベースセンター

デルに導入されており[4][5][6], 今後より広範な問題に応用されることが期待されている.

ベイズ学習モデルにおいて学習の対象である個々のパラメータは確率変数で表現されているため, 従来の最尤法とは異なりしばしば計算コストの高い高次元の期待値計算を伴う方法でパラメータ推定を行う. 特に高い精度が要求されるような場合においては高精度なパラメータ推定手法として知られるMCMC法が利用されるが, 膨大なサンプルを必要とするため計算コストが高く, ベイズモデルの応用における問題となっている.

PCFG[9]は計算機科学分野だけでなくバイオインフォマティクス分野など[10]幅広い分野で活用されている一般性の高い確率モデルであり, そのベイズモデルである Bayesian PCFG についても期待が高い. 近年, 変分ベイズ法[7][8]のような近似に基づく高速パラメータ推定手法が提案され, さかんに研究されているものの, 高精度な推定を必要とする場合は依然として計算コストの高いMCMCのような手法が用いられている.

Bayesian PCFG に対する MCMC 法を構成する場合, 単純には Gibbs Sampler のような手法の適用が考えられるが, PCFG のパラメータ間に強い相互依存がありパラメータ空間も大きいことから Gibbs Sampler では収束が遅く効率が悪い. これに対し Johnson らは構文木のみのサンプリングによる Metropolis-Hastings Sampler を構築し, より効率的なサンプリング手法を提案している[1]. しかし Johnson らの手法は依然として計算コストが大きく, 大規模データに適用するには困難を伴う.

近年 Bayesian HMM のための高速なサンプリング手法として Beam Sampler が提案された. この手法は Slice Sampling の手法を応用し, 隠れ状態間の遷移確率分布を補助変数(スライサー)で“間引く”ことで計算の高速化を図っている. この“間引き”の処理は近似等の処理ではなくサンプリングの処理の一部であるため, Beam Sampler は精度を損なわず Bayesian HMM のパラメータ推定を高速に行うことができる. HMM における遷移確率分布は PCFG における書き換え規則の確率分布に対応するため, 正しいスライス手法が与えられれば, 同様な高速化を達成できることが期待される.

そこで我々は, Johnson らの手法に Beam Sampler の枠組みを導入することで, 精度を保ったまま高速

に動作する Bayesian PCFG のための MCMC 法を構築することを目指す.

## 2 Background

### 2.1. PCFG

文脈自由文法(CFG)は文の生成過程を明らかにするモデルであり  $G = (V_N, V_T, S, R)$  の 4-タプルで定義される.  $V_T$  は終端記号集合,  $V_N$  は非終端記号集合,  $R$  は生成規則の集合,  $S \in V_N$  は開始記号である. ここでは生成規則についてチョムスキー標準形, すなわち  $A \rightarrow BC$  または  $A \rightarrow s$  の形式の生成規則のみを考える. ここで  $A, B, C \in V_N$ ,  $s \in V_T$  である.

確率文脈自由文法(PCFG)は, CFG の各々のルールに対して確率値を割り振ったものであり,  $(G, \theta)$  と定義される.  $\theta$  は  $|R|$  次元の実数ベクトルであり,  $\theta$  のそれぞれの要素は  $\theta_r$  で表され, これは  $r \in R$  の生成規則の確率値であることを表す. たとえば,

$\theta_{A \rightarrow BC}$  ならば  $A \rightarrow BC$  の規則の確率,  $\theta_{A \rightarrow s}$  ならば  $A \rightarrow s$  の規則の確率を表す. ここで, 確率の定義から  $\theta_r \geq 0$  かつ  $\sum_{a \in V_T \cup V_N \times V_N} \theta_{A \rightarrow a} = 1$  である必要がある.

### 2.2. Bayesian PCFG

Bayesian PCFG では, 生成規則のパラメータ  $\theta$  を確率変数  $P(\theta)$  として扱う. ベイズアプローチにおいて離散パラメータは, その扱いやすさから Dirichlet 分布のような共役な確率分布によって表され, 本論文における Bayesian PCFG についても Dirichlet 分布を導入した PCFG を仮定している. ここで Dirichlet 分布を  $P_D(\cdot)$ ,  $\Gamma$  をガンマ関数,  $A \rightarrow *$  の形式の文法規則を  $R_A$ ,  $\theta$  における  $A \rightarrow *$  の形式の文法規則についてのパラメータベクトル  $\theta_A$ , Dirichlet 分布のハイパーパラメータベクトルを  $\alpha$ ,  $\alpha$  における  $A \rightarrow *$  の形式の文法規則についてのサブベクトルを  $\alpha_A$  で表すと,  $\theta$  の事前分布  $P_D(\theta | \alpha)$  は

$$P_D(\theta | \alpha) = \prod_{A \in V_N} P_D(\theta_A | \alpha_A) \quad (1)$$

のような Dirichlet 分布で表現される. ここで,

$$P_D(\theta_A | \alpha_A) = \frac{1}{C(\alpha_A)} \prod_{r \in R_A} \theta_r^{\alpha_r - 1} \quad (2)$$

$$C(\alpha_A) = \frac{\prod_{r \in R_A} \Gamma(\alpha_r)}{\Gamma(\sum_{r \in R_A} \alpha_r)} \quad (3)$$

である。一般にモデルパラメータは未知でありデータから推定する必要がある。ベイズアプローチにおける推定の対象はデータ観測後の事後確率分布であり、

$$\begin{aligned} P(\boldsymbol{\theta} | \mathbf{w}) &\propto P(\mathbf{w} | \boldsymbol{\theta}) P(\boldsymbol{\theta} | \boldsymbol{\alpha}) \\ &= \prod_{i=1}^n P(w_i | \boldsymbol{\theta}) P(\boldsymbol{\theta} | \boldsymbol{\alpha}) \end{aligned} \quad (4)$$

の左項を求めることである。 $\mathbf{w} = (w_1, \dots, w_n)$  は終端記号列からなるデータ集合であり、 $w_i$  が個別の終端記号列を意味する。

### 2.3. Gibbs Sampler

Bayesian PCFG の事後分布は解析的積分が不可能であり、その推定には MCMC 法や変分ベイズ法などの手法が用いられる。MCMC 法のもっとも単純な手法の 1 つである Gibbs Sampler を Bayesian PCFG のパラメータ推定法として用いる場合、パラメータ  $\boldsymbol{\theta}$  と構文木  $\mathbf{t}$  を交互にサンプリングすることでアルゴリズムが構築される。ここで、 $\mathbf{t} = (t_1, \dots, t_n)$  であり、それぞれの  $t_i$  は終端記号列  $w_i$  に対する構文木を意味する。詳しい導出は Johnson らの論文に譲り、ここでは具体的な Gibbs Sampler のアルゴリズムを直接示す。

構文木  $t$  において文法規則  $r \in R$  が使われた回数を  $f_r(t)$ 、 $\mathbf{t}$  における  $A \rightarrow *$  の形式の文法規則が使われた回数のベクトルを  $\mathbf{f}_A(\mathbf{t})$  とすると、Bayesian PCFG の Gibbs Sampler は、 $\mathbf{t}$  を固定して  $\boldsymbol{\theta}$  についてのサンプリング、

$$P(\boldsymbol{\theta} | \mathbf{t}, \mathbf{w}, \boldsymbol{\alpha}) = \prod_{A \in V_N} P_D(\boldsymbol{\theta}_A | \mathbf{f}_A(\mathbf{t}) + \boldsymbol{\alpha}_A) \quad (5)$$

$\boldsymbol{\theta}$  を固定して  $\mathbf{t}$  についてのサンプリング、

$$P(\mathbf{t} | \boldsymbol{\theta}, \mathbf{w}, \boldsymbol{\alpha}) = \prod_{i=1}^n P(t_i | w_i, \boldsymbol{\theta}) \quad (6)$$

を交互に繰り返すことになる。パラメータ  $\boldsymbol{\theta}$  についてのサンプリングは、構文木  $\mathbf{t}$  における各文法規則の使用回数から容易に計算できるが、構文木  $\mathbf{t}$  の具体的なサンプリングについては自明ではない。

Johnson らは構文木のサンプリングについて、動的計画法を用いた効率的な手法を提案している。まず、

ある終端記号列  $\mathbf{s} = (s_1, \dots, s_n)$  が与えられたとき、その内側確率を

$$\begin{aligned} p_{k,k}^A &= \theta_{A \rightarrow s_k} \\ p_{i,k}^A &= \sum_{A \rightarrow BC} \theta_{A \rightarrow BC} p_{i,j}^B p_{j+1,k}^C \end{aligned} \quad (7)$$

のように計算する。ここで、 $p_{i,k}^A$  は非終端記号  $A$  が終端記号列  $s_i, \dots, s_k$  を生成する確率を意味する。次に、作られた内側確率表の確率にしたがい構文木を再帰的にサンプリングする。これを疑似コードで示すと、

```
Function SAMPLE(A, i, k)
  If i = k
    return TREE(A, s_k)
  Else
    (j, B, C) = MULTI(A, i, k)
    return TREE(A, SAMPLE(B, i, j), SAMPLE(C, j, k))
```

となる。ここで、関数  $\text{SAMPLE}(A, i, k)$  はある構文木のノードにおいて非終端記号  $A$  が最終的に終端記号列  $s_i, \dots, s_k$  を生成するとき、そのノード以下の構文木を確率的にサンプリングする関数である。また、関数  $\text{MULTI}(A, i, k)$  はどの分割点で構文木が枝分かれするのか、その分割位置  $j$  と子ノードの非終端記号  $B, C$  を確率的に返す関数であり、その確率は

$$P(j, B, C) = \frac{\theta_{A \rightarrow BC} p_{i,j}^B p_{j+1,k}^C}{p_{i,k}^A} \quad (9)$$

と表すことができる。

### 2.4. Metropolis-Hastings Sampler

PCFG における各々のパラメータ間には強い依存があるため、パラメータ  $\boldsymbol{\theta}$  を個別にサンプリングする Gibbs Sampler は十分なサンプルを得るために多大な時間が必要になる。この問題に対し Johnson らは Bayesian PCFG のパラメータ  $\boldsymbol{\theta}$  を積分消去し、構文木のサンプリングと Metropolis-Hastings の枠組みによるサンプルの確率的な受理から構成されるサンプリングアルゴリズムを提案している。この手法では、依存の強い  $\boldsymbol{\theta}$  のサンプリングは回避されるため、Gibbs Sampler に比べ速い収束が見込まれる。さらに  $\boldsymbol{\theta}$  のサンプリングにかかる計算コストも不要であるため、効率的な Bayesian PCFG のサンプリングアルゴリズムであるといえる。

このアルゴリズムでは、Gibbs Sampler とは違いパラメータ  $\theta$  のサンプリングを行わず、構文木をサンプリングしたのち、そのサンプルを確率的に受諾もしくは拒否する。いま、 $\mathbf{w} = (w_1, \dots, w_n)$  を終端記号列の  $\mathbf{s}$  の集合とし、それぞれの終端記号列  $w_i$  に対応する構文木を  $\mathbf{t} = (t_1, \dots, t_n)$  で表す。  $t'_i$  が新しくサンプリングされた構文木、  $\mathbf{t}_{-i}$  が  $t_i$  を除いた構文木サンプル群とすると、そのサンプルの受理確率は

$$A(t_i, t'_i) = \min \left\{ 1, \frac{P(t'_i | w_i, \mathbf{t}_{-i}, \alpha) P(t_i | w_i, \theta')}{P(t_i | w_i, \mathbf{t}_{-i}, \alpha) P(t'_i | w_i, \theta')} \right\} \quad (7)$$

$$= \min \left\{ 1, \frac{P(t'_i | \mathbf{t}_{-i}, \alpha) P(t_i | w_i, \theta')}{P(t_i | \mathbf{t}_{-i}, \alpha) P(t'_i | w_i, \theta')} \right\}$$

となる。ここで、

$$P(t_i | \mathbf{t}_{-i}, \alpha) = \prod_{A \in \mathcal{F}_N} \frac{C(\mathbf{a}_A + \mathbf{f}_A(\mathbf{t}))}{C(\mathbf{a}_A + \mathbf{f}_A(\mathbf{t}_{-i}))} \quad (9)$$

である。  $\theta'$  は  $\theta_r'$  を要素に持つベクトルである。  $\theta_r'$  は  $\mathbf{t}_{-i}$  と  $\alpha$  が与えられたときの  $\theta_r$  の期待値で、

$$\theta_r' = \frac{f_r(\mathbf{t}_{-i}) + \alpha_r}{\sum_{r' \in R_A} f_{r'}(\mathbf{t}_{-i}) + \alpha_{r'}} \quad (10)$$

として計算される。

以上が Johnson らの Bayesian PCFG のための構文木サンプリングの枠組みである。この手法はパラメータ  $\theta$  と構文木  $\mathbf{t}$  を交互にサンプリングする Gibbs Sampler に比べ効率的ではあるが、パラメータ空間が大きい場合、内側確率の計算に多大な計算コストがかかり依然として大規模データには不向きである。

## 2.5. Beam Sampler

Beam Sampler は、Bayesian HMM のための高速 MCMC アルゴリズムであり、HMM のパラメータ推定にしばしば用いられる動的計画法(この場合 Forward-Backward)と Slice Sampling の手法を応用することで高速化をしている。HMM のための Gibbs Sampler は、各時刻の隠れ状態を交互にサンプリングする形式で行われるが、Beam Sampler は各時刻の隠れ状態を個別にサンプリングするのではなく、動的計画法を利用し隠れ状態列をひとつのサンプルとする形式で行われるため、Johnson らの手法同様パラ

メータ間の依存の問題が解決される。また Slice Sampling に基づく補助変数 (スライサー) を各時刻について導入し、

1. 前回サンプルされた状態遷移列に従いスライサーをサンプリング
2. 状態遷移分布をスライスすることで、各時刻における可能な状態遷移を間引き
3. 間引かれた分布に従い隠れ状態列をサンプル

という手順でサンプリングを行う。ここでスライスとは、対象となる確率分布においてスライサーの値以上のもののみを考慮するよう分布を“間引く”操作であり、スライスされた分布は等確率でサンプリングされる。この手順により各時刻で考慮すべき状態遷移数が減るために計算が高速に行われる。Beam Sampler における分布のスライスおよびスライスされた分布のサンプリングという処理は近似処理ではなく、Slice Sampling を援用したもので、それ自体がサンプリングの手続きである。そのため Beam Sampler は高精度なパラメータ推定というマルコフ連鎖モンテカルロ法の性質を維持したまま高速化に成功した手法であるといえる。

## 3 Method

### 3.1. PCFG へのスライスサンプリングの導入

我々のアプローチは Johnson らの手法に Beam Sampler の枠組みを導入することで高速化を図るものである。Beam Sampler でサンプルとする隠れ状態列は PCFG における構文木にあたり、Beam Sampler で用いている動的計画法は Johnson らの手法では内側確率を計算し、構文木をサンプリングすることと等価である。ここから Slice Sampling を内側確率の計算に利用すれば Beam Sampler が PCFG へと適用可能である。Johnson らのアルゴリズムが

1. 内側確率表を計算
2. 構文木をサンプリング
3. 構文木を Accept/Reject

という手順で構成されていたのに対し、提案手法は

1. スライサーをサンプリング
2. 生成確率をスライスして、内側確率表の各セルにおける文法規則適用分布を間引き
3. 構文木をサンプリング
4. 構文木を Accept/Reject

という手順で構成されることになる。しかし PCFG と HMM のモデルの違いから Beam Sampler をそのまま適用しようとしてもアルゴリズムは構築できない。

問題となるのは、Beam Sampler における前回のサンプルにおける各時刻の状態遷移確率からスライサーの値をサンプリングするという手続きである。

HMM では各時刻についてスライサーがそれぞれ割り当てられており (Figure 1), これを PCFG にそのまま適用しようとする内側確率表の各セル, つまり考える全ての構文木のノード位置についてスライサーを導入するという形式となる。この場合前回のサンプルにあたるのが構文木であり, スライサーのサンプルに必要なのは, 構文木の各ノードで用いられた文法規則の確率値である。しかしこの場合, 前回のサンプル, すなわちひとつの構文木からすべてのセルのスライサーを作ることができない。なぜなら, サンプルとして得られた構文木は, 内側確率表全てのセルに一一対応するだけノード数を保持しておらず, 従ってこのような配置でスライサーを導入しようとした場合, 全てのスライサーをサンプリングするには情報が足りないからである (Figure 2)。これは  $n$  を文中の単語数とした場合, 構文木は  $2n-1$  個のノードしか持たず, 一方内側確率表には  $n(n+1)/2$  個のセルがあるということからも分かる。

### 3.2. Split Position Slice Sampler

そのため我々はスライサーの配置に関する新たなルールとして, Split Position Slice Sampler という形式

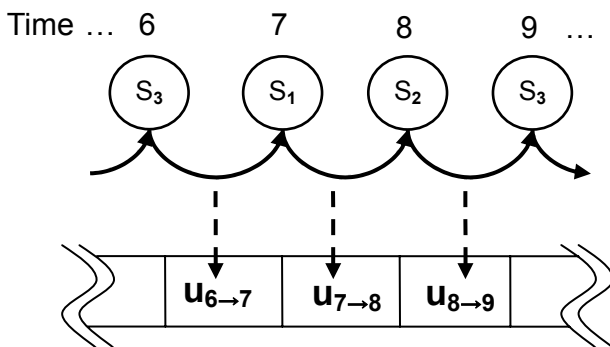


Figure 1. HMM のための Beam Sampler におけるスライサー  
上段が 1 つの隠れ状態列サンプル, 下段が動的計画法に導入されるスライサー列を表し, 破線がその対応を示している。各時刻間の状態遷移確率を用い, 動的計画法の各セルに対応するスライサー  $u_i$  をサンプリングする。観測列は省略している。

を提案する。この形式におけるスライサーは内側確率表において, 最終的に終端記号を出力する“終端セル”を担当するスライサー  $u$  と, 構文木の“分割点”を担当するスライサー  $v$  の 2 種類に分けて扱われる (Figure 3)。終端セルに対応する構文木上のノードは常に存在するため,  $u$  については Beam Sampler 同様に, 前回の構文木サンプル中の対応するノードにおいて用いられた文法規則の確率値をもとに値を決める。一方, それぞれの終端記号の間を“分割点”と定めると, 構文木の終端ノード以外のノード, すなわち分岐ノードは必ず 1 つの“分割点”を担当することがわかる。ここから, “分割点”を担当するスライサー  $v$  は, 構文木中の“分割点”を担当する分岐ノードの文法規則の確率値から作られる (Figure 4)。この手法では内側確率計算の際, 終端セル以外のセルにおける計算に複数の  $v$  が利用されることになるが, サンプルされた構文木における各ノードでは, そのいずれかのスライサーにより作られた文法規則が用いられ, 構文木全体で見るときにスライサーの重複はない。また, 導入すべきスライサーは最下列セル数  $n$ , 分割点数  $n-1$  であり, 構文木の持つノード数と一致するため, 必ず一一対応が出来る。この手法を我々は Split Position Slice Sampler と呼ぶ。以下では Split Position Slice Sampler を利用した構文木のサンプリングアルゴリズムの具体的な手順を示す。

終端記号列  $s$  に対する前回の構文木サンプルにおいて,  $s_k$  を導出した文法規則の確率値を  $\pi_k$ , 構文木の分割点  $j$  における文法規則の確率値を  $\rho_j$  と定義する。関数  $I(C)$  は, 条件  $C$  が真のとき  $I(C) = 1$ , 偽の時  $I(C) = 0$  となる関数である。

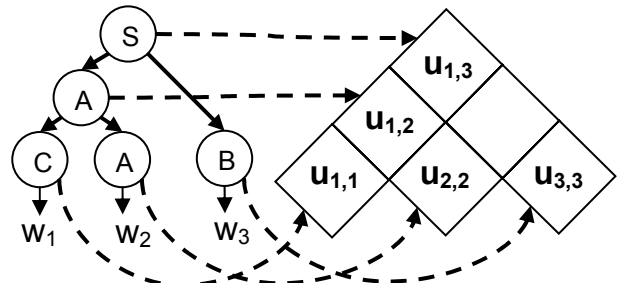


Figure 2. PCFG へ Beam Sampler を単純に適用した場合  
左が 1 つの構文木サンプル, 右が内側確率表に導入されるスライサーを表し, 破線がその対応を示している。情報が足りずスライサーを導入できないセルが存在する

### 1. Sampling $\mathbf{u}, \mathbf{v}$ step

$$u_k \sim \text{Uniform}(0, \pi_k) \quad (11)$$

$$v_j \sim \text{Uniform}(0, \rho_j) \quad (12)$$

### 2. Inside-Filtering Step

$$p_{k,k}^A = \mathbf{I}(u_k < \theta_{A \rightarrow w_k}) \quad (13)$$

$$p_{i,k}^A = \sum_{A \rightarrow BC \in R} \sum_{i \leq j < k} \mathbf{I}(v_j < \theta_{A \rightarrow BC}) p_{i,j}^B p_{j+1,k}^C \quad (14)$$

### 3. Tree-Sampling Step

Function SAMPLE( $A, i, k$ )

If  $i = k$

return TREE( $A, s_k$ )

Else

( $j, B, C$ ) = MULTI( $A, i, k$ )

return TREE( $A, \text{SAMPLE}(B, i, j), \text{SAMPLE}(C, j, k)$ )

関数 MULTIにおける確率計算は

$$P(j, B, C) = \frac{\mathbf{I}(\theta_{A \rightarrow BC}) p_{i,j}^B p_{j+1,k}^C}{p_{i,k}^A} \quad (15)$$

と置き換えられる。

このアルゴリズムは Johnson らの手法における構

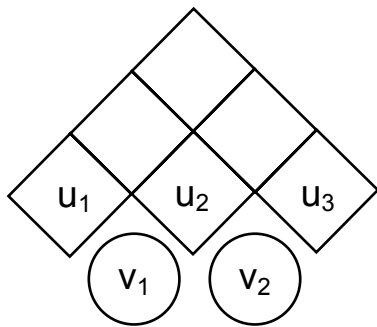


Figure 3. 終端セルにスライサー $\mathbf{u}$ , 分割点にスライサー $\mathbf{v}$ を導入. 結果として導入するスライサーの総数は構文木のノード数になる.

文木のサンプリングステップについての置き換えであり, その他の変更は無い. つまり本手法は Gibbs Sampler に組み込んで  $\theta$  のサンプリングと組みあわせるか, Metropolis-Hastings Sampler と組み合わせてサンプリングされた構文木の確率的な受理をすることで PCFG のサンプラーとして機能する.

以上で Beam Sampling 法の枠組みを PCFG へと適用した高速な MCMC 法である Split Position Slice Sampler が構築された. この手法は, 単純に考えた場合に対応のつかない, 形の違う構文木同士に対応をつけることで Beam Sampler の枠組みを拡張した. これにより, Beam Sampler の枠組みは Bayesian PCFG のためのサンプリングアルゴリズムに導入可能となり, 高速なサンプラーが構築された.

## 4 Experiments

人工的な CFG 文法から小規模なコーパスを生成し, ノンパラメトリック PCFG モデル[11]による教師無し文法学習により, 本手法と従来手法の比較を行った. コーパスは文法規則  $S \rightarrow xSy$ ,  $S \rightarrow zSw$ ,  $S \rightarrow SS$ ,  $S \rightarrow \epsilon$  から生成され, 平均文長 9 単語の 100 文を訓練データ, 平均文長 13 単語の 20 文をテストデータとして用い, 従来 of Johnson の手法と提案手法についてそれぞれ 50 回の実験を行い, 計算速度とテストデータの対数尤度について評価を行った.

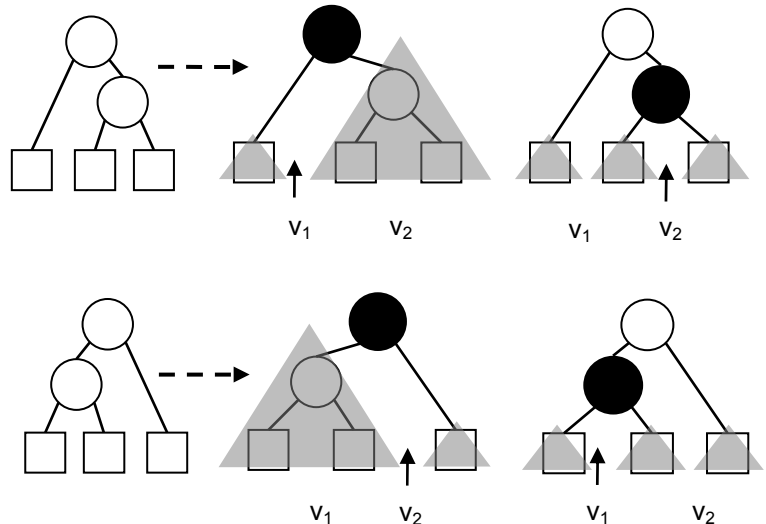


Figure 4. 分割点と分岐ノード, スライサーの関係

丸が分岐ノード, 四角が終端ノードを表し, 3 文字の終端記号列に対応する 2 種類の構文木を上下に配し, 黒いノードが担当する分割点を実線矢印で示している. どのような形の構文木でも必ず分割点を担当する分岐ノードが 1 つ存在し, 分割点を担当するスライサーを作ることができる.

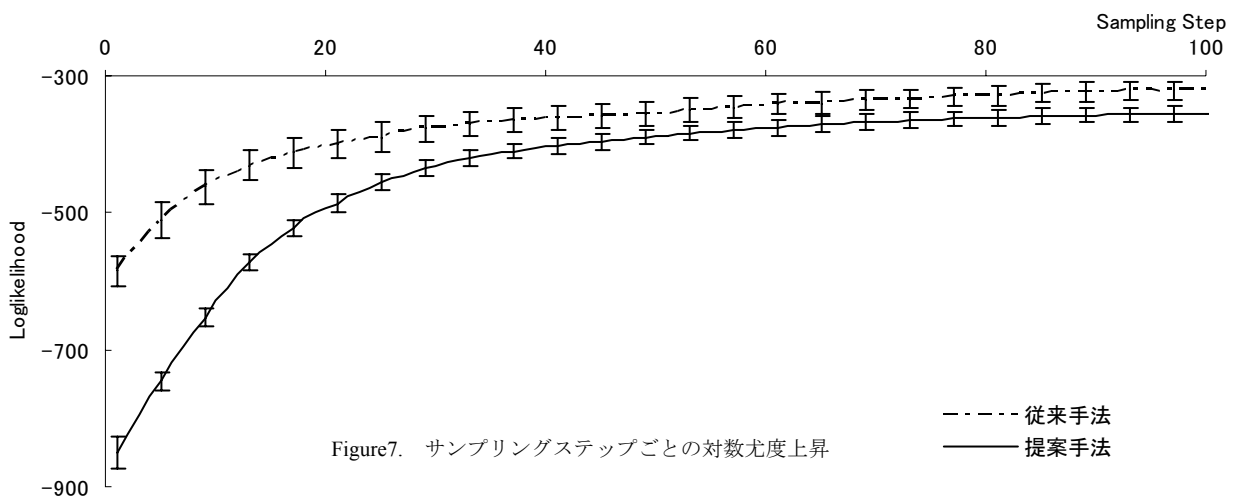
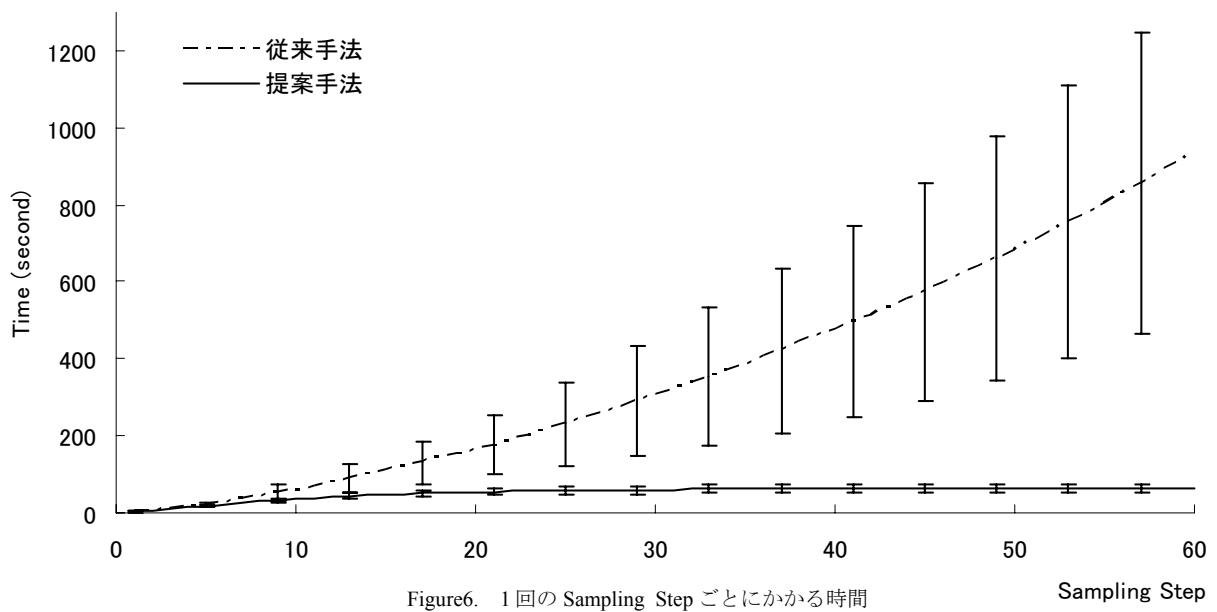
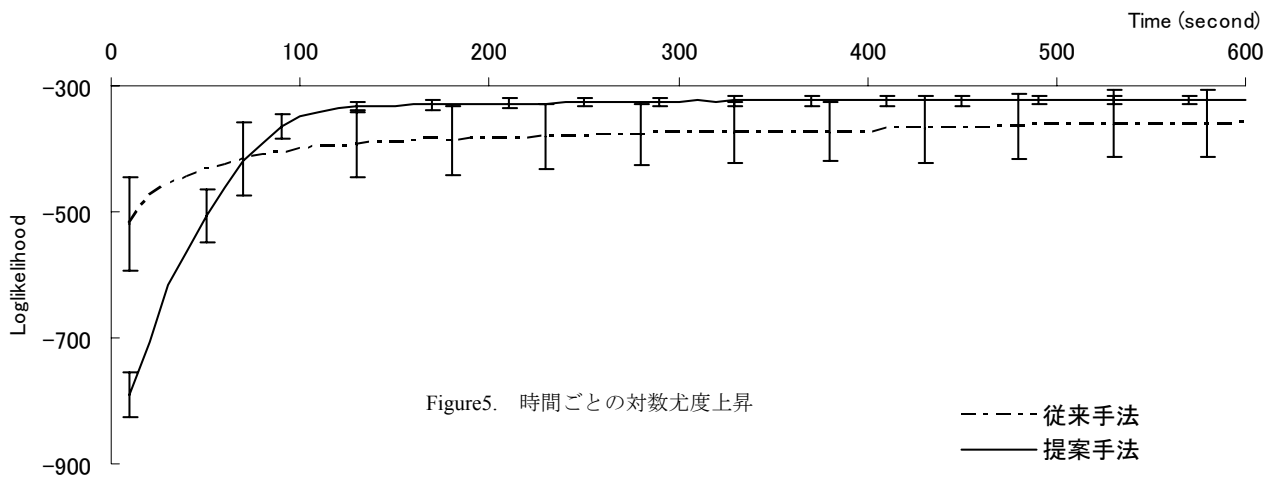


Figure5は時間ごとに見たコーパスの対数尤度, Figure6はサンプリング1ステップごとにかかる時間, Figure7はサンプリング1ステップごとに見たテストデータの対数尤度である. グラフは50回の実験による結果の平均値であり, 誤差棒は95%の信頼区間を示している. Figure5より, 提案手法は時間ごとの対数尤度が急速に上昇しており, 従来手法に比べより高速に学習していることがわかる. Figure5とFigure7においてサンプリングの初期に対数尤度が低いのは, スライサーという, 新たに推定すべきパラメータが加わったため, つまりモデルが複雑になったためであると考えられる. しかし, Figure6からは, 1ステップの計算速度が従来手法に比べて提案手法が大幅に高速化されているために, その影響を補って余りある高速化効果が得られ, その結果, 従来手法に比べ計算速度の向上がなされていることがわかった.

## 5 まとめ

本論文では Beam Sampler の枠組みを Split Position Slice Sampler へと拡張し, 従来手法に組み込むことで Bayesian PCFG のための高速なパラメータ推定手法を構築し, それを小規模な実験によって評価した. 我々の手法では一切の近似をせず, 高い精度でのパラメータ推定が可能であるという MCMC 法の性質を保っているため, 高速かつ高精度な手法であるといえる. このアルゴリズムは CKY アルゴリズムを伴えばどのような問題にも適用可能であり, ベイズ拡張された PCFG のスーパーセットのようなモデルにも適用可能であることが期待できる.

また, Beam Sampler は元々ノンパラメトリック拡張された HMM において提案されたものであるということから, 実験でも示したとおり, ノンパラメトリック拡張された PCFG のためのサンプリング法としても用いることができる. ただし, PCFG ではノンパラメトリック HMM と動的計画法の向きが逆であることなどが原因で, いくつかの変更が必要となる. しかしそれらは提案手法とは直接関係が無く, 性能評価にも影響しないため, 本論文では省略した.

今後は PCFG のスーパーセットのためのサンプリングアルゴリズムへの拡張といったアルゴリズム的な応用, 実際の大規模データへの適用などの実用面での応用の可能性を探るつもりである.

## 参考文献

- [1] Johnson, M., Griffiths, T. L. & Goldwater, S. (2007). Bayesian Inference for PCFGs via Markov Chain Monte Carlo. In *Proceedings of North American Chapter of Association for Computational Linguistics – Human Language Technologies*.
- [2] Gael, J. V., Saatchi, Y. Teh, Y. W. & Gharamani, Z. (2008) Beam Sampling for the Infinite Hidden Markov Model. In *Proceedings of the 25<sup>th</sup> International Conference on Machine Learning*.
- [3] Neal, R. M. (2003). Slice sampling. *The annals of Statistics*.
- [4] Teh, Y. W., Jordan, M. I., Beal, M. & Blei, D. (2006). Hierarchical Dirichlet processes. *Journal of the American Statistical Association*.
- [5] Blei, M., Gharamani, Z. & Rasmussen, C. (2002). The infinite hidden Markov model. In *Advances in Neural Information Processing Systems*.
- [6] Liang, P., Petrov, S., Jordan, M. I. & Klein, D. (2007) The infinite PCFG using hierarchical Dirichlet processes. In *Proceedings of the 2007 joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*.
- [7] Attias, H. (1999). Inferring Parameters and Structure of Latent Variable Models by Variational Bayes. In *Proceedings of Uncertainty in Artificial Intelligence*.
- [8] Kenichi, K., Yoshitaka, K., Taisuke, S. (2004). Variational Bayesian Approach to Probabilistic Context-Free Grammar based on Dynamic Programming. *Journal of Information Processing Society, Japan*.
- [9] Chaniak, E. (1996). Treebank grammars. *Association for the Advancement of Artificial Intelligence*.
- [10] Sakakibara, Y.; Brown, M.; Hughey, R.; Mian, I. S.; Sjölander, K.; Underwood, R. C.; and Haussler, D. (1994). Stochastic Context-Free Grammars for tRNA Modeling. *Nuc. Acids Res*.
- [11] Liang, P., Petrov, S., Jordan, M. I. & Klein, D. (2007) The infinite PCFG using hierarchical Dirichlet processes. In *Proceedings of the 2007 joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*.